PHOTOREALISTIC IMAGE SYNTHESIS USING RAY-BUNDLES

A dissertation submitted to the Hungarian Academy of Sciences in fulfillment of the requirements for the degree of Doctor of Science (*Doctor Academiae Scientiarum Hungaricae*)

by

Szirmay-Kalos László

Department of Control Engineering and Information Technology Technical University of Budapest

Year 2000

Contents

1	Introduction 1							
	1.1	Global pass						
	1.2	Local pass						
	1.3	Tone mapping 3						
2	Global illumination problem 5							
	2.1	The rendering equation						
	2.2	Measuring the radiance						
	2.3	3 The potential equation						
	2.4	4 Measuring the potential						
	2.5	The rendering problem						
		2.5.1 Geometry of the surfaces						
		2.5.2 Bi-directional Reflection Distribution Functions						
		2.5.3 Lightsources						
		2.5.4 Measuring devices						
	2.6	Numerical solution of the rendering equation						
		2.6.1 Error measures for numeric techniques						
		2.6.2 Properties of the rendering equation						
	2.7	Classification of the solution techniques						
3	Solution strategies for the global illumination problem 20							
	3.1	1 Inversion						
	3.2	Expansion						
		3.2.1 Expansion of the rendering equation: gathering walks						
		3.2.2 Expansion of the potential equation: shooting walks						
		3.2.3 Merits and disadvantages of expansion methods						
	3.3	Iteration						
		3.3.1 Analysis of the iteration						
	3.4	Analytical solution of the rendering equation 26						
	0	3.4.1 Scenes with constant radiance 26						
		3.4.2 Scenes with constant reflected radiance						
4	Finite-element methods for the global illumination problem 28							
	4.1	Galerkin's method						
	4.2	Point collocation method						
	4.3	Finite-element methods for the diffuse global illumination problem						
		4.3.1 Geometric methods for form factor computation						
5	Numerical quadrature for high dimensional integrals 34							
	5.1	Monte-Carlo quadrature						
	5.2	Quasi-Monte Carlo quadrature						
		5.2.1 Error analysis for integrands of finite variation: Koksma-Hlawka inequality 37						

		5.2.2 Generation of the sample points 40 5.2.3 Generation of low-discrepancy sequences 41
	5.3	Importance sampling 43
	0.0	5.3.1 Generation of a random variable with a prescribed probability density 43
		5.3.2 Importance sampling in quasi-Monte Carlo integration
		5.3.3 Metropolis sampling
6	Ran	dom walk solution of the global illumination problem 47
	6.1	Why should we use Monte-Carlo expansion methods?
	6.2	Quasi-Monte Carlo quadrature for the rendering equation
		6.2.1 Integrating functions of unbounded variation
	6.3	Importance sampling for the rendering equation
		6.3.1 BRDF sampling 53
		6.3.2 Lightsource sampling
		6.3.3 Sampling the lightsources in gathering random walks
		6.3.4 Importance sampling in colored scenes
		6.3.5 Multiple importance sampling
	6.4	Handling infinite-dimensional integrals
		6.4.1 Russian-roulette
	~ -	6.4.2 Russian-roulette in quasi-Monte Carlo quadrature
	6.5	Review of random walk algorithms
		6.5.1 Gathering-type random walk algorithms
		6.5.2 Shooting-type walk methods
7	Stoc 7 1	hastic iteration solution of the global illumination problem74Why should we use Monte-Carlo iteration methods?74
	7.1	Formal definition of stochastic iteration
	1.2	7.2.1 Other averaging techniques 77
	73	Stochastic iteration for the diffuse radiosity 77
	1.5	7.3.1 Stochastic radiosity 78
		7.3.2 Transillumination radiosity
		7.3.3 Randomly placed hemicubes
		7.3.4 Stochastic ray-radiosity
7.4 Definition of the random transport operator for the non-diffuse finite- ϵ		Definition of the random transport operator for the non-diffuse finite-element case 79
		7.4.1 Single ray based transport operator
8	Sim	ulating light transport using ray-bundles 83
	8.1	Reformulation of the rendering equation using finite-elements
	8.2	Stochastic expansion using ray bundles
		8.2.1 Generating uniformly distributed points on the sphere
		8.2.2 Simple Monte-Carlo, or quasi-Monte Carlo walks
		8.2.3 Calculation of the image estimate
		8.2.4 Improved walking techniques
		8.2.5 <i>D</i> -step iteration
	8.3	Importance sampling for the evaluation of directional integrals
		8.3.1 Application of the VEGAS algorithm
		8.3.2 Application of Metropolis sampling
		8.3.3 Evaluation of the performance of the Metropolis method
	8.4	Stochastic iteration using ray-bundles
	0 -	8.4.1 Can we use quasi-Monte Carlo techniques in iteration?
	8.5	Calculation of the radiance transport in a single direction
		o.s.i Galerkin's method with piece-wise constant basis functions 102

		8.5.2	Analysis of the finite resolution problem of discrete methods	109			
		8.5.3	Point collocation method with piece-wise linear basis functions	110			
	8.6	Handling sky-light illumination					
	8.7	7 Improving the efficiency					
		8.7.1	Self-correcting iteration	112			
		8.7.2	Preprocessing the small lightsources	113			
		8.7.3	Adaptive importance sampling and resolution control in iteration	120			
		8.7.4	Reducing the power defect of the iteration	120			
		8.7.5	Constant radiance step	121			
9	Simulation results						
	9.1	Testing	the walk method	122			
	9.2	Testing	stochastic iteration and self-correcting iteration	125			
		9.2.1	Self-correcting stochastic iteration	125			
		9.2.2	Self-correcting stochastic iteration with incoming first-shot	125			
10	Conclusions						
	10.1	Future	improvements	134			
BI	BIBLIOGRAPHY						
SU	SUBJECT INDEX						

Chapter 1

Introduction

The ultimate objective of **image synthesis** or **rendering** is to provide the user with the illusion of watching real objects on the computer screen (figure 1.1). The image is generated from an internal model which is called the **virtual world**. To provide the illusion of watching the real world, the color sensation of an observer looking at the artificial image generated by the graphics system must be approximately equivalent to the color perception which would be obtained in the real world. The color perception of humans depends on the **light power** reaching the eye from a given direction and on the operation of the eye. The power, in turn, is determined from the **radiance** of the visible points. The radiance depends on the shape and optical properties of the objects and on the intensity of the lightsources. In order to model this complex phenomenon, both the physical-mathematical structure of the light-object interaction and the operation of the eye must be understood.



Figure 1.1: Tasks of rendering

The image synthesis uses an internal model consisting of the **geometry of the virtual world**, **optical material properties** and the description of the **lighting** in the scene (figure 1.2). From these, applying the laws of physics (e.g. Maxwell equations) the real world optical phenomena can be simulated to find the light distribution in the scene. This step is called the view-independent step or the **global pass** of

rendering. Then a **measurement device**, called the **eye** or **camera**, is placed into the scene and the light distribution is measured from a given location and orientation. This is called the view-dependent step or the **local pass**. Note that not all rendering algorithms make a clear distinction between the determination of the view-independent light distribution and the measurement of this distribution by the camera, but simultaneously compute the light distribution and its effect on the camera.

Rendering results in a representation of the perceived image, which is usually the collection of pixel colors or some discrete sampling of the radiance function. The exact simulation of the light perceived by the eye is impossible, since it would require endless computational process. On the other hand, it is not even worth doing since the possible distributions which can be produced by computer screens are limited in contrast to the infinite variety of real world light distributions. Consequently, color perception is approximated instead of having a completely accurate simulation. The accuracy of this approximation is determined by the ability of the eye to make a distinction between two light distributions.

Computer screens can produce controllable electromagnetic waves, or colored light, mixed from three separate wavelengths for their observers. Thus in the final step of image synthesis **tone mapping** is needed which converts the computed color or radiance information to the R, G, B intensities that can be produced by the color monitor.



Figure 1.2: Dataflow of rendering

1.1 Global pass

The **global pass** determines the light reflected off the surface points at different directions. Since light is an electromagnetic wave, light distribution in a point and at a given direction can be represented by a wavelength-dependent function [Ábr97, Ant80]. Rendering algorithms usually evaluate this functions at a few representative wavelengths. On a given wavelength the intensity of the light is described by the **radiance**. In scenes not incorporating **participating media** it is enough to calculate the radiance at surface points. The radiance reflected off a surface point is affected by the emission of this point (**lighting**), the illumination provided by other surface points and the optical properties of the material at this point (**material properties**). Formally this dependence is characterized by a Fredholm type integral equation of the second kind, which is called the **rendering equation**. From mathematical point of view, global pass is the solution of this integral equation for the representative wavelengths.

1.2 Local pass

The **local pass** means the measurement of the global radiance function by a camera. A camera is a collection of light measuring devices which usually correspond to pixels in the image. A certain measuring device is characterized by a sensitivity function that describes which points and directions may affect the device.

1.3 Tone mapping

Light is an electromagnetic wave, and its **color** is determined by the eye's perception of its spectral energy distribution. Due to its internal structure, the eye is a very poor spectrometer since it actually samples and integrates the energy in three overlapping frequency ranges by three types of photopigments according to a widely accepted (but also argued) model. As a consequence of this, any color perception can be represented by three scalars (called **tristimulus** values) instead of complete functions [Ábr97, Ant80, Nem90].

A convenient way to define the axes of a coordinate system in the three-dimensional space of color sensations is to select three wavelengths where one type of photopigments is significantly more sensitive than the other two [SK99d]. This method has been devised by Grassmann, who also specified a criterion for separating the three representative wavelengths. **Grassmann laws** state that the representative wavelengths should be selected such that no one of them can be matched by the mixture of the other two in terms of color sensation (this criterion is similar to the concept of linear independence.) An appropriate collection of the representative wavelengths is:

$$\lambda_{\text{red}} = 645 \text{ nm}, \quad \lambda_{\text{green}} = 526 \text{ nm}, \quad \lambda_{\text{blue}} = 444 \text{ nm}.$$
 (1.1)

Now let us suppose that monochromatic light of wavelength λ is perceived by the eye. The equivalent portions of red, green and blue light, or (r, g, b) tristimulus values, can be generated by three **color matching functions** $(r(\lambda), g(\lambda) \text{ and } b(\lambda))$ which are based on physiological measurements. Note the negative section of $r(\lambda)$ (and to a less extent in $g(\lambda)$) in figure 1.3. It means that not all colors can be represented by positive (r, g, b) values.



Figure 1.3: Mean 10-deg color matching functions of Stiles and Burch: $r(\lambda)$, $g(\lambda)$, $b(\lambda)$.

If the perceived color is not monochromatic, but is described by an $L(\lambda)$ distribution, the tristimulus coordinates are computed using the assumption that the sensation is produced by an additive mixture of the perceptions of elemental monochromatic components:

$$r = \int_{\lambda} L(\lambda) \cdot r(\lambda) \, d\lambda, \quad g = \int_{\lambda} L(\lambda) \cdot g(\lambda) \, d\lambda, \quad b = \int_{\lambda} L(\lambda) \cdot b(\lambda) \, d\lambda. \tag{1.2}$$

For computer generated images, the color sensation of an observer watching a virtual world on the screen must be approximately equivalent to the color sensation obtained in the real world. Since color sensations are represented by (r, g, b), it means that the tristimulus values should be identical. If two

energy distributions are associated with the same tristimulus coordinates, they produce the same color sensation, and are called **metamers**.

In computer monitors and in television screens three phosphor layers can be stimulated to produce red, green and blue light. The objective, then, is to find the necessary stimulus to produce a metamer of the real energy distribution of the light [Sch96, BS95]. This stimulus can be controlled by the (R, G, B) values of the actual pixel. The (r, g, b) matching functions of figure 1.3 depend on the wavelength of the selected primaries, which are not necessarily identical to the wavelengths on which our monitor can emit light. This requires the conversion of tristimulus values by a linear transformation.

The calculation of pixel R, G, B values thus consists of the following steps. First the spectrum associated with the pixel is computed. Then the spectrum is matched by three standard color matching functions defined by three primaries. Finally, the standard color coordinates are transformed to the monitor color coordinates taking into account the monitor properties. In practice, the standard color system is usually the **CIE XYZ** system [WS82] which uses three hypothetical primaries to allow the definition of any color by positive weights.



Figure 1.4: Mean 10-deg color XYZ matching functions of Stiles and Burch: $X(\lambda), Y(\lambda), Z(\lambda)$

The linear transformation that converts from the **XYZ** system to the monitor **RGB** system can be obtained from the X, Y, Z coordinates of the emissions of the three phosphors and of the white point of the monitor. For a monitor with **standard NTSC phosphors** and **white point**, the following transformation can be used [Gla95]:

$$\begin{bmatrix} R\\G\\B \end{bmatrix} = \begin{bmatrix} 1.967 & -0.548 & -0.297\\-0.955 & 1.938 & -0.027\\0.064 & -0.130 & 0.982 \end{bmatrix} \cdot \begin{bmatrix} X\\Y\\Z \end{bmatrix}.$$
 (1.3)

The whole computation of the (R, G, B) values in order for the monitor color to be a metamer of the calculated spectrum is called **tone mapping**. The (R, G, B) values are positive numbers usually in the range of [0...255] if 8 bits are available to represent them. Unfortunately, not all colors can be reproduced on the computer screen, because of the negative sections of the color matching functions and due to the fact that the number of available intensity levels is usually much less than can be perceived in the real world. Thus tone mapping is also responsible for optimally selecting from the available intensity levels for color reproduction. The mapping from the computed levels to the available ones can be either linear or logarithmic. The latter takes advantage of the logarithmic characteristics of the human perception system [PP98].

Chapter 2

Global illumination problem

In this chapter the mathematical model of the light-surface interaction is presented. This mathematical model is an integral equation, which has to be solved to obtain physically accurate images.

2.1 The rendering equation

Hereinafter, monochromatic light of a representative wavelength λ will be assumed, since the complete color calculation can be broken down to these representative wavelengths. The parameters of the equations usually depend on the wavelength, but for notational simplicity, we do not always include the λ variable in them.

In this section, we briefly review the measures of the light transport and the mathematical formulation that can compute them.



Figure 2.1: Definition of directions in a spherical coordinate system (left) and calculation of differential solid angles (right)

The directional property of the light energy emission is described in a so-called **illumination sphere** Ω or in **illumination hemisphere** Ω_H which contain those solid angles to where the surface point can emit energy. The surface of transparent materials can emit in any directions of a sphere, while the surface of opaque materials can transfer energy only to the hemisphere that is "above" the surface.

Setting up a **spherical coordinate system** (figure 2.1), a direction ω can be defined by two angles θ, ϕ , where θ is the angle between the given direction and the *z*-axis, and ϕ is the angle between the projection of the given direction onto the *x*, *y* plane and the *x*-axis.

Sets of directions are defined by solid angles. By definition, a **solid angle** is a cone or a pyramid, with its size determined by its subtended area of a unit sphere centered around the apex (figure 2.2). A differential (infinitesimal) solid angle can also be given by a vector $d\vec{\omega}$, where the vector equals to a direction of the differential set.

A differential solid angle can also be expressed by the θ , ϕ angles. Suppose that θ is modified by $d\theta$ and ϕ is by $d\phi$. During this the directional vector scans a differential rectangle having $d\theta$ vertical and



Figure 2.2: Definition of the solid angle

 $\sin\theta \cdot d\phi$ horizontal sizes (right of figure 2.1), thus the size of the solid angle is

$$d\omega = \sin\theta \cdot d\phi d\theta. \tag{2.1}$$

The solid angle, in which a differential dA surface can be seen from point \vec{p} , is the projected (visible) area per the square of the distance of the surface (figure 2.2). If the angle between the surface normal of dA and the directional vector from dA to \vec{p} is θ , and the distance from dA to \vec{p} is r, then this solid angle is:

$$d\omega = \frac{dA \cdot \cos\theta}{r^2}.$$
(2.2)

The intensity of the energy transfer is characterized by several metrics in computer graphics depending on whether or not the directional and positional properties are taken into account.

The light **power** or **flux** Φ is the energy radiated through a boundary per unit time over a given range of the spectrum (say $[\lambda, \lambda + d\lambda]$). Since a **photon** has \hbar/λ energy where \hbar is the Planck-constant, the power can be supposed to be proportional to the number of photons that go through the boundary in a unit time. The power is not always a convenient measure since it also needs the definition of a boundary. We can get rid of this additional information if the boundary is defined in a differential way focusing on a single surface point and a single direction. The resulting measure is called the radiance.

The **radiance** or **intensity** $L(\vec{x}, \omega)$ is the differential light flux $\Phi(\vec{x}, dA, \omega, d\omega)$ leaving a surface element dA around \vec{x} in a differential solid angle $d\omega$ around ω per the projected area of the surface element dA and the size of the solid angle $d\omega$. If the angle of the surface normal and the direction of interest is θ , then the projected area is $dA \cdot \cos \theta$, hence the radiance is:

$$L(\vec{x},\omega) = \frac{d\Phi(\vec{x},dA,\omega,d\omega)}{dA \cdot d\omega \cdot \cos\theta}.$$
(2.3)

Since a differential surface area can also be seen as a vector $d\vec{A}$, which is parallel to the surface normal at the given point, the radiance can also be obtained in a simpler form

$$L(\vec{x},\omega) = \frac{d\Phi(\vec{x},dA,\omega,d\omega)}{d\vec{A}\cdot d\vec{\omega}}.$$
(2.4)

Having introduced the most important metrics, we turn to their determination in the simplest case, where there are only two differential surface elements in the 3D space, one (dA) emits light energy and the other (dA') absorbs it (figure 2.3). According to the definition of the radiance (equation (2.3)), if dA' is visible from dA in solid angle $d\omega$ and the radiance of the surface element dA is L in this direction, then the flux leaving dA and reaching dA' is:

$$d\Phi = L \cdot dA \cdot d\omega \cdot \cos \theta. \tag{2.5}$$

Using equation (2.2), the solid angle can be expressed from the projected area of dA', thus we get:

$$d\Phi = L \cdot \frac{dA \cdot \cos\theta \cdot dA' \cdot \cos\theta'}{r^2}.$$
 (2.6)



Figure 2.3: Energy transfer between two differential surface elements

This formula is called the fundamental law of photometry.

Note that if equation (2.2) is used again for the emitter, the transferred power can also be written in the following form:

$$d\Phi = L \cdot \frac{dA \cdot \cos \theta \cdot dA' \cdot \cos \theta'}{r^2} = L \cdot dA' \cdot d\omega' \cdot \cos \theta'.$$
(2.7)

Thus similar formula applies for the patch that receives the power as for the patch that emits it.

In light-surface interaction the surface illuminated by an incident beam may reflect a portion of the incoming energy in various directions or it may absorb the rest. It has to be emphasized that a physically correct model must maintain energy equilibrium, that is, the reflected and the transmitted (or absorbed) energy must be equal to the incident energy.

Optically perfect or smooth surfaces will reflect or transmit only **coherent components** governed by the laws of geometric optics, including the law of reflection and the Snellius–Descartes law of refraction. The surface irregularities, however, reflect or refract the incident light incoherently in any direction. Since the exact nature of these irregularities is not known, light-surface interaction is modeled by means of probability theory.

Assume that a photon comes from the direction denoted by ω' to point \vec{x} . The probability of reflection or refraction at \vec{x} into solid angle $d\omega$ around ω is expressed by the following probability density function, also called as **transfer probability density**:

 $w(\omega', \vec{x}, \omega) \cdot d\omega = \Pr\{\text{photon is reflected or refracted to } d\omega \text{ around } \omega \mid \text{coming from } \omega'\}.$ (2.8)

Note that this probability distribution is a mixed, discrete-continuous distribution, since the probability that the light is reflected exactly to the ideal reflection direction may be non-zero.



Figure 2.4: Geometry of the rendering equation

The light flux (Φ^{out}) leaving the surface at solid angle $d\omega$ around ω consists of the emission and reflected (or refracted) components.

In order to obtain the **reflected/refracted component**, let us assume that photons are coming to area dA around \vec{x} from solid angle $d\omega'$ around ω' and their total represented power is $\Phi^{\text{in}}(\vec{x}, dA, \omega', d\omega')$. The probability that a photon is reflected to $d\omega$ is $w(\omega', \vec{x}, \omega) \cdot d\omega$ thus the expected power leaving the surface element after reflection or refration is:

$$w(\omega', \vec{x}, \omega) \ d\omega \cdot \Phi^{\mathrm{in}}(\vec{x}, \omega', d\omega').$$

To obtain the total reflected/refracted power, all the possible incoming directions ω' should be taken into account and their contributions should be integrated:

$$\int_{\Omega} (w(\omega', \vec{x}, \omega) \ d\omega) \Phi^{\text{in}}(\vec{x}, \omega', d\omega').$$
(2.9)

If the surface itself emits energy, i.e. if it is a lightsource, then the **emission** also contributes to the output flux:

$$\Phi^e(\vec{x},\omega) = L^e(\vec{x},\omega) \cdot dA \cdot \cos\theta \cdot d\omega.$$
(2.10)

Adding the possible contributions we obtain:

$$\Phi^{\text{out}}(\vec{x},\omega) = \Phi^{e}(\vec{x},\omega) + \int_{\Omega} (w(\omega',\vec{x},\omega) \ d\omega) \Phi^{\text{in}}(\vec{x},\omega',d\omega').$$
(2.11)

The fluxes can be expressed by the radiant intensities according to equations (2.5) and (2.7), thus:

$$\Phi^{\text{in}}(\vec{x},\omega',d\omega') = L^{\text{in}}(\vec{x},\omega') \cdot dA \cdot \cos\theta' \cdot d\omega',$$

$$\Phi^{\text{out}}(\vec{x},\omega,d\omega) = L(\vec{x},\omega) \cdot dA \cdot \cos\theta \cdot d\omega.$$
(2.12)

Substituting these terms into equation 2.11 and dividing both sides by $dA \cdot d\omega \cdot \cos \theta$ we get:

$$L(\vec{x},\omega) = L^{e}(\vec{x},\omega) + \int_{\Omega} L^{in}(\vec{x},\omega') \cdot \cos\theta' \cdot \frac{w(\omega',\vec{x},\omega)}{\cos\theta} d\omega'.$$
(2.13)

Let us realize that $L^{\text{in}}(\vec{x}, \omega')$ is equal to $L(\vec{y}, \omega')$ if \vec{y} is the point that is visible from \vec{x} at direction $-\omega'$, usually expressed by the so called **visibility function** h (i.e. $\vec{y} = h(\vec{x}, -\omega')$).

Using these equations and introducing the **bi-directional reflection/refraction function** — or **BRDF** for short — as defined by the following formula

$$f_r(\omega', \vec{x}, \omega) = \frac{w(\omega', \vec{x}, \omega)}{\cos \theta},$$
(2.14)

we obtain the following fundamental law, called the rendering equation:

$$L(\vec{x},\omega) = L^{e}(\vec{x},\omega) + \int_{\Omega} L(h(\vec{x},-\omega'),\omega') \cdot f_{r}(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega'.$$
(2.15)

Introducing the following notation for the integral operator

$$(\mathcal{T}L)(\vec{x},\omega) = \int_{\Omega} L(h(\vec{x},-\omega'),\omega') \cdot f_r(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega', \qquad (2.16)$$

the short form of the rendering equation can also be established:

$$L = L^e + \mathcal{T}L. \tag{2.17}$$

In fact, every color calculation problem consists of several rendering equations, one for each representative wavelength. Optical parameters (L^e, f_r) obviously vary in the different equations, but the geometry represented by function h does not.

2.2 Measuring the radiance

Having solved the rendering equation, the radiance can be determined for any point and direction. To obtain an image, the power that affect different parts of light sensitive material (retina or the film of a camera) must be determined. Mathematically each distinct part is associated with a measuring device, thus the collection of these devices is, in fact, the model of the human eye or cameras.

A measuring device is characterized by a **sensitivity function** $W^e(\vec{y}, \omega')$, which gives a **scaling** value C if a light-beam of unit power leaving \vec{y} in direction ω' contributes to the measured value and 0 otherwise (for example, this device can measure the power going through a single pixel of the image and landing at the eye, or leaving a surface element at any direction).

Since the power leaving surface area $d\vec{y}$ around \vec{y} in solid angle $d\omega$ around ω is $L(\vec{y}, \omega)d\vec{y} \cos\theta d\omega$, the measured value by a certain device is:

$$\int_{\Omega} \int_{S} d\Phi(\vec{y},\omega) \cdot W^{e}(\vec{y},\omega) = \int_{\Omega} \int_{S} L(\vec{y},\omega) \cos \theta \cdot W^{e}(\vec{y},\omega) \, d\vec{y} \, d\omega = \mathcal{M}L.$$
(2.18)

Operator \mathcal{M} is called the **radiance measuring operator**.

2.3 The potential equation

Looking at the light propagation as an interaction between an emitter and a receiver, the radiance describes this interaction from the point of view of the emitter. The same phenomenon can also be explained from the point of view of the receiver, when the appropriate measure is called the **potential**. By definition, the potential $W(\vec{y}, \omega')$ expresses the effect of that portion of the unit power light-beam emitted by point \vec{y} in direction ω' , which actually lands at a given measuring device either directly or indirectly after some reflections or refractions. Using probabilistic terms the potential is the product of the scaling value C and the probability that a light-beam emitted by a point in a given direction reaches the measuring device.



Figure 2.5: Geometry of the potential equation

If only **direct contribution** is considered, then $W(\vec{y}, \omega') = W^e(\vec{y}, \omega')$. To take into account light reflections, we can establish the **potential equation** [PM95]:

$$W = W^e + \mathcal{T}'W. \tag{2.19}$$

In this equation integral operator \mathcal{T}' — which is the adjoint of \mathcal{T} — describes the potential transport

$$(\mathcal{T}'W)(\vec{y},\omega') = \int_{\Omega} W(h(\vec{y},\omega'),\omega) \cdot f_r(\omega',h(\vec{y},\omega'),\omega) \cdot \cos\theta \, d\omega, \qquad (2.20)$$

where θ is the angle between the surface normal and the outgoing direction ω .

To prove this equation, we can use the probabilistic interpretation of the potential. Denoting the "*unit power light-beam that is leaving* \vec{x} *at direction* ω " by $\Pi(\vec{x}, \omega)$, we can write:

$$W(\vec{y}, \omega') = C \cdot \Pr{\Pi(\vec{y}, \omega') \text{ lands at the device}} =$$

 $C \cdot \Pr{\{\Pi(\vec{y}, \omega') \text{ goes directly to the device}\} + C \cdot \Pr{\{\Pi(\vec{y}, \omega') \text{ lands at the device after at least one reflection}\}}.$

The first term is equal to the sensitivity function. Considering the second term, we can obtain:

 $C \cdot \Pr{\{\Pi(\vec{y}, \omega') \text{ lands at the device after at least one reflection}\}} =$

$$\int_{\Omega} C \cdot \Pr\{\Pi(h(\vec{y}, \omega'), \omega) \text{ lands at the device}\} \cdot \Pr\{\Pi(\vec{y}, \omega') \text{ is reflected to } d\omega \text{ around } \omega \text{ at } h(\vec{y}, \omega')\} = \int_{\Omega} W(h(\vec{y}, \omega'), \omega) \cdot f_r(\omega', h(\vec{y}, \omega'), \omega) \cdot \cos \theta \, d\omega.$$
(2.22)

2.4 Measuring the potential

Alternatively to the radiance, the power arriving at the measuring device can also be computed from the potential. Since

$$d\Phi^e(\vec{y},\omega') = L^e(\vec{y},\omega') \cdot \cos\theta \ d\vec{y} \ d\omega'$$

is proportional to the power of the light-beam emitted by $d\vec{y}$ in $d\omega'$ and the probability that the photons of this beam really go either directly or indirectly to the measuring device is $W(\vec{y}, \omega')/C$, the total measured value that takes into account all possible emission points and directions is

$$C \cdot \int_{\Omega} \int_{S} d\Phi^{e}(\vec{y}, \omega') \cdot \frac{W(\vec{y}, \omega')}{C} = \int_{\Omega} \int_{S} W(\vec{y}, \omega') \cdot L^{e}(\vec{y}, \omega') \cdot \cos\theta \, d\vec{y} \, d\omega' = \mathcal{M}'W, \tag{2.23}$$

where \mathcal{M}' is the **potential measuring operator**. Note that unlike the radiance measuring operator, the potential measuring operator integrates on the lightsource.

2.5 The rendering problem

Generally, the **rendering problem** is a quadruple [Kel97]:

$$\langle S, f_r(\omega', \vec{x}, \omega), L^e(\vec{x}, \omega), \mathbf{W}^e(\vec{x}, \omega) \rangle$$

where S is the geometry of surfaces, f_r is the BRDF of surface points, L^e is the emitted radiance of surface points at different directions and \mathbf{W}^e is a collection of measuring functions. Rendering algorithms aim at modeling and simulation of light-surface interactions to find out the power emitted by the surfaces and landing at the measuring devices. The light-surface interaction can be formulated by the *rendering equation* or alternatively by the *potential equation*. Since according to the definition of the radiance the total power of the scene is the integral of the radiance function

$$\Phi = \int_{\Omega} \int_{S} L(\vec{y}, \omega) \ d\vec{y} \ d\vec{\omega},$$

we are interested only in those solutions where the total integral of the radiance function is finite. Formally the solution is searched in the L_1 function space.

Recall that the measured value can be computed by the measuring function from the radiance

$$\int_{\Omega} \int_{S} L(\vec{y}, \omega) \cos \theta \cdot W^{e}(\vec{y}, \omega) \, d\vec{y} \, d\omega = \mathcal{M}L,$$
(2.24)

(2.21)

where \mathcal{M} is the radiance measuring operator. Let us introduce the scalar product $\langle u, v \rangle$

$$\langle u, v \rangle = \int_{\Omega} \int_{S} u(\vec{y}, \omega) \cdot v(\vec{y}, \omega) \, d\vec{y} \, d\vec{\omega} = \int_{\Omega} \int_{S} u(\vec{y}, \omega) \cdot v(\vec{y}, \omega) \, d\vec{y} \cos \theta \, d\omega$$

where θ is the angle between the surface normal and direction ω and thus $d\vec{y} \cos \theta$ is the infinitesimal visible area from direction ω . Using this scalar product, we can obtain an alternative form of the measuring operator:

$$\mathcal{M}L = \langle L, W^e \rangle.$$

The **potential measuring operator** can also be given in a scalar product form:

$$\mathcal{M}'W = \langle L^e, W \rangle. \tag{2.25}$$

Let us examine a single reflection of the light. The measured value taking into account a single reflection can be expressed in two ways:

$$\langle \mathcal{T}L^e, W^e \rangle = \langle L^e, \mathcal{T}'W^e \rangle. \tag{2.26}$$

Thus the radiance and potential transport operators are adjoint operators [Mát81].

2.5.1 Geometry of the surfaces

A surface is a set of 3D points which are defined by an equation [SK99d]. Those points are said to be in this set, which satisfy the definition equation. The equation can produce the points on the surface either implicitly, when the general form is

$$F(x, y, z) = 0$$

or explicitly, which is preferred in rendering algorithms. The general form of an explicit surface definition is:

$$\vec{r} = \vec{r}(u, v), \quad u \in [0, 1], \ v \in [0, 1].$$
 (2.27)

Points on the surface can be generated by selecting u, v parameters from the unit interval and substituting their values into the function $\vec{r}(u, v)$. For example, a sphere of radius R and of center (x_0, y_0, z_0) can be defined either by the following explicit equation:

 $x = x_0 + R \cdot \cos 2\pi u \cdot \sin \pi v, \quad y = y_0 + R \cdot \sin 2\pi u \cdot \sin \pi v, \quad z = z_0 + R \cdot \cos \pi v, \quad u, v \in [0, 1],$ or by an **implicit equation**:

$$(x - x_0)^2 + (y - y_0)^2 + (z - z_0)^2 - R^2 = 0.$$



Figure 2.6: Tessellation process

Some rendering algorithms do not work with the original geometry but approximate the surfaces by triangle meshes. This approximation — also called the tessellation — can be realized by selecting $n \times m$ points in parameter space $u \in [0, 1], v \in [0, 1]$ and adding the

$$[\vec{r}(u_i, v_j), \vec{r}(u_{i+1}, v_j), \vec{r}(u_{i+1}, v_{j+1})]$$
 and $[\vec{r}(u_i, v_j), \vec{r}(u_{i+1}, v_{j+1}), \vec{r}(u_i, v_{j+1})]$

triangles for all $i = 1 \dots n - 1$ and $j = 1 \dots m - 1$ indices to the resulting list (figure 2.6). For the discussion of surface definition methods using forward and reverse engineering and of transformations refer to [SK99d, Chi88, VMC97, RVW98] and to [SKe95, Kra89, Her91, Lan91], respectively.

2.5.2 Bi-directional Reflection Distribution Functions

The **Bi-directional Reflection Distribution Functions** (or **BRDF**s) characterize the optical material properties. Photorealistic rendering algorithms require the BRDFs not to violate physics. Such BRDF models must satisfy both reciprocity and energy balance, and are called **physically plausible** [Lew93].

Reciprocity that was recognized by Helmholtz is the symmetry property of the BRDF characterizing reflections, which is defined by the following equation [Min41]:

$$f_r(\omega, \vec{x}, \omega') = f_r(\omega', \vec{x}, \omega), \qquad (2.28)$$

where ω' is the vector pointing towards the incoming light and vector ω defines the viewing direction. Reciprocity is important because it allows for the backward tracing of the light as happens in visibility ray-tracing algorithms. Note that reciprocity is valid if the incoming and outgoing beams are in the same material, that is, reflection BRDFs are expected to be reciprocal but refraction BRDFs are not.

Suppose that the surface is illuminated by a beam from direction ω' . Energy balance means that the number of outgoing photons cannot be more than how many were received. To examine it formally, we can introduce the **albedo** [Lew93] that is the ratio of the total reflected power and incoming power, or equivalently, the probability of the reflection to any direction. Energy balance means that the albedo cannot be greater than 1:

$$a(\vec{x}, \omega') = \Pr\{\text{photon is reflected} \mid \text{coming from } \omega'\} = \int_{\Omega_H} w(\omega', \vec{x}, \omega) \, d\omega = \int_{\Omega_H} f_r(\omega', \vec{x}, \omega) \cdot \cos\theta \, d\omega \le 1.$$
(2.29)

If the BRDF is reciprocal, then the albedo can also be interpreted as the optical response of the surface to homogeneous sky-light illumination of unit intensity:

$$\mathcal{T}1 = \int_{\Omega_H} 1 \cdot f_r(\omega', \vec{x}, \omega) \cdot \cos \theta' \, d\omega' = a(\vec{x}, \omega).$$
(2.30)

Using the definition of the transfer probability density w in equation (2.8), we can obtain the following probabilistic interpretation of the BRDF:

$$f_r(\omega', \vec{x}, \omega) \cdot \cos \theta = \frac{1}{d\omega} \cdot \Pr\{\text{photon is reflected or refracted to } d\omega \text{ around } \omega \mid \text{it comes from } \omega'\}.$$

Different BRDF models are discussed in [SKe95, SK99d, SK99b, Pho75, Bli77, CT81, War92, HTSG91, ON94, Sch93, NNSK99c, NNSK99a, CSK98].

2.5.3 Lightsources

The **lightsource**s are surfaces that have non-zero L^e emission. Rendering algorithms also use abstract lightsources that are not associated with surfaces [SKe95]. The most important types of these **abstract lightsource models** are the following:

- a **point-lightsource** is located at a given point of the virtual world and its size is zero. The direction of the illumination provided by this lightsource at a point is the vector from the point to the location of the lightsource and the intensity of the illumination decreases with the square of the distance. An example of this category is an electric bulb.
- a **directional-lightsource** is located at infinity in a given direction. The direction and the intensity of its illumination are constant everywhere. A directional-lightsource is, in fact, equivalent to an infinitely distant plane emitter. The sun can be considered as a directional-lightsource.
- sky-light illumination provides constant illumination at any point and at any direction.

2.5.4 Measuring devices

In order to establish models for the **camera**, the operation of the **human eye** can be analyzed when the human is looking at the real world and when he is watching the computer screen.



Figure 2.7: A model of the human eye

The human eye contains a lens called the **pupil** of size Δe (figure 2.7). We shall assume that the pupil is small, which is often referred as the **pinhole camera model**. When the human is watching the computer screen, a pixel p is visible in a small solid angle Ω_p . To provide the illusion that the screen is a replacement of the real world, the Φ_p power emitted by the pixel towards the eye should be similar to that Φ power which would come from the real world and would land at the pupil from solid angle Ω_p . If the emission intensity of the pixel is L_p , then the power landing at the eye from this pixel is

$$\Phi_p = L_p \cdot \Delta e \cdot \cos \theta_e \cdot \Omega_p,$$

where θ_e is the angle between the surface normal on the pupil and the direction of the pixel.

We have to find a camera model that computes a measured value P which is stored in the **image-buffer** to control the monitor. The response of **monitors** can be characterized by a function $B \cdot \mathcal{R}$, where B represents a scaling according to the brightness settings and \mathcal{R} might be non-linear. In order to compensate the monitor's non-linearity, the **look-up table** is set to distort the values of the image-buffer with the inverse of this non-linear mapping, that is by \mathcal{R}^{-1} . This distortion is called **gamma-correction**. Using gamma-correction, the radiant intensity of the pixel is:

$$L_p = B \cdot \mathcal{R}(\mathcal{R}^{-1}(P)) = B \cdot P.$$

Since we require that $\Phi_p = \Phi$, our model camera is expected to provide the following measured value:

$$P = \mathcal{R}\left(\mathcal{R}^{-1}\left(\frac{L_p}{B}\right)\right) = \frac{L_p}{B} = \frac{\Phi}{\Delta e \cdot \cos\theta_e \cdot \Omega_p \cdot B}.$$
(2.31)

Let us assign a measuring device for this pixel. This device is sensitive to those surface points that are visible in this pixel — i.e. in Ω_p — and to those directions that point from the surface points to the pupil. Mathematically, this can be represented by the following measuring function

$$W^{e}(\vec{y},\omega) = \begin{cases} C & \text{if } \vec{y} \text{ is visible in } \Omega_{p} \text{ and } \omega \text{ points from } \vec{y} \text{ through } \Delta e, \\ 0 & \text{otherwise,} \end{cases}$$
(2.32)

where

$$C = \frac{1}{\Delta e \cdot \cos \theta_e \cdot \Omega_p \cdot B}.$$

The measuring device provides the following value:

$$P = \mathcal{M}L = \int_{\Omega} \int_{S} L(\vec{y}, \omega) \cdot W^{e}(\vec{y}, \omega) \cdot \cos \theta \, d\vec{y} d\omega.$$
(2.33)

Applying equation (2.2), we can use the following substitutions:

$$d\omega = d\vec{e} \cdot \frac{\cos \theta_e}{|\vec{y} - \vec{e}|^2}, \quad d\vec{y} = \frac{|\vec{y} - \vec{e}|^2}{\cos \theta} \cdot d\omega_p,$$

where $d\vec{e}$ is a differential area on the pupil and $|\vec{y} - \vec{e}|$ is the distance between the pupil and the visible surface. Substituting these and taking advantage that the pupil Δe is small, we obtain

$$P = \int_{\Omega_p} \int_{\Delta_e} L(h(\vec{e}, -\omega_p), \omega_p) \cdot C \cdot \frac{\cos \theta_e}{|\vec{y} - \vec{e}|^2} \cdot \frac{|\vec{y} - \vec{e}|^2}{\cos \theta} \cdot \cos \theta \, d\vec{e} d\omega_p = \int_{\Omega_p} \int_{\Delta_e} L(h(\vec{e}, -\omega_p), \omega_p) \cdot C \cdot \cos \theta_e \, d\vec{e} d\omega_p \approx \int_{\Omega_p} L(h(\vec{e}, -\omega_p), \omega_p) \cdot C \cdot \cos \theta_e \cdot \Delta e \, d\omega_p = \int_{\Omega_p} L(h(\vec{e}, -\omega_p), \omega_p) \cdot \frac{1}{\Omega_p \cdot B} \, d\omega_p,$$
(2.34)

where \vec{vy} is the position of the very small pupil, and ω_p is the direction which points from \vec{y} to the eye position.

Note that the measured value is independent of both the distance and the orientation of the visible surface! This is explained by the fact that as we move farther from a surface, although the power coming from a unit surface area decreases with the square of the distance, the area that can be visible in a given solid angle increases with the same speed. The two factors compensate each other. Similarly, when we look at a surface from a non perpendicular direction, the radiance is weighted by the cosine of this angle in the power, but the surface area that can be seen is inversely proportional to the cosine of the same angle.



Figure 2.8: Evaluation of the measured value as an integral on the pixel (left) and on the surface (right)

Since Ω_p is the collection of those directions that go through the pixel, the measured value can also be expressed as an integral over the pixel area S_p (figure 2.8). Let \vec{p} be the running point on the pixel, θ_p be the angle between the pixel normal and vector $\vec{p} - e\vec{y}e$, and θ_{pix} be equal to θ_p in the center of the pixel. Since $|\vec{p} - e\vec{y}e| = f/\cos\theta_p$ where f is the **focal distance**, i.e. the distance between the eye and the plane of the window, we can establish the following correspondance:

$$d\omega_p = \frac{d\vec{p} \cdot \cos\theta_p}{|\vec{p} - e\vec{y}e|^2} = \frac{d\vec{p} \cdot \cos^3\theta_p}{f^2}.$$
(2.35)

Substituting this to equation (2.34), the measured value becomes

$$P = \int_{S_p} L(h(\vec{p}, -\omega_{\vec{p}}), \omega_{\vec{p}}) \cdot \frac{1}{\Omega_p \cdot B} \cdot \frac{\cos^3 \theta_p}{f^2} d\vec{p}.$$
(2.36)

Let us introduce **camera parameter** $c(\vec{p})$ by

$$c(\vec{p}) = \frac{S_p}{\Omega_p \cdot B} \cdot \frac{\cos^3 \theta_p}{f^2}.$$
(2.37)

Since

$$\Omega_p = \int_{\Omega_p} d\omega_p = \int_{S_p} \frac{\cos^3 \theta_p}{f^2} d\vec{p} \approx \frac{S_p \cdot \cos^3 \theta_{\text{pix}}}{f^2}, \qquad (2.38)$$

the camera parameter can also be expressed in the following exact and approximating forms:

$$c(\vec{p}) = \frac{S_p \cdot \cos^3 \theta_p}{B \cdot \int\limits_{S_p} \cos^3 \theta_p \, d\vec{p}} \approx \frac{\cos^3 \theta_p}{B \cdot \cos \theta_{\text{pix}}^3} \approx \frac{1}{B}.$$
(2.39)

The approximations are accurate if the pixel is small compared to the focal distance. If image space **filtering** is also applied, then S_p may be greater than the pixel and the camera parameter is not necessarily constant but follows the shape of a reconstruction filter [SKe95, SK95]. Summarizing the measured value is an integral on the pixel:

$$P = \int_{S_p} L(h(\vec{p}, -\omega_{\vec{p}}), \omega_{\vec{p}}) \cdot \frac{c(\vec{p})}{S_p} d\vec{p}.$$
(2.40)

Equation (2.34) can also be evaluated on object surfaces. Using the notations of figure 2.8, we can obtain:

$$d\omega_p = \frac{d\vec{y} \cdot \cos\theta}{|\vec{y} - \vec{eye}|^2} \tag{2.41}$$

where \vec{y} is a surface point visible in the pixel. Let the indicator function of those points that are visible in Ω_p be $V_{\Omega_p}(\vec{y})$. The measured value is then:

$$P = \int_{S} L(\vec{y}, \omega_{\vec{y} \to e\vec{y}e}) \cdot V_{\Omega_p}(\vec{y}) \cdot \frac{1}{\Omega_p \cdot B} \cdot \frac{\cos \theta}{|\vec{y} - e\vec{y}e|^2} d\vec{y}.$$
 (2.42)

Let us introduce surface dependent camera parameter $g(\vec{y})$ by

$$g(\vec{y}) = \frac{1}{\Omega_p \cdot B \cdot |\vec{y} - \vec{\text{eye}}|^2} = \frac{f^2}{B \cdot |\vec{y} - \vec{\text{eye}}|^2 \cdot \int_{S_p} \cos^3 \theta_p \, d\vec{p}} \approx \frac{f^2}{B \cdot |\vec{y} - \vec{\text{eye}}|^2 \cdot S_p \cdot \cos^3 \theta_{\text{pix}}} \quad (2.43)$$

using formula (2.38) that expresses Ω_p . Thus the final form of the measured value is:

$$P = \int_{S} L(\vec{y}, \omega_{\vec{y} \to e\vec{y}e}) \cdot V_{\Omega_p}(\vec{y}) \cdot g(\vec{y}) \cdot \cos \theta \ d\vec{y}.$$
(2.44)

Comparing this formula with equation (2.33), the following sensitivity function can be derived:

$$W^{e}(\vec{y},\omega) = \begin{cases} \delta(\omega - \omega_{\vec{y} \to e\vec{y}e}) \cdot g(\vec{y}) & \text{if } \vec{y} \text{ is visible through the pixel,} \\ 0 & \text{otherwise.} \end{cases}$$
(2.45)

2.6 Numerical solution of the rendering equation

2.6.1 Error measures for numeric techniques

When solving the rendering problem numerically, the required radiance, potential or measured power can only be approximated. To characterize the accuracy of the approximation, **error measures** are needed. The errors of radiance or potential are called **object-space error measures**. The error of the measured power is called **image-space error measure**. A good error measure should be invariant to changes that modify neither the object space nor the camera. Particularly, the error should not change when a surface is subdivided to two surfaces or the complete scene is scaled — i.e. it should be tessellation and

scaling independent —, or when a pixel is subdivided to increase the image resolution — i.e. it should be resolution independent.

Tessellation and scaling independent measures can be defined if the norm of the error function is divided by the total surface area. The **norm** [Mát81, ST93] is selected from the following possibilities:

$$||f||_{1} = \int_{\Omega} \int_{S} |f(\vec{x},\omega)| \, d\vec{x} \, d\vec{\omega}, \quad ||f||_{2} = \sqrt{\int_{\Omega} \int_{S} (f(\vec{x},\omega))^{2} \, d\vec{x} \, d\vec{\omega}}, \quad ||f||_{\infty} = \max_{\vec{x},\omega} |f(\vec{x},\omega)|.$$
(2.46)

Using any of these, if L is the calculated radiance and \tilde{L} is the real radiance, then an absolute object-space error ϵ_a and a relative object-space error ϵ_r are

$$\epsilon_a = \frac{||L - \tilde{L}||}{S}, \quad \epsilon_r = \frac{||L - \tilde{L}||}{||\tilde{L}||}.$$
(2.47)

Image space norms are based on the powers that go through different pixels $p = 1, 2, ..., N_P$.

$$||\Phi||_{1} = \sum_{p=1}^{N_{P}} |\Phi_{p}|, \quad ||\Phi||_{2} = \sqrt{\sum_{p=1}^{N_{P}} |\Phi_{p}|^{2}}, \quad ||\Phi||_{\infty} = \max_{p} |\Phi_{p}|.$$
(2.48)

Resolution independent measures divide the total error of pixels by the number of pixels (P). If the calculated and real powers that go through pixel p are Φ_p and $\tilde{\Phi}_p$, respectively, then an absolute and a relative image space error measures are

$$\epsilon_a = \frac{||\Phi - \tilde{\Phi}||}{N_P}, \quad \epsilon_r = \frac{||\Phi - \tilde{\Phi}||}{||\tilde{\Phi}||}.$$
(2.49)

2.6.2 Properties of the rendering equation

The integral operators of both the rendering and the potential equations are **contractions**. This statement can be proven from the fact that for physically plausible models, the albedo is less than 1. Here, the ∞ -norm is used:

$$\|\mathcal{T}L\|_{\infty} = \max|\mathcal{T}L| \le \max|L| \cdot \max|\mathcal{T}1| = \max|L| \cdot \max|a(\vec{x},\omega)| = \|L\|_{\infty} \cdot \max|a(\vec{x},\omega)|.$$

For the potential, similar results can be obtained:

$$\begin{aligned} \|\mathcal{T}'W\|_{\infty} &= \max |\mathcal{T}'W| \le \max |W| \cdot \max |\mathcal{T}'1| = \max |W| \cdot \max |a(h(\vec{y}, \omega'), \omega')| = \\ \|W\|_{\infty} \cdot \max |a(\vec{x}, \omega)|. \end{aligned}$$

2.7 Classification of the solution techniques

In order to find the color of a pixel, the radiance has to be evaluated for that surface which is visible through the given pixel. The determination of this surface is called the **hidden surface problem** or **visibility calculation**. Having solved the visibility problem, the surface visible in the given pixel is known, and the radiance may be calculated on the representative wavelengths by the rendering equation.

Due to multiple reflections of light beams, the calculation of the radiance of the light leaving a point on a surface in a given direction requires the radiances of other surfaces visible from this point, which generates new visibility and shading problems to solve (figure 2.9). To calculate those radiances, other surfaces should be evaluated, and our original point on the given surface might contribute to those radiances. As a consequence of that, the rendering equation has complicated coupling between its left and right sides, making the evaluation difficult.

There are three general and widely accepted approaches to attack this coupling:



Figure 2.9: Multiple reflections of light

1. Local illumination methods

Local illumination methods take a drastic approach and simplify or disregard completely all the terms causing problems. The unknown radiance inside the integral of the rendering equation is replaced by some approximation of the emission function. Formally, these methods evaluate the following simplified rendering equation to obtain the radiance:

$$L(\vec{x},\omega) = L^{e}(\vec{x},\omega) + \int_{\Omega} L_{\text{lightsource}}(h(\vec{x},-\omega'),\omega') \cdot f_{r}(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega', \qquad (2.50)$$

where $L_{\text{lightsource}}$ may be a simplification of the emission function L^e . Abstract lightsources, such as point or directional lightsources are preferred here, since their radiance is a Dirac-delta like function, which simplifies the integral of equation (2.50) to a sum.

These methods take into account only a single reflection of the light coming from the abstract lightsources. Ideal mirrors and refracting objects cannot be rendered with these methods.

2. Recursive ray-tracing

Another alternative is to eliminate from the rendering equation those energy contributions which cause the difficulties, and thus give ourselves a simpler problem to solve. For example, if limited level, say n, coupling caused by ideal reflection and refraction were allowed, and we were to ignore the other non-ideal components coming from non-abstract lightsources, then the number of surface points which would need to be evaluated to calculate a pixel color can be kept under control. Since the illumination formula contains two terms regarding the coherent components (reflective and refracting lights), the maximum number of surfaces involved in the color calculation of a pixel is two to the power of the given depth, i.e. 2^n . An implementation of this approach is called **visibility ray-tracing** which uses the following illumination formula:

$$L(\vec{x},\omega) = L^{e}(\vec{x},\omega) + \int_{\Omega} L_{\text{light source}}(h(\vec{x},-\omega'),\omega') \cdot f_{r}(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega' + k_{r}(\omega_{r},\vec{x},\omega) \cdot L(h(\vec{x},-\omega_{r}),\omega_{r}) + k_{t}(\omega_{t},\vec{x},\omega) \cdot L(h(\vec{x},-\omega_{t}),\omega_{t}),$$
(2.51)

where ω_r and ω_t are the ideal reflection and refraction directions, and k_r and k_t are the integrated BRDF components representing the ideal reflection and refraction, respectively.

3. Global illumination solution

Local illumination and recursive ray-tracing apply brutal simplifications and thus provide physically inaccurate images. Thus these methods cannot be reliably used in engineering applications, such as in lighting design. These application areas require the rendering equation to be solved without relevant simplifications, that leads to the family of **global illumination** methods. Global illumination methods rely on numerical techniques to resolve the coupling in the rendering or potential equations and to solve the integral equation without unacceptable simplifications.



local illumination method

local illumination method with shadow computation



recursive ray-tracing

global illumination method

Figure 2.10: Comparison of local illumination method, recursive ray-tracing and global illumination method. Ambient light was also added when the local illumination method and recursive ray-tracing were computed. The images have been generated by a path tracing program [SK99d]. Note that local illumination methods cannot render ideal reflection or refraction, thus there are no transparent and mirror like objects. Recursive ray-tracing, on the other hand, is unable to follow multiple non-ideal reflections, thus the illumination diffusely reflected from the back wall is missing on the spheres. The running times are 90 seconds, 95 seconds, 135 seconds, 9 hours, respectively, which clearly shows the price we should pay for a physically accurate simulation. The three different approaches represent three levels of the compromise between image generation speed and quality. By ignoring more and more terms in the illumination formula, its calculation can be speeded up, but the result inevitably becomes more and more artificial.

The ignoration of the terms, however, violates the energy equilibrium, and causes portions of objects to come out extremely dark, sometimes unexpectedly so. These artifacts can be reduced by reintroducing the ignored terms in simplified form, called **ambient light**. The ambient light represents the ignored energy contribution in such a way as to satisfy the energy equilibrium. Since this ignored part is not calculated, nothing can be said of its positional and directional variation, hence it is supposed to be constant in all directions and everywhere in the 3D space. From this aspect, the role of ambient light also shows the quality of the shading algorithm. The more important a role it has, the poorer quality picture it will generate.

In order to formally express the capabilities of a certain algorithm, Heckbert has developed a notation that is based on the regular expressions originally proposed for the syntax of formal languages [Hec91]. The elements of the notation are as follows:

- E is the eye,
- L is the lightsource,
- D is a non-ideal i.e. incoherent reflection or refraction,
- S is an ideal reflection or refraction,
- * is the sign of iteration, that is, it can mean $0, 1, 2, \ldots$ applications of the given operator,
- [] represents optionality,
- means selection.

A global illumination algorithm is expected to model all types of lightpaths, that is, it must have $L[D|S]^*E$ type. Visibility ray-tracing allows multiple steps only for the ideal reflection or refraction, thus it can model only $L[DS^*]E$ path. Finally, local illumination models simulate only a single non-ideal reflection and fall into the L[D]E category.

Chapter 3

Solution strategies for the global illumination problem

Since the rendering or the potential equations contain the unknown radiance function both inside and outside the integral, in order to express the solution, this coupling should be resolved. The possible solution techniques fall into one of the following three categories: **inversion**, **expansion** and **iteration**.

Operator \mathcal{T} represents light-surface interaction, thus each of its application generates a higher-bounce estimate of the light transport (or alternatively \mathcal{T}' represents potential-surface interaction). For physically plausible optical material models, a reflection or refraction always decreases the total energy, thus the integral operator is always a contraction. However, when the transport is evaluated numerically, computation errors may pose instability problems if the scene is highly reflective. As we shall see, expansion and iteration exploit the contractive property of the transport operator, but inversion does not.

3.1 Inversion

Inversion groups the terms that contain the unknown function on the same side of the equation and applies formally an inversion operation:

$$(1-\mathcal{T})L = L^e \implies L = (1-\mathcal{T})^{-1}L^e.$$
(3.1)

Thus the measured power is

$$\mathcal{M}L = \mathcal{M}(1 - \mathcal{T})^{-1}L^e.$$
(3.2)

However, since \mathcal{T} is infinite dimensional, it cannot be inverted in closed form. Thus it should be approximated by a finite dimensional mapping, that is usually given as a matrix. This kind of approximation is provided by **finite-element methods** that project the problem into a finite dimensional function space, and approximate the solution here. This projection converts the original integral equation to a system of linear equations, which can be inverted, for example, by Gaussian elimination method. This approach was used in early radiosity methods, but have been dropped due to the cubic time complexity and the numerical instability of the matrix inversion.

Inversion has a unique property that is missing in the other two methods. Its efficiency does not depend on the contractivity of the integral operator, neither does it even require the integral operator to be a contraction.

3.2 Expansion

Expansion techniques eliminate the coupling by obtaining the solution in the form of an infinite Neumann series.

3.2.1 Expansion of the rendering equation: gathering walks

Substituting the right side's L by $L^e + TL$, which is obviously L according to the equation, we get:

$$L = L^e + \mathcal{T}L = L^e + \mathcal{T}(L^e + \mathcal{T}L) = L^e + \mathcal{T}L^e + \mathcal{T}^2L.$$
(3.3)

Repeating this step n times, the original equation can be expanded into a Neumann series:

$$L = \sum_{i=0}^{n} \mathcal{T}^{i} L^{e} + \mathcal{T}^{n+1} L.$$
 (3.4)

If integral operator \mathcal{T} is a contraction, then $\lim_{n\to\infty} \mathcal{T}^{n+1}L = 0$, thus

$$L = \sum_{i=0}^{\infty} \mathcal{T}^i L^e.$$
(3.5)

The measured power is

$$\mathcal{M}L = \sum_{i=0}^{\infty} \mathcal{M}\mathcal{T}^i L^e.$$
(3.6)

The terms of this infinite Neumann series have intuitive meaning as well: $\mathcal{MT}^0 L^e = \mathcal{M}L^e$ comes from the emission, $\mathcal{MT}^1 L^e$ comes from a single reflection, $\mathcal{MT}^2 L^e$ from two reflections, etc.



Figure 3.1: The integrand of \mathcal{MT}^2L^e *is a two-step gathering walk*

In order to understand how this can be used to determine the power going through a single pixel, let us examine the structure of $\mathcal{MT}^i L^e$ as a single multi-dimensional integral for the i = 2 case:

$$\mathcal{M}(\mathcal{T}^2 L^e) = \int_{S_p} \int_{\Omega_1'} \int_{\Omega_2'} \frac{c(\vec{p})}{S_p} \cdot w_1(\vec{x}_1) \cdot w_2(\vec{x}_2) \cdot L^e(\vec{x}_3, \omega_2') \ d\omega_2' d\omega_1' d\vec{p}.$$
(3.7)

where S_p is the pixel area, $c(\vec{p})$ is the camera parameter, \vec{p} is a running point on this pixel, ω'_1 and ω'_2 are the directions of the first and second reflections, respectively (figure 3.1) and

$$\vec{x}_{1} = h(\vec{p}, -\omega_{\vec{p}}),
\vec{x}_{2} = h(\vec{x}_{1}, -\omega'_{1}),
\vec{x}_{3} = h(\vec{x}_{2}, -\omega'_{2}) = h(h(\vec{x}_{1}, -\omega'_{1}), -\omega'_{2}),$$
(3.8)

and the weights are

Thus to evaluate the integrand at point $(\vec{p}, \omega'_1, \omega'_2)$, the following algorithm must be executed:

- 1. Point $\vec{x}_1 = h(\vec{p}, -\omega_{\vec{p}})$ that is visible through the point \vec{p} of the pixel from the eye should be found. This can be done by sending a ray from the eye into the direction of \vec{p} and identifying the surface that is first intersected.
- 2. Surface point $\vec{x}_2 = h(\vec{x}_1, -\omega'_1)$ that is the point which is visible from \vec{x}_1 at direction $-\omega'_1$ must be determined. This can be done by sending another ray from \vec{x}_1 into direction $-\omega'_1$ and identifying the surface that is first intersected.
- 3. Surface point $\vec{x}_3 = h(h(\vec{x}_1, -\omega'_1), -\omega'_2)$ that is the point visible from \vec{x}_2 at direction $-\omega'_2$ is identified. This means the continuation of the ray tracing at direction $-\omega'_2$.
- 4. The emission intensity of the surface at \vec{x}_3 in the direction of ω'_2 is obtained and is multiplied with the cosine terms and the BRDFs of the two reflections.

This algorithm can easily be generalized for arbitrary number of reflections. A ray is emanated recursively from the visible point at direction ω'_1 then from the found surface at ω'_2 , etc. until ω'_n . The emission intensity at the end of the walk is read and multiplied by the BRDFs and the cosine terms of the stages of the walk. These walks provide the value of the integrand at "point" $\vec{p}, \omega'_1, \omega'_2, \ldots, \omega'_n$. Note that a single walk of length n can be used to estimate the 1-bounce, 2-bounce, etc. n-bounce transfer simultaneously, if the emission is transferred not only from the last visited point but from all visited points.

The presented walking technique starts at the eye and **gathers** the illumination encountered during the walk. The gathered illumination is attenuated according to the cosine weighted BRDFs of the path.

So far, we have examined the structure of the terms of the Neumann series as a single multidimensional integral. Alternatively, it can also be considered as recursively evaluating many directional integrals. For example, the two-bounce transfer is:

$$\mathcal{MT}^{2}L^{e} = \int_{S_{p}} w_{0} \cdot \left[\int_{\Omega'_{1}} w_{1} \cdot \left[\int_{\Omega'_{2}} w_{2} \cdot L^{e} d\omega'_{2} \right] d\omega'_{1} \right] d\vec{p}.$$
(3.10)

In order to estimate the outer integral of variable \vec{p} , the integrand is needed in sample point \vec{p} . This, in turn, requires the estimation of the integral of variable ω'_1 at \vec{p} , which recursively needs again the approximation of the integral of variable ω'_2 at (\vec{p}, ω'_1) .

If the same number — say m — of sample points are used for each integral quadrature, then this recursive approach will use m points for the 1-bounce transfer, m^2 for the two-bounces, m^3 for the three-bounces, etc. This kind of subdivision of paths is called **splitting** [AK90]. Splitting becomes prohibitive for high-order reflections and is not even worth doing because of the contractive property of the integral operator. Due to the contraction, the contribution of higher-order bounces is less, thus it is not very wise to compute them as accurately as low-order bounces.

3.2.2 Expansion of the potential equation: shooting walks

The potential equation can also be expanded into a Neumann series similarly to the rendering equation:

$$W = \sum_{i=0}^{\infty} \mathcal{T}^{i} W^e, \qquad (3.11)$$

which results in the following measured power

$$\mathcal{M}'W = \sum_{i=0}^{\infty} \mathcal{M}'\mathcal{T}'^i W^e.$$
(3.12)



Figure 3.2: The integrand of $\mathcal{M}'\mathcal{T}'^2W^e$ is a two-step shooting walk

 $\mathcal{M}'W^e$ is the power measured by the device from direct emission, $\mathcal{M}'\mathcal{T}'W^e$ is the power after a single reflection, $\mathcal{M}'\mathcal{T}'^2W^e$ is after two reflections, etc.

Let us again consider the structure of $\mathcal{M}'\mathcal{T}'^2W^e$:

$$\mathcal{M}'\mathcal{T}'^{2}W^{e} = \int_{S} \int_{\Omega_{1}} \int_{\Omega_{2}} \int_{\Omega_{3}} \int_{\Omega_{3}} L^{e}(\vec{y}_{1}, \omega_{1}) \cdot w_{0}(\vec{y}_{1}) \cdot w_{1}(\vec{y}_{2}) \cdot w_{2}(\vec{y}_{3}) \cdot W^{e}(\vec{y}_{3}, \omega_{3}) \ d\omega_{3}d\omega_{2}d\omega_{1}d\vec{y}_{1}, \quad (3.13)$$

where

$$\vec{y}_2 = h(\vec{y}_1, \omega_1), \vec{y}_3 = h(\vec{y}_2, \omega_2) = h(h(\vec{y}_1, \omega_1), \omega_2)$$
(3.14)

and the weights are

$$w_0 = \cos \theta_1,$$

$$w_1 = f_r(\omega_1, \vec{y}_2, \omega_2) \cdot \cos \theta_2,$$

$$w_2 = f_r(\omega_2, \vec{y}_3, \omega_3) \cdot \cos \theta_3.$$
(3.15)

Substituting the measuring function of the pin-hole camera (equation (2.45)), we can conclude that $\mathcal{M}'\mathcal{T}'^2W^e$ can be non-zero only if \vec{y}_3 is visible through the given pixel, and the integral over Ω_3 is simplified to a single value by the Dirac-delta function, thus the final form of the two-bounce reflection is:

$$\int_{S} \int_{\Omega_{1}} \int_{\Omega_{2}} L^{e}(\vec{y}_{1}, \omega_{1}) \cdot w_{0}(\vec{y}_{1}) \cdot w_{1}(\vec{y}_{2}) \cdot w_{2}(\vec{y}_{3}) \cdot g(\vec{y}_{3}) \, d\omega_{2} d\omega_{1} d\vec{y}_{1}.$$
(3.16)

if \vec{y}_3 is visible though the given pixel and 0 otherwise, where $g(\vec{y}_3)$ is the surface dependent camera parameter.

Thus to evaluate the integrand at point $(\vec{y}_1, \omega_1, \omega_2)$, the following algorithm must be executed:

- 1. The cosine weighted emission of point \vec{y}_1 in direction ω_1 is computed. Surface point $\vec{y}_2 = h(\vec{y}_1, \omega_1)$ that is the point which is visible from \vec{y}_1 at direction ω_1 must be determined. This can be done by sending a ray from \vec{y}_1 into direction ω_1 and identifying the surface that is first intersected. This point "receives" the computed cosine weighted emission.
- 2. Surface point $\vec{y}_3 = h(h(\vec{y}_1, \omega_1), \omega_2)$ that is the point visible from \vec{y}_2 at direction ω_2 is identified. This means the continuation of the ray tracing at direction ω_2 . The emission is weighted again by the local BRDF and the cosine of the normal and the outgoing direction.

3. It is determined whether or not this point \vec{y}_3 is visible from the eye, and through which pixel. The contribution to this pixel is obtained by weighting the transferred emission by the local BRDF, cosine angle and the surface dependent camera parameter.

This type of walk, called **shooting**, starts at a known point \vec{y}_1 of a lightsource and simulates the photon reflection for a few times and finally arrives at a pixel whose radiance this walk contributes to.

Note that in gathering walks the BRDF is multiplied with the cosine of the angle between the normal and the incoming direction, while in shooting walks with the cosine of the angle between the normal and the outgoing direction.

3.2.3 Merits and disadvantages of expansion methods

The main problem of expansion techniques is that they require the evaluation of very high dimensional integrals that appear as terms in the infinite series. Practical implementations usually truncate the infinite Neumann series, which introduces some bias, or stop the walks randomly, which significantly reduces the samples of higher order interreflections. These can result in visible artifacts for highly reflective scenes.

On the other hand, expansion methods also have an important advantage. Namely, they do not require temporary representations of the complete radiance function, thus do not necessitate finite-element approximations. Consequently, these algorithms can work with the original geometry without tessellating the surfaces to planar polygons.

Expansion techniques generate random walks independently. It can be an advantage, since these algorithms are suitable for **parallel computing**. However, it also means that these methods "forget" the previous history of walks, and they cannot reuse the visibility information gathered when computing the previous walks, thus they are not as fast as they could be.

3.3 Iteration

Iteration techniques realize that the solution of the rendering equation is the fixed point of the following iteration scheme

$$L_n = L^e + \mathcal{T}L_{n-1},\tag{3.17}$$

thus if operator \mathcal{T} is a contraction, then this scheme will converge to the solution from any initial function L_0 . The measured power can be obtained as a limiting value

$$\mathcal{M}L = \lim_{n \to \infty} \mathcal{M}L_n. \tag{3.18}$$

In order to store the approximating functions L_n , usually **finite-element methods** are applied, as for example, in **diffuse radiosity** [SC94], or in non-diffuse radiosity using **partitioned hemisphere** [ICG86], **directional distributions** [SAWG91] or **illumination networks** [BF89].

There are two critical problems here. On the one hand, since the domain of L_n is 4 dimensional and L_n has usually high variation, an accurate finite-element approximation requires very many basis functions, which, in turn, need a lot of storage space. Although **hierarchical methods** [HSA91, AH93], **wavelet** or **multiresolution methods** [CSSD96, SGCH94] and **clustering** [SDS95, CLSS97, SPS98] can help, the memory requirements are still prohibitive for complex scenes. This problem is less painful for the diffuse case since here the domain of the radiance is only 2 dimensional.

On the other hand, when finite element techniques are applied, operator \mathcal{T} is only approximated, which introduces some non-negligible error in each step. If the contraction ratio of the operator is λ , then the total accumulated error will be approximately $1/(1 - \lambda)$ times the error of a single step [SKFNC97]. For highly reflective scenes — i.e. when λ is close to 1 — the iteration is slow and the result is inaccurate if the approximation of the operator is not very precise. Very accurate approximations of the transport operator, however, require a lot of computation time and storage space.

In the diffuse radiosity setting several methods have been proposed to improve the quality of the iteration. For example, we can use Chebyshev iteration instead of the Jacobi or the Gauss-Seidel method

for such ill conditioned systems [BBS96]. On the other hand, realizing that the crucial part of designing such an the algorithm is finding a good and "small" approximation of the transport operator, the method called **well-distributed ray-sets** [NNB97, BNN⁺98] proposes the adaptive approximation of the transport operator. This approximation is a set of rays selected carefully taking into account the important patches and directions. In [BNN⁺98], the adaptation of the discrete transport operator is extended to include patch subdivision as well, to incorporate the concepts of **hierarchical radiosity** [HSA91]. The adaptation strategy is to refine the discrete approximation (by adding more rays to the set), when the iteration with the coarse approximation is already stabilized. Since the discrete approximation of the transport operator is not constant but gets finer in subsequent phases, the error accumulation problem can be controlled but is not eliminated.

This thesis proposes a new method called the **stochastic iteration** to attack both the problem of prohibitive memory requirements and the problem of error accumulation.

Compared to expansion techniques, iteration has both advantages and disadvantages. Its important advantage is that it can potentially reuse all the information gained in previous computation steps and can exploit the coherence of the radiance function, thus iteration is expected to be faster than expansion. Iteration can also be seen as a single infinite length random walk. If implemented carefully, iteration does not reduce the number of estimates for higher order interreflections, thus it is more robust when rendering highly reflective scenes than expansion.

The property that iteration requires tessellation and finite-element representation is usually considered as a disadvantage. And indeed, sharp shadows and highlights on highly specular materials can be incorrectly rendered and light-leaks may appear, not to mention the unnecessary increase of the complexity of the scene description (think about, for example, the definition of an original and a tessellated sphere). However, finite-element representation can also provide smoothing during all stages of rendering, which results in more visually pleasing and dot-noise free images. Summarizing, iteration is the better option if the scene is not highly specular.

3.3.1 Analysis of the iteration

In order to find necessary conditions for the convergence, let us examine two subsequent iteration steps:

$$L_{n} = L^{e} + \mathcal{T}L_{n-1},$$

$$L_{n-1} = L^{e} + \mathcal{T}L_{n-2}.$$
(3.19)

Substracting the two equations and assuming that $L_0 = 0$, we can obtain:

$$L_n - L_{n-1} = \mathcal{T}(L_{n-1} - L_{n-2}) = \mathcal{T}^{n-1}(L_1 - L_0) = \mathcal{T}^{n-1}L^e.$$
(3.20)

If operator \mathcal{T} is a contraction, that is if

$$||\mathcal{T}L|| < \lambda \cdot ||L||, \quad \lambda < 1, \tag{3.21}$$

with some function norm, then

$$||L_n - L_{n-1}|| = ||\mathcal{T}^{n-1}L^e|| < \lambda^{n-1} \cdot ||L^e||.$$
(3.22)

Thus iteration converges to the real solution at least with the speed of a geometric series. The contraction ratio λ depends on both the geometry and the average reflectivity of the scene. From the point of view of the geometry, the contraction ratio is maximum if the environment is closed, when λ corresponds to the average albedo.

Error caused by the approximation of the transport operator

In practice operator \mathcal{T} cannot be evaluated exactly, which introduces some error in each step of the iteration. The errors may accumulate in the final solution. This section examines this phenomenon.

Assume that an operator \mathcal{T}^* which is only an approximation of \mathcal{T} is used in the iteration. Suppose that both operators are contractions, thus both iterations will converge from any initial function.

Let the radiance after n iterations of operator \mathcal{T}^* be L_n and let us assume that the iteration starts at the solution of the exact equation, that is $L_0 = L$ (since the iteration converges to the same limit from any initial function, this assumption can be made). The solution of the approximated equation is $L^* = L_{\infty}$. The error at step n is

$$||L_n - L|| = ||L_n - L_{n-1} + L_{n-1} - \dots + L_1 - L|| \le \sum_{i=1}^n ||L_i - L_{i-1}||.$$
(3.23)

Since if i > 1, then

$$||L_i - L_{i-1}|| = ||\mathcal{T}^* L_{i-1} - \mathcal{T}^* L_{i-2}|| = ||\mathcal{T}^* (L_{i-1} - L_{i-2})|| \le \lambda \cdot ||L_{i-1} - L_{i-2}|| \le \lambda^{i-1} \cdot ||L_1 - L_0||$$

we have

we obtain

$$\sum_{i=1}^{n} ||L_i - L_{i-1}|| \le ||L_1 - L_0|| \cdot (1 + \lambda + \lambda^2 + \dots \lambda^{n-1}).$$
(3.24)

Letting n go to infinity, we obtain the error between the fixed points of the original and the approximated equations

$$||L^* - L|| \le ||L_1 - L_0|| \cdot (1 + \lambda + \lambda^2 + \ldots) = \frac{||L_1 - L_0||}{1 - \lambda}.$$
(3.25)

On the other hand, subtracting the real solution defined as the fixed point of the exact equation $L = L^e + TL$ from the first iteration of the approximated operator

 $L_1 = L^e + \mathcal{T}^* L,$

$$L_1 - L = \mathcal{T}^* L - \mathcal{T} L. \tag{3.26}$$

Thus the final error formula is

$$||L^* - L|| \le \frac{||\mathcal{T}^*L - \mathcal{T}L||}{1 - \lambda}.$$
 (3.27)

3.4 Analytical solution of the rendering equation

In order to test the error and convergence of global illumination algorithms, we need scenes for which the exact solution is known. Generally, there are two alternatives. Either we use simple scenes for which analytical solution is possible, or numerical methods are applied using so many samples that the approximate solution can be accepted as a reference.

In order to find analytically solvable scenes, we use a reverse approach. Normally, a scene is analyzed to find the radiance distribution. Now, we start with a prescribed radiance distribution and search for scenes where the radiance would be identical to the given radiance. Two special cases are examined. In the first case we are looking for scenes where the radiance is constant everywhere and at every direction. In the second case we establish criteria for making only the reflected radiance constant.

3.4.1 Scenes with constant radiance

Suppose that the constant radiance is L and also that the scene is a **closed environment**, i.e. looking at any direction we can see a surface, thus the incoming radiance is also constant. Substituting this to the rendering equation we get:

$$\tilde{L} = L^e(\vec{x},\omega) + (\mathcal{T}\tilde{L})(\vec{x},\omega) = L^e(\vec{x},\omega) + \tilde{L} \cdot (\mathcal{T}1)(\vec{x},\omega) = L^e(\vec{x},\omega) + \tilde{L} \cdot a(\vec{x},\omega),$$
(3.28)

since the albedo is the response to homogeneous unit illumination. According to this equation, the radiance will be constant for scenes of arbitrary geometry and of arbitrary emission function if the albedo function is: $\mathbf{r}_{a}(\mathbf{r}_{a})$

$$a(\vec{x},\omega) = 1 - \frac{L^e(\vec{x},\omega)}{\tilde{L}}.$$
(3.29)

If L^e is constant, then the required albedo will also be constant. In this special case — called the homogeneous diffuse environment — the geometry is arbitrary, but all surfaces have the same diffuse reflection and emission [Shi91b, Kel95].

3.4.2 Scenes with constant reflected radiance

Let us assume that the reflected radiance — i.e. the radiance without the emission — is L_r and assume again that the scene is closed. Substituting this into the rendering equation, we obtain:

$$L(\vec{x},\omega) = L^{e}(\vec{x},\omega) + \dot{L}_{r} = L^{e}(\vec{x},\omega) + (\mathcal{T}(L^{e} + \dot{L}_{r}))(\vec{x},\omega) = L^{e}(\vec{x},\omega) + (\mathcal{T}L^{e})(\vec{x},\omega) + \dot{L}_{r} \cdot a(\vec{x},\omega).$$
(3.30)

This imposes the following constraint on the albedo function:

$$a(\vec{x},\omega) = 1 - \frac{(\mathcal{T}L^e)(\vec{x},\omega)}{\tilde{L}_r}.$$
(3.31)

Note that if the reflected radiance is constant, then the albedo is also constant, which is guaranteed if the BRDF is diffuse and has the same f_r value everywhere. An appropriate geometry, which makes $\mathcal{T}L^e$ constant for arbitrary diffuse emission if the BRDF is diffuse and constant, is the internal surface of the sphere as proposed originally by [HMF98]. Formally, let the scene be an inner surface S of a sphere of radius R, the BRDF be constant $f_r = a/\pi$ and the emission be diffuse and defined by $L^e(\vec{x})$. Using the

$$d\omega' = \frac{\cos\theta_{\vec{y}} \cdot d\vec{y}}{|\vec{y} - \vec{x}|^2}$$

substitution for the solid angle (equation (2.2)), we obtain the following form of the light transport operator:

$$(\mathcal{T}L^e)(\vec{x},\omega) = \int_S f_r \cdot \cos\theta_{\vec{x}} \cdot L^e(\vec{y}) \cdot \frac{\cos\theta_{\vec{y}}}{|\vec{y} - \vec{x}|^2} \, d\vec{y}.$$
(3.32)



Figure 3.3: Geometry of the reference scene

Looking at figure 3.3, we can see that inside a sphere $\cos \theta_{\vec{x}} = \cos \theta_{\vec{y}} = \frac{|\vec{y} - \vec{x}|}{2R}$, thus we can obtain the following final form:

$$(\mathcal{T}L^{e})(\vec{x},\omega) = \int_{S} f_{r} \cdot L^{e}(\vec{y}) \cdot \frac{\cos\theta_{\vec{x}} \cdot \cos\theta_{\vec{y}}}{|\vec{y} - \vec{x}|^{2}} \, d\vec{y} = \frac{f_{r}}{4R^{2}} \cdot \int_{S} L^{e}(\vec{y}) \, d\vec{y} = a \cdot \frac{\int_{S} L^{e}(\vec{y}) \, d\vec{y}}{4R^{2}\pi}.$$
 (3.33)

The response is equal to the product of the albedo and the average emission of the total spherical surface.

Chapter 4

Finite-element methods for the global illumination problem

Iteration requires the representation of the temporary radiance function L_n . So does expansion if viewindependent solution is needed since the final radiance distribution must be represented in a continuous domain. The exact representation of such functions might need infinite data, which is not available for practical algorithms. Thus finite approximations are required.

To represent a function over a continuous domain, **finite-element methods** [Pop87] can be used, which approximate the function in a finite function series form:

$$L(\vec{x},\omega) \approx L^{(n)}(\vec{x},\omega) = \sum_{j=1}^{n} \mathbf{L}_{j} \cdot b_{j}(\vec{x},\omega) = \mathbf{b}^{T}(\vec{x},\omega) \cdot \mathbf{L}$$
(4.1)

where $b_i(\vec{x}, \omega)$ is a system of predefined basis functions, and \mathbf{L}_i factors are unknown coefficients.

This representation can also be seen as projecting the infinite dimensional space of the possible radiance functions into a finite-dimensional function space defined by the basis functions.



Figure 4.1: Finite-element approximation



Figure 4.2: Projection to the adjoint base

Substituting this approximation into the rendering equation we can obtain:

$$\sum_{j=1}^{n} \mathbf{L}_{j} \cdot b_{j}(\vec{x},\omega) \approx \sum_{j=1}^{n} \mathbf{L}_{j}^{e} \cdot b_{j}(\vec{x},\omega) + \mathcal{T} \sum_{j=1}^{n} \mathbf{L}_{j} \cdot b_{j}(\vec{x},\omega) = \sum_{j=1}^{n} \mathbf{L}_{j}^{e} \cdot b_{j}(\vec{x},\omega) + \sum_{j=1}^{n} \mathbf{L}_{j} \cdot \mathcal{T} b_{j}(\vec{x},\omega).$$
(4.2)

Note that equality cannot be guaranteed, since even if $\sum_{j=1}^{n} \mathbf{L}_{j} \cdot b_{j}(\vec{x}, \omega)$ is in the subspace defined by the basis functions, the integral operator \mathcal{T} may result in a function that is out of this space. Instead, equality is required in an appropriate subspace defined by **adjoint basis** functions $\tilde{b}_{1}(\vec{x}), \tilde{b}_{2}(\vec{x}), \ldots, \tilde{b}_{n}(\vec{x})$ (figure 4.2). This set is required to be orthogonal to the original base $b_{1}(\vec{x}), b_{2}(\vec{x}), \ldots, b_{n}(\vec{x})$ in the following sense:

$$\langle b_i(\vec{x}), \tilde{b}_j(\vec{x}) \rangle = \begin{cases} 1 \text{ if } i = j, \\ 0 \text{ otherwise.} \end{cases}$$
(4.3)

Having projected equation 4.2 to the adjoint base — i.e. multiplying it by each adjoint basis functions b_i — we obtain the following system of linear equations:

$$\mathbf{L}_{i} = \mathbf{L}_{i}^{e} + \sum_{j=1}^{n} \langle \mathcal{T}b_{j}, b_{i} \rangle \cdot \mathbf{L}_{j}.$$
(4.4)

This system of linear equations can also be presented in a vector form:

$$\mathbf{L} = \mathbf{L}^e + \mathbf{R} \cdot \mathbf{L}, \qquad \mathbf{R}_{ij} = \langle \mathcal{T}b_j, b_i \rangle.$$
(4.5)

An adjoint of this linear equation can be derived by selecting the adjoint base as the normalization of the measuring functions. Suppose that each basis function b_i is associated with a measuring device W_i^e that measures the power \mathbf{P}_i leaving the support of the basis function, thus the appropriate adjoint basis function is $\tilde{b}_i = W_i^e / \langle b_i, W_i^e \rangle$. Thus the measured radiance is

$$\langle \sum_{j=1}^{n} \mathbf{L}_{j} \cdot b_{j}, W_{i}^{e} \rangle = \langle b_{i}, W_{i}^{e} \rangle \cdot \mathbf{L}_{i} = \mathbf{P}_{i}.$$

Similarly, the measured emission is

$$\langle \sum_{j=1}^{n} \mathbf{L}_{j}^{e} \cdot b_{j}, W_{i}^{e} \rangle = \langle b_{i}, W_{i}^{e} \rangle \cdot \mathbf{L}_{i}^{e} = \mathbf{P}_{i}^{e}.$$

Applying measuring operator W_i^e for equation (4.2), we can obtain the following equation:

$$\langle b_i, W_i^e \rangle \cdot \mathbf{L}_i = \langle b_i, W_i^e \rangle \cdot \mathbf{L}_i^e + \sum_{j=1}^n \langle \mathcal{T}b_j, W_i^e \rangle \cdot \mathbf{L}_j.$$
 (4.6)

Replacing the radiances by measured values, we get

$$\mathbf{P}_{i} = \mathbf{P}_{i}^{e} + \sum_{j=1}^{n} \frac{\langle \mathcal{T}b_{j}, W_{i}^{e} \rangle}{\langle b_{j}, W_{j}^{e} \rangle} \cdot \mathbf{P}_{j}.$$
(4.7)

This can also be presented in a matrix form

$$\mathbf{P} = \mathbf{P}^e + \mathbf{H} \cdot \mathbf{P},\tag{4.8}$$

where

$$\mathbf{H}_{ij} = \frac{\langle \mathcal{T}b_j, W_i^e \rangle}{\langle b_j, W_j^e \rangle} = \langle \mathcal{T}b_j, \tilde{b}_i \rangle \cdot \frac{\langle b_i, W_i^e \rangle}{\langle b_j, W_j^e \rangle} = \mathbf{R}_{ij} \cdot \frac{\langle b_i, W_i^e \rangle}{\langle b_j, W_j^e \rangle}.$$
(4.9)

When finite-element techniques are used together with expansion, finite-element representation can either be used to represent the final result [Kel95], or even be involved in the random walk [PM95].

The latter case may correspond either to the random-walk solution of the linear equation derived by projecting the integral equation, or to the Monte-Carlo evaluation of the multi-dimensional integrals containing both the transport and the projection operators. The second case is preferred, because it does not require matrix \mathbf{R} to be explicitly computed and stored.

The main problem of finite-element representations is that they require a lot of basis functions to accurately approximate high-variation, high-dimensional functions. Not surprisingly, finite-element methods have become really popular only for the diffuse case, where the radiance depends on 2 scalars and is relatively smooth. For solving the non-diffuse case, they are good only if the surfaces are not very specular.

4.1 Galerkin's method

Galerkin's method finds an approximation of the solution by making the error orthogonal to the set of basis functions. It means that the projection of error to the original base is zero. Formally, in Galerkin's method the set of basis functions is the same — except for normalization constants — as the set of adjoint basis functions. A particularly important case is the piece-wise constant approximation, when the surfaces and the directions are partitioned into patches A_1, A_2, \ldots, A_n and solid angles $\Omega_1, \Omega_2, \ldots, \Omega_m$, respectively. The basis functions are then:

$$b_{ij}(\vec{x},\omega) = \begin{cases} 1 \text{ if } \vec{x} \in A_i \land \omega \in \Omega_j, \\ 0 \text{ otherwise.} \end{cases}$$
(4.10)

The adjoint basis functions are:

$$\tilde{b}_{ij}(\vec{x},\omega) = \begin{cases} 1/(|A_i| \cdot |\Omega_j|) \text{ if } \vec{x} \in A_i \land \omega \in \Omega_j, \\\\ 0 \text{ otherwise.} \end{cases}$$
(4.11)

The system of linear equations determining the unknown L_{ij} values is

$$\mathbf{L}_{ij} = \mathbf{L}_{ij}^e + \sum_{k=1}^n \sum_{l=1}^m \mathbf{L}_{kl} \cdot \mathbf{R}_{ijkl},$$
(4.12)

where

$$\mathbf{R}_{ijkl} = \langle \mathcal{T}b_{kl}(\vec{x},\omega), \tilde{b}_{ij}(\vec{x},\omega) \rangle = \frac{1}{|A_i| \cdot |\Omega_j|} \cdot \int_{\Omega_j} \int_{A_i} \int_{\Omega} b_{kl}(h(\vec{x},-\omega'),\omega') \cdot f_r(\omega',\vec{x},\omega) \cos \theta'_{\vec{x}} \, d\omega' \, d\vec{x} d\omega.$$
(4.13)

4.2 Point collocation method

The **point collocation** method finds the unknown coefficients by requiring the finite-element approximation to be exact at the predefined knot points only. Formally, it uses Dirac-delta type adjoint basis functions which emphasize these knot-points:

$$b_{ij}(\vec{x},\omega) = \delta_{ij}(\vec{x}-\vec{x}_i,\omega-\omega_j).$$
(4.14)

The coefficients of the linear equation can be expressed as follows:

$$\mathbf{R}_{ijkl} = \langle \mathcal{T}b_{kl}(\vec{x},\omega), \tilde{b}_{ij}(\vec{x},\omega) \rangle = \int_{\Omega} b_{kl}(h(\vec{x}_i,-\omega'),\omega') \cdot f_r(\omega',\vec{x}_i,\omega_j) \cos \theta'_{\vec{x}_i} \, d\omega'.$$
(4.15)

4.3 Finite-element methods for the diffuse global illumination problem

If the surfaces have only diffuse reflection and emission — which is a general assumption of the **radiosity method** [CG85] — then the rendering (or the potential) equation has a simplified form:

$$L(\vec{x}) = L^e(\vec{x}) + \int_{\Omega_H} L(h(\vec{x}, -\omega')) \cdot f_r(\vec{x}) \cdot \cos \theta'_{\vec{x}} \, d\omega'.$$
(4.16)

In this case, the BRDF and the radiance depend on the surface point, but are independent of the direction, which reduces the inherent dimensionality of the problem and simplifies the finite-element representation:

$$L(\vec{x},\omega) \approx \sum_{j=1}^{n} L_j \cdot b_j(\vec{x}).$$
(4.17)

A widely used approach is the application of piece-wise constant basis functions for which $b_j(\vec{x}) = 1$ if \vec{x} is on surface element A_j and 0 otherwise. An appropriate adjoint basis is $\tilde{b}_j(\vec{x}) = 1/A_j$ if \vec{x} is on surface element A_j and 0 otherwise. Using this set of basis functions, the original rendering equation is projected to the following linear equation:

$$\mathbf{L} = \mathbf{L}^e + \mathbf{R} \cdot \mathbf{L} \tag{4.18}$$

where

$$\mathbf{R}_{ij} = \langle \mathcal{T}b_j, \tilde{b}_i \rangle = \frac{1}{A_i} \cdot \int_{A_i} \int_{\Omega} b_j (h(\vec{x}, -\omega')) \cdot f_r(\vec{x}) \cdot \cos \theta'_{\vec{x}} \, d\omega' \, d\vec{x}.$$
(4.19)

Let us replace the directional integral by a surface integral using equation (2.2):

$$d\omega' = \frac{d\vec{y} \cdot \cos\theta_{\vec{y}}}{|\vec{x} - \vec{y}|^2}.$$

This is true only when \vec{x} and \vec{y} are visible from each other, otherwise the solid angle of the visible surface is obviously zero. In order to extent the formula to this case, a **visibility indicator** $v(\vec{x}, \vec{y})$ is introduced, which is 1 if the two points are visible from each other and zero otherwise. Using this substitution we obtain

$$\mathbf{R}_{ij} = \frac{1}{A_i} \cdot \int\limits_{A_i} \int\limits_{S} b_j(\vec{y}) \cdot f_r(\vec{x}) \cdot \frac{\cos \theta'_{\vec{x}} \cdot \cos \theta_{\vec{y}}}{|\vec{x} - \vec{y}|^2} \cdot v(\vec{x}, \vec{y}) \, d\vec{y} \, d\vec{x}. \tag{4.20}$$

Taking advantage that the base functions are zero except for their specific domain, $\cos \theta_{\vec{y}} = \cos \theta_i$ and $\cos \theta'_{\vec{x}} = \cos \theta_i$ are constant on these patches, and assuming that the BRDF on patch *i* is f_i , we get

$$\mathbf{R}_{ij} = \frac{f_i}{A_i} \cdot \int\limits_{A_i} \int\limits_{A_j} \frac{\cos \theta_i \cdot \cos \theta_j}{|\vec{x} - \vec{y}|^2} \cdot v(\vec{x}, \vec{y}) \, d\vec{y} \, d\vec{x}.$$
(4.21)

Note that \mathbf{R}_{ij} is a product of two factors, the albedo of patch *i* — that is $a_i = f_i \cdot \pi$ —, and a so called **form factor** \mathbf{F}_{ij} which describes the geometry of the scene:

$$\mathbf{F}_{ij} = \frac{1}{A_i} \cdot \int\limits_{A_i} \int\limits_{A_j} \frac{\cos \theta_i \cdot \cos \theta_j}{\pi \cdot |\vec{x} - \vec{y}|^2} \cdot v(\vec{x}, \vec{y}) \, d\vec{y} \, d\vec{x}.$$
(4.22)

So far we have discussed the light propagation problem from the point of view of gathering. For shooting, similar formulae can be produced if incoming direction ω' is replaced by the outgoing direction $\omega = -\omega'$ in equation (4.19):

$$\mathbf{R}_{ij} = \frac{1}{A_i} \cdot \int_{A_i} \int_{\Omega} b_j(h(\vec{x},\omega)) \cdot f_r(\vec{x}) \cdot \cos\theta_{\vec{x}} \, d\omega \, d\vec{x}.$$
(4.23)

An adjoint equation can also be derived as a special case of equation (4.8). Let W_i^e be 1 in points of A_i and at the directions of the hemisphere of A_i . The **power equation** is then

$$\mathbf{P} = \mathbf{P}^e + \mathbf{H} \cdot \mathbf{P},\tag{4.24}$$

where the different forms of \mathbf{H}_{ij} are taken from equation (4.9):

$$\mathbf{H}_{ij} = \mathbf{R}_{ij} \cdot \frac{A_i}{A_j} = \mathbf{R}_{ji} \cdot \frac{f_i}{f_j} = \frac{f_i}{A_j} \cdot \int_{A_j} \int_{\Omega} b_i(h(\vec{y},\omega)) \cdot \cos\theta_j \, d\omega \, d\vec{y}.$$
(4.25)

In order to solve the projected integral equation or the power equation, basically the same techniques can be applied as for the original integral equation: inversion, expansion and iteration. Both the number of unknown variables and the number of equations are equal to the number of surfaces (*n*). The calculated L_i radiances represent the light density of the surface on a given wavelength. To generate color pictures at least three independent wavelengths should be selected (say red, green and blue), and the color information will come from the results of the three different calculations.

Thus, to sum up, the basic steps are these:

- 1. \mathbf{F}_{ij} form factor calculation.
- 2. Describe the light emission (\mathbf{L}_{i}^{e}) on the representative wavelengths, or in the simplified case on the wavelength of red, green and blue colors.
- 3. Solve the linear equation for the representative wavelengths.
- 4. Generate the picture taking into account the camera parameters by any known hidden surface algorithm.

In practical circumstances the number of elemental surface patches can be very large, making the form factor computation and the solution of the linear equation rather time consuming.

4.3.1 Geometric methods for form factor computation

Geometric form factor calculation methods approximate the outer integral of the double integral of the form factor from a single value, thus they apply the following approximation:

$$\mathbf{F}_{ij} = \frac{1}{A_i} \cdot \int\limits_{A_i} \int\limits_{A_j} \frac{\cos \theta_i \cdot \cos \theta_j}{\pi \cdot |\vec{x} - \vec{y}|^2} \cdot v(\vec{x}, \vec{y}) \, d\vec{y} \, d\vec{x} \approx \int\limits_{A_j} \frac{\cos \theta_i \cdot \cos \theta_j}{\pi \cdot |\vec{x}_i - \vec{y}|^2} \cdot v(\vec{x}_i, \vec{y}) \, d\vec{y}. \tag{4.26}$$

where \vec{x}_i is the center of patch *i*.

Nusselt [SH81] has realized that this formula can be interpreted as projecting the visible parts of A_j onto the unit hemisphere centered above \vec{x}_i , then projecting the result orthographically onto the base circle of this hemisphere in the plane of \vec{x}_i , and finally calculating the ratio of the doubly projected area and the area of the unit circle (π). Due to the central role of the unit hemisphere, this method is called the **hemisphere algorithm**.

This means that a single row of the form factor matrix can be determined by solving a visibility problem, where the "window" is a half sphere and the "eye" is in the center of surface i.

Discrete hemisphere algorithm and its variations

Discrete methods subdivide the hemisphere or its base circle into finite number of areas called "pixels" and assume that what can be seen through these small areas is homogeneous (figure 4.3). Thus it is enough to determine the visibility through a single point in each pixel.

The complicated form of the "hemispherical window" can be simplified if the hemisphere is replaced by other immediate surfaces, such as a **hemicube** [CG85] or a **cubic tetrahedron** [BKP91]. In these cases the modification of the geometry must be compensated by appropriate weighting functions in area computation. For hemicube and hemishpere algorithms, the window surface consists of normal rectangles, which can be exploited by the built in transformation and scan-conversion hardware of graphics workstations.


Figure 4.3: Geometric interpretation of hemisphere form factor algorithm and the discrete algorithm for form factor computation

Chapter 5

Numerical quadrature for high dimensional integrals

The solution of the rendering equation requires the numerical evaluation of high-dimensional integrals. Numerical quadrature formulae take finite samples from the domain, evaluate the integrand for these samples and generate the integral as a weighted sum of the integrand values. The general form for the solution is

$$I = \int_{V} f(\mathbf{z}) \, d\mathbf{z} \approx \sum_{i=1}^{N} f(\mathbf{z}_{i}) \cdot w(\mathbf{z}_{i})$$
(5.1)

where \mathbf{z}_i is a sample point from the *s*-dimensional domain *V*, and *w* is the weighting function of the given quadrature rule.

Classical quadrature rules usually trace back the computation of multi-dimensional integrals to a series of one-dimensional integrals. Thus they select the coordinates of the sample points independently, and the high-dimensional points are defined as the Cartesian product of the 1-dimensional sample sets. The simplest techniques, including the **brick-rule**, the **trapezoidal-rule**, **Simpson-rule**, etc. space the abscissas equally, which results in a regular grid in the domain. More sophisticated techniques, such as the **Gaussian quadrature**, select the sample points non uniformly along each abscissa, giving more freedom to decrease the error of the integral.



Figure 5.1: Error of the brick rule in one and two dimensional spaces

Let us evaluate the error of the brick rule in one and two dimensional spaces (figures 5.1). Suppose that N sample points are used in the domains of [0, 1] and $[0, 1]^2$, respectively. In the one dimensional space the error is the sum of the areas of triangles that are between the function and the series of bricks. The average height of a triangle is $\Delta f/(2N)$ where Δf is the variation of the integrand and N is the number of bricks. Since the base of a triangle is 1/N and the number of triangles is N, the error is:

$$\epsilon_1 = \frac{\Delta f}{2N} \cdot \frac{1}{N} \cdot N = \frac{\Delta f}{2N} = \mathcal{O}\left(\frac{1}{N}\right).$$
(5.2)

In the two dimensional space the N sample points should be distributed along two coordinates, thus we have just $n = \sqrt{N}$ points in a single row and column. Now the error is the sum of the volume of the roof-like objects that are between the function and the series of bricks. The average height of a volume element is $\Delta f/(2n)$, its base has $1/n^2$ area and there are n^2 elements. Thus the error is

$$\epsilon_2 = \frac{\Delta f}{2n} \cdot \frac{1}{n^2} \cdot n^2 = \frac{\Delta f}{2\sqrt{N}} = \mathcal{O}\left(\frac{1}{\sqrt{N}}\right).$$
(5.3)

This idea can be generalized to any dimensions and we can conclude that the integration error of the brick rule in an *s*-dimensional space is proportional to $\Delta f/N^{1/s}$ (Δf is the total change, i.e. the variation, of the integrand). From another point of view, it means that the required number of sample points to find an estimate with error ϵ is proportional to $(\Delta f/\epsilon)^s$, thus the computational efforts grow exponentially with the dimension of the domain. This phenomenon is called as the **dimensional explosion** or **dimensional core**. Concerning other classical quadrature rules, Δf measures higher order changes, such as the distance from the piece-wise linear or polynomial functions, but they also exhibit dimensional core.

The dimensional explosion can be avoided by **Monte-Carlo** [Sob91] or **quasi-Monte Carlo** [Nie92] integration methods. Monte-Carlo methods trace back the estimation of an integral to the calculation of an expected value, which is estimated by averaging random samples. Quasi-Monte Carlo techniques, on the other hand, use deterministic samples that are uniformly distributed in the integration domain.

5.1 Monte-Carlo quadrature

Monte-Carlo integration converts the calculation of an integral to an equivalent expected value problem [Sob91]. Assume that a random vector variable \mathbf{z} is uniformly distributed in V, thus its probability density is $p(\mathbf{z}) = 1/V$. The expected value of random variable $f(\mathbf{z})$ is

$$E[f(\mathbf{z})] = \int_{V} f(\mathbf{z}) \cdot p(\mathbf{z}) \, d\mathbf{z} = \int_{V} f(\mathbf{z}) \cdot \frac{1}{V} \, d\mathbf{z} = \frac{1}{V} \cdot I, \tag{5.4}$$

thus the required integral can easily be found if this expected value is available. According to the **theo**rems of large numbers, if independent random samples $\mathbf{z}_1, \mathbf{z}_2, \ldots, \mathbf{z}_N$ are generated using probability density p, then the expected value can be estimated by the average \hat{f} :

$$E[f(\mathbf{z})] \approx \hat{f} = \frac{1}{N} \sum_{i=1}^{N} f(\mathbf{z}_i)$$
(5.5)

thus we obtain the following quadrature formula:

$$\int_{V} f(\mathbf{z}) \, d\mathbf{z} \approx \frac{V}{N} \sum_{i=1}^{N} f(\mathbf{z}_{i}).$$
(5.6)

Estimator \hat{f} is also a random variable whose expected value is equal to $E[f(\mathbf{z})] = I/V$. Suppose that the variance of $f(\mathbf{z})$ is σ^2 . If the samples are independent random variables, then the variance of estimator \hat{f} is:

$$D^{2}\left[\hat{f}\right] = \frac{1}{N^{2}} \sum_{i=1}^{N} \sigma^{2} = \frac{\sigma^{2}}{N}.$$
(5.7)

Thus the standard deviation of the estimator is $D[\hat{f}] = \sigma/\sqrt{N}$. According to the **central limit theorem**, estimator \hat{f} will have normal distribution asymptotically, with mean $E[f(\mathbf{z})] = I/V$ and standard deviation σ/\sqrt{N} . Examining the shape of the probability density of the normal distribution, we can conclude that the probability that the distance between the variable and the mean is less than 3 times the standard deviation is 0.997. Thus with 0.997 confidence level we can say that the (probabilistic) error bound is:

$$\left|\int_{V} f(\mathbf{z}) \, d\mathbf{z} - \frac{V}{N} \sum_{i=1}^{N} f(\mathbf{z}_{i})\right| < \frac{3V\sigma}{\sqrt{N}}.$$
(5.8)

Let us realize that this bound is independent of the dimension of the domain! This property means that by Monte-Carlo quadrature the dimensional explosion can be avoided.

5.2 Quasi-Monte Carlo quadrature

In the previous section we concluded that randomly distributing the sample points solves the dimensional core problem. We show that it is also possible with deterministically selected sample points. The resulting method is called **quasi-Monte Carlo quadrature** [Sob91, Nie92]. For the sake of simplicity, assume that an *s*-dimensional function $f(\mathbf{z})$ needs to be integrated over the domain $[0, 1]^s$. In order to guarantee that *f* is Riemann-integrable, *f* is assumed to be piece-wise continuous. This integral is approximated by a finite sum as follows:

$$\int_{\mathbf{z}\in[0,1]^s} f(\mathbf{z}) \, d\mathbf{z} \approx \frac{1}{N} \sum_{i=1}^N f(\mathbf{z}_i).$$
(5.9)

The question is how the sequence of sample points $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N$ should be selected to minimize the error of the integral quadrature. A minimal requirement is the stability of the sequence, which means that in asymptotic sense the error is zero for any Riemann integrable function:

$$\int_{\mathbf{z}\in[0,1]^s} f(\mathbf{z}) d\mathbf{z} = \lim_{N\to\infty} \frac{1}{N} \sum_{i=1}^N f(\mathbf{z}_i).$$
(5.10)

Sequences meeting this requirement are called **uniform** or **equidistribution** (i.e. **1-equidistribution**) sequences [Nie92].

In order to find other necessary requirements for uniform sequences, let us consider the integration of a very simple function which is 1 inside a *d*-dimensional "brick" originating at the center of the coordinate system and 0 elsewhere:

$$\mathcal{L}(\mathbf{z}) = \begin{cases} 1 \text{ if } 0 \le \mathbf{z}|_1 \le v_1, 0 \le \mathbf{z}|_2 \le v_2, \dots, 0 \le \mathbf{z}|_s \le v_s, \\ 0 \text{ otherwise.} \end{cases}$$
(5.11)

Let us denote the volume of this brick by $V(A) = \prod_{j=1}^{s} v_j$. Integrating this function we have:

$$\int_{\mathbf{z}\in[0,1]^s} \mathcal{L} \, d\mathbf{z} = \prod_{j=1}^s v_j = V(A).$$
(5.12)

If the number of sample points that are inside the s-dimensional "brick" A is m(A), then the finite approximation sum is

$$\frac{1}{N}\sum_{i=1}^{N}f(\mathbf{z}_{i}) = \frac{m(A)}{N}.$$
(5.13)

Since the exact value of the integral is V(A) now, for uniform sequences, the average number of sample points that are located inside a volume should be proportional to the size of this volume:

$$\lim_{N \to \infty} \frac{m(A)}{N} = V(A).$$
(5.14)

If the size of the sequence of the sample points is not infinite, then the proportionality requirement can only be approximately met. The maximal error in this approximation is called the **discrepancy** (or the **star-discrepancy**) of the sequence:

$$\mathcal{D}^*(\mathbf{z}_1, \mathbf{z}_2, \dots \mathbf{z}_N) = \sup_A \left| \frac{m(A)}{N} - V(A) \right|.$$
(5.15)

If a sequence is appropriate for integral-quadrature, then the approximation sum is expected to converge to the real value of the integral, which requires the discrepancy to converge to 0. This is another necessary requirement which is derived from the examination of a very simple function. Note that this requirement also emphasizes the uniformness of those sequences that can be used for numerical integration.

However, this requirement is not only necessary, but also sufficient, since any Riemann-integrable function can be approximated by piece-wise constant step-functions with arbitrary precision. To show how step-functions can do this, an example of a 2-dimensional function is shown in figure 5.2.



Figure 5.2: Definition of functions from bricks originating at the center of the coordinate system

5.2.1 Error analysis for integrands of finite variation: Koksma-Hlawka inequality

In this section, an upper-bound is given to the error of the quasi-Monte Carlo quadrature for functions that have bounded and piece-wise continuous mixed derivatives.

The error of the quadrature is shown below:

$$\left|\int_{\mathbf{z}\in[0,1]^s} f(\mathbf{z}) \, d\mathbf{z} - \frac{1}{N} \sum_{i=1}^N f(\mathbf{z}_i)\right|.$$
(5.16)

Intuitively this error must depend on two independent factors. On the one hand, if the discrepancy of the sequence of sample points is high, then there are large regions where there are no sample point at all, which increases the error. This means that the error is expected to be proportional to the discrepancy of the sequence of sample locations.

On the other hand, the error also depends on how quickly the function changes between the sample points. If the function can change significantly in a small domain, then the error can be quite large. However, if the slope of the function is small, then nothing dramatic happens between the sample points thus the error will be small.

Measures describing how rapidly a function can change are called **variations**. For a 1-dimensional function the **variation in the sense of Vitali** [Deá89] is defined as:

$$\mathcal{V}_{\rm V}(f(x)) = \limsup \sum_{i=1}^{N} |f(x_{i+1}) - f(x_i)|.$$
(5.17)

For a 2-dimensional function, the definition is analogous:

$$\mathcal{V}_{\mathcal{V}}(f(x,y)) = \limsup \sum_{i=1}^{N} \sum_{j=1}^{M} |f(x_{i+1}, y_{j+1}) - f(x_{i+1}, y_i) - f(x_i, y_{i+1}) + f(x_i, y_i)|.$$
(5.18)

Note that for higher dimensions, the variation of Vitali does not provide a good measure: if the function is constant in x, for instance, then the variation is zero, regardless how the function changes depending on y.

Thus, it is worth introducing a somehow more stronger variation type, called the **Hardy-Krause variation**. The variation in the sense of Hardy-Krause is the sum of the variations of the function and of its restrictions to the end of the domain. For the two-dimensional case, the new variation is:

$$\mathcal{V}_{\mathrm{HK}}(f(x,y)) = \mathcal{V}_{\mathrm{V}}f(x,y) + \mathcal{V}_{\mathrm{V}}f(x,1) + \mathcal{V}_{\mathrm{V}}f(1,y).$$
(5.19)

If a function has bounded and piece-wise continuous mixed derivatives, then its variation is finite. For a 2-dimensional function meeting this requirement, the variation can be given by the following formula:

$$\mathcal{V}_{\mathrm{HK}}(f(u,v)) = \int_{0}^{1} \int_{0}^{1} \left| \frac{\partial^{2} f(u,v)}{\partial u \partial v} \right| \, du \, dv + \int_{0}^{1} \left| \frac{\partial f(u,1)}{\partial u} \right| \, du + \int_{0}^{1} \left| \frac{\partial f(1,v)}{\partial v} \right| \, dv. \tag{5.20}$$

The property that a function is not continuous does not necessarily mean that the variation is infinite. If at most finite or countable infinite discontinuities occur at hyper-planes parallel to the coordinate axes, then the variation is still finite. An example of a discontinuous function that have finite variation is

$$f(x,y) = \begin{cases} 1 \text{ if } x > x_0, \\ 0 \text{ otherwise.} \end{cases}$$
(5.21)

However, when the discontinuity is not parallel to the coordinate axes, then the variation is infinite. A simple function of infinite variation is [Deá89]:

$$f(x,y) = \begin{cases} 1 \text{ if } x > y, \\ 0 \text{ otherwise.} \end{cases}$$
(5.22)

Now, let us turn to the error of quasi-Monte Carlo integration. The following formula, which expresses the previous intuition that the error depends on the uniformness of the sample points and on the variation of the integrand, is called the **Koksma-Hlawka inequality**:

$$\left|\int_{\mathbf{z}\in[0,1]^s} f(\mathbf{z}) \, d\mathbf{z} - \frac{1}{N} \sum_{i=1}^N f(\mathbf{z}_i)\right| \le \mathcal{V}_{\mathrm{HK}} \cdot \mathcal{D}^*(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N).$$
(5.23)

For the sake of notational simplicity, the Koksma-Hlawka inequality is proven here only for the two-dimensional case (s = 2). The outline of the proof is as follows. We start with the assumption that the function has piece-wise continuous and bounded mixed derivatives, thus the mixed derivatives are Riemann-integrable and their integration gives back the original function. First the function is expressed as the integral of its derivatives, and the function values in the approximation sum are converted accordingly. On the other hand, the integral of f is converted to a similar form using partial integration. Finally the difference of the approximating sum and the integral is examined and upperbounded.

First, the value of function f(u, v) is expressed by its derivatives at point x, y:

$$\begin{aligned} f(x,y) &= f(1,1) + [f(1,1) - f(x,1) - f(1,y) + f(x,y)] - [f(1,1) - f(x,1)] - [f(1,1) - f(1,y)] = \\ f(1,1) &+ \int_{x}^{1} \int_{y}^{1} f_{uv}(u,v) \, du \, dv - \int_{x}^{1} f_{u}(u,1) \, du - \int_{y}^{1} f_{v}(1,v) \, dv. \end{aligned}$$

where

$$f_{uv}(u,v) = \frac{\partial^2 f(u,v)}{\partial u \partial v}, \quad f_u(u,v) = \frac{\partial f(u,v)}{\partial u}, \quad f_v(u,v) = \frac{\partial f(u,v)}{\partial v}.$$

Let us introduce the step function ϵ :

$$\epsilon(u, v) = \begin{cases} 1 \text{ if } u \ge 0, v \ge 0, \\ 0 \text{ otherwise.} \end{cases}$$

Using this step function, the domains of the integrals can be extended to [0, 1], as follows:

$$f(x,y) = f(1,1) + \int_{0}^{1} \int_{0}^{1} f_{uv}(u,v) \cdot \epsilon(u-x,v-y) \, du \, dv - \int_{0}^{1} f_{u}(u,1) \cdot \epsilon(u-x,1) \, du - \int_{0}^{1} f_{v}(1,v) \cdot \epsilon(1,v-y) \, dv.$$

Substituting $\mathbf{z}_i = (x_i, y_i)$ into this formula we have:

$$f(\mathbf{z}_i) = f(1,1) + \int_0^1 \int_0^1 f_{uv}(u,v) \cdot \epsilon(u-x_i,v-y_i) \, du \, dv - \int_0^1 f_u(u,1) \cdot \epsilon(x_i,1) \, du - \int_0^1 f_v(1,v) \cdot \epsilon(1,y_i) \, dv.$$

Averaging these formulae for i = 1, 2, ..., N, the approximating sum has the following form:

$$\frac{1}{N}\sum_{i=1}^{N}f(\mathbf{z}_{i}) = f(1,1) + \frac{1}{N}\int_{0}^{1}\int_{0}^{1}f_{uv}(u,v)\cdot m(u,v)\,du\,dv - \frac{1}{N}\int_{0}^{1}f_{u}(u,1)\cdot m(u,1)\,du - \frac{1}{N}\int_{0}^{1}f_{v}(1,v)\cdot m(1,v)\,dv.$$

where

$$m(u,v) = \sum_{i=1}^{N} \epsilon(u - x_i, v - y_i),$$

which is the number of points located in the rectangle [(0,0), (u, v)]. The integral

$$\int_{\mathbf{z} \in [0,1]^2} f(\mathbf{z}) \, d\mathbf{z} = \int_{v=0}^1 \int_{u=0}^1 f(u,v) \, du \, dv$$

can also be converted to a similar form, if partial integration is applied. First the inner integral is considered:

$$\int_{u=0}^{1} f(u,v) \, du = \int_{u=0}^{1} f(u,v) \cdot 1 \, du = f(1,v) - \int_{u=0}^{1} f_u(u,v) \cdot u \, du = F(v)$$

Then the outer integral is processed in a similar way:

$$\int_{v=0}^{1} \int_{u=0}^{1} f(u,v) \, du \, dv = \int_{v=0}^{1} F(v) \, dv = \int_{v=0}^{1} F(v) \cdot 1 \, dv = F(1) - \int_{v=0}^{1} F_v(v) \cdot v \, dv$$

Substituting the definition of F(v) we get:

$$\int_{v=0}^{1} \int_{u=0}^{1} f(u,v) \, du \, dv = f(1,1) + \int_{0}^{1} \int_{0}^{1} f_{uv}(u,v) \cdot uv \, du \, dv - \int_{0}^{1} f_{u}(u,1) \cdot u \, du - \int_{0}^{1} f_{v}(1,v) \cdot v \, dv.$$

Thus the error of quadrature is then:

$$\begin{split} |\int_{v=0}^{1} \int_{u=0}^{1} f(u,v) \, du \, dv - \frac{1}{N} \sum_{i=1}^{N} f(\mathbf{z}_{i})| = \\ |\int_{0}^{1} \int_{0}^{1} f_{uv}(u,v) \cdot \left(\frac{m(u,v)}{N} - uv\right) \, du \, dv - \int_{0}^{1} f_{u}(u,1) \cdot \left(\frac{m(u,1)}{N} - u\right) \, du - \int_{0}^{1} f_{v}(1,v) \cdot \left(\frac{m(1,v)}{N} - v\right) \, dv| \leq \\ \left(\int_{0}^{1} \int_{0}^{1} |f_{uv}(u,v)| \cdot \, du \, dv + \int_{0}^{1} |f_{u}(u,1)| \, du + \int_{0}^{1} |f_{v}(1,v)| \, dv\right) \cdot \sup_{u,v} \left| \left(\frac{m(u,v)}{N} - uv\right) \right| = \\ \mathcal{V}_{\mathrm{HK}} \cdot \mathcal{D}^{*}(\mathbf{z}_{1}, \mathbf{z}_{2}, \dots \mathbf{z}_{N}). \end{split}$$

This is exactly what we wanted to prove.

According to this inequality, the error can be upperbounded by the product of two independent factors, the variation of the integrand and the discrepancy of the used sample set. The discrepancy shows how uniformly the set is distributed [Shi91a]. This immediately presents two orthogonal strategies to improve the quality of quadratures. Either we try to make the function flat by appropriate variable transformations, or use very uniformly distributed sample sets. The first technique is called **importance sampling** [Sob91], while the second involves the **stratification** [Sob91, Mit96, Arv95] of random points or the application of **low-discrepancy series** [Nie92, War95, PFTV92, Knu81, Sob91].

If the integrand has finite variation, then the error is proportional to the discrepancy of the sequence of sample locations. For carefully selected sample points the discrepancy can converge to zero with almost linear speed. Thus, quadratures having almost linear convergence can be obtained in this way, which is better than the $O(1/\sqrt{N})$ speed of Monte-Carlo quadratures.

Note, on the other hand, that functions having infinite variation can also be integrated by quasi-Monte Carlo quadrature. The quadrature will be asymptotically exact for any uniform sequence and for any Riemann integrable function. The fact that the Koksma-Hlawka inequality cannot be applied means that it cannot be used to provide a qualitative measure for the speed of the convergence. Practical experience shows that quasi-Monte Carlo integration outperforms the classical Monte-Carlo integration even for discontinuous functions [SKDP99]. However, the difference in the effectiveness becomes significantly less when the integrand has infinite variation. This phenomenon will be investigated in the subsequent sections. Since the integrand of the rendering and potential equations is usually discontinuous, this case is very important for computer graphics applications.

5.2.2 Generation of the sample points

As a conclusion of error analysis we can state that we need very uniform sequences for quasi-Monte Carlo quadrature. Regular grids should be rejected because of their $O(1/N^{1/s})$ discrepancy which results in dimensional explosion.



Figure 5.3: 100 points distributed by a regular grid (left), random distribution (middle) and Halton low-discrepancy sequence (right)

Monte-Carlo method proposed the application of random samples to avoid dimensional explosion. This can also be justified by analyzing the discrepancy. In fact, the discrepancy of a uniformly distributed random series is

$$\mathcal{O}\left(\sqrt{\frac{\log\log N}{2N}}\right)$$

asymptotically with probability 1.

In order to prove this, let us define a new random variable ξ_i from uniformly distributed random sample z_i in the following way:

$$\xi_i = \begin{cases} 1 \text{ if } \mathbf{z}_i \text{ is in an } s \text{-dimensional brick } A \text{ that originates at the center,} \\ 0 \text{ otherwise.} \end{cases}$$

Note that if z_i is uniformly distributed in $[0, 1]^s$, then the expected value and the variance of ξ_i are

$$E[\xi_i] = V(A), \qquad D^2[\xi_i] = V(A) - V^2(A) \le \frac{1}{4},$$

where V(A) is the volume of brick A. If the samples are independent random variables, the generated $\xi_1, \xi_2, \ldots, \xi_N$ random variables will also be independent, and of the same distribution having mean $E[\xi]$ and variance $D^2[\xi]$.

According to the **theorem of iterated logarithm** [Rén62], the difference between the average of independent random variables $\xi_1, \xi_2, \ldots, \xi_N$ of the same distribution and their mean $E[\xi]$ can be upperbounded in the following way:

$$\Pr\left\{\lim\sup\left|\sum\frac{\xi_i}{N} - E[\xi]\right| \le \sqrt{D^2[\xi] \cdot \frac{2\log\log N}{N}}\right\} = 1.$$

In our case the variance $D^{2}[\xi]$ cannot exceed 1/4 and

$$\sup\left|\sum \frac{\xi_i}{N} - E[\xi]\right| = \sup_A \left|\frac{m(A)}{N} - V(A)\right| = \mathcal{D}^*(\mathbf{z}_1, \mathbf{z}_2, \dots \mathbf{z}_N),$$

thus the theorem of iterated logarithm becomes what we want to prove:

$$\Pr\left\{\lim \mathcal{D}^*(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N) \le \sqrt{\frac{\log \log N}{2N}}\right\} = 1.$$

5.2.3 Generation of low-discrepancy sequences

Beyond random sequences, however, there are deterministic sequences that have even better discrepancy. The discrepancy of the best sequences known is in the order of $O(\log^s N/N)$ or even in the order of $O(\log^{s-1} N/N)$ if N is known before starting the sequence. These sequences are called **lowdiscrepancy** sequences. There are many sequences published in the literature [Nie92, War95, Dea89, Knu81]. The most famous one is probably the **Halton-sequence** (its one-dimensional version is also called **Van der Corput sequence**).

The element i of the one-dimensional Halton sequence of base b is defined as the radical inverse of the expansion of i in base b. This means that number i is expanded in radix b, then the number is mirrored onto the "radix" point. The first few points in base 2 are shown in table 5.1.

i	binary form of <i>i</i>	radical inverse	H_i
1	1	0.1	0.5
2	10	0.01	0.25
3	11	0.11	0.75
4	100	0.001	0.125
5	101	0.101	0.625
6	110	0.011	0.375
7	111	0.111	0.875

Table 5.1: The calculation of the *i*th Halton point H_i in base 2

Why is this sequence uniform? Note that the construction algorithm generates as binary form of *i* all binary combinations of length *k* before producing a combination of length k + 1. This means that after the radical inverse the sequence H_i will visit all intervals of length 2^{-k} before putting a second point in an interval already visited. Thus the sequence is really uniform.

On the other hand, as k increases, the algorithm produces a single point in each interval of length 1/2, then in each interval of length 1/4, etc. thus the sequence is not only asymptotically uniform, but also the first N points are fairly uniformly distributed (this is guaranteed by the property that the radical inverse makes the most significant bit the most rapidly changing bit). This is also true in other radix



Figure 5.4: The distribution of the first 10,100 and 1000 Halton points in 2 dimensions

systems as well. If the base is b, then the Halton sequence will place a sample point in all intervals of size b^{-k} before putting a new point into any intervals.

A Halton sequence is able to generate points that are uniformly distributed in the 1-dimensional [0, 1] interval. If higher dimensional regions, such as rectangles, cubes, etc. should be filled uniformly, then different coordinates of the sample vectors can be generated from Halton sequences of different base numbers. In order for these vectors to uniformly fill the whole region, the "interdependence" of different coordinates should be as little as possible. To examine this, let us assume that a two-dimensional Halton sequence is generated with base numbers b_1 and b_2 . According to the behavior of the Halton sequence, the generation algorithm would visit all columns of width $b_1^{-k_1}$ before visiting a column again, and similarly it would visit all rows of height $b_2^{-k_2}$ before putting a new point into a row. The columns and rows form $b_1^{k_1} \cdot b_2^{k_2}$ cells. Uniformness in the two-dimensional space means that the algorithm is expected to put a point in each cell before putting a second point in any cells. Since the periodicity of the columns and rows are $b_1^{k_1}$ and $b_2^{k_2}$, respectively, the periodicity of the cells is the smallest common multiple of $b_1^{k_1}$ and $b_2^{k_2}$. This equals to the their product, that is total number of cells, if b_1 and b_2 are relative primes.

This can also be stated in a general way. If a multi-dimensional Halton sequence is to be constructed, then the base numbers of different coordinates must be relative primes.

A C++ class that can initialize an arbitrary Halton point and then it can generate incrementally all subsequent points using a very fast algorithm [Kel96b] is presented in the following:

```
class Halton {
        double value, inv_base;
        Number( long i, int base ) {
                 double f = inv_base = 1.0/base;
                 value = 0.0;
                 while ( i > 0 ) {
                         value += f * (double)(i % base);
                         i /= base; f *= inv_base;
                 }
        double Next( ) {
                  double r = 1.0 - value - 0.000000001;
                  if (inv_base < r) value += inv_base;</pre>
                  else {
                         double h = inv_base, hh;
                         do {
                                 hh = h; h *= inv_base;
                         } while ( h >= r );
                         value += hh + h - 1.0;
                 }
                return value;
        }
};
```

5.3 Importance sampling

Importance sampling is a well-known method of Monte-Carlo integration to reduce variance. The basic idea is to use non-uniform distribution to find sample points, which places more samples where the function is large. More specifically it means that

$$I = \int_{V} f(\mathbf{z}) \, d\mathbf{z} = \int_{V} \frac{f(\mathbf{z})}{p(\mathbf{z})} \cdot p(\mathbf{z}) \, d\mathbf{z} = E\left[\frac{f(\mathbf{z})}{p(\mathbf{z})}\right] \approx \frac{1}{N} \cdot \sum_{i=1}^{N} \frac{f(\mathbf{z}_{i})}{p(\mathbf{z}_{i})} \pm \frac{3V\sigma}{\sqrt{N}},\tag{5.24}$$

where $p(\mathbf{z})$ is a probability density in V, the \mathbf{z}_i points are selected according to this probability density, and the variance σ^2 is defined by

$$\sigma^{2} = D^{2} \left[\frac{f(\mathbf{z})}{p(\mathbf{z})} \right] = E \left[\left(\frac{f(\mathbf{z})}{p(\mathbf{z})} - I \right)^{2} \right] = \int_{V} \left(\frac{f(\mathbf{z})}{p(\mathbf{z})} - I \right)^{2} \cdot p(\mathbf{z}) \, d\mathbf{z}.$$
(5.25)

The probability density $p(\mathbf{z})$ should be selected to minimize the variance. As can be shown easily, the variance can be minimized if $p(\mathbf{z})$ is proportional to the integrand $f(\mathbf{z})$. In order to demonstrate this, let us express the ratio of the integrand and the probability density in the following way:

$$\frac{f(\mathbf{z})}{p(\mathbf{z})} = I + \beta \cdot \delta(\mathbf{z}), \tag{5.26}$$

where $I = E[\frac{f(\mathbf{z})}{p(\mathbf{z})}]$ and β is a normalization constant to allow $\int_{V} (\delta(\mathbf{z}))^2 \cdot p(\mathbf{z}) d\mathbf{z} = 1$. The variance of the integral quadrature is then:

$$\sigma^2 = E[(I + \beta \cdot \delta(\mathbf{z}) - E[I + \beta \cdot \delta(\mathbf{z})])^2] = \beta^2 \cdot E[(\delta(\mathbf{z}))^2] = \beta^2.$$
(5.27)

This is obviously minimal if $\beta = 0$, when the variance is also zero.

Thus in Monte-Carlo integration it is worth applying probability distributions that are large where the integrand is large and small where the integrand is negligible.

5.3.1 Generation of a random variable with a prescribed probability density

We concluded that importance sampling requires random samples generated from a probability density which is proportional — at least approximately — to the integrand. This section examines how such samples can be found. First, let us consider the 1-dimensional case and assume that the domain of the integration is an interval [a, b].

Suppose that we want samples with a probability density that is proportional to a function g(z). This function is an approximation of the integrand f. If this function is different from the integrand, then the importance sampling is not optimal. The needed probability density p(z) can be obtained by scaling function g to integrate to 1 as probability densities do:

$$p(z) = \frac{g(z)}{\int\limits_{a}^{b} g(z)dz}.$$
(5.28)

Note that this gives us a clear explanation why we should use non-optimal densities. If g were equal to f, then the construction of the probability density would require the integral of f.

From the probability density *p*, probability distribution function *P* is obtained:

$$P(z) = \int_{a}^{z} p(Z) \, dZ.$$
 (5.29)

A random variable ξ having probability distribution P can be constructed by transforming another random variable r, which is uniformly distributed in the [0, 1] interval, with the $\xi = P^{-1}(r)$ transformation.

To prove this, the probability distribution of ξ is examined:

$$\Pr\{\xi < z\} = \Pr\{P^{-1}(r) < z\} = \Pr\{r < P(z)\} = P(z).$$
(5.30)

since P(z) is not decreasing and the probability that r is in an interval of [0, 1] equals to the size of this interval.

The multi-dimensional case can be traced back to a series of 1-dimensional constructions. After normalization we have

$$p(\mathbf{z}) = \frac{g(\mathbf{z})}{\int\limits_{V} g(\mathbf{z}) d\mathbf{z}}.$$
(5.31)

The probability density is expressed as a product of 1-dimensional conditional densities:

$$p(z_1, z_2, \dots, z_s) = p_1(z_1 | z_2, \dots, z_s) \cdot p_2(z_2 | z_3, \dots, z_s) \cdot \dots \cdot p_s(z_s).$$
(5.32)

For these conditional densities the same method can be used recursively as for the 1-dimensional case. If the different coordinates are independent random variables, then we obtain:

$$p(z_1, z_2, \dots, z_s) = p_1(z_1) \cdot p_2(z_2) \cdot \dots \cdot p_s(z_s).$$
 (5.33)

5.3.2 Importance sampling in quasi-Monte Carlo integration

Classical Monte-Carlo integration places more samples where the integrand is large. The same basic idea also makes sense in quasi-Monte Carlo integration. However, for formal analysis, we have to find another approach since terms like probability density or variance cannot be applied in the deterministic quasi-Monte Carlo framework.

The alternative formulation is the integration using variable transformation. Suppose that a function f needs to be integrated in domain V and we have a monotonously increasing mapping T that maps this domain onto V'. The integral can also be evaluated in the new domain using the following formula:

$$\int_{V} f(\mathbf{z}) \, d\mathbf{z} = \int_{V'} f(T^{-1}(\mathbf{y})) \left| \frac{\partial T^{-1}(\mathbf{y})}{\partial \mathbf{y}} \right| \, d\mathbf{y},\tag{5.34}$$

where $\left|\frac{\partial T^{-1}(\mathbf{y})}{\partial \mathbf{y}}\right|$ is the Jacobi determinant of the inverse transformation.

If quasi-Monte Carlo integration is used, then domain V' is $[0, 1]^s$. In order to quantify the error of the quadrature, we can use the Koksma-Hlawka inequality which expresses the error-bound as a product of the discrepancy of the sample points and the variation of the function. Since the discrepancy of the sample points is independent of the function to be integrated, the error-bound can be controlled by the appropriate transformation of the integrand to reduce variation. In order to reduce the variation, the function should be flattened. In the ideal case when the integrand is constant, the variation is 0.

To make the transformed integrand constant, the Jacobi determinant should be inversely proportional to f. Since the intuitive meaning of the Jacobi determinant is the compression or expansion ratio of the two corresponding sub-cubes in V and V' respectively, this criterion states that if the sample points are uniformly distributed in V', then the transformation will concentrate them around regions where f is high.

For the sake of simplicity, the details are discussed only for the 1-dimensional case, when the variable transformation has the following form

$$\int_{z_{\min}}^{z_{\max}} f(z) \, dz = \int_{0}^{1} f(T^{-1}(y)) \cdot \frac{dT^{-1}(y)}{dy} \, dy.$$
(5.35)

In the ideal case mapping T makes the integrand have zero variation, that is constant C:

$$f(T^{-1}(y)) \cdot \frac{dT^{-1}(y)}{dy} = \mathcal{C}.$$

From this we can have

$$T(z) = \frac{1}{\mathcal{C}} \cdot \int_{z_{\min}}^{z} f(Z) \, dZ.$$

Since mapping T is expected to map to [0, 1], we require that $T(z_{\max}) = 1$. Thus the constant C should be equal to $\int_{z_{\min}}^{z_{\max}} f(Z) dZ$. Summarizing, the uniformly distributed point y should be transformed by the inverse of the following function

$$T(z) = \frac{\int\limits_{z_{\min}}^{z} f(Z) dZ}{\int\limits_{z_{\min}}^{z_{\max}} f(Z) dZ}.$$
(5.36)

Note that this is the same transformation as has been derived for the random samples. It means that the method of importance sampling is independent of whether random samples are transformed in Monte-Carlo quadrature, or deterministic samples are transformed in quasi-Monte Carlo integration.

5.3.3 Metropolis sampling

Importance sampling requires a probability density that is proportional to a function g which, in turn, should mimic the integrand. In the previous section we discussed a random variable transformation method that can generate samples with the required probability density. This method expects function g to be symbolically integrable and its primitive function to be invertible. This requirement often contradicts the requirement of the good approximation of the original integrand.

This section discusses a different sampling strategy, proposed by Metropolis et. al [MRR⁺53], that has less expectations towards function g. In fact, it only supposes that the integral of g can be determined either symbolically or numerically.

The Metropolis method carries out sampling by establishing a discrete time Markov process \mathbf{z}_i , i = 1, 2, ... in the space of samples, whose limiting distribution is proportional to the selected function. A discrete time Markov process visits states which correspond to samples. The Metropolis algorithm constructs this process such that the probability of being in a given state converges to a limiting value and in the limiting case $g(\mathbf{z}) = b \cdot p(\mathbf{z})$, where $b = \int_V g(\mathbf{z}) d\mathbf{z}$.

A Markov process can be defined by the state transition probabilities, that is by the conditional probability of the next state provided that the current state is known. In Metropolis method, the next state \mathbf{z}_{i+1} of this process is found by letting an almost arbitrary **tentative transition function** $T(\mathbf{z}_i \rightarrow \mathbf{z}_t)$ generate a **tentative sample** \mathbf{z}_t which is either accepted as the real next state or rejected making the next state equal to the actual state using an **acceptance probability** $a(\mathbf{z}_i \rightarrow \mathbf{z}_t)$. Thus the state transition probability density from state \mathbf{x} to a different state \mathbf{y} is:

$$P(\mathbf{x} \to \mathbf{y}) \, d\mathbf{y} = T(\mathbf{x} \to \mathbf{y}) \, d\mathbf{y} \cdot a(\mathbf{x} \to \mathbf{y}). \tag{5.37}$$

The event that the process remains in the same state is the complement of moving to any other state, thus the probability of no state transition happening is:

$$1 - \int_{\mathbf{x} \in V, \mathbf{x} \neq \mathbf{y}} P(\mathbf{x} \to \mathbf{y}) \, d\mathbf{y} = 1 - \int_{\mathbf{x} \in V, \mathbf{x} \neq \mathbf{y}} T(\mathbf{x} \to \mathbf{y}) \cdot a(\mathbf{x} \to \mathbf{y}) \, d\mathbf{y}.$$
 (5.38)

Let us denote the probability density of being in state \mathbf{x} at step n by $p_n(\mathbf{x})$. Using the total probability theorem and taking advantage that in Markov processes the future depends only on the present state and is independent of the past, the following recursion can be established for these state probabilities:

$$p_{n+1}(\mathbf{y}) = \int_{\mathbf{x}\in V, \mathbf{x}\neq \mathbf{y}} p_n(\mathbf{x}) \cdot P(\mathbf{x} \to \mathbf{y}) \, d\mathbf{x} + \left(1 - \int_{\mathbf{x}\neq \mathbf{y}} P(\mathbf{y} \to \mathbf{x}) \, d\mathbf{x}\right) \cdot p_n(\mathbf{y}). \tag{5.39}$$

If the limiting probability $p(\mathbf{y}) = \lim_{n \to \infty} p_n(\mathbf{y})$ exists, then it should be the fixed point of this recursion, that is:

$$p(\mathbf{y}) = p(\mathbf{y}) + \int_{\mathbf{x} \in V, \mathbf{x} \neq \mathbf{y}} p(\mathbf{x}) \cdot P(\mathbf{x} \to \mathbf{y}) - p(\mathbf{y}) \cdot P(\mathbf{y} \to \mathbf{x}) \, d\mathbf{x}.$$
 (5.40)

If the Markov process is ergodic, then the limiting distribution is unambigous and is independent of the initial state of the process. The process is ergodic if after a given number of steps any state can be reached from any other state with non-zero probability.

The core idea of Metropolis sampling is to construct acceptance probability $a(\mathbf{x} \rightarrow \mathbf{y})$ in such a way that the limiting probability of the process will be $p(\mathbf{z}) = g(\mathbf{z})/b$. Substituting this goal into equation (5.40) we get:

$$g(\mathbf{y}) = g(\mathbf{y}) + \int_{\mathbf{x} \in V, \mathbf{x} \neq \mathbf{y}} g(\mathbf{x}) \cdot P(\mathbf{x} \to \mathbf{y}) - g(\mathbf{y}) \cdot P(\mathbf{y} \to \mathbf{x}) \, d\mathbf{x}.$$
 (5.41)

This holds when the total incoming and outgoing flows of state \mathbf{x} are balanced. One way of ensuring this is to require that:

$$g(\mathbf{x}) \cdot P(\mathbf{x} \to \mathbf{y}) = g(\mathbf{y}) \cdot P(\mathbf{y} \to \mathbf{x}).$$
(5.42)

This condition — that is also called as the **detailed balance** — means that the transitions between any two states are balanced. Using the formulae (5.37) for state transition probabilities, we can further obtain:

$$g(\mathbf{x}) \cdot T(\mathbf{x} \to \mathbf{y}) \cdot a(\mathbf{x} \to \mathbf{y}) = g(\mathbf{y}) \cdot T(\mathbf{y} \to \mathbf{x}) \cdot a(\mathbf{y} \to \mathbf{x}).$$
(5.43)

Thus the required ratio of the two acceptance probabilities is:

$$\frac{a(\mathbf{y} \to \mathbf{x})}{a(\mathbf{x} \to \mathbf{y})} = \frac{g(\mathbf{x}) \cdot T(\mathbf{x} \to \mathbf{y})}{g(\mathbf{y}) \cdot T(\mathbf{y} \to \mathbf{x})}.$$
(5.44)

Any acceptance probability satisfying this requirement makes the limiting distribution proportional to g. Considering the speed of convergence, on the other hand, the state transition probabilities and consequently the acceptance probabilities should be large. Since a probability cannot exceed 1, the optimal acceptance probability is:

$$a(\mathbf{x} \to \mathbf{y}) = \min\left\{\frac{g(\mathbf{y}) \cdot T(\mathbf{y} \to \mathbf{x})}{g(\mathbf{x}) \cdot T(\mathbf{x} \to \mathbf{y})}, 1\right\}.$$

The algorithm generating a trajectory of the Markov process to obtain samples $\{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N\}$ is as follows:

for i = 1 to N do Based on the actual state \mathbf{z}_i , choose another random, tentative state \mathbf{z}_t using $T(\mathbf{z}_i \to \mathbf{z}_t)$ $a(\mathbf{z}_i \to \mathbf{z}_t) = (g(\mathbf{z}_t) \cdot T(\mathbf{z}_t \to \mathbf{z}_i))/(g(\mathbf{z}_i) \cdot T(\mathbf{z}_i \to \mathbf{z}_t))$ if $a(\mathbf{z}_i \to \mathbf{z}_t) \ge 1$ then $\mathbf{z}_{i+1} = \mathbf{z}_t$ else // accept with probability $a(\mathbf{z}_i \to \mathbf{z}_t)$ Generate uniformly distributed random number r in [0, 1]. if $r < a(\mathbf{z}_i \to \mathbf{z}_t)$ then $\mathbf{z}_{i+1} = \mathbf{z}_t$ else $\mathbf{z}_{i+1} = \mathbf{z}_i$ endif

endfor

Chapter 6

Random walk solution of the global illumination problem

Recall that expansion obtains the measured power as a Neumann series: $\mathcal{M}L = \sum_{i=0}^{\infty} \mathcal{M}\mathcal{T}^i L^e$. The terms of this infinite Neumann series have intuitive meaning as well: $\mathcal{M}\mathcal{T}^0 L^e = \mathcal{M}L^e$ comes from the emission, $\mathcal{M}\mathcal{T}^1 L^e$ comes from a single reflection, $\mathcal{M}\mathcal{T}^2 L^e$ from two reflections, etc.

6.1 Why should we use Monte-Carlo expansion methods?

Expansion techniques require the evaluation of very high-dimensional — in fact, infinite dimensional — integrals. When using classical quadrature rules for multi-dimensional integrals [PFTV92], such as for example the trapezoidal rule, in order to provide a result with a given accuracy, the number of sample points is in the order of $\mathcal{O}(N^s)$, where s is the dimension of the domain. This phenomenon is called the **dimensional core** or **dimensional explosion** and makes classical quadrature rules prohibitively expensive for higher dimensions.

However, Monte-Carlo or quasi-Monte Carlo techniques distribute the sample points simultaneously in all dimensions, thus they can avoid dimensional explosion. For example, the probabilistic error bound of Monte-Carlo integration is $\mathcal{O}(N^{-0.5})$, independently of the dimension of the domain. *s*-dimensional low discrepancy series [Nie92] can even achieve $\mathcal{O}(\log^s N/N) = \mathcal{O}(N^{-(1-\epsilon)})$ convergence rates for finite variation integrands.

Furthermore, classical quadrature cannot be used for infinite dimensional integrals, thus the Neumann series should be truncated after D terms. This truncation introduces a bias of order $\lambda^{D+1} \cdot ||L^e||/(1-\lambda)$, where λ is the contraction of the light transport operator. Using a Russian-roulette based technique, on the other hand, Monte-Carlo methods are appropriate for even infinite dimensional integrals. Thus we can conclude that the stochastic approach is indispensable for expansion methods.

In computer graphics the first Monte-Carlo random walk algorithm — called **distributed ray-tracing** — was proposed by Cook et al. [CPC84], which spawned to a set of variations, including **path-tracing** [Kaj86], **light-tracing** [DLW93], **bi-directional path-tracing** [LW93, VG95], **Monte-Carlo radiosity** [Shi91b, Neu95, PM95], and **two-pass methods** which combine radiosity and ray-tracing [Shi90, ZS95, WCG87].

The problem of naive generation of walks is that the probability that a shooting path finds the eye is zero for a pin-hole camera or very small if a non-zero aperture camera model is used, while the probability that a gathering random path ends in a lightsource may be very little if the lightsources are small, thus the majority of the paths do not contribute to the image at all, and their computation is simply waste of time. Note that shooting is always superior for view-independent algorithms since they do not have to face the problem of small aperture. Thus, on the one hand, random walk must be combined with a deterministic step that forces the walk to go to the eye and to find a lightsource. On the other hand, **importance sampling** [Sob91] should be incorporated to prefer useful paths along which significant radiance is transferred.

Steps of the walk transfer the radiance or the potential in the scene. The source and destination of the transfer can be points in the case of continuous methods or patches in the case of finite-element methods. If the algorithm is such that it always selects a single source for shooting or single destination for gathering, then the method is called **local method**. On the other hand, if many sources and destinations are taken into consideration simultaneously in each transfer, then the method is called **global method** or **multi-path method** [Sbe96].

6.2 Quasi-Monte Carlo quadrature for the rendering equation

Quasi-Monte Carlo walk techniques mean that instead of generating the next direction randomly, the direction is sampled from a low-discrepancy point set. Since the low-discrepancy sequences have better asymptotic disrepancy than random sequences do, quasi-Monte Carlo methods are expected to provide more accurate results. However, the integrand of the rendering equation is discontinuous where the discontinuity is not aligned with the coordinate axes, thus its variation is infinite. These discontinuities are usually produced by the projected object boundaries. This property makes the Koksma-Hlawka inequality not appropriate for the error analysis and for the prediction of the convergence rates.

6.2.1 Integrating functions of unbounded variation

In this section the convergence speed is examined for functions which are generally smooth but have general discontinuities of finite "length". First the domain of the integration is assumed to be 2-dimensional, then the results will be generalized to arbitrary dimensions.



Figure 6.1: A typical integrand of the rendering equation

Suppose that N samples have been generated to estimate the integral of a function such as in figure 6.1 using a low-discrepancy sequence. In order to overcome the difficulty that the integrand f has infinite variation, the function is decomposed into two functions, one is smooth \tilde{f} having continuous mixed derivatives and the other \hat{f} inherits the discontinuity of f (figure 6.2).



Figure 6.2: Decomposition of f into a smooth (\tilde{f}) and a discontinuous (\hat{f}) function

Low-discrepancy sequences generate points in a cellular grid in a way that the difference of the number of points in two cells is at most 1. If there are already N number of points, the size of a cell on the finest, filled level is approximately $1/\sqrt{N} \times 1/\sqrt{N}$. Let us define the domain of \hat{f} as the set of those cells that are intersected by the discontinuity. This domain is called the **domain of discontinuity**. The number of such cells is in the order of the "digital length" of the discontinuity curve, which is the product of the maximum extent l and the resolution of the grid \sqrt{N} . Since each cell has at least 1 (and at most 2) points, the number of points in this domain is at least $l\sqrt{N}$.

The error of quadrature is as follows:

$$\left| \int_{\mathbf{z} \in [0,1]^2} f(\mathbf{z}) \, d\mathbf{z} - \frac{1}{N} \sum_{i=1}^N f(\mathbf{z}_i) \right| \le \left| \int_{\mathbf{z} \in [0,1]^2} \tilde{f}(\mathbf{z}) \, d\mathbf{z} - \frac{1}{N} \sum_{i=1}^N \tilde{f}(\mathbf{z}_i) \right| + \left| \int_{\mathbf{z} \in [0,1]^2} \hat{f}(\mathbf{z}) \, d\mathbf{z} - \frac{1}{N} \sum_{i=1}^N \hat{f}(\mathbf{z}_i) \right|.$$
(6.1)

Since \tilde{f} has finite variation, the first term in the error is bounded by $\mathcal{V}_{HK}(\tilde{f}) \cdot \mathcal{D}^*(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N)$.

Concerning the second term, the integration of \hat{f} is estimated taking $l\sqrt{N}$ uniformly distributed samples and averaging the result. Since the samples and the discontinuity are not related in any way, we can suppose that this is a normal Monte-Carlo integration [PFTV92]. The uniform property of lowdiscrepancy sequence guarantees that this pseudo-random set can be assumed to have uniform distribution. If Δf is the difference between the maximum and minimum values in the domain of discontinuity, then $\sigma^2 \leq (\Delta f)^2$. In our case the number of sample points N' is $l\sqrt{N}$ and the size of the domain V is l/\sqrt{N} , thus we obtain with 0.997 confidence level:

$$\int_{V} \hat{f}(\mathbf{z}) \, d\mathbf{z} = \frac{V}{N'} \cdot \sum_{i=1}^{N'} \hat{f}(\mathbf{z}_i) \pm 3 \cdot V \cdot \frac{\Delta f}{\sqrt{N'}} = \frac{1}{N'} \cdot \sum_{i=1}^{N'} \hat{f}(\mathbf{z}_i) \pm 3 \cdot \Delta f \cdot \sqrt{l} \cdot N^{-3/4}. \tag{6.2}$$

Taking into account that \hat{f} is zero outside the domain of discontinuity, equality 6.1 can be expressed as:

$$\left|\int_{\mathbf{z}\in[0,1]^2} f(\mathbf{z}) \, d\mathbf{z} - \frac{1}{N} \sum_{i=1}^N f(\mathbf{z}_i)\right| \le \mathcal{V}_{\mathrm{HK}}(\tilde{f}) \cdot \mathcal{D}^*(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N) + 3 \cdot \Delta f \cdot \sqrt{l} \cdot N^{-3/4}.$$
(6.3)

For large N values the second term will be dominant, which results in $\mathcal{O}(N^{-3/4})$ error bound. This is poorer than the $\mathcal{O}(\log^2 N/N)$ bound suggested by the Koksma-Hlawka inequality assuming, for example, the application of the Halton sequence. Note that the point from where the second term dominates the first one depends on "intensity" Δf and size of the discontinuity \sqrt{l} . The same analysis can be carried out in higher dimensions as well. In *s* dimensions a discontinuity of size *l* would intersect $V = l \cdot N^{-1/s}$ volume of cells which would contain $N' = l \cdot N^{(s-1)/s}$ sample points. Thus the general error bound is:

$$\left|\int_{\mathbf{z}\in[0,1]^s} f(\mathbf{z}) \, d\mathbf{z} - \frac{1}{N} \sum_{i=1}^N f(\mathbf{z}_i)\right| \le \mathcal{V}_{\mathrm{HK}}(\tilde{f}) \cdot \mathcal{D}^*(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N) + 3 \cdot \Delta f \cdot \sqrt{l} \cdot N^{-\frac{(s+1)}{2s}}.$$
 (6.4)

Thus, for quasi-Monte Carlo integration of discontinuous functions, the order of the error bound will be in between the $\mathcal{O}(N^{-(1-\epsilon)})$ bound of finite variation functions and the $\mathcal{O}(N^{-0.5})$ bound of Monte-Carlo quadrature. The higher the dimension of the integral, the closer the Monte-Carlo and quasi-Monte Carlo techniques get in convergence speed. Thus it is still worth using quasi-Monte Carlo quadrature if the size of the discontinuity l is not very large, since in this case the error could be significantly less then for Monte-Carlo quadrature.

Numerical evidence using simple functions

In order to demonstrate the previous results, the convergences of a 2-dimensional and a 3-dimensional functions are examined, that are simple enough to analytically compute their integrals.

The 2-dimensional function is:

$$f_2(x,y) = \begin{cases} (x+y) \cdot a + 1 - 2 \cdot a & \text{if } x + y > 1, \\ (x+y) \cdot a & \text{otherwise}, \end{cases}$$
(6.5)



Figure 6.3: Error of integrating f_2 (*left*) *and* f_3 (*right*)



Figure 6.4: Error measurements for 1 and 10 bounces in the spherical reference scene (section 3.4.2) where the BRDF is diffuse, the albedo is 0.5, and 25 percents of the area is a diffuse lightsource

where a is a free parameter in the range of [0, 0.5]. Note that by setting a appropriately, the intensity of the discontinuity can be controlled without altering either the value of the integral or the variation of the continuous part. If a = 0.5, then the function has finite variation, otherwise it has infinite variation. The results of the simulation are shown in the left of figure 6.3. This figure shows the maximum error after a given number of samples.

The 3-dimensional function is:

$$f_3(x, y, z) = \begin{cases} (x + y + z) \cdot a + 0.6 - 1.8 \cdot a & \text{if } x + y + z > 1, \\ (x + y + z) \cdot a & \text{otherwise}, \end{cases}$$
(6.6)

where a is a free parameter in the range of [0, 1/3]. If a = 1/3, then f_3 has finite variation, otherwise it has not. The error of integration of f_3 is summarized in the right of figure 6.3.

Numerical evidence for the rendering equation

The efficiency of Monte-Carlo and quasi-Monte Carlo quadratures have been tested for the presented spherical scene (section 3.4.2) assuming a single pixel camera. The error has been measured separately for the different bounces.

Looking at the error measurements of figure 6.4, we can see that even for integrands of infinite variation, quasi-Monte Carlo methods are still better but they lose their advantage when computing higher bounces as predicted by the theoretical results. The other important problem in higher dimensions is that although a low-discrepancy series has almost linearly decreasing discrepancy in the asymptotic sense, this discrepancy can still be high for not very many points (in the solution of the rendering equation we rarely use more than 1000 samples for the estimation of a single pixel). In the case of the Halton series, for example, the **base** of the series strongly affects the initial behavior of the discrepancy. These base numbers are different prime numbers for different dimensions, thus for high-dimensional integrals the base numbers can be quite high, which results in degraded performance.

6.3 Importance sampling for the rendering equation

When solving the rendering equation, usually directional integrals (or surface integrals in other formulation) should be evaluated. These directional integrals have the following form:

$$\mathcal{T}L^{\mathrm{in}}(\vec{x},\omega) = \int_{\Omega} L^{\mathrm{in}}(\vec{x},\omega') \cdot f_r(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega'.$$
(6.7)

To allow the application of random or low-discrepancy point sets, the integration domain should be transformed to the unit cube or square. To establish such a mapping, first direction ω' is expressed by spherical coordinates ϕ , θ' , which converts the directional integral to the following form:

$$\int_{\Omega} L^{\text{in}}(\vec{x},\omega') \cdot f_r(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega' = \int_{\phi=0}^{2\pi} \int_{\theta'=0}^{\pi} L^{\text{in}}(\phi,\theta') \cdot f_r(\phi,\theta') \cdot \cos\theta' \cdot \sin\theta' \, d\theta' d\phi \quad (6.8)$$

since $d\omega = \sin \theta' d\theta' d\phi$. Let us denote $f_r(\phi, \theta') \cdot \cos \theta' \cdot \sin \theta'$ by $w(\phi, \theta')$, which is the transfer probability density. Now we find a mapping $T(\phi, \theta') = \mathbf{z}$ that emphasizes those directions where the integrand is large and projects the domain of the spherical coordinates onto the unit square:

$$\int_{\phi=0}^{2\pi} \int_{\theta'=0}^{\pi} L^{\text{in}}(\phi, \theta') \cdot w(\phi, \theta') \, d\theta' d\phi = \int_{[0,1]^2} L^{\text{in}}(T^{-1}(\mathbf{z})) \cdot w(T^{-1}(\mathbf{z})) \cdot \left| \frac{dT^{-1}(\mathbf{z})}{d\mathbf{z}} \right| \, d\mathbf{z} = \int_{[0,1]^2} L^{\text{in}}(T^{-1}(\mathbf{z})) \cdot \frac{w(T^{-1}(\mathbf{z}))}{t(\mathbf{z})} \, d\mathbf{z},$$
(6.9)

where

$$\left|\frac{dT^{-1}(\mathbf{z})}{d\mathbf{z}}\right| = \frac{1}{t(\mathbf{z})}$$

is the Jacobi determinant of the inverse mapping. If the Jacobi determinant is large, then a small portion of the unit square is mapped onto a large region. Thus sample points that are uniformly distributed in the unit square will be quite rare in these regions. Alternatively, where the Jacobi determinant is small, the sample points will be dense. Considering this, the meaning of $t(\mathbf{z})$ is the **density** of sample points in the neighborhood of $\omega = (\phi, \theta') = T^{-1}(\mathbf{z})$. This has an illustrative content for the random case. If \mathbf{z} is uniformly distributed random variable, then the probability density of $\omega = T^{-1}(\mathbf{z})$ will be $t(\mathbf{z})$.

The solution of the rendering equation for a given point (\vec{x}, ω) requires the evaluation of the following multi-dimensional integral (equation (3.5)):

$$L(\vec{x},\omega) = L^e + \mathcal{T}L^e + \mathcal{T}^2L^e + \dots = \int_{[0,1]^2} \dots \int_{[0,1]^2} L^e + \frac{w_1}{t_1} \cdot L^e + \frac{w_1}{t_1} \cdot \frac{w_2}{t_2} \cdot L^e + \dots d\mathbf{z}_1 d\mathbf{z}_2 \dots$$
(6.10)

This can be estimated by Monte-Carlo or quasi-Monte Carlo quadratures which evaluate the integrand in sample points and average the results. A crucial design decision of such an algorithm is the selection of mappings T_i to have good importance sampling. Using probabilistic approach, it means that the probability of selecting a walk is proportional to its contribution. Following the directions concluded from the Koksma-Hlawka inequality, the mappings should make the integrand flat — that is of low variation, or constant in the ideal case.

Looking at formula (6.10), which is the single multi-dimensional solution of the rendering equation, this decision seems to be hard to make, since there are too many free parameters to control simultaneously. Fortunately, the solution can also be presented in the following recursive form:

$$L(\vec{x},\omega) = L^e + \int_{[0,1]^2} \frac{w_1}{t_1} \cdot [L^e + \int_{[0,1]^2} \frac{w_2}{t_2} \cdot [L^e + \ldots] \ldots] d\mathbf{z}_1 d\mathbf{z}_2 \ldots$$
(6.11)

If we could ensure that each of the integrands of the form

$$\int_{[0,1]^2} \frac{w_i}{t_i} \cdot \left[L^e + \int_{[0,1]^2} \ldots\right] d\mathbf{z}_i$$

is constant (at least approximately), then the integrand of the single multi-dimensional integral will also be constant [SKCP99]. An optimal importance sampling strategy thus requires density t_i to be proportional to the product of the incoming illumination $L^e + \int \ldots$ and the cosine weighted BRDF w_i . Unfortunately, during random walks the incoming non-direct illumination is not known (the random walk is just being done to estimate it).

Thus we have to use approximations for which we have three alternatives. Firstly, information about the illumination in the space can be gathered in a preprocessing phase, then this information can be used to obtain probability densities for importance sampling. This is called the **global importance sampling**. These methods can be classified according to the data structure built in the preprocessing phase. Since the ray-space is 5-dimensional, it is straightforward to apply a **5D adaptive tree** [LW96] that is similar to the well-known octree to store radiance information. Jensen proposed the application of the **photonmap** as the basis of importance sampling [Jen95]. We assigned the power computed in the preprocessing phase to **links** established between two interacting patches [SKCP98].

The second alternative is using the information gained during previous walks to approximate the illumination. This strategy is called **adaptive importance sampling**. Adaptive importance sampling methods neither require the non-uniform probability densities to be constructed in advance, nor simplify them to take into account only local properties, but converge to a desired probability density using the knowledge of previous samples. Three techniques are particularly important, which have also been used in rendering: **genetic algorithms** [LB94] the **Metropolis sampling** [MRR⁺53, VG97] and the **VEGAS method** [Lep80, SK98a]. The first use of Metropolis sampling in rendering aimed at speeding up bidirectional path tracing [VG97].

In the third alternative the problem is simplified and the indirect illumination is not considered in importance sampling. When the directions are generated, we use only transfer probability density w_i and L^e representing the direct illumination of the actual point. This is called the **local importance sampling**.

It turns out that we have to encounter severe problems when we have to find a mapping which has a density that is proportional to the product of the effects of the transfer probability density and the direct lighting. Consequently, local importance sampling strategies usually use only either w_i or L^e to identify important directions. The first alternative is called the **BRDF sampling**, while the second is called the **lightsource sampling**.

6.3.1 BRDF sampling

BRDF based importance sampling means that at step i the density t_i is proportional to the transfer probability density w_i , that is

$$t_i \propto w_i = f_r(\omega_{\rm in}, \vec{x}, \omega_{\rm out}) \cdot \cos\theta \sin\theta.$$
(6.12)

In gathering algorithms ω_{out} is known, θ is the angle between ω_{in} and the surface normal, and ω_{in} should be determined. In shooting algorithms, on the other hand, ω_{in} is known, θ is the angle between ω_{out} and the surface normal, and ω_{out} should be determined.

Due to the fact that t_i represents density (probability density for Monte-Carlo methods), its integral is 1. Thus for gathering walks and for non-transparent materials, the ratio of proportionality in equation (6.12) is

$$\int_{\Omega_H} w \, d\omega_{\rm in} = \int_{\Omega_H} f_r(\omega_{\rm in}, \vec{x}, \omega_{\rm out}) \cdot \cos \theta_{\rm in} \, d\omega_{\rm in} = a(\vec{x}, \omega_{\rm out})$$

where $a(\vec{x}, \omega_{out})$ is the **albedo** of the surface at point \vec{x} in the outgoing direction. Similarly, the proportionality ratio for shooting walks is

$$\int_{\Omega_H} w \, d\omega_{\text{out}} = \int_{\Omega_H} f_r(\omega_{\text{in}}, \vec{x}, \omega_{\text{out}}) \cdot \cos \theta_{\text{out}} \, d\omega_{\text{out}} = a(\vec{x}, \omega_{\text{in}}).$$

Thus the weights $w^* = w_i/t_i$ are the albedos at the visited points.

BRDF sampling for diffuse materials

Diffuse materials have constant BRDF, that is

$$t_i(\phi, \theta) \propto w_i = f_r \cdot \cos \theta \sin \theta.$$

The proportionality ratio is found to make t_i to integrate to 1:

$$t_i(\phi,\theta) = \frac{w_i}{\int\limits_{\Omega_H} w_i \, d\omega} = \frac{f_r \cdot \cos\theta \sin\theta}{\int\limits_{\phi=0}^{2\pi} \int\limits_{\theta=0}^{\pi/2} f_r \cdot \cos\theta \sin\theta \, d\theta d\phi} = \frac{\cos\theta \sin\theta}{\pi}$$

Assuming that the random variables used to produce the two coordinate samples are independent, this density is obtained in a product form:

$$t_i(\phi,\theta) = \frac{1}{2\pi} \cdot \left[2\cos\theta\sin\theta\right],\tag{6.13}$$

where $1/(2\pi)$ is the probability density of ϕ and $2\cos\theta\sin\theta = \sin 2\theta$ is the probability density of θ .

The corresponding probability distribution functions are:

$$P(\phi) = \int_{0}^{\phi} \frac{1}{2\pi} d\Phi = \frac{\phi}{2\pi}, \qquad P(\theta) = \int_{0}^{\theta} \sin 2\Theta \, d\Theta = \sin^2 \theta.$$

Thus the required ϕ and θ random variables can be found by the following transformations of u, v variables that are uniformly distributed in [0, 1] (section 5.3.1):

$$\phi = 2\pi \cdot u, \quad \theta = \arcsin\sqrt{v}.$$

The transformed weight after importance sampling is the albedo

$$w_i^* = \frac{w_i}{t_i} = f_r \cdot \pi = a.$$
(6.14)

BRDF sampling for specular materials

Specular materials can be characterized by the reciprocal version of the Phong's BRDF model, that is

$$f_r(\omega_{\rm in}, \vec{x}, \omega_{\rm out}) = k_s \cdot \cos^n \psi \cdot \epsilon(\pi/2 - \theta),$$

where ψ is the angle between ω_{out} and the mirror direction of ω_{in} onto the surface normal, which will be referred to as ω_r , and $\epsilon(\pi/2 - \theta)$ indicates that the outgoing direction cannot point into the object, i.e. the angle θ between the surface normal and the outgoing direction should be less than 90 degrees.



Figure 6.5: Parameterization for the calculation of the albedo

In order to appropriately parameterize the directional sphere, now the north pole is selected by the reflection direction ω_r (figure 6.5). Let us identify a direction by an angle from ω_r , that is by ψ , and by another angle ϕ between its projection onto a plane perpendicular to ω_r and an arbitrary vector on this plane.

BRDF sampling requires a density which satisfies the following criterion:

$$t_i(\phi,\psi) \propto w_i = k_s \cdot \cos^n \psi \cdot \cos \theta(\psi,\phi) \cdot \epsilon(\pi/2 - \theta(\psi,\phi)) \cdot \sin \psi$$

Unfortunately, the $\cos \theta \cdot \epsilon(\pi/2 - \theta)$ factor forbids the symbolic integration of this formula, thus we will use a density that is proportional only to $\tilde{w}_i = k_s \cdot \cos^n \psi \sin \psi$. The proportionality ratio is found to make t_i to integrate to 1:

$$t_i(\phi,\psi) = \frac{k_s \cdot \cos^n \psi \sin \psi}{\int\limits_{\phi=0}^{2\pi} \int\limits_{\psi=0}^{\pi/2} k_s \cdot \cos^n \psi \sin \psi \, d\psi d\phi} = \frac{n+1}{2\pi} \cos^n \psi \sin \psi.$$

Assuming that the random variables used for producing the two coordinate samples are independent, this density is obtained in a product form:

$$t_i(\phi, \psi) = \frac{1}{2\pi} \cdot [(n+1)\cos^n \psi \sin \psi],$$
 (6.15)

where $1/(2\pi)$ is the probability density of ϕ and $(n+1)\cos^n\psi\sin\psi$ is the probability density of ψ .

The corresponding probability distribution functions are:

$$P(\phi) = \frac{\phi}{2\pi}, \qquad P(\psi) = \int_{0}^{\psi} (n+1)\cos^{n}\Psi\sin\Psi \, d\Psi = 1 - \cos^{n+1}\psi.$$

Thus the required ϕ and θ random variables can be found by the following transformations of u, v variables that are uniformly distributed in [0, 1]:

 $\phi = 2\pi \cdot u, \quad \psi = \arccos(1-v)^{1/(n+1)}.$

The transformed weight after importance sampling is

$$w_i^* = \frac{w_i}{t_i} = \frac{2\pi k_s}{n+1} \cdot \cos\theta(\psi, \phi) \cdot \epsilon(\pi/2 - \theta(\psi, \phi)).$$
(6.16)

Different other specular BRDF models are presented and their BRDF sampling is discussed in [War92, NNSK98b, NNSK98a, NNSK99a].

6.3.2 Lightsource sampling

Lightsource sampling is used in **direct lightsource calculations** [SWZ96] and as a complementary sampling strategy to BRDF sampling in random walks.

Since in this case, the samples are selected from a lightsource instead of the directional sphere, the surface integral form of the transport operator is needed:

$$(\mathcal{T}L^{e})(\vec{x},\omega) = \int_{\Omega} L^{e}(h(\vec{x},-\omega'),\omega') \cdot f_{r}(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega' = \int_{S_{e}} L^{e}(\vec{y},\omega_{\vec{y}\to\vec{x}}) \cdot f_{r}(\omega_{\vec{y}\to\vec{x}},\vec{x},\omega) \cdot \frac{\cos\theta'_{\vec{x}}\cdot\cos\theta_{\vec{y}}}{|\vec{x}-\vec{y}|^{2}} \cdot v(\vec{y},\vec{x}) \, d\vec{y},$$
(6.17)

where $v(\vec{y}, \vec{x})$ is 1 if points \vec{x} and \vec{y} are not occluded from each other and 0 otherwise, and S_e is the surface of non-zero emission. To obtain a Monte-Carlo estimate for this integral, N points $\vec{y}_1, \ldots, \vec{y}_N$ are sampled uniformly on the lightsource and the following formula is used:

$$(\mathcal{T}L^e)(\vec{x},\omega) \approx \frac{|S_e|}{N} \cdot \sum_{i=1}^N L^e(\vec{y}_i,\omega_{\vec{y}_i\to\vec{x}}) \cdot v(\vec{y}_i,\vec{x}) \cdot f_r(\omega_{\vec{y}_i\to\vec{x}},\vec{x},\omega) \cdot \frac{\cos\theta'_i \cdot \cos\theta_{\vec{y}_i}}{|\vec{x}-\vec{y}_i|^2}.$$
 (6.18)

If the scene has a single homogeneous lightsource which is relatively small and is far from the considered point, then the integrand will be approximately constant on the lightsource surface, thus this estimator has low variance.

6.3.3 Sampling the lightsources in gathering random walks

Since lightsource sampling generates samples only on the direct lightsources, it completely ignores indirect illumination. Thus it cannot be used alone in global illumination algorithms, but only as a complementary part of, for example, BRDF sampling.

The simplest way to combine the two strategies is to generate all but the last directions of the gathering walk by sampling the BRDF and to compute the last direction by sampling the lightsource. Note that when stopping the walk, the indirect illumination is assumed to be zero, thus following the directions of the lightsources is a reasonable approach.

Another combination strategy is to trace one or more shadow rays from each visited point of the walk towards the lightsources, not only from the last of them.

Formally, this approach can be presented as a restructuring of the Neumann series

$$L = L^e + \mathcal{T}L^e + \mathcal{T}^2L^e + \mathcal{T}^3L^e \dots = L^e + (\mathcal{T}L^e) + \mathcal{T}(\mathcal{T}L^e) + \mathcal{T}^2(\mathcal{T}L^e) \dots$$
(6.19)

and using lightsource sampling for the $L^{e*} = (\mathcal{T}L^e)$ integral while sampling the BRDFs when evaluating the $\mathcal{T}^i L^{e*}$ integrals. Practically it means that having hit a surface, one or more shadow rays are traced towards the lightsources and the reflection of the illumination of this point is estimated. This reflection is used as if it were the emission of the surface. This method is particularly efficient if the scene consists of point lightsources. Tracing a single ray to each point lightsource, the illumination due to the point lightsources can be determined exactly (with zero variance).

6.3.4 Importance sampling in colored scenes

So far, we have assumed that the weights containing the BRDFs and the emissions are scalars thus the densities can be made proportional to them. This is only true if the rendering equation is solved on a single wavelength.

However, if color images are needed, the rendering equation should be solved on several (at least 3) different wavelengths. If the different wavelengths are handled completely independently, then the proposed importance sampling strategy can be used without any modifications. However, this approach carries out geometric calculations, such as tracing rays, independently and redundantly for different wavelengths, thus it cannot be recommended. A better approach is using rays that transport light on all wavelengths simultaneously. In this case the emission and the BRDF can be represented by vectors, thus to allow importance sampling, we need a scalar **importance function** \mathcal{I} that is large when the elements in the vector are large and small when the elements are small. The importance is a functional of the spectrum. A straightforward way is using the **luminance** of the spectrum since it emphasizes those wavelengths to which the eye is more sensitive.

6.3.5 Multiple importance sampling

So far, we mentioned two basic importance sampling strategies, the BRDF sampling and the lightsource sampling, which are local in the sense that they focus on a single reflection. It is easy to imagine that if the sampling considers simultaneously many reflections, then the number of possible strategies increases dramatically.

Obviously, we desire to use the best sampling strategy. Unfortunately the performance of a sampling strategy depends on the properties of the scene, which is usually not known a-priori, thus the best strategy cannot be selected. Instead of selecting the best, Veach and Guibas [VG95] proposed to combine several strategies in a way that the strengths of the individual sampling methods are preserved.

Suppose that we can use *n* different sampling techniques for generating random paths, where the distribution of the samples is constructed from several $p_1, ..., p_n$ importance sampling distributions. The number of samples taken from p_i is denoted by M_i , and the total number of samples by $M = \sum_i M_i$. The M_i values are fixed in advance before any samples are taken. The "average probability density" of selecting the sample \mathbf{z} is then

$$\hat{p}(\mathbf{z}) = \sum_{i=1}^{n} \frac{M_i}{M} \cdot p_i(\mathbf{z}).$$
(6.20)

Thus the integral quadrature using these samples is

$$\int_{[0,1]^s} f(\mathbf{z}) \, d\mathbf{z} = \int_{[0,1]^s} \frac{f(\mathbf{z})}{\hat{p}(\mathbf{z})} \cdot \hat{p}(\mathbf{z}) \, d\mathbf{z} \approx \frac{1}{M} \sum_{i=1}^n \sum_{j=1}^{M_i} \frac{f(\mathbf{z}_{i,j})}{\hat{p}(\mathbf{z}_{i,j})} = \sum_{i=1}^n \frac{1}{M_i} \sum_{j=1}^{M_i} w_i(\mathbf{z}_{i,j}) \cdot \frac{f(\mathbf{z}_{i,j})}{p_i(\mathbf{z}_{i,j})}$$
(6.21)

where $\mathbf{z}_{i,j}$ is the *j*th sample taken from the *i*th distribution, and the weights are

$$w_i(\mathbf{z}) = \frac{M_i \cdot p_i(\mathbf{z})}{\sum_{k=1}^n M_k \cdot p_k(\mathbf{z})}.$$
(6.22)

Let us interpret this result when all methods use the same number of samples. $p_i(\mathbf{z})$ is the probability that a sample \mathbf{z} is generated by method *i*. The samples are combined with this weight, which guarantees that no sample will be accounted for twice. In order to have an unbiased estimation, $\sum_i w_i(\mathbf{z}) = 1$ should hold for all \mathbf{z} .

6.4 Handling infinite-dimensional integrals

Walk methods require the evaluation of a series of integrals where the dimension goes to infinity. One way of attacking the problem is truncating the Neumann series, but this introduces some bias which can be quite high if the scene is highly reflective. Fortunately, there is another approach that solves the infinite-dimensional integration problem through randomization. In the context of Monte-Carlo integration, this approach is called the **Russian-roulette** [AK90], but here a somewhat more general treatment is also given that can also justify this approach for quasi-Monte Carlo quadratures.

6.4.1 Russian-roulette

Note that the Neumann series contains a sequence of the following integrals:

$$L = \int_{[0,1]^2} \frac{w}{t} \cdot [L^e + \ldots] d\mathbf{z} = \int_{[0,1]^2} w^*(\mathbf{z}) \cdot L^{\mathrm{in}}(\mathbf{z}) d\mathbf{z} = E\left[w^* \cdot L^{\mathrm{in}}\right]$$

if z is uniformly distributed in $[0, 1]^2$. A Monte-Carlo quadrature would generate random samples in the domain and estimate the integral as an average of the integrand at these samples. Let us further randomize this computation and before each sample let us decide randomly with probability s whether we really evaluate the integrand at the sample point or simply assume that the integrand is zero without any calculations. In order to compensate the not computed terms, when the integrand is really computed, it is divided by probability s. This randomization introduces a new random variable L^{ref} which is equal to $w^* \cdot L^{\text{in}}/s$ if the integrand is evaluated and zero otherwise. The Monte-Carlo quadrature which provides the estimate as an expected value will still be correct:

$$E[L^{\text{ref}}] = s \cdot E\left[L^{\text{ref}} \mid \text{sample is used}\right] + (1-s) \cdot E\left[L^{\text{ref}} \mid \text{sample is not used}\right] = s \cdot E\left[\frac{w^* \cdot L^{\text{in}}}{s}\right] + (1-s) \cdot 0 = E\left[w^* \cdot L^{\text{in}}\right] = L.$$
(6.23)

The variance of the new estimator, on the other hand, is increased:

$$D^{2}[L^{\text{ref}}] = E[(L^{\text{ref}})^{2}] - E^{2}[L^{\text{ref}}] = s \cdot E\left[\left(\frac{L^{\text{ref}}}{s}\right)^{2}\right] + (1-s) \cdot 0 - E^{2}[L^{\text{ref}}] = \left(\frac{1}{s} - 1\right) \cdot E[(w^{*} \cdot L^{\text{in}})^{2}] + D^{2}[w^{*} \cdot L^{\text{in}}].$$
(6.24)

6.4.2 Russian-roulette in quasi-Monte Carlo quadrature

In the context of quasi-Monte Carlo integration Russian-roulette should be explained differently because there is no "randomization" in deterministic techniques. This discussion is also used to generalize the basic concept and to show that termination decision can be made using the complete previous path not only the actual state.

Instead of randomization, we can suppose that the domain of integration is extended by additional variables on which a contribution indicator function is defined that will determine whether or not a higher order term is included in the quadrature (interestingly, in order to get rid of high-dimensional integrals, we increase the dimension of the integration). In order to compensate the missing terms in the integral quadrature, the really computed terms are multiplied by an appropriate factor. If the used contribution indicator is such that the domain where it is non-zero shrinks quickly, then the possibility of evaluating samples of high-dimensional functions is rather low, which saves computation time. However, the integral quadrature will still be correct asymptotically.

A term of the Neumann series has generally the following form

$$L = \int \dots \int W(\mathbf{z}_1, \dots, \mathbf{z}_n) \cdot L^e(\mathbf{z}_1, \dots, \mathbf{z}_n) \, d\mathbf{z}_1 \dots d\mathbf{z}_n, \tag{6.25}$$

where $W(\mathbf{z}_1, \ldots, \mathbf{z}_n)$ is the product of the weights $w_1(\mathbf{z}_1) \cdot \ldots \cdot w_n(\mathbf{z}_n)$ including the cosine functions of the angles and the BRDFs, or the product of the ratios of weights and densities $w_1^* \cdot \ldots \cdot w_n^* = w_1/t_1 \cdot \ldots \cdot w_n/t_n$, depending whether or not BRDF sampling has already been applied.

Let us extend this by a contribution indicator function $C(\mathbf{z}_1, r_1, \dots, \mathbf{z}_n, r_n)$ that is constant 1 if a sample $\mathbf{z}_1, \dots, \mathbf{z}_n$ should be taken into account in the integral quadrature and 0 if it should not and its contribution is assumed to be zero. Usually this function is in separable form

$$C(\mathbf{z}_1, r_1, \dots \mathbf{z}_n, r_n) = \prod_{i=1}^n c_i(\mathbf{z}_i, r_i),$$

where $c_i(\mathbf{z}_i, r_i) = 1$ means that the walk must be continued at step *i* and $c_i(\mathbf{z}_i, r_i) = 0$ forces the walk to stop. Function $c_i(\mathbf{z}_i, r_i)$ can be defined, for instance, by a new weight function $\xi(\mathbf{z}_i)$ in the following way:

$$c_i(\mathbf{z}_i, r_i) = \begin{cases} 1 & \text{if } \xi(\mathbf{z}_i) > r_i, \\ 0 & \text{otherwise.} \end{cases}$$

The "possibility" of really computing a walk $\mathbf{z}_1, \ldots \mathbf{z}_n$ is

$$P(\mathbf{z}_1,\ldots,\mathbf{z}_n) = \int_{r_1=0}^1 \ldots \int_{r_n=0}^1 C(\mathbf{z}_1,r_1,\ldots,\mathbf{z}_n,r_n) dr_1\ldots dr_n$$

We can define the following function of variables r_1, \ldots, r_n ,

$$L_r(r_1,\ldots,r_n) = \int \ldots \int C(\mathbf{z}_1,r_1,\ldots,\mathbf{z}_n,r_n) \cdot \tilde{W} \cdot \tilde{L}^e \, d\mathbf{z}_1\ldots d\mathbf{z}_n, \tag{6.26}$$

where \tilde{W} and \tilde{L}^e are appropriate modifications of W and L^e , which can compensate the missing terms.

The integral of this function is

$$\int_{r_1=0}^{1} \dots \int_{r_n=0}^{1} L_r(r_1, \dots, r_n) dr_1 \dots dr_n =$$

$$\int_{r_1=0}^{1} \dots \int_{r_n=0}^{1} \int \dots \int C(\mathbf{z}_1, r_1, \dots, \mathbf{z}_n, r_n) \cdot \tilde{W} \cdot \tilde{L}^e d\mathbf{z}_1 \dots d\mathbf{z}_n dr_1 \dots dr_n =$$

$$\int \dots \int P(\mathbf{z}_1, \dots, \mathbf{z}_n) \cdot \tilde{W} \cdot \tilde{L}^e d\mathbf{z}_1 \dots d\mathbf{z}_n.$$
(6.27)

A sufficient requirement for this integral to be equal to the original integral L is

$$P(\mathbf{z}_1, \dots, \mathbf{z}_n) \cdot \tilde{W} \cdot \tilde{L}^e = W \cdot L^e.$$
(6.28)

There are many possible selections of the contribution indicator and the \hat{W} and \hat{L}^e functions, that can satisfy this requirement, thus there are many different unbiased estimators. A widely used selection is letting

$$W = 1$$
, $L^e = L^e$ and $\xi_i(\mathbf{z}_i) = w_i(\mathbf{z}_i)$.

which corresponds to continuing the walk after step *i* only if $w(\mathbf{z}_i) > r_i$. If importance sampling has already been applied, then the walk is continued if $w^* > r_i$. Since w^* approximates the albedo, this strategy continues the walk with the probability of the albedo.

BRDF sampling for materials of multiple reflection type

Practical reflection models incorporate different simple BRDFs. For example, a lot of materials can be well modeled by a sum of diffuse and specular reflections. So far, methods have been presented that are good for either the diffuse or the specular reflection, but not for the sum of them.

Fortunately, Russian-roulette can also be extended to handle these cases. If the reflection model is a sum of different BRDFs, then a random selection can be made from the different components. Suppose that the transfer probability density is available in the form of a sum of the weights corresponding to elementary BRDFs:

$$w = w_1 + w_2 + \ldots + w_n.$$

Thus the radiance of a single reflection is:

$$L = \int_{\Omega} w \cdot L^{\text{in}} d\omega = \int_{\Omega} w_1 \cdot L^{\text{in}} d\omega + \ldots + \int_{\Omega} w_n \cdot L^{\text{in}} d\omega.$$

Suppose that mappings t_i can be found for each integral, that also mimics important directions:

$$L = \int_{[0,1]^2} \frac{w_1}{t_1} \cdot L^{\text{in}} \, d\mathbf{z} + \dots + \int_{[0,1]^2} \frac{w_n}{t_n} \cdot L^{\text{in}} \, d\mathbf{z} = E\left[w_1^* \cdot L^{\text{in}}\right] + \dots + E\left[w_n^* \cdot L^{\text{in}}\right].$$

Let us select the *i*th BRDF with probability p_i and weight the resulting radiance by $1/p_i$ or stop the walk with probability $p_0 = 1 - p_1 - \ldots - p_n$. Thus the new random variable L^{ref} is $w_i^* \cdot L^{\text{in}}/p_i$ if the *i*th model is used, and 0 if no model is selected. The expected value of L^{ref} will still be correct:

$$E[L^{\text{ref}}] = p_1 \cdot E\left[\frac{w_1^* \cdot L^{\text{in}}}{p_1}\right] + \dots + p_n \cdot E\left[\frac{w_n^* \cdot L^{\text{in}}}{p_n}\right] + (1 - p_1 - \dots - p_n) \cdot 0 = E\left[(w_1^* + \dots + w_n^*)L^{\text{in}}\right] = L.$$
(6.29)

This is a Monte-Carlo estimation of the sum. According to importance sampling, the variance will be small if $w_i^* L^{in}/p_i$ can be made nearly constant. Since we usually do not have a-priori information about L^{in} , w_i^*/p_i can be made a constant number. Thus to obtain a low-variance estimator, an elementary BRDF should be selected with the probability of its transformed weight w_i^* . Note that the weight w_i^* may either be equal or approximate the albedo, thus a low-variance estimator selects an elementary BRDF with the probability of its albedo.

In order to generate an "out" direction from the "in" direction and the surface normal, the following general BRDF sampling algorithm can be used:

```
BRDFSampling(in, normal, out)
      prob = SelectBRDFModel(normal, in)
      if prob = 0 then return 0
      prob *= Reflection(in, normal, out)
      if prob = 0 then return 0
      return prob
```

end

In this program "SelectBRDFModel" randomly selects an elementary BRDF from those that compose the given BRDF with a probability which approximates the albedo and also returns the selection probability. If it returns 0, then it has decided that the walk has to be stopped because of Russian-roulette. "Reflection" generates a new direction "out" with a probability density that is approximately proportional to the transfer probability density of the selected reflection model and returns the selection probability.



Figure 6.6: Metallic spheres generated by gathering walks with 50 samples per pixel (top: importance sampling according to the cosine angle only (51 min); middle: BRDF sampling with Russian-roulette (46 min); bottom: BRDF sampling with truncating the Neumann series to 5 terms (54 min)



Figure 6.7: Left: Scene with metallic objects rendered by gathering walks using BRDF sampling with Russian roulette. Right: Scene with metallic objects rendered by gathering walks using BRDF sampling without Russian-roulette but truncating the Neumann series to 5 terms.

6.5 Review of random walk algorithms

In this chapter a number of practical random walk algorithms are reviewed and analyzed. For completeness, non-random and non global illumination methods, such as ray-casting and recursive ray-tracing are also included. The primary classification of random walk algorithms is based on the direction of generated rays. If the walks are started at the eye and go opposite to the light, then the algorithm is called **gathering**. On the other hand, if the walks originate at the lightsources and go in the direction of the light, then the algorithm is called **shooting**.

6.5.1 Gathering-type random walk algorithms

Gathering type random walks correspond to the Monte-Carlo solution of the rendering equations. They start at the eye position and gather the emission of the visited points.

The general structure of these algorithms is as follows:

```
for each pixel p do

color = 0

for i = 1 to N do

ray = sample ray randomly from the eye through pixel p

samplecolor = c \cdot \text{Trace(ray)}

color += samplecolor/N

endfor

write(p, color)

endfor
```



Figure 6.8: Gathering-type random walks

In different gathering algorithms the "Trace" function is implemented differently. This function returns the radiance carried by this ray to the eye. The radiance is then multiplied by value $c = (c/S_p) \cdot S_p$ where c/S_p scaling is required by the measuring function (equation (2.39)) and S_p is the size of the integration domain (equation 5.6).

Ray-casting

Ray-casting is a local-illumination algorithm of type LDE, which replaces the unknown radiance inside the integral of the rendering equation by an approximation of the emission function. In its trace function the following simplification is used to determine the radiance of a point:

$$L(\vec{x},\omega) = L^{e}(\vec{x},\omega) + \int_{\Omega} L_{\text{lightsource}}(h(\vec{x},-\omega'),\omega') \cdot f_{r}(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega', \qquad (6.30)$$

where $L_{\text{light source}}$ may be a simplification of the emission function L^e . The "Trace" function is:

```
\begin{array}{l} {\rm Trace(ray)} \\ ({\rm object},\,\vec{x}) = {\rm FirstIntersect(ray)} \\ {\rm if \ no \ intersection \ then \ return \ } L_{\rm sky} \\ {\rm return \ } L^e(\vec{x},\, {\rm -ray.direction}\,) + {\rm DirectLightsource}(\vec{x},\, {\rm -ray.direction}\,) \\ {\rm end} \end{array}
```



Figure 6.9: Ray-casting

In this algorithm L_{sky} is the radiance of background illumination (e.g. sky), "FirstIntersect" is responsible for finding that object which is first intersected by the ray and also the intersection point. "DirectLightsource", on the other hand, computes an estimate of the single reflection of the light from the lightsources, which happens at point \vec{x} into the given direction.

Visibility ray-tracing

Visibility ray-tracing is a recursive ray-tracing like algorithm which can follow multiple light bounces only for ideal reflection and refraction (it is of $L[D]S^*E$ type).



Figure 6.10: Visibility ray-tracing

Formally, it simplifies the rendering equation to the following form:

$$L(\vec{x},\omega) = L^{e}(\vec{x},\omega) + \int_{\Omega} L_{\text{light source}}(h(\vec{x},-\omega'),\omega') \cdot f_{r}(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega' + k_{r}(\omega_{r},\vec{x},\omega) \cdot L(h(\vec{x},-\omega_{r}),\omega_{r}) + k_{t}(\omega_{t},\vec{x},\omega) \cdot L(h(\vec{x},-\omega_{t}),\omega_{t}),$$
(6.31)

where ω_r and ω_t are the ideal reflection and refraction directions, and k_r and k_t are the reflection and refraction coefficients. The implementation of the "Trace" function of the visibility ray-tracing is:

```
Trace(ray)

(object, \vec{x}) = FirstIntersect(ray)

if no intersection then return L_{sky}

color = L^e(\vec{x}, -ray.direction) + DirectLightsource(\vec{x}, -ray.direction)

if k_r > 0 then color += k_r \cdot Trace(reflected ray)

if k_t > 0 then color += k_t \cdot Trace(refracted ray)

return color

end
```

This subroutine calls itself recursively to find the radiance of the illumination at the reflection and refraction directions. In order to avoid infinite recursion, the algorithm is usually extended to limit the maximum level of recursion.

Distributed ray-tracing

Distributed ray-tracing suggested by Cook [CPC84] is a global illumination algorithm, which can model all the possible paths.



Figure 6.11: Distributed ray-tracing

In this method the ray tracing is not terminated when reaching a surface having neither ideal reflection nor ideal refraction. After a ray has hit a surface, child rays are generated randomly according to the BRDF characterizing the surface. For the appropriate estimation of the general interreflection, child rays have to be traced and the average of their contributions have to be computed. This approach is based on the recursive formulation of the integrals in the Neumann series (equation (3.10)). The implementation of the "Trace" function of distributed ray-tracing is:

```
Trace(ray)

(object, \vec{x}) = FirstIntersect(ray)

if no intersection then return L_{sky}

color = L^e(\vec{x}, -ray.direction) + DirectLightsource(\vec{x}, -ray.direction)

for sample = 1 to N do

prob = BRDFSampling(-ray.direction, normal, newray)

if prob > 0 then color += Trace( newray ) · w(newray.direction, normal, -ray.direction) / prob /N

endfor

return color

end
```

In this program "BRDFSampling" — as defined in section 6.4.2 — finds a new ray to follow, which is then traced recursively.

Path-tracing

Another Monte-Carlo approach proposed by Kajiya is **path-tracing** [Kaj86], which is based on the multi-dimensional integral formulation of the terms of the Neumann series (equation (3.7)).



Figure 6.12: Path tracing without (left) and with (right) direct lightsource computation

This method creates a path history for a single particle interacting with the environment until absorption using BRDF sampling and Russian-roulette. Rather than spawning new rays at each intersection, a random direction is chosen according to a density t_i which is approximately proportional to w_i . The walk is continued with a probability $a_i = w_i/t_i$ which is equal to the approximation of the albedo (Russian-roulette). The measured value of a single path is

$$P = c \cdot (L_1^e + L_2^e \cdot \frac{w_1}{t_1 \cdot a_1} + L_3^e \cdot \frac{w_2}{t_2 \cdot a_2} \cdot \frac{w_1}{t_1 \cdot a_1} + \dots)$$

where L_i^e is the emission of the point visited at step *i* of the path and w_i is the transfer density of this point, and *c* is the scaling factor of the measurement device. Note that if ideal BRDF sampling is used, then w_i is proportional to t_i and both w_i/t_i and a_i are equal to the albedo, which results in the following estimate:

$$P = c \cdot (L_1^e + L_2^e + L_3^e + \ldots).$$

This estimate has very high variation if the lightsources are small. This problem can be solved if lightsource sampling is combined with the gathering walk, which means that at each visited point the effects of the lightsources are estimated. The implementation of the "Trace" function of path-tracing is:

```
Trace(ray)

(object, \vec{x}) = FirstIntersect(ray)

if no intersection then return L_{sky}

color = L^e(\vec{x}, -ray.direction) + DirectLightsource(\vec{x}, -ray.direction)

prob = BRDFSampling(-ray.direction, normal, newray)

if prob = 0 then return color

color += Trace( newray ) · w(newray.direction, normal, -ray.direction) / prob

return color

end
```

In this program "BRDFSampling" finds a new direction or if it returns 0, then it has decided that the walk has to be stopped because of Russian-roulette. Note that this algorithm generates all but the last directions of the path by BRDF sampling and the last is obtained by lightsource sampling. Thus if the surface at the light reflection is shiny (close to ideal mirror or ideal refractor), then the quality of importance sampling can be quite bad. Since almost ideal surfaces close to the lightsources are responsible for **caustics**, path tracing as other gathering algorithms are poor in rendering caustics effects.

6.5.2 Shooting-type walk methods

Shooting walks are based on the Monte-Carlo solution of the potential equation. They start at the eye, go through the scene and try to find the eye.

The general structure of shooting algorithms is as follows:

```
Clear Image
DirectCamera
for i = 1 to N do
ray = Sample randomly from a lightsource with selection probability p^e
power = L^e \cdot \cos \theta / p^e / N
Shoot(ray, power)
endfor
```

Function "DirectCamera" calculates the direct contribution of the lightsources on the image by determining the points visible in each pixel and integrating their L^e emission. In different shooting algorithms the "Shoot" function is implemented differently. This function is responsible for determining the power carried to the eye by the complete path, for the identification of the pixel through which the path arrives at the eye, and adding this contribution to the pixel color.



Figure 6.13: Shooting-type walks

Photon tracing

Photon tracing (forward ray-tracing) is the inverse of visibility ray-tracing and uses similar simplifying assumptions.



Figure 6.14: Photon tracing

It also stops tracing when hitting a surface that does not have coherent reflection or refraction. In photon tracing the rays are emitted from the lightsources, and at each hit it is examined whether the surface has ideal reflection, refraction and incoherent reflection or refraction. In the directions of ideal reflection or refraction, the tracing is continued by starting new child rays. The implementation of its Shoot function is:

```
Shoot(ray, power)

(object, \vec{x}) = FirstIntersect(ray)

if no intersection then return

if \vec{x} is visible from pixel p then

color[p] += power \cdot w(ray.direction, <math>\vec{x}, eye direction) \cdot g(\vec{x})

endif

if k_r > 0 then Shoot( reflected ray, k_r \cdot power )

if k_t > 0 then Shoot( refracted ray, k_t \cdot power )

return

end
```

The "eye direction" is a vector pointing from \vec{x} to the eye position. The algorithm is capable of handling LS^*DE paths.

Light-tracing

In **light-tracing** [DLW93] photons perform random walk through the scene starting at the lightsources. Whenever a surface is hit, a ray is traced from the intersection point to the eye and the contribution is added to the selected pixel (if any).



Figure 6.15: Light tracing

Light tracing is the direct implementation of the Monte-Carlo quadrature of the multi-dimensional formulation of the potential equation. When the next direction is determined, BRDF based importance sampling can be applied and combined with Russian-roulette. It chooses a random direction according to a density t_i which is approximately proportional to w_i (importance sampling). The walk is continued with a probability a_i that is equal to the approximation of the albedo (Russian-roulette). The measured value of a single step of the path is

$$P = \frac{L^e \cos \theta}{N \cdot p^e} \cdot \frac{w_1}{t_1 \cdot a_1} \cdot \frac{w_2}{t_2 \cdot a_2} \cdot \ldots \cdot w(eye) \cdot g,$$

if this point is visible at the pixel and zero otherwise. Here L^e is the emission of the starting point, θ is the angle between the surface normal of the lightsource and the first direction, p^e is the probability of selecting this lightsource point and starting direction, w(eye) is the cosine weighted BRDF at the given point from the last direction to the eye, and g is the surface dependent camera parameter. Note that if ideal BRDF sampling is used, i.e. w_i is proportional to t_i and both w_i/t_i and a_i are equal to the albedo, and ideal lightsource sampling is used, i.e. p^e is proportional to $L^e \cos \theta$, thus $L^e \cos \theta / N \cdot p^e = \Phi / N$, then the following estimate can be obtained:

$$P = \frac{\Phi}{N} \cdot w(eye) \cdot g.$$

This estimate has high variation if the camera is often hidden since if the point is not visible from the camera, the contribution is zero. The implementation of the "Shoot" function of light tracing is:

```
Shoot(ray, power)

(object, \vec{x}) = FirstIntersect(ray)

if no intersection then return

if \vec{x} is visible from pixel p then

color[p] += power \cdot w(ray.direction, <math>\vec{x}, eye direction) \cdot g(\vec{x})

endif

prob = BRDFSampling(-ray.direction, normal, newray)

if prob = 0 then return

newpower = power * w( -ray.direction, normal, newray.direction) / prob

Shoot( newray, newpower )

return

end
```

This algorithm also applies BRDF sampling for all but the last steps. The direction of the last visibility ray might be far from the direction preferred by the BRDF. This degrades the performance of importance sampling if the visible surface is very shiny. Thus visible mirrors or refractors (glass) pose difficulties to shooting algorithms.

Random walks for the radiosity setting

Expansion expands the solution of equation $\mathbf{L} = \mathbf{L}^e + \mathbf{R} \cdot \mathbf{L}$ (equation (4.5)) into a discrete Neumann series

$$\mathbf{L} = \mathbf{L}^e + \mathbf{R} \cdot \mathbf{L}^e + \mathbf{R}^2 \cdot \mathbf{L}^e + \mathbf{R}^3 \cdot \mathbf{L}^e + \dots$$
(6.32)

Let us again examine the $\mathbf{R}^2 \cdot \mathbf{L}^e$ term. Using the definition of the matrix \mathbf{R} , this can also be expressed as a multi-dimensional integral:

$$(\mathbf{R}^2 \cdot \mathbf{L}^e)|_i = \sum_{j=1}^n \sum_{k=1}^n \mathbf{R}_{ij} \cdot \mathbf{R}_{jk} \cdot \mathbf{L}_k^e = \int_{S} \int_{\Omega} \int_{S} \int_{\Omega} \tilde{b}_i(\vec{x}_1) \cdot w_1(i) \cdot \sum_{j=1}^n b_j(h(\vec{x}_1, -\omega_1')) \cdot \tilde{b}_j(\vec{x}_2) \cdot w_2(j) \cdot \sum_{k=1}^n b_k(h(\vec{x}_2, -\omega_2')) \cdot \mathbf{L}_k^e \ d\omega_2' d\vec{x}_2 d\omega_1' d\vec{x}_1,$$

where

$$w_1(i) = f_i \cdot \cos \theta'_1, \quad w_2(j) = f_j \cdot \cos \theta'_2.$$

Considering the integrand, \vec{x}_1 should be in patch *i* for \hat{b}_i to be non zero. Then, only a single b_j will give non-zero value for the $\vec{y}_1 = h(\vec{x}_1, -\omega'_1)$ point. To select this, a ray has to be traced from \vec{x}_1 in direction $-\omega'_1$ and the visible patch should be identified. Following this, another point on the identified patch *i* should be selected, which is denoted by \vec{x}_2 , and a ray is traced in direction $-\omega'_2$ to obtain an index *k* of a patch whose emission should be propagated back on the walk. During propagation, the emission is multiplied by the BRDFs (f_i, f_j) and the cosine $(\cos \theta'_2, \cos \theta'_1)$ factors of the visited patches (figure 6.16).

Note that this is basically the same walking scheme, as used to solve the original integral equation. The fundamental difference is that when a patch is hit by the ray, the walk is not continued from the found point but from another uniformly sampled point of the patch.



Figure 6.16: Random walk solution of the projected rendering equation

The power equation $\mathbf{P} = \mathbf{P}^e + \mathbf{H} \cdot \mathbf{P}$ (equation (4.8)) can be treated similarly. Again, let us examine the two-bounce case

$$(\mathbf{H}^{2} \cdot \mathbf{P}^{e})|_{i} = \sum_{j=1}^{n} \sum_{k=1}^{n} \mathbf{R}_{ji} \cdot \frac{f_{i}}{f_{j}} \cdot \mathbf{R}_{kj} \cdot \frac{f_{j}}{f_{k}} \cdot \mathbf{P}_{k}^{e} = \int_{S} \int_{\Omega} \int_{S} \int_{\Omega} \mathbf{P}_{k}^{e} \cdot \tilde{b}_{k}(\vec{y}_{1}) \cdot w_{1}(k) \cdot \sum_{j=1}^{n} b_{j}(h(\vec{y}_{1},\omega_{1})) \cdot \tilde{b}_{j}(\vec{y}_{2}) \cdot w_{2}(j) \cdot \sum_{k=1}^{n} b_{i}(h(\vec{y}_{2},\omega_{2})) \cdot w_{3}(i) \ d\omega_{2}d\vec{y}_{2}d\omega_{1}d\vec{y}_{1},$$

where

$$w_1(k)=\cos heta_1, \hspace{1em} w_2(j)=f_j\cdot\cos heta_2, \hspace{1em} w_3(i)=f_i.$$

It means that the integrand in a single point can be obtained by selecting a point \vec{y}_1 on patch k, then tracing a ray in direction ω_1 . Having identified the intersected patch j a new point \vec{y}_2 is selected on this patch and the ray-tracing is continued at direction ω_2 . The patch which is hit by this ray receives the power of patch k attenuated by the BRDFs and the cosine factors of the steps.

Thus, the projected rendering equation can also be solved by random walks [Shi91b, Sbe96]. The basic difference is that when a patch is hit by a ray, then instead of initiating the next ray from this point, another independent point is selected on the same patch.

Considering the concept of importance sampling and Russian-roulette, many different strategies can be elaborated by appropriately defining the P, \tilde{W} and \tilde{L}^e functions (recall that according to equation (6.28) the requirement of an unbiased estimate is $P(\mathbf{z}_1, \ldots, \mathbf{z}_n) \cdot \tilde{W} \cdot \tilde{L}^e = W \cdot L^e$). For example, let us use the following simulation [Shi91b, Sbe96] to obtain a radiance estimate of patch i_1 :

First a ray is found that starts on this patch. The starting point \vec{x}_1 is sampled from a uniform distribution, while the direction ω'_1 is sampled from a cosine distribution, thus the probability density is $1/A_{i_1} \cdot \cos \theta'_1/\pi$. This ray is traced and the next patch is identified. Let it be patch i_2 . At patch i_2 it is decided whether or not the walk should be stopped with probability of the albedo of the patch. Note that for diffuse surfaces the albedo is $a = f \cdot \pi$. If the walk has to be continued, then a new starting point \vec{x}_2 is found on patch i_2 , and the same procedure is repeated recursively.

With this strategy, the probability density of completing an n step walk is

$$p(\vec{x}_{1}, \omega_{1}', \vec{x}_{2}, \omega_{2}', \dots, \vec{x}_{n-1}, \omega_{n-1}') = \frac{1}{A_{i_{1}}} \cdot \frac{\cos \theta_{1}'}{\pi} \cdot \frac{a_{i_{2}}}{A_{i_{2}}} \cdot \frac{\cos \theta_{2}'}{\pi} \dots \frac{a_{i_{n-1}}}{A_{i_{n-1}}} \cdot \frac{\cos \theta_{n-1}'}{\pi} \cdot (1 - a_{i_{n}}) = \frac{f_{i_{1}}}{A_{i_{1}}} \cdot \cos \theta_{1}' \cdot \frac{f_{i_{2}}}{A_{i_{2}}} \cdot \cos \theta_{2}' \dots \frac{f_{i_{n-1}}}{A_{i_{n-1}}} \cdot \cos \theta_{n-1}' \cdot \frac{1 - a_{i_{n}}}{a_{i_{1}}} = W \cdot \frac{1 - a_{i_{n}}}{a_{i_{1}}}.$$
(6.33)

Thus the required weight \hat{W} of the walk is

$$\tilde{W} = \frac{a_{i_1}}{1 - a_{i_n}}.$$
(6.34)

Thus if the patch on which the walk is terminated is a source having emission \mathbf{L}_n^e , then the estimator of the radiance of patch *i* is

$$\mathbf{L}_n^e \cdot \frac{a_{i_1}}{1 - a_{i_n}}$$
6.5.3 Bi-directional random walk algorithms

Bi-directional algorithms are based on the combination of shooting and gathering walks, thus they can combine the advantages of both techniques. Namely, they can effectively handle small lightsources and small aperture cameras, and can render caustics and ideally refracting or reflecting visible objects.

Bi-directional path-tracing

Bi-directional path-tracing [LW93, VG95] initiates paths at the same time from a selected lightsource and from the viewpoint. After some steps, either a single deterministic shadow ray is used to connect the two types of walks [VG95], or all points of the gathering walk are connected to all points of the shooting walk using deterministic rays [LW93]. If the deterministic shadow ray detects that the two points are occluded from each other, then the contribution of this path is zero.

Note that gathering and shooting walks use different integration variables, namely a gathering walk is specified by a point on the pixel area and a sequence of incoming directions, while a shooting walk is defined by a point on the lightsource and a sequence of the outgoing directions. Thus when the two walks are connected, appropriate transformations should take place.



Figure 6.17: Correspondence between the solid angles of incoming and outgoing directions

Let us first consider a walk of a single bounce (figure 6.17). According to the definition of the solid angle, we obtain

$$\frac{d\omega_1'}{d\omega_2} = \frac{dA \cdot \cos\theta_{\rm out}/r_1^2}{dA \cdot \cos\theta_{\rm in}/r_2^2} = \frac{r_2^2}{r_1^2} \cdot \frac{\cos\theta_{\rm out}}{\cos\theta_{\rm in}},\tag{6.35}$$

and for the substitution of the surface integral on the lightsource

$$d\omega_2' = \frac{d\vec{y} \cdot \cos\theta}{r_2^2}.$$
(6.36)

Thus the transformation rule is

$$\cos\theta_1' \cdot \cos\theta_{\rm in} \ d\omega_1' d\omega_2' = \frac{\cos\theta_1' \cdot \cos\theta_{\rm out}}{r_1^2} \cdot \cos\theta \ d\omega_2 d\vec{y},$$

which means that when converting a shooting type walk to a gathering type walk, then the radiance should be multiplied by

$$\frac{\cos\theta_1'\cdot\cos\theta_{\rm out}}{r_1^2}.$$

When the shooting walk consists of more than 1 steps, then formula (6.35) should be applied to each of them, but formula (6.36) only to the last step. This conversion replaces the incoming directions by the outgoing directions and the subsequent steps compensate r_{k+1}^2/r_k^2 scaling. Finally, we end up with a formula which is similar to the 1-step case:

$$\cos \theta'_k \cdot \cos \theta'_{k+1} \cdot \ldots \cos \theta'_n \, d\omega'_k \ldots d\omega'_n = \frac{\cos \theta'_k \cdot \cos \theta_{n-k+1}}{r_k^2} \cdot \cos \theta_{n-k} \cdot \ldots \cos \theta_1 \, d\omega_{n-k} \ldots d\omega_1 d\vec{y}.$$



Figure 6.18: Bi-directional path tracing with a single deterministic step

Figure 6.18 shows an example when k = 2 and n = 4. This formula means that we can use the rules of sections 3.2.1 and 3.2.2 to generate the shooting and gathering walks — gathering walks use the cosine of the incoming angle, while shooting walks use the cosine of the outgoing angle — and the transformation of the combined walk to a single gathering walk requires a multiplication by

$$\frac{\cos\theta_k'\cdot\cos\theta_{n-k+1}}{r_k^2}$$

Formally, if the endpoints of the shooting and gathering walks are visible from each-other, then the measured value of a single path is

$$P = \frac{L^e \cos \theta}{N \cdot p^e} \cdot \frac{w_1^s}{t_1^s \cdot a_1^s} \cdot \frac{w_2^s}{t_2^s \cdot a_2^s} \cdot \ldots \cdot f_r^s \cdot \frac{\cos \theta^s \cdot \cos \theta^{\prime g}}{r^2} \cdot f_r^g \cdot \ldots \cdot \frac{w_2^g}{t_2^g \cdot a_2^g} \cdot \frac{w_1^g}{t_1^g \cdot a_1^g} \cdot dr$$

where superscripts s and g stand for shooting and gathering steps, respectively. Note that if ideal BRDF sampling is used, then the estimate is:

$$P = \frac{\Phi}{N} \cdot f_r^s \cdot \frac{\cos \theta^s \cdot \cos \theta'^g}{r^2} f_r^g \cdot c.$$



Figure 6.19: Bi-directional path tracing with multiple deterministic steps

In Lafortune's version of the bi-directional path tracing [LW93] not only the endpoints of the shooting and gathering walks are connected, but all intersection points are linked by shadow rays.

Note that in bi-directional path-tracing a path of given length n can be generated by many ways, for instance, by a gathering walk of length i combined with a shooting walk of length n - i, where i = 0, ..., n. This creates a danger of accounting a walk more than one time. To eliminate this danger, the

result can be estimated by a weighted sum of the different walks as suggested by the concept of multiple importance sampling. Other heuristic weight factors can also provide good results [LW93, Vea97]. For example, when a gathering walk is computed, at each step the radiance coming from the continuation of the walk is weighted by W and the radiance coming from deterministic connections to shooting walks is weighted by 1 - W. Since for shiny surfaces the continuation of the walk by BRDF sampling is supposed to be better, W may express how shiny the surface is.

Metropolis light transport

Random walk methods usually generate the ray-paths independently. Thus when a difficult path is found, it is thrown away right after its application. Metropolis method, on the other hand, generates samples by perturbing the previous path, thus they are expected to be better for difficult lighting conditions.

Recall that Metropolis method [MRR⁺53] offers samples with a probability density that is proportional to a given "importance function" in an asymptotic case. Let this importance function \mathcal{I} be the luminance of the light carried by a ray-path to the eye through any pixel. This selection is justified by the fact that the eye is particularly sensitive to luminance variations and different pixels have equal importance in the image. The integral of the importance function on the whole domain is denoted by *b*.



Figure 6.20: Generating walks by mutations in Metropolis light transport

In order to generate samples according to probability density $\mathcal{I}(\mathbf{z})/b$, a Markov process is constructed whose stationary distribution is just this (here \mathbf{z} denotes a ray-path). Scalar b can be estimated in a preprocessing phase using a normal random walk algorithm. The Metropolis algorithm generates a sequence of ray-paths starting from an initial path. Veach and Guibas proposed bi-directional path tracing [VG97] to find an initial path, although any random walk algorithm can be used for this. Generating a new path \mathbf{z}_{i+1} after path \mathbf{z}_i consists of two steps. First the **tentative transition function** $T(\mathbf{z}_i \rightarrow \mathbf{z}_t)$ produces a **tentative path** \mathbf{z}_t by mutating path \mathbf{z}_i a little, i.e. by changing directions and adding or deleting steps. Then the tentative path is either accepted or rejected using an **acceptance probability**

$$a(\mathbf{z}_i \to \mathbf{z}_t) = \min\left\{\frac{\mathcal{I}(\mathbf{z}_t) \cdot T(\mathbf{z}_t \to \mathbf{z}_i)}{\mathcal{I}(\mathbf{z}_i) \cdot T(\mathbf{z}_i \to \mathbf{z}_t)}, 1\right\}$$

that expresses the increase of the importance. The definition of mutations is almost arbitrary if they make the Markov process ergodic. Ergodic processes have stationary distribution, and this distribution is independent of the starting state of the process. Practically it requires that any path of positive power could be generated from any other path after a certain number of perturbations. Thus mutations should modify all features, including directions (or visited points), starting point, length, etc. Furthermore, in order to avoid situations when the process is stuck into a region which is surrounded by regions of zero importance, the tentative transition should take large enough steps to jump over these zero importance regions. Veach [VG97], for example, proposed the generation of a completely new path when the contribution of the tentative path is zero. Summarizing, the **Metropolis light-transport** algorithm is:

Generate an initial ray-path \mathbf{z}_1 using random walk, e.g. bi-directional path tracing for i = 1 to N do Based on the actual ray-path, find another, tentative path \mathbf{z}_t mutating \mathbf{z}_i with $T(\mathbf{z}_i \to \mathbf{z}_t)$ if $\mathcal{I}(\mathbf{z}_t) = 0$ then Generate a completely new path \mathbf{z}_{i+1} from scratch using random walk else $a(\mathbf{z}_i \to \mathbf{z}_t) = (\mathcal{I}(\mathbf{z}_t) \cdot T(\mathbf{z}_t \to \mathbf{z}_i))/(\mathcal{I}(\mathbf{z}_i) \cdot T(\mathbf{z}_i \to \mathbf{z}_t))$ Generate uniformly distributed random number r in [0, 1] // accept with "probability" $a(\mathbf{z}_i \to \mathbf{z}_t)$ if $r < a(\mathbf{z}_i \to \mathbf{z}_t)$ then $\mathbf{z}_{i+1} = \mathbf{z}_t$ else $\mathbf{z}_{i+1} = \mathbf{z}_i$ endif Compute the contribution of the ray-path \mathbf{z}_{i+1} to the affected pixel Multiply this contribution by $b/(\mathcal{I}(\mathbf{z}_{i+1}) \cdot N)$ and accumulate to the pixel endfor

Photon-map

Bi-directional path tracing connects a single gathering walk to a single shooting walk. However, if the effects of all shooting walks could be stored, then when a new gathering walk is computed, it could be connected to all of the shooting walks simultaneously, which can significantly increase the number of samples in the integral quadrature. This is exactly what Jensen [JC95, Jen96, JC98] proposed, also giving the definition of a data structure, called the **photon-map** which can efficiently store the effects of many shooting walks.

A photon map is a collection of photon hits at the end of the paths generated in the shooting phase of the algorithm. The photon-map is organized in a **kd-tree** to support efficient retrieval. A photon hit is stored with the power of the photon on different wavelengths, position, direction of arrival and with the surface normal.



Figure 6.21: Rendering with photon-map

The gathering phase is based on the following approximation of the transport operator:

$$L(\vec{x},\omega') = \int_{\Omega} L(h(\vec{x},-\omega'),\omega') \cdot f_r(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega' = \int_{\Omega} \frac{d\Phi(\omega')}{dA\cos\theta' d\omega'} \cdot f_r(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega' \approx \sum_{i=1}^n \frac{\Delta\Phi(\omega'_i)}{\Delta A} \cdot f_r(\omega'_i,\vec{x},\omega),$$
(6.37)

where $\Delta \Phi(\omega'_i)$ is the power of a photon landing at the surface ΔA from direction ω'_i .

The $\Delta \Phi$ and ΔA quantities are approximated from the photons in the neighborhood of \vec{x} in the following way. A sphere centered around \vec{x} is extended until it contains *n* photons. If at this point the radius of the sphere is *r*, then the intersected surface area is $\Delta A = \pi r^2$ (figure 6.22).



Figure 6.22: Retrieving data from the photon-map

Instant radiosity

Instant radiosity [Kel97] elegantly subdivides the shooting walks into a view-independent walk and to the last bounce which reflects the contribution to the eye. The view-independent walks are calculated in the first phase of the algorithm and each walk results in a point-lightsource that represents the power at the end of the walk. The view-independent walk is quite similar to the light-tracing algorithm, but the new directions are sampled from the Halton sequence instead of a random distribution. In the radiosity setting the reflection at the end of the walks is also diffuse, thus it can be stored as a point-lightsource of homogeneous directional radiation.



Figure 6.23: Instant radiosity

In the second phase of the algorithm the power of the diffuse lightsources are reflected towards the eye. Since this is a simple, local illumination problem, the algorithm can take advantage of the rendering hardware of advanced workstations which can render the effect of these lightsources on the scene and also to compute shadows. If the number of lightsources is more than what can be handled by the hardware, the computation is broken into different phases. Each phase uses just a few of the lightsources and the final image is obtained as the average of the estimates of the phases, which is computed by the hardware accumulation buffer.

Instant radiosity is quite similar to photon-map based techniques. However, instead of using raytracing for final gather, the photons in the photon map are used as lightsources and fast and hardware supported visibility and shadow algorithms are applied. The other fundamental difference is that instant radiosity allows just a relatively low number of photons which therefore should be very well distributed. The optimal distribution is provided by quasi-Monte Carlo light walks.

Chapter 7

Stochastic iteration solution of the global illumination problem

7.1 Why should we use Monte-Carlo iteration methods?

Deterministic iteration has two critical problems. On the one hand, since the domain of the radiance function is 4 dimensional and has usually high variation, an accurate finite-element approximation often requires very many basis functions, which, in turn, need a lot of storage space. On the other hand, when finite element techniques are applied, the transport operator is only approximated, which introduces some non-negligible error in each step. As concluded in section 3.3.1, if the contraction ratio of the operator is λ , then the total accumulated error will be approximately $1/(1 - \lambda)$ times the error of a single step [SKFNC97]. For highly reflective scenes, the iteration is slow and the result is inaccurate if the approximation of the operator is not very precise. Very accurate approximations of the transport operator, however, require a lot of computation time and storage space.

Both problems can be successfully attacked by randomizing the iteration, which is called the **stochastic iteration**. The basic idea of stochastic iteration is that instead of approximating operator \mathcal{T} in a deterministic way, a much simpler random operator is used during the iteration, which "behaves" as the real operator just in the "average" case. When stochastic iteration is applied, the transport operator should be like the real operator just in the average case. As we shall see, it is relatively easy to find random operators whose expected case behavior matches exactly to that of the real operator. Thus the error accumulation problem can be avoided.

On the other hand, if the operator is carefully randomized, it does not require the integrand everywhere in the domain, which allows us not to store the complete radiance function, thus a lot of storage space can be saved. Compared to the astronomical storage requirements of non-diffuse radiosity methods, for example, with stochastic iteration we can achieve the same goal with one variable per patch [SKP98a]. This argument loses some of its importance when view-independent solution is also required, since the final solution should be stored anyway. This is not a problem if only the diffuse case is considered, since using a single radiosity value per patch the image can be generated from any viewpoint. For the non-diffuse case, the reduced storage gets particularly useful when the image is to be calculated in only a single, or in a few eye positions.

Summarizing, the advantages of stochastic iteration are the simplicity, speed, affordable storage requirements and numerical stability even for very large systems containing highly reflective materials.

7.2 Formal definition of stochastic iteration

The concept of stochastic iteration has been proposed and applied for the diffuse radiosity problem in [Neu95, NFKP94, NPT⁺95, SKFNC97], that is for the solution of finite-dimensional linear equations. In this section we generalize the fundamental concepts to solve integral equations [SK98c, SK99c], then the generalized method will be used for attacking non-diffuse global illumination problems.

Suppose that we have a random linear operator \mathcal{T}^* so that

$$E[\mathcal{T}^*L] = \mathcal{T}L \tag{7.1}$$

for any Riemann-integrable function L.

During stochastic iteration a random sequence of operators $\mathcal{T}_1^*, \mathcal{T}_2^*, \ldots, \mathcal{T}_i^*, \ldots$ is generated, which are instantiations of \mathcal{T}^* , and this sequence is used in the iteration:

$$L_m = L^e + \mathcal{T}_m^* L_{m-1}.$$
 (7.2)

Since in computer implementations the calculation of a random operator may invoke finite number of random number generator calls, we are particularly interested in those random operators which have the following construction scheme:

- 1. Random "point" p_i is found from a finite dimensional set Π using probability density $\operatorname{prob}(p)$. This probability density may or may not depend on function L.
- 2. Using p_i a "deterministic" operator $\mathcal{T}^*(p_i)$ is applied to radiance L.

Point p_i is called the **randomization point** since it is responsible for the random nature of operator \mathcal{T}^* . Using a sequence of random transport operators, the measured value

$$P_m = \mathcal{M}L_m \tag{7.3}$$

will also be a random variable which does not converge but fluctuates around the real solution. However, the solution can be found by averaging the estimates of the subsequent iteration steps:

$$P = \lim_{M \to \infty} \sum_{m=1}^{M} P_m.$$
(7.4)

Formally the sequence of the measured values during the iteration is the following:

$$P_{1} = \mathcal{M}L_{1} = \mathcal{M}(L^{e} + \mathcal{T}_{1}^{*}L^{e}),$$

$$P_{2} = \mathcal{M}L_{2} = \mathcal{M}(L^{e} + \mathcal{T}_{2}^{*}L^{e} + \mathcal{T}_{2}^{*}\mathcal{T}_{1}^{*}L^{e}),$$

$$\vdots$$

$$P_{M} = \mathcal{M}L_{M} = \mathcal{M}(L^{e} + \mathcal{T}_{M}^{*}L^{e} + \mathcal{T}_{M}^{*}\mathcal{T}_{M-1}^{*}L^{e} + \mathcal{T}_{M}^{*}\mathcal{T}_{M-2}^{*}L^{e} + \dots).$$

Averaging the first M steps, we obtain:

$$\tilde{P}_{M} = \frac{1}{M} \sum_{i=1}^{M} \mathcal{M}L_{i} = \mathcal{M}(L^{e} + \frac{1}{M} \sum_{i=1}^{M} \mathcal{T}_{i}^{*}L^{e} + \frac{1}{M} \sum_{i=1}^{M-1} \mathcal{T}_{i+1}^{*}\mathcal{T}_{i}^{*}L^{e} + \frac{1}{M} \sum_{i=1}^{M-2} \mathcal{T}_{i+2}^{*}\mathcal{T}_{i+1}^{*}\mathcal{T}_{i}^{*}L^{e} + \dots) = \mathcal{M}(L^{e} + \frac{1}{M} \sum_{i=1}^{M} \mathcal{T}_{i}^{*}L^{e} + \frac{M-1}{M} \cdot \frac{1}{M-1} \sum_{i=1}^{M-1} \mathcal{T}_{i+1}^{*}\mathcal{T}_{i}^{*}L^{e} + \frac{M-2}{M} \cdot \frac{1}{M-2} \sum_{i=1}^{M-2} \mathcal{T}_{i+2}^{*}\mathcal{T}_{i+1}^{*}\mathcal{T}_{i}^{*}L^{e} + \dots).$$
(7.5)

In order to prove that \tilde{P}_M really converges to the solution of the integral equation, first it is shown that the expected value of $\mathcal{T}_{i+k}^* \mathcal{T}_{i+k-1}^* \dots \mathcal{T}_{i+1}^* \mathcal{T}_i^* L^e$ is $\mathcal{T}^{k+1} L^e$. For k = 0, it comes directly from the requirement of equation (7.1). For k = 1, the **total expected value theorem** [Rén62] can be applied:

$$E[\mathcal{T}_{i+1}^*\mathcal{T}_i^*L^e] = \int_{\Pi} E[\mathcal{T}_{i+1}^*\mathcal{T}_i^*L^e | p_{i+1} = p] \cdot \operatorname{prob}(p) \, dp.$$
(7.6)

Since for a fixed $p_{i+1} = p$, operator \mathcal{T}_{i+1}^* becomes a deterministic linear operator, its order can be exchanged with that of the expected value operator:

$$E[\mathcal{T}_{i+1}^*\mathcal{T}_i^*L^e | p_{i+1} = p] = \mathcal{T}_{i+1}^*(p) \left(E[\mathcal{T}_i^*L^e] \right).$$
(7.7)

Using requirement (7.1), the expected value $E[\mathcal{T}_i^*L^e]$ is $\mathcal{T}L^e$, thus we further obtain

$$E[\mathcal{T}_{i+1}^*\mathcal{T}_i^*L^e|p_{i+1}=p] = \mathcal{T}_{i+1}^*(p)(\mathcal{T}L^e).$$
(7.8)

Substituting this to equation (7.6), we get

$$E[\mathcal{T}_{i+1}^*\mathcal{T}_i^*L^e] = \int_{\Pi} \mathcal{T}_{i+1}^*(p)(\mathcal{T}L^e) \cdot \operatorname{prob}(p) \, dp = E[\mathcal{T}_{i+1}^*(\mathcal{T}L^e)] = \mathcal{T}(\mathcal{T}L^e) = \mathcal{T}^2L^e, \tag{7.9}$$

which concludes our proof for the k = 1 case. The very same idea can be used recursively for more than two terms.

Returning to the averaged solution \tilde{P}_M , its expected value is then

$$E[\tilde{P}_M] = \mathcal{M}(L^e + \mathcal{T}L^e + \frac{M-1}{M}\mathcal{T}^2L^e + \frac{M-2}{M}\mathcal{T}^3L^e + \dots + \frac{1}{M}\mathcal{T}^ML^e).$$
(7.10)

Note also that there is some power "defect" ΔP between this expected value and the real solution

$$P = \mathcal{M}(L^e + \mathcal{T}L^e + \mathcal{T}^2L^e + \mathcal{T}^3L^e + \ldots)$$

because of the missing higher order terms for finite M values. Denoting the contraction ratio of the integral operator \mathcal{T} by λ , and assuming that the measuring device is calibrated to show value c for unit homogeneous radiance, this defect can be estimated as follows:

$$\Delta P| = \left| \mathcal{M} \left(\frac{1}{M} \mathcal{T}^2 L^e + \frac{2}{M} \mathcal{T}^3 L^e + \dots \frac{M-1}{M} \mathcal{T}^M L^e + \mathcal{T}^{M+1} L^e + \mathcal{T}^{M+2} L^e + \dots \right) \right| \leq \frac{c\lambda^2}{M} \cdot ||L^e|| \cdot (1 + 2\lambda + 3\lambda^2 + \dots + (M-1)\lambda^{M-2} + M\lambda^{M-1} + M\lambda^M + \dots) = \frac{c\lambda^2}{M} \cdot ||L^e|| \cdot \left[\frac{d}{d\lambda} \left(\sum_{i=1}^{M-1} \lambda^i \right) + M \cdot \frac{\lambda^{M-1}}{1-\lambda} \right] \leq \frac{c}{M} \cdot \frac{\lambda^2}{(1-\lambda)^2} \cdot ||L^e||.$$
(7.11)

This defect converges to zero if the operator is a contraction and M goes to infinity, thus we have

$$\lim_{M \to \infty} E[\tilde{P}_M] = \mathcal{M}(L^e + \mathcal{T}L^e + \mathcal{T}^2L^e + \mathcal{T}^3L^e + \dots).$$
(7.12)

Finally, it must be explained why and when random variable \tilde{P}_M converges to its expected value. Looking at formula (7.5) we can realize that it consists of sums of the following form:

$$\frac{1}{M-k} \cdot \sum_{i=1}^{M-k} \mathcal{T}_{i+k}^* \mathcal{T}_{i+k-1}^* \dots \mathcal{T}_{i+1}^* \mathcal{T}_i^* L^e.$$

According to the theorems of large numbers, and particularly to the Bernstein [Rén62] theorem, these averages really converge to the expected value if the terms in the average are not highly correlated (note that here the terms are not statistically independent as assumed by most of the laws of large numbers). It means that random variables $\mathcal{T}_{i+k}^* \mathcal{T}_{i+k-1}^* \dots \mathcal{T}_i^* L^e$ and $\mathcal{T}_{j+k}^* \mathcal{T}_{j+k-1}^* \dots \mathcal{T}_j^* L^e$ should not have strong correlation if $i \neq j$. This is always true if the operators are generated from independent random variables. To show a negative example, let us assume that there is a very strong correlation between the random operators, namely that different \mathcal{T}_i^* random operators use exactly the same randomization point. When this randomization point is known, the iteration is fully deterministic using operator $\mathcal{T}_i^* = \mathcal{T}^*$. The limiting value of this iteration will be a random variable

$$\mathcal{M}(L^{e} + (\mathcal{T}^{*})L^{e} + (\mathcal{T}^{*})^{2}L^{e} + (\mathcal{T}^{*})^{3}L^{e} + \ldots),$$

which usually differs from the expected solution.

There is another danger that should be considered. Random variable \tilde{P}_M is a sum of random variables whose number goes to infinity. Even if the variances of all single terms in this sum converge to zero, the variance of the sum of these terms might still converge to a non-zero value. In the context of random walks, this phenomenon is called "non-existing variance" [Sbe99, Erm75]. To avoid this case, random operators should not be "*overrandomized*", thus its variance must not exceed a certain limit which depends on the contraction ratio.

7.2.1 Other averaging techniques

In the previous section we solved the problem that stochastic iteration is not convergent by simply averaging the values generated during iteration. There are other averaging schemes, on the other hand, that use even more combinations of the preceding random operators.

Self-correcting iteration

Self-correcting iteration [Neu95] uses the following formulae to derive a new value from the previous one:

$$L'_{m} = L^{e} + \mathcal{T}_{m}^{*}L_{m-1},$$

$$L_{m} = \tau_{m} \cdot L'_{m} + (1 - \tau_{m}) \cdot L_{m-1},$$

$$\tilde{P}_{m} = \mathcal{M}L_{m},$$
(7.13)

where τ_m is an appropriate sequence that converges to 0, as for example, $\tau_m = 1/m$.

To allow comparison, the corresponding formulae of the normal iteration are also presented here:

$$L_{m} = L^{e} + \mathcal{T}_{m}^{*} L_{m-1},$$

$$\tilde{P}_{m} = \tau_{m} \cdot \mathcal{M} L_{m} + (1 - \tau_{m}) \cdot \tilde{P}_{m-1}.$$
(7.14)

Note that the fundamental difference is that self-correction iteration uses the average of the previous samples not only in the final estimate but also to continue iteration. Self-correcting iteration thus can use all combinations of the preceding random operators to compute the actual result. However, it also has energy defect.

7.3 Stochastic iteration for the diffuse radiosity

In the gathering type radiosity algorithms the projected rendering equation (formula (4.5)) has the following form

$$\mathbf{L} = \mathbf{L}^e + \mathbf{R} \cdot \mathbf{L}.$$

Alternatively, shooting radiosity algorithms are based on the projected potential equation (formula (4.8)):

$$\mathbf{P} = \mathbf{P}^e + \mathbf{H} \cdot \mathbf{P}.$$

According to the basic requirement of stochastic iteration we need to find random operators \mathcal{T}_F^* or $\mathcal{T}_F'^*$ that behave as the real operator in average, that is

$$E[\mathcal{T}_F^* \mathbf{L}] = \mathbf{R} \cdot \mathbf{L},$$

$$E[\mathcal{T}_F'^* \mathbf{P}] = \mathbf{H} \cdot \mathbf{P}.$$
(7.15)

The evaluation of $(\mathbf{R} \cdot \mathbf{L})|_i$ or alternatively of $(\mathbf{H} \cdot \mathbf{P})|_i$ requires a surface and a directional integration (or in other formulations two surface integrations).

The possible alternatives for the random transport operator are as follows:

- 1. Both integrals are explicitly computed but only for a randomly selected subset of the patches (stochastic radiosity).
- 2. The surface integral explicitly computed but the directional integral implicitly (stochastic iteration version of the transillumination radiosity).
- 3. Compute the surface integral implicitly but the directional integral explicitly (randomly placed hemicubes).
- 4. Both integrals are computed implicitly (stochastic ray-radiosity).

7.3.1 Stochastic radiosity

In **stochastic radiosity** [NFKP94], the randomized operator is simplified in a sense that it first selects a single (or a few) patches with probability proportional to their power and then calculates the transfer only from this important patch as if it had all the power $\Phi = \sum_{k=1}^{n} \mathbf{P}_{k}$. Thus both integrals are explicitly computed but only for a subset of patches.

To prove that it meets the requirement stated by equation (7.15), let us suppose that patch j has been selected and let us examine the new power of patch i:

$$(\mathcal{T}_{F}^{'*}\mathbf{P})|_{i} = \mathbf{H}_{ij} \cdot \Phi.$$
(7.16)

Since the probability of selecting patch j is \mathbf{P}_{i}/Φ , the expectation of the new power is

$$E[(\mathcal{T}_{F}^{'*}\mathbf{P})|_{i}] = \sum_{j=1}^{n} \mathbf{H}_{ij} \cdot \Phi \cdot \frac{\mathbf{P}_{j}}{\Phi} = \sum_{j=1}^{n} \mathbf{H}_{ij} \cdot \mathbf{P}_{j}$$
(7.17)

which we wanted to prove.

7.3.2 Transillumination radiosity

The transillumination radiosity method [Neu95, SKFNC97] has also a stochastic iteration version. It defines the random transport operator by uniformly selecting D transillumination directions $\omega'_1, \ldots, \omega'_D$ and allowing patches to interact only in these transillumination directions. In order to calculate these interactions, a large discretized window is placed perpendicularly to each transillumination direction and the radiance transfer to a patch is approximated by elementary transfers going through the pixels covering the projection of the patch.

Let us consider a single transillumination direction. Projecting patch A_i onto a plane that is perpendicular to the transillumination direction and then approximating the integral of the incoming radiance here by a discrete sum, we get

$$\int_{A_i} L(h(\vec{x}, -\omega'_d)) \cdot \cos \theta'_d \, d\vec{x} = \int_{A_i^p} L(h(\vec{x}', -\omega'_d)) \cdot \, d\vec{x}' \approx \sum_{P \in A_i^p} \mathbf{L}_{\text{buffer}_d[P]} \cdot \delta P. \tag{7.18}$$

where $\operatorname{buffer}_{d}[P]$ stores the index of that patch which is visible in pixel P in the transillumination direction ω'_{d} from patch *i*, and δP is the size of a pixel of the buffer (figure 7.1).



Figure 7.1: Integration on the transillumination plane

Thus the random transfer operator is

$$(\mathcal{T}_F^* \mathbf{L})|_i = \frac{4\pi \cdot f_i \cdot \delta P}{D} \sum_{d=1}^D \sum_{P \in A_i^p} \mathbf{L}_{\text{buffer}_d[P]}.$$
(7.19)

If the transillumination directions are uniformly distributed and the buffer is uniformly jittered, then the expected value of this operator is

$$E[(\mathcal{T}_F^*\mathbf{L})|_i] = \int_{\Omega} \int_{P} \frac{4\pi \cdot f_i \cdot \delta P}{D} \sum_{d=1}^{D} \sum_{P \in A_i^p} \mathbf{L}_{\text{buffer}_d[P]} \frac{dp}{\delta P} \frac{d\omega_d'}{4\pi} = \frac{1}{D} \sum_{d=1}^{D} \int_{\Omega} \int_{P} \sum_{P \in A_i^p} \mathbf{L}_{\text{buffer}_d[P]} \cdot f_i \, dp \, d\omega_d'.$$

If uniform jittering is applied, then we can usually assume that the discrete approximation of the positional radiance distribution gives back the real distribution in the average case, that is

$$\int_{P} \sum_{P \in A_i^p} \mathbf{L}_{\text{buffer}_d[P]} \, dp = \int_{A_i^p} L(h(\vec{x}', -\omega_d')) \, d\vec{x}'.$$
(7.20)

However, this statement is not always true if incremental polygon filling algorithms are applied [SK98a, SK98c]. Note, for example, that an incremental polygon filling algorithm always generates an approximation whose width and height are at least 1. Thus this assumption is correct if the resolution of the discrete buffer is high enough to project each polygon onto at least a few pixels. Substituting this to the expected value integral, we get

$$E[(\mathcal{T}_F^*\mathbf{L})|_i] = \frac{1}{D} \sum_{d=1}^D \int_{\Omega} \int_{A_i^p} L(h(\vec{x}', -\omega_d')) d\vec{x}' \cdot f_i d\omega_d' = \int_{\Omega} \int_{A_i} L(h(\vec{x}', -\omega')) \cdot f_i \cdot \cos\theta' d\vec{x} d\omega'.$$
(7.21)

Using $L(h(\vec{x}', -\omega')) = \sum_{j=1}^{n} b_j (h(\vec{x}', -\omega')) \cdot \mathbf{L}_j$ and equation (4.21), we can prove that the expectation really gives back the real projected transport operator:

$$E[(\mathcal{T}_F^*\mathbf{L})|_i] = \sum_{j=1}^n \int_{\Omega} \int_{A_i} b_j(h(\vec{x}', -\omega')) \cdot f_i \cdot \cos\theta'_d \, d\vec{x} d\omega' \cdot \mathbf{L}_j = \sum_{j=1}^n \mathbf{R}_{ij} \cdot \mathbf{L}_j.$$
(7.22)

7.3.3 Randomly placed hemicubes

This method can evaluate the directional integral explicitly by hemicubes, for instance. To simulate Monte-Carlo surface integration, the center of the hemicube is selected randomly on the patch using a uniform distribution.

7.3.4 Stochastic ray-radiosity

Stochastic ray-radiosity [NPT⁺95] approximates the transport operator by M random rays that are sampled proportionally to the power of the patches. On a patch the starting point of the ray is sampled using a uniform distribution, while the direction follows a cosine distribution. A single ray carries Φ/M power. Thus this method approximates both integrals implicitly.

Let us examine the case when a single ray is selected (since different rays are sampled from the same distribution, the effect of M rays will be M times the effect of a single ray in the expected value). Suppose that patch j is selected as a shooting patch. The probability of the selection event is \mathbf{P}_j/Φ . Thus the probability density of selecting a point \vec{x} of a patch and a direction ω is

$$\frac{\mathbf{P}_j}{\Phi} \cdot \frac{1}{A_j} \cdot \cos \theta.$$

This transfers Φ/M power to the patch that is hit by the ray where the reflected power is computed. Thus the random transport operator for a single ray is

$$E[(\mathcal{T}_{F}^{'*}\mathbf{P})|_{i}] = M \cdot f_{i} \cdot \sum_{j=1}^{n} \iint_{A_{j}} \iint_{\Omega} b_{i}(h(\vec{y},\omega)) \cdot \frac{\Phi}{M} \cdot \frac{1}{A_{j}} \cdot \cos\theta \ d\vec{y}d\omega \cdot \frac{\mathbf{P}_{j}}{\Phi} = \sum_{j=1}^{n} \mathbf{H}_{ij} \cdot \mathbf{P}_{j}.$$
(7.23)

7.4 Definition of the random transport operator for the non-diffuse finite-element case

When moving towards the non-diffuse case, another requirement should be imposed upon the random transport operator. It must not only meet the requirement of equation (7.1) and be easy to compute, but it must also allow the compact representation of the \mathcal{T}_i^*L functions. This extra requirement is evident

if we take into account that unlike in the diffuse case, the domain of L is a 4-dimensional continuous space, so is the domain of \mathcal{T}_i^*L . From the point of view of compactness, what we have to avoid is the representation of these functions over the complete domain.

Thus those transport operators are preferred, which require the value of L just in a few "domain points" (e.g. in a single "domain point"). Note that the evaluation of \mathcal{T}_i^*L now consists of the following steps: first a randomization point p_i is found to define random operator \mathcal{T}_i^* , which in turn determines at which domain point the value of L is required. Up to now, we have had complete freedom to define the set of randomization points. One straightforward way is defining this set to be the same as (or a superset of) the domain of the radiance function and using random transport operators that require the value of the radiance function at their randomization points. Although this equivalence is not obligatory, it can significantly simplify the computations, since when the randomization point is generated, the required domain point is also known.

Using random operators that evaluate the radiance in a single point is not enough in itself, since even a single "point" can result in a continuous \mathcal{T}_i^*L function, which must be stored and re-sampled in the subsequent iteration step and also by the measurement. The solution is the postponing of the complete calculation of \mathcal{T}_i^*L until it is known where its value is needed in the next iteration step and by the measuring device. In this way, the random operator should be evaluated twice but just for two points. Once for the actual and the previous "points" resulting in $[\mathcal{T}^*(p_i)L(p_i)](p_{i+1})$, and once for p_{eye} which is needed by the measuring device and for previous point providing $[\mathcal{T}^*(p_i)L(p_i)](p_{eye})$. The complete iteration goes as follows:

$$\begin{split} P &= 0\\ \text{Find } p_1 \text{ randomly}\\ L(p_1) &= L^e(p_1)\\ \text{for } i &= 1 \text{ to } M \text{ do}\\ P^{\text{new}} &= L^e(p_{\text{eye}}) + [\mathcal{T}^*(p_i)L(p_i)](p_{\text{eye}})\\ P &= \mathcal{M}P^{\text{new}} \cdot 1/i + (1 - 1/i) \cdot P\\ \text{Find } p_{i+1} \text{ randomly}\\ L(p_{i+1}) &= L^e(p_{i+1}) + [\mathcal{T}^*(p_i)L(p_i)](p_{i+1})\\ \text{endfor}\\ \text{Display final image} \end{split}$$

// initialize the measured value to zero // select the randomization point of the first iteration

// measure the radiance // average the measured value // select the randomization point of the next iteration // a single step of stochastic iteration

7.4.1 Single ray based transport operator

The continuous formulation has just a single directional integral, thus a random transport operator can evaluate this single integral implicitly. This results in a method that uses a "single" random walk to obtain the solution.

More specifically, let the random transport operator use a single ray having random origin \vec{y}_i and direction ω_i generated with a probability that is proportional to the cosine weighted radiance of this point at the given direction.

This ray transports the whole power

$$\Phi = \int\limits_{S} \int\limits_{\Omega} L(ec{y},\omega') \cos heta_{ec{y}} \, d\omega' \, dec{y}$$

to that point \vec{x} which is hit by the ray. Formally, the random transport operator is

$$(\mathcal{T}^*L)(\vec{x},\omega) = \Phi \cdot \delta(\vec{x} - h(\vec{y}_i,\omega_i)) \cdot f_r(\omega_i,\vec{x},\omega).$$
(7.24)

Let us prove that this random operator meets the requirement of equation (7.1). The probability density of selecting surface point \vec{y} and direction ω' is

$$\frac{d\Pr\{\vec{y},\omega'\}}{d\vec{y}\,d\omega_{\vec{y}}} = \frac{L(\vec{y},\omega')\cdot\cos\theta_{\vec{y}}}{\Phi}$$
(7.25)



Figure 7.2: Symmetry of solid angles of shooting and gathering

Using the definition of the solid angle $d\omega_{\vec{y}} = d\vec{x} \cdot \cos \theta'_{\vec{x}} / |\vec{y} - \vec{x}|^2$ we can obtain a symmetry relation (figure 7.2) for the shooting and gathering solid angles:

$$d\vec{y} \cdot d\omega_{\vec{y}} \cdot \cos\theta_{\vec{y}} = d\vec{y} \cdot \frac{d\vec{x} \cdot \cos\theta'_{\vec{x}}}{|\vec{y} - \vec{x}|^2} \cdot \cos\theta_{\vec{y}} = d\vec{x} \cdot \frac{d\vec{y} \cdot \cos\theta_{\vec{y}}}{|\vec{y} - \vec{x}|^2} \cdot \cos\theta'_{\vec{x}} = d\vec{x} \cdot d\omega'_{\vec{x}} \cdot \cos\theta'_{\vec{x}}.$$
 (7.26)

Thus the probability of selecting \vec{y}, ω' can also be expressed in the following way:

$$d\Pr\{\vec{y},\omega'\} = \frac{L(\vec{y},\omega')\cdot\cos\theta_{\vec{y}}}{\Phi} \cdot d\vec{y} \, d\omega_{\vec{y}} = \frac{L(h(\vec{x},-\omega'),\omega')\cdot\cos\theta_{\vec{x}}}{\Phi} \cdot d\vec{x} \, d\omega'_{\vec{x}}.$$
(7.27)

Now we can easily prove that the random transport operator meets requirement (7.1) since

$$E[(\mathcal{T}^*L)(\vec{x},\omega)] = \int_{S} \int_{\Omega} \Phi \cdot \delta(\vec{x} - h(\vec{y},\omega')) \cdot f_r(\omega',\vec{x},\omega) \, d\Pr\{\vec{y},\omega'\} = \int_{\Omega} L(h(\vec{x},-\omega'),\omega') \cdot \cos\theta'_{\vec{x}} \cdot f_r(\omega',\vec{x},\omega) \, d\omega'_{\vec{x}} = (\mathcal{T}L)(\vec{x},\omega).$$
(7.28)

Note that the result of the application of the random operator can be a single point that receives all the power and reflects some part of it or the result can be no point at all if the ray leaves the scene.

Suppose that the first random operator \mathcal{T}_1^* is applied to L^e which may transfer all power

$$\Phi_1 = \int_{S} \int_{\Omega} L^e(\vec{y}_1, \omega_1) \cos \theta_{\vec{y}_1} \ d\omega_1 \ d\vec{y}_1$$

to a single point $\vec{x}_1 = h(\vec{y}_1, \omega_1)$ using probability density

$$\frac{d\Pr_1\{\vec{y}_1,\omega_1\}}{d\vec{y}_1d\omega_1} = \frac{L^e(\vec{y}_1,\omega_1)\cdot\cos\theta_{\vec{y}_1}}{\Phi}$$

Before continuing with the second step of the iteration, the radiance should be measured, that is, an image estimate should be computed from $L^e + \mathcal{T}_1^* L^e$. We can separately calculate the effect of the lightsources on the image and then add the effect of $\mathcal{T}_1^* L^e$. Note that $\mathcal{T}_1^* L^e$ is concentrated in a single point, thus its contribution can be computed by tracing a ray from the eye to this point, and if this point is not occluded, then adding $f_r(\omega_1, \vec{x}, \omega_{\text{eye}}) \cdot \Phi$ to that pixel in which \vec{x} is visible.

The second operator \mathcal{T}_2^* should be applied to $L_1 = L^e + \mathcal{T}_1^* L^e$, thus both the total power and the probability density have been modified:

$$\Phi_2 = \int_{S} \int_{\Omega} L_1(\vec{y}_2, \omega_2) \cos \theta_{\vec{y}_2} \ d\omega_2 \ d\vec{y}_2 = \Phi_1 \cdot (1 + a_{\vec{x}_1}(\omega_1))$$

where $a_{\vec{x}_1}$ is the **albedo** at point \vec{x}_1 defined by $a_{\vec{x}}(\omega) = \int_{\Omega} f_r(\omega, \vec{x}, \omega') \cos \theta'_{\vec{x}} d\omega'$, and the new probability density is

$$\frac{d\Pr_2\{\vec{y}_2,\omega_2\}}{d\vec{y}_2d\omega_2} = \frac{L_1(\vec{y}_2,\omega_2)\cdot\cos\theta_{\vec{y}_2}}{\Phi} = \frac{L^e(\vec{y}_2,\omega_2)\cdot\cos\theta_{\vec{y}_2} + f_r(\omega_1,\vec{y}_2,\omega_2)\cos\theta_{\vec{y}_2}\cdot\delta(\vec{y}_2-\vec{x}_1)}{\Phi_1(1+a_{\vec{x}_1}(\omega_1))}.$$

Sampling according to this mixed, discrete-continuous probability density can be realized in the following way. First it is decided randomly whether we sample L^e or the newly generated point using probabilities $1/(1 + a_{\vec{x}_1}(\omega_1))$ and $a_{\vec{x}_1}(\omega_1)/(1 + a_{\vec{x}_1}(\omega_1))$, respectively. If L^e is selected, then the sampling process is the same as before, i.e. a random point and a random direction are found with probability density

$$rac{L^e(ec{y}_2,\omega_2)\cos heta_{ec{y}_2}}{\Phi_1}$$

However, if the new point is chosen, then the direction of the next transfer is found with probability density

$$\frac{f_r(\omega_1, \vec{y}_2, \omega_2) \cos \theta_{\vec{y}_2}}{a_{\vec{x}_1}(\omega_1)}$$

In either case, a ray defined by the selected point and direction is traced, and the complete power $\Phi_2 = \Phi_1 \cdot (1 + a_{\vec{x}_1}(\omega'_1))$ is transferred to that point which is hit by the ray. The subsequent steps of the iteration are similar.

Interestingly this iteration is a sequence of variable length random walks, since at each step the point that is last hit by the ray is only selected with a given probability as the starting point of the next ray. To understand how it works, first let us assume that the environment is closed, that is, rays always hit objects. The algorithm selects a point from a lightsource and then starts a random walk. A step is computed by sending a ray to the randomly selected direction. At the first hit the image contribution is computed. To prepare for the next step, the power to be transferred is set to $(1 + a_1)\Phi$. The walk finishes after the first step according to the powers, that is with probability $1/(1 + a_1)$ and continues with probability $a_1/(1 + a_1)$. If a walk finishes, another walk is initiated from the lightsource, if not, then the ray is emanated from the previous hit point. The new ray is traced, the image contribution is computed, and the power is set again to $(1 + a_2(1 + a_1))\Phi$. Thus now the walk is continued with

$$\frac{a_2+a_2a_1}{\mathsf{l}+a_2+a_2a_1}$$

probability. Note that the limiting case of this continuation probability is a special average of the albedos:

$$\Pr\{\text{continue at step } n\} = \tilde{a} = \frac{a_n + a_n a_{n-1} + a_n a_{n-1} a_{n-2} + \dots + a_n \dots + a_1}{1 + a_n + a_n a_{n-1} + a_n a_{n-1} a_{n-2} + \dots + a_n \dots + a_1}$$

If all points have the same albedo a, then in the limiting case $\tilde{a} = a$, thus the continuation probability is similar to what is usually used in expansion. However, a difference still remains. When expansion continues a walk, it scales the power by $1/a_i$ according to Russian-roulette and resets the scaling factor to 1 when a new walk is initiated. In iteration, however, the power is scaled only by a_i , but the emission power is added at each reflection, and this is not reinitialized when a new lightsource point is sampled.

Now let us also consider open environments where a ray can leave the space with positive probability. When this happens, the power will be the emission power and the algorithm starts to build up the walk from scratch again. Generally, the continuation probability will be in between $\tilde{a}/(1+\tilde{a})$ and \tilde{a} depending on how open the scene is.

Note that even this simple example highlighted the fundamental difference between iteration and expansion. Expansion — i.e. random walk — methods usually continue the walk with the probability of the local albedo, that is, they make decisions based on local features. Iteration, on the other hand, takes into account the available knowledge about the radiance function — it makes decision using global information — which eventually results in a continuation probability which depends on the average albedo and on how open the scene is. Since these factors determine the contraction of the light transport operator, we can state that iteration uses a more precise approximation of the contraction to randomly terminate the walk.

Chapter 8

Simulating light transport using ray-bundles

This chapter proposes expansion and iteration solutions of the global illumination problem, in which ray-bundles are used to simulate the radiance transport.

Walk methods proposed so far use individual ray-paths as samples of the integrand of the rendering equation. However, ray-shooting may waste a lot of computation by ignoring all the intersections but the one closest to the eye. Thus it seems worth using a set of **global directions** [Sbe96, Neu95, SMP98, Sbe97] for the complete scene instead of solving the visibility problem independently for different points \vec{x} . Moreover, ray-shooting is a simple but by no means the most effective visibility algorithm since it is unable to take advantage of image or object coherence. Other methods based on the exploitation of image coherence, such as the z-buffer, painter's, Warnock's, etc. algorithms can be considered as handling a bundle of parallel rays and solving the visibility problem for all of them simultaneously. Continuous (also called object-precision) methods can even determine the visibility problem independently of the resolution, which corresponds to tracing infinitely many parallel rays simultaneously.

The set of parallel global rays is called the ray-bundle.

Using ray-bundles, we have to realize that even single-ray techniques use recursive ray-tracing to simulate multiple interreflections. Thus ray-bundles should also be traced several times in different directions to model multiple interreflections. This tracing composes a **walk** using ray-bundles in each step.

In order to make this method work, three problems must be solved. Firstly, the radiance on the wavefront of the ray-bundle must be represented, for which finite-elements are used. Secondly, the directions should be selected in a way that all the possible light-paths are covered and the integral quadrature should be accurately approximated. The application of random or low-discrepancy series on the directional sphere is proposed to solve this problem. Thirdly, efficient algorithms are needed that can compute the radiance transfer of all patches in a single direction, for which the generalization of discrete and continuous visibility algorithms are applied.

8.1 **Reformulation of the rendering equation using finite-elements**

In order to represent the radiance that is transferred by a ray-bundle, a 2-dimensional function is needed, which is defined by finite-elements. Using finite-element concepts, the radiance function is searched in the following form:

$$L(\vec{x},\omega) \approx L^{(n)}(\vec{x},\omega) = \sum_{j=1}^{n} L_j(\omega) \cdot b_j(\vec{x})$$
(8.1)

where $L^{(n)}(\vec{x}, \omega)$ is the approximating radiance and $b_j(\vec{x})$ is a complete function system. In this function space, the scalar product of two functions f, g is defined as the integral of their products on the total

surface S:

$$\langle f,g\rangle = \int_{S} f(\vec{x}) \cdot g(\vec{x}) \, d\vec{x}. \tag{8.2}$$

Since the radiance is approximated in a subspace, we cannot expect the radiance approximation to satisfy the original rendering equation everywhere. Instead, equality is required in an appropriate subspace defined by **adjoint basis** functions $\tilde{b}_1(\vec{x}), \tilde{b}_2(\vec{x}), \dots, \tilde{b}_n(\vec{x})$ (figure 4.1). This set is called adjoint of $b_1(\vec{x}), b_2(\vec{x}), \dots, b_n(\vec{x})$ since we require that $\langle b_i(\vec{x}), \tilde{b}_j(\vec{x}) \rangle = 1$ if i = j and 0 otherwise. Projecting the rendering equation into the subspace of $\tilde{b}_1(\vec{x}), \tilde{b}_2(\vec{x}), \dots, \tilde{b}_n(\vec{x})$ we obtain

$$\langle L^{(n)}(\vec{x},\omega), \tilde{b}_i(\vec{x}) \rangle = \langle L^e(\vec{x},\omega), \tilde{b}_i(\vec{x}) \rangle + \langle \int_{\Omega} L^{(n)}(h(\vec{x},-\omega'),\omega') \cdot f_r(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega', \tilde{b}_i(\vec{x}) \rangle.$$
(8.3)

Using the orthogonal property of the basis and adjoint basis functions, we get

$$L_{i}(\omega) = L_{i}^{e}(\omega) + \sum_{j=1}^{n} \int_{\Omega} L_{j}(\omega') \cdot \langle b_{j}(h(\vec{x}, -\omega')) \cdot f_{r}(\omega', \vec{x}, \omega) \cdot \cos \theta', \tilde{b}_{i}(\vec{x}) \rangle \, d\omega'.$$
(8.4)

The same equation can also be presented in a matrix form:

$$\mathbf{L}(\omega) = \mathbf{L}^{e}(\omega) + \int_{\Omega} \mathbf{T}(\omega', \omega) \cdot \mathbf{L}(\omega') \, d\omega', \tag{8.5}$$

where $\mathbf{L}(\omega)|_i = L_i(\omega)$ is the vector of radiance values, and

$$\mathbf{T}(\omega',\omega)|_{ij} = \langle f_r(\omega',\vec{x},\omega) \cdot b_j(h(\vec{x},-\omega')) \cdot \cos\theta', \tilde{b}_i(\vec{x}) \rangle$$

is the **bi-directional transport matrix**.

Assume that the BRDF function $f_r(\omega', \vec{x}, \omega)$ can be well approximated by $\tilde{f}_i(\omega', \omega)$ inside the support of \tilde{b}_i (if the support of these basis functions is small, this is always possible). This allows for the separation of the transport matrix to a diagonal matrix $\mathbf{F}(\omega', \omega)$ of BRDF functions

$$\mathbf{F}(\omega',\omega)|_{ii} = \tilde{f}_i(\omega',\omega),$$

and to a **geometry matrix** $\mathbf{A}(\omega')$ that is independent of direction ω :

$$\mathbf{A}(\omega')|_{ij} = \langle b_j(h(\vec{x}, -\omega')) \cdot \cos\theta', \tilde{b}_i(\vec{x}) \rangle.$$
(8.6)

The geometry matrix contains a scalar product of basis functions at points that are visible from each-other in direction ω' . Thus it expresses the strength of coupling as the degree of visibility.

Using the geometry matrix, equation (8.5) can also be written as

$$\mathbf{L}(\omega) = \mathbf{L}^{e}(\omega) + \int_{\Omega} \mathbf{F}(\omega', \omega) \cdot \mathbf{A}(\omega') \cdot \mathbf{L}(\omega') \, d\omega'.$$
(8.7)

This equation — that can be called as the projected rendering equation — is highly intuitive. The radiance of a patch is the sum of the emission and the reflection of all incoming radiances. The role of the patch-direction-patch "form-factor matrix" is played by $\mathbf{A}(\omega')$.

8.2 Stochastic expansion using ray bundles

This section discusses an expansion algorithm for the solution of equation (8.7). The algorithm takes samples of the global ray-bundle walks and uses them in the quadrature. A single walk starts by selecting a direction either randomly or quasi-randomly, and the emission transfer of all patches is calculated into this direction (figure 8.1). Then a new direction is found, and the emission is transferred and the irradiance generated by the previous transfer is reflected from all patches into this new direction. The



Figure 8.1: A path of ray-bundles

algorithm keeps doing this for a few times depending on how many bounces should be considered, then the emission is sent and the irradiance caused by the last transfer is reflected towards the eye. Averaging these contributions results in the final image. When the radiance reflection is calculated from the previous direction to the current direction or to the direction of the eye, the radiance is attenuated by the BRDF of the corresponding surface element.

Concerning the memory requirements of the method, each patch holds the irradiance of the last step of the walk and the accumulated radiance towards the eye. Since the selected directions are the same for all surfaces, they must be stored only once. Consequently the memory requirement is comparable to that of the diffuse radiosity algorithms although it is also capable to handle specular reflections.

Now, let us formally examine the steps of the method. In order to simplify the notations, the integral operator of the projected rendering equation is denoted by T_F :

$$\int_{\Omega} \mathbf{T}(\omega',\omega) \cdot \mathbf{L}(\omega') \, d\omega' = \mathcal{T}_F \mathbf{L}(\omega).$$
(8.8)

Thus the short form of the projected rendering equation is:

$$\mathbf{L}(\omega) = \mathbf{L}^{e}(\omega) + \mathcal{T}_{F}\mathbf{L}(\omega).$$
(8.9)

In equation (8.9) the unknown radiance function $\mathbf{L}(\omega)$ appears on both sides. Substituting the right side's $\mathbf{L}(\omega)$ by the complete right side, which is obviously $\mathbf{L}(\omega)$ according to the equation, we get:

$$\mathbf{L}(\omega) = \mathbf{L}^{e}(\omega) + \mathcal{T}_{F}(\mathbf{L}^{e}(\omega) + \mathcal{T}_{F}\mathbf{L}(\omega)) = \mathbf{L}^{e}(\omega) + \mathcal{T}_{F}\mathbf{L}^{e}(\omega) + \mathcal{T}_{F}^{2}\mathbf{L}(\omega).$$
(8.10)

Repeating this step m times, the original equation can be expanded into a Neumann series:

$$\mathbf{L}(\omega) = \sum_{i=0}^{m} \mathcal{T}_{F}^{i} \mathbf{L}^{e}(\omega) + \mathcal{T}_{F}^{m+1} \mathbf{L}(\omega).$$
(8.11)

If integral operator \mathcal{T}_F is a contraction, that is if $||\mathcal{T}_F \mathbf{L}(\omega)|| \leq \lambda \cdot ||\mathbf{L}(\omega)||$, for some $\lambda < 1$, then $\mathcal{T}_F^{m+1}\mathbf{L}$ converges to zero, thus

$$\mathbf{L}(\omega) = \sum_{i=0}^{\infty} \mathcal{T}_F^i \mathbf{L}^e(\omega).$$
(8.12)

The contractive property of T_F comes from the fact that a reflection or refraction always decreases the energy. Using, for example, the infinite norm, we obtain

$$||\mathcal{T}_F \mathbf{L}(\omega)||_{\infty} \le \max_{\vec{x}} \int_{\Omega} f_r(\omega', \vec{x}, \omega) \cdot \cos \theta' \, d\omega' \cdot ||\mathbf{L}(\omega)||_{\infty} = \max_{\vec{x}} a_{\vec{x}}(\omega) \cdot ||\mathbf{L}(\omega)||_{\infty},$$

where $a_{\vec{x}}(\omega)$ is the **albedo** of the material at point \vec{x} . For physically plausible material models, the albedo must be less than 1.

The terms of this infinite Neumann series have intuitive meaning as well: $\mathcal{T}_F^0 \mathbf{L}^e(\omega) = \mathbf{L}^e(\omega)$ comes from the emission, $\mathcal{T}_F^1 \mathbf{L}^e(\omega)$ comes from a single reflection (called 1-bounce), $\mathcal{T}_F^2 \mathbf{L}^e(\omega)$ from two reflections (called 2-bounces), etc. Using the definition of integral operator \mathcal{T}_F , the full form of the Neumann series is:

$$\mathbf{L}(\omega) = \mathbf{L}^{e}(\omega) + \int_{\Omega} \mathbf{T}(\omega_{1}', \omega) \cdot \mathbf{L}^{e}(\omega_{1}') \, d\omega_{1}' + \int_{\Omega} \int_{\Omega} \mathbf{T}(\omega_{1}', \omega) \cdot \mathbf{T}(\omega_{2}', \omega_{1}') \cdot \mathbf{L}^{e}(\omega_{2}') \, d\omega_{2}' d\omega_{1}' + \dots$$
(8.13)

In practice the infinite sum of the Neumann series is approximated by a finite sum. The number of required terms is determined by the contraction λ of operator \mathcal{T}_F — that is the overall reflectivity and the openness of the scene. Let us denote the maximum number of calculated bounces by D. The truncation of the Neumann series introduces a bias in the estimation, which can be tolerated if D is high enough. On the other hand, this bias can even be eliminated by randomly terminating the walk in the sense of Russian-roulette.

In order to simplify the notations, we introduce the *d*-bounce irradiance J_d for d = 1, 2, ... as follows:

$$\begin{aligned} \mathbf{J}_0 &= \mathbf{A}(\omega'_D) \cdot \mathbf{L}^e(\omega'_D), \\ \mathbf{J}_d &= 4\pi \cdot \mathbf{A}(\omega'_{D-d}) \cdot \mathbf{F}(\omega'_{D-d+1}, \omega'_{D-d}) \cdot \mathbf{J}_{d-1} \end{aligned}$$

where \mathbf{J}_d is a d + 1 dimensional function of directions $(\omega'_{D-d}, \omega'_{D-d+1}, \dots, \omega'_D)$. The *d*-bounce irradiance represents the irradiance arriving at each patch, that is emitted from a patch and is bounced exactly *d* times.

Similarly, we can define the **max** d-bounce irradiance I_d for d = 1, 2, ... as follows:

$$\mathbf{I}_0 = \mathbf{A}(\omega'_D) \cdot \mathbf{L}^e(\omega'_D), \mathbf{I}_d = \mathbf{A}(\omega'_{D-d}) \cdot (\mathbf{L}^e(\omega'_{D-d}) + 4\pi \cdot \mathbf{F}(\omega'_{D-d+1}, \omega'_{D-d}) \cdot \mathbf{I}_{d-1})$$

where I_d is a d + 1 dimensional function of directions $(\omega'_{D-d}, \omega'_{D-d+1}, \dots, \omega'_D)$. The max *d*-bounce irradiance represents the irradiance arriving at each patch after completing a path of length *d* following the given directions, and gathering and then reflecting the emission of the patches visited during the path.

Limiting the analysis to at most D+1 bounces, the solution of the rendering equation can be obtained as a 2D-dimensional integral:

$$\mathbf{L}(\omega) = \left(\frac{1}{4\pi}\right)^{D} \cdot \int_{\Omega} \dots \int_{\Omega} \left[\mathbf{L}^{e}(\omega) + 4\pi \cdot \mathbf{F}(\omega_{1}', \omega) \cdot \mathbf{I}_{D}\right] d\omega_{D}' \dots d\omega_{1}'.$$
(8.14)

This high-dimensional integral (2*D* is 10 to 20 in practical cases) can be evaluated by numerical quadrature. The properties of the integrand and the dimension determine which quadrature rules are applicable for the calculation. Since classical quadrature rules, such as the trapezoidal rule or Gaussian quadrature, are not appropriate for the evaluation of high-dimensional integrals due to their dimensional explosion, Monte-Carlo techniques are proposed. Since the integrand is square integrable (in L_2), the variance of the Monte-Carlo estimate will converge to zero. Furthermore, it is also continuous and is of finite variation in the sense of Hardy and Krause, thus quasi-Monte Carlo techniques are also appropriate [SKFNC97].

8.2.1 Generating uniformly distributed points on the sphere

Since random or low-discrepancy sequences are defined in the unit interval, the directional sphere must be mapped into the unit square, or alternatively, the points in the unit square should be projected onto the directional hemisphere. In order to generate a uniform distribution on the surface of the directional sphere, the following transformation can be used [Sob91]:

$$\phi = 2\pi u, \quad \theta = \arccos(1 - 2v). \tag{8.15}$$

In this case

$$d\omega = det \begin{bmatrix} \frac{\partial \phi}{\partial u} & \frac{\partial \phi}{\partial v} \\ \\ \frac{\partial \theta}{\partial u} & \frac{\partial \theta}{\partial v} \end{bmatrix} \cdot \sin \theta \ du \ dv = 4\pi \cdot du \ dv. \tag{8.16}$$

thus the spatial distribution is uniform on the sphere if it was in the unit square.

8.2.2 Simple Monte-Carlo, or quasi-Monte Carlo walks

In order to evaluate formula (8.14), M random or quasi-random walks should be generated. When the D-bounce irradiance is available, it is multiplied by the BRDF defined by last direction ω_1 and viewing direction ω to find a Monte-Carlo estimate of the radiance that is visible from the eye position. Note that this step makes the algorithm view-dependent.

The final image is the average of these estimates. The complete algorithm — which requires just one variable for each patch *i*, the max *d*-bounce irradiance I[i] — is summarized in the following:

for m = 1 to M do // samples of global walks Generate $(\omega_1^{(m)}, \omega_2^{(m)}, \dots, \omega_D^{(m)})$ // $\omega_i^{(m)}$ is a uniformly distributed direction $\mathbf{I} = 0$ for d = 0 to D - 1 do $\mathbf{I} = \mathbf{A}(\omega'_{D-d}) \cdot (\mathbf{L}^e(\omega'_{D-d}) + 4\pi \cdot \mathbf{F}(\omega'_{D-d+1}, \omega'_{D-d}) \cdot \mathbf{I})$ // a walk Calculate the image estimate from the irradiance \mathbf{I} Divide the estimate by M and add to the Image endfor Display Image

8.2.3 Calculation of the image estimate

A walk generates a sample of the irradiance field on each patch of the scene. From this irradiance field an image estimation can be obtained taking into account the position of the camera. There are basically two different methods to calculate the image estimate. On the one hand, evaluating the BRDF once for each patch, a radiance value is assigned to them, then in order to avoid "blocky" appearance, bi-linear smoothing is applied. Using Phong interpolation, on the other hand, the radiance is evaluated at each point visible through a given pixel using the irradiance field, the surface normal and the BRDF of the found point. Phong interpolation is more time consuming but the generated image is not only numerically precise, but is also visually pleasing.

Bi-linear interpolation

A simple way of the image generation is to assign a single radiance value to each patch using the normal vector of the patch and the viewing direction that is computed for the center of the patch. Thus the image estimate is calculated as follows:

// do it for each patch

for each patch p do Calculate viewing direction ω from patch p $\mathbf{L}[p] += \mathbf{L}^{e}[p](\omega) + 4\pi \cdot \tilde{f}_{p}(\omega_{1}^{(m)}, \omega) \cdot \mathbf{I}[p]$ endfor

However, the constant color of the surfaces results in the annoying effect of faceted objects, since the eye psychologically accentuates the discontinuities of the color distribution. To create the appearance of smooth surfaces, the tricks of **Gouraud shading** can be applied to replace the jumps of color by linear changes. In contrast to Gouraud shading as used in incremental methods, in this case vertex

colors are not available to form a set of knot points for interpolation. These vertex colors, however, can be approximated by averaging the colors of adjacent polygons (see figure 8.2). Assuming triangular patches, the algorithm of image creation is:

```
for each vertex v do

L_v[v] = 0

adjacent_patches[v] = 0

endfor

for each patch p do

for each vertex v of patch p do

L_v[v] += L[p]

adjacent_patches[v]++

endfor

for each vertex v do L_v[v] /= adjacent_patches[v]

for each patch p do

Find vertices of patch p: v_1, v_2, v_3

GouraudShading(v_1, v_2, v_3, L_v[v_1], L_v[v_2], L_v[v_3])

endfor
```



Figure 8.2: Bi-linear interpolation

Although bi-linear interpolation results in continuous color, but the derivative of the color is still discontinuous. Unfortunately, the human eye is sensitive to these discontinuous derivatives, which leads to **Mach-banding**. On the other hand, highly specular surfaces may have strongly non-linear color, which can hardly be approximated by linear functions. This can produce noticeable artifacts and can miss important highlights. The improved interpolation scheme proposed in the following section can solve these problems.

Phong integpolation

In Phong interpolation, the surface normal inside the patches are interpolated from the surface normals at the vertices and the BRDF is evaluated for every pixel using the irradiance information. If the normal vectors are not available in the vertices, then they can be computed as an average normal vectors of patches adjacent to this vertex. The algorithm of the image estimation is:

```
for each pixel P of the image

Find point \vec{x} which is visible in this pixel

Determine the viewing direction \omega, the visible patch s and the surface normal \vec{N} at this point \vec{x}

\mathbf{L}[P] = \mathbf{L}^{e}[s](\omega) + 4\pi \cdot f_{r}(\omega_{1}^{(m)}, \vec{x}, \vec{N}, \omega) \cdot \mathbf{I}[s]

endfor
```

This procedure can be speeded up if in the preprocessing phase the surface visible at each pixel and its surface normal are determined and stored in a map.

8.2.4 Improved walking techniques

The walk algorithm that has been directly derived from the quadrature formulae uses direction ω_1 to evaluate the contribution of 1-bounces, directions (ω_1, ω_2) for the 2-bounces, $(\omega_1, \omega_2, \omega_3)$ for the 3-bounces, etc. This is just a little fraction of information what can be gathered during the complete walk. Why do not we use the samples of $\omega_1, \omega_2, \omega_3$, etc. to calculate the 1-bounce contribution, (ω_1, ω_2) , $(\omega_1, \omega_3), \ldots, (\omega_2, \omega_3)$, etc. combinations of directions for 2-bounces, etc. as shown in figure 8.3. This is possible, since if the samples of $(\omega_1, \omega_2, \ldots, \omega_D)$ are taken from a uniform sequence, then all combinations of its elements also form uniform sequences in lower dimensional spaces.



Figure 8.3: Combined and bi-directional walks

If all possible combinations are used, then each walk generates $\binom{D}{k}$ number of samples for the kbounces, which can be used to increase the accuracy. Note that the increased accuracy is for free in terms of additional visibility computation. However, due to the dependence of the BRDF functions onto two directions and due to the fact that different bounces will be estimated by different number of samples and thus cannot be summed, the required storage per patch is increased to D(D + 1)/2 + 1 number of variables. Since D is 5–8 in practical cases, this storage overhead is affordable.

Now each patch is represented by a triangle matrix \mathcal{I} , where the (i, j) element stores the sum of those *i*-bounce irradiances where the last direction is ω_j . Table 8.1 shows an example for D = 3. The complete algorithm of **combined walking** is shown below:

for
$$m = 1$$
 to M do
Generate $(\omega_1^{(m)}, \omega_2^{(m)}, \dots, \omega_D^{(m)})$ // $\omega_i^{(m)}$ is a uniformly distributed direction
 $\mathcal{I} = 0$ // a walk
 $\mathcal{I}[1][D - d] = \mathbf{A}(\omega_{D-d}^{(m)}) \cdot \mathbf{L}^e(\omega_{D-d}^m)$
for $b = 2$ to $d + 1$ do
 $\mathcal{I}[b][D - d] = 0$
for $pd = 0$ to $d - 1$ do $\mathcal{I}[b][D - d] += 4\pi \cdot \mathbf{A}(\omega_{D-d}^{(m)}) \cdot \mathbf{F}(\omega_{D-pd}^{(m)}, \omega_{D-d}^{(m)}) \cdot \mathcal{I}[b - 1][D - pd]$
endfor
endfor
for each patch p do // accumulation of the contributions
Calculate viewing direction ω
for $d = 0$ to $D - 1$ do
for $b = 1$ to $d + 1$ do $\mathbf{L}[p] += (\mathbf{L}^e(\omega)[p] + 4\pi \cdot \tilde{f}_p(\omega_{D-d}^{(m)}, \omega) \cdot \mathcal{I}[b][D - d]/{\binom{D}{b}})/M$
endfor
endfor
Render Image

last direction	1-bounce	2-bounce	3-bounce
3	$\mathbf{J}(\omega_3)$		
2	$\mathbf{J}(\omega_2)$	$\mathbf{J}(\omega_3,\omega_2)$	
1	$\mathbf{J}(\omega_1)$	$\mathbf{J}(\omega_3,\omega_1)+\mathbf{J}(\omega_2,\omega_1)$	$\mathbf{J}(\omega_3,\omega_2,\omega_1)$

Table 8.1: Irradiance matrix of a patch for D = 3

Bi-directional walking techniques are based on the recognition that when computing geometry matrix $\mathbf{A}(\omega')$ for some direction, the matrix for the reverse direction $\mathbf{A}(-\omega')$ is usually also available paying very little or no additional effort.

On the other hand, since $\mathbf{A}(\omega')_{ij}$ represents that ratio of the radiance emitted or reflected in direction ω' by patch *i* which is actually received by patch *j*, $\mathbf{A}(\omega')_{ij}$ can be non zero only if patch *i* and patch *j* are facing towards each other and patch *i* is facing in direction ω' . Thus at most one value from the $\mathbf{A}(\omega')_{ij}$ and $\mathbf{A}(-\omega')_{ij}$ can be non zero. It means that bi-directional techniques do not even require additional storage and a single geometry matrix can be used to store the values for both ω' and $-\omega'$. It can be decided whether a matrix element is valid for ω' or $-\omega'$ by inspecting the angle between its normal vector and the given directions.



Figure 8.4: Error of different walking techniques. The test scene was the homogeneous, closed Cornell-box, where all surfaces have diffuse reflection of albedo 0.5 and diffuse emission of intensity 0.5 (section 3.4)

Formally, in bi-directional technique $\omega_1, -\omega_1, \omega_2, -\omega_2$, etc. are used to calculate the 1-bounce contribution, $(\omega_1, \omega_2), (\omega_1, -\omega_2), (-\omega_1, \omega_2), (-\omega_1, -\omega_2), (\omega_1, \omega_3), (\omega_1, -\omega_3)$, etc. combinations of directions for 2-bounces, etc. This multiplies the number of samples used for the computation of 1-bounces by 2, for the 2-bounces by 4 and generally for the *D* bounces by 2^D , which can be quite significant. As for the combined walking technique, these additional samples can be used in the numerical quadrature, since if the samples of $(\omega_1, \omega_2, \dots, \omega_D)$ are taken from a uniform sequence, then all combinations of its elements and its inversions also form uniform sequences.

The improved walking techniques use more samples to estimate the 1-bounce, 2-bounce, etc. contributions than the normal technique, thus we can expect an increase of the accuracy (figure 8.4). However, if we use quasi-Monte Carlo integration, this increase is less than what could be observed for Monte-Carlo quadrature. The reason is that low-discrepancy sequences are so "well-designed" to be as uniform as possible, that mixing two different sequences cannot improve the discrepancy much further.

8.2.5 *D*-step iteration

The presented approach truncates the Neumann series after D steps, which introduces some bias. This bias can be eliminated using a simple correction of the emission function L^e when calculating higher order interreflections.

Note that a global walk of length D provides the following terms:

$$L^e + \mathcal{T}_1^* L^e + \mathcal{T}_{(1,2)}^* L^e \dots + \mathcal{T}_{(1,D)}^* L^e,$$

where

$$\mathcal{T}^*_{(i,j)} = \mathcal{T}^*_j \mathcal{T}^*_{j-1} \dots \mathcal{T}^*_{i+1} \mathcal{T}^*_i.$$

Thus having computed the first walk, we also have an estimate for $\mathcal{T}^*_{(1,D)}L^e = \mathcal{T}^*_D\mathcal{T}^*_{D-1}\dots\mathcal{T}^*_2\mathcal{T}^*_1L^e$. Let us use this estimate to correct the emission function in the higher order terms when the second walk is computed:

$$L^{e} + \mathcal{T}^{*}_{D+1}(L^{e} + \mathcal{T}^{*}_{(1,D)}L^{e}) + \dots + \mathcal{T}^{*}_{(D+1,2D)}(L^{e} + \mathcal{T}^{*}_{(1,D)}L^{e}) =$$

$$L^{e} + \mathcal{T}^{*}_{D+1}L^{e} + \dots + \mathcal{T}^{*}_{(D+1,2D)}L^{e} + \mathcal{T}^{*}_{(1,D+1)}L^{e} + \dots + \mathcal{T}^{*}_{(1,2D)}L^{e}.$$
(8.17)

This gives us estimates not only for the bounces from 0 to D but also for the bounces from D + 1 to 2D. Again the last-bounce will store $\mathcal{T}^*_{(1,D)}L^e + \mathcal{T}^*_{(1,2D)}L^e$, which can be used to compensate the emission. Thus after the second step we have estimates for the 0 to 3D bounces. Asymptotically, this method will generate estimates for all bounces. However, if M global walks are generated, then the number of estimates for bounces of 0 to D is M, for bounces of D + 1 to 2D is M - 1, for bounces 2D + 1 to 3Dis M - 2 etc., which still results in some small energy defect.

8.3 Importance sampling for the evaluation of directional integrals

As for other Monte-Carlo methods, importance sampling is a primary candidate to improve the presented methods. Importance sampling means that samples of the integral quadrature are selected from a non-uniform probability density which emphasizes those samples that make significant contribution to the integral.

Let us consider walks of length at most D. In order to generate the resulting image, the radiance of the visible surfaces should be evaluated for each pixel p. A single walk can be characterized by the vector of the directions of individual steps, that is by

$$\mathbf{z} = (\omega'_1, \omega'_2, \dots, \omega'_D).$$

This vector can be regarded as a point in a space of walks and also as a sample in the domain of the integration. A single walk can contribute to the powers of all pixels and at all representative wavelengths. Thus the function is not scalar, but rather a vector. In order to establish an importance-driven algorithm, first of all an importance function should be defined, which can rank different contributions made by different walks.

Let the **importance function** \mathcal{I} be the sum of luminances of all pixels of the image estimate, resulting from the walk. This importance function will really concentrate on those walks that have significant influence on a few pixels or some influence on a lot of pixels, and does not waste computation time on those walks that have little contribution to just a few pixels. Using the **luminance** information is justified by the fact that the human eye is more sensitive to luminance variations than to color variations.

8.3.1 Application of the VEGAS algorithm

The **VEGAS algorithm** [Lep80] can be applied to gather information about the importance sequences of directions. Let us approximate the probability density in the following product form:

$$p(\omega_1, \omega_2, \dots, \omega_D) \propto g_1(\omega_1) \cdot g_2(\omega_2) \cdot \dots \cdot g_D(\omega_D).$$
 (8.18)

It can be shown [Lep80] that the optimal selection of g_1 is

$$g_1(\omega_1) = \sqrt{\int \dots \int \frac{\mathcal{I}^2(\omega_1, \dots, \omega_D)}{g_2(\omega_2) \dots g_D(\omega_D)}} \, d\omega_2 \dots d\omega_D, \tag{8.19}$$

and similar formulae apply to g_2, \ldots, g_D . These g_1, \ldots, g_D functions can be tabulated as 2-dimensional arrays (note that a single direction is defined by 2 scalars, ϕ and θ). The (i, j) element of this matrix represents the importance of the directional region where

$$\phi \in [\frac{2i\pi}{N}, \frac{2(i+1)\pi}{N}), \quad \theta \in [\frac{j\pi}{N}, \frac{(j+1)\pi}{N}).$$

This immediately presents a recursive importance sampling strategy. The algorithm is decomposed into phases consisting of a number of samples. At the end of each phase weights g_1, \ldots, g_D are refined, to provide a better probability density for the subsequent phase. Assuming that g_1, \ldots, g_D are initially constant, a standard Monte-Carlo method is initiated, but in addition to accumulating to compute the integral, g_1, \ldots, g_D are also estimated using equation (8.19). Then for the following phase, the samples are selected according to the g functions. In order to calculate a sample for ω_i , for instance, a single random value is generated in the range of 0 and the sum of all elements in the array defining g_i . Then the elements of the array is retained one by one and summed to a running variable. When this running variable exceeds the random sample, then the searched directional region is found. The direction in this region is then found by uniformly selecting a single point from the region.

VEGAS method is not optimal in the sense that the probability density can only be approximately proportional to the importance even in the limiting case since only product form densities are considered.

We have to mention that the original VEGAS method used 1-dimensional g functions, but in our case, the 2 scalars defining a single direction are so strongly correlated, thus it is better not to separate them. In theory, higher dimensional tables could also be used, but this would pose unacceptable memory requirements.

8.3.2 Application of Metropolis sampling

The Metropolis method achieves the sampling according to a probability that is proportional to the importance function \mathcal{I} by establishing a Markov process on the space of global walks, whose limiting distribution is proportional to the selected importance function. The integration problem is formulated as follows:

$$\mathbf{L}(\omega) = \left(\frac{1}{4\pi}\right)^D \int_{\Omega} \int_{\Omega} \int_{\Omega} \dots \int_{\Omega} \mathbf{L}^e(\omega) + 4\pi \cdot \mathbf{F}(\omega'_1, \omega) \cdot \mathbf{I}_D(\omega'_1, \omega'_2, \dots, \omega'_D) \, d\omega'_D \dots d\omega'_1 = \int_V f(\mathbf{z}) \, d\mathbf{z}.$$

The Metropolis algorithm [MRR⁺53] converges to the optimal probability density that is proportional to the importance, that is in the limiting case $\mathcal{I}(\mathbf{z}) = b \cdot p(\mathbf{z})$. However, this probability density cannot be stored, thus in the Monte-Carlo formula the importance should be used instead, in the following way:

$$\mathbf{L}(\omega) = \int_{V} f(\mathbf{z}) \, d\mathbf{z} = \int_{V} \frac{f(\mathbf{z})}{\mathcal{I}(\mathbf{z})} \cdot \mathcal{I}(\mathbf{z}) \, d\mathbf{z} = b \cdot \int_{V} \frac{f(\mathbf{z})}{\mathcal{I}(\mathbf{z})} \cdot p(\mathbf{z}) \, d\mathbf{z} = b \cdot E\left[\frac{f(\mathbf{z})}{\mathcal{I}(\mathbf{z})}\right] \approx \frac{b}{N} \cdot \sum_{i=1}^{N} \frac{f(\mathbf{z}_{i})}{\mathcal{I}(\mathbf{z}_{i})}.$$
 (8.20)

Scalar b of equation (8.20) has also an intuitive meaning. It is actually the integral of the importance function on the whole domain since

$$\int_{V} \mathcal{I}(\mathbf{z}) \, d\mathbf{z} = \int_{V} \frac{\mathcal{I}(\mathbf{z})}{p(\mathbf{z})} \cdot p(\mathbf{z}) \, d\mathbf{z} = \int_{V} b \cdot p(\mathbf{z}) \, d\mathbf{z} = b.$$
(8.21)

In order to generate samples according to $p(\mathbf{z}) = 1/b \cdot \mathcal{I}(\mathbf{z})$, a Markov process is constructed whose stationary distribution is just $p(\mathbf{z})$. Informally, the next state \mathbf{z}_{i+1} of this process is found by letting an

almost arbitrary **tentative transition function** $T(\mathbf{z}_i \to \mathbf{z}_t)$ generate a **tentative sample** \mathbf{z}_t by mutating the actual sample \mathbf{z}_i . Then, the tentative sample is either accepted as the real next state or rejected making the next state equal to the actual state using an **acceptance probability** $a(\mathbf{z}_i \to \mathbf{z}_t)$ that expresses the increase of the importance:

$$a(\mathbf{z}_i \to \mathbf{z}_t) = \min \left\{ \frac{\mathcal{I}(\mathbf{z}_t) \cdot T(\mathbf{z}_t \to \mathbf{z}_i)}{\mathcal{I}(\mathbf{z}_i) \cdot T(\mathbf{z}_i \to \mathbf{z}_t)}, 1 \right\}.$$

When we use Metropolis sampling in the solution of the global illumination problem, the "state" z corresponds to a complete walk. Mutation strategies are responsible for changing the walk a "little", by perturbing one or more directions.

Variance reduction

Metropolis method may ignore calculated function values if their importance is low. However, these values can be used to reduce variance. Suppose that the importance degrades at step *i*. Thus the process is in \mathbf{z}_t with probability $a(\mathbf{z}_i \rightarrow \mathbf{z}_t)$ and in \mathbf{z}_i with probability $1 - a(\mathbf{z}_i \rightarrow \mathbf{z}_t)$. In order to compute the integral quadrature, random variable $f(\mathbf{z}_{i+1})/\mathcal{I}(\mathbf{z}_{i+1})$ is needed. A common variance reduction technique is to replace a random variable by its mean, thus we can use

$$(1 - a(\mathbf{z}_i \to \mathbf{z}_t)) \cdot \frac{f(\mathbf{z}_i)}{\mathcal{I}(\mathbf{z}_i)} + a(\mathbf{z}_i \to \mathbf{z}_t) \cdot \frac{f(\mathbf{z}_t)}{\mathcal{I}(\mathbf{z}_t)}$$

in the integral quadrature instead of $f(\mathbf{z}_{i+1})/\mathcal{I}(\mathbf{z}_{i+1})$.

Mutation strategies

The statespace of the Markov process consists of D-dimensional vectors of directions that define the sequence of directions in the global walks. Thus the tentative transition function is allowed to modify one or more directions in these sequences.



Figure 8.5: Mutation strategy

The set of possible sequences of directions can be represented by a 2D-dimensional unit cube (each direction is defined by two angles). In the actual implementation random mutations are used that are uniformly distributed in a 2D dimensional cube of edge-size s. In order to find the optimal extent of the random perturbation, several, contradicting requirements must be taken into consideration. First of all, in order to cover the whole statespace of unit size, mutations cannot be very small. Small mutations also emphasize the start-up bias problem which is a consequence of the fact that the Markov process only converges to the desired probability density (this phenomenon will be examined in detail later). On the other hand, if the mutations are large, then the Markov process "forgets" which regions are important, thus the quality of importance sampling will decrease. Finally, another argument against small mutations is that it makes the subsequent samples are statistically independent, which guarantees that if the standard deviation of random variable $f(\mathbf{z})$ is σ , then the standard deviation of the Monte-Carlo quadrature will

be σ/\sqrt{M} after evaluating M samples. Since Metropolis method uses statistically correlated samples, the standard deviation of the quadrature can be determined using the Bernstein theorem [Rén62], which states that it can be upperbounded by

$$\sigma \cdot \sqrt{\frac{1+2\sum_{k=1}^{M} R(k)}{M}}$$
(8.22)

where R(k) is an upperbound of the correlation between $f(\mathbf{z}_i)$ and $f(\mathbf{z}_{i+k})$. It means that strong correlation also increases the variance of the integral estimate.

Generating an initial distribution and automatic exposure

The Metropolis method promises to generate samples with probabilities proportional to their importance in the stationary state. Although the process converges to this probability from any initial distribution as shown in figure 8.8, the samples generated until the process is in the stationary state should be ignored.

The initial bias problem can be reduced if several Metropolis sequences are generated starting from an initial distribution that mimics the stationary distribution and their estimates are averaged. Selecting samples with probabilities proportional to the importance can be approximated in the following way. A given number of seed points are found in the set of sequences of global directions. The importances of these seed points are evaluated, then, to simulate the distribution following this importance, the given number of initial points are selected randomly from these seed points using the discrete distribution determined by their importance.

Equation (8.20) also contains an unknown b constant that expresses the luminance of the whole image. The initial seed generation can also be used to determine this constant. Then at a given point of the algorithm the total luminance of the current image — that is the sum of the importances of the previous samples — is calculated and an effective scaling factor is found that maps this luminance to the expected one.

8.3.3 Evaluation of the performance of the Metropolis method

To evaluate the efficiency of Metropolis sampling for ray-bundle tracing, the scene of figure 8.6 has been selected. The surfaces have both diffuse and specular reflections and the lightsource is well hidden from the camera (figure 8.6), thus this scene is rather difficult to render by other algorithms.



Figure 8.6: A test scene

The error measurements of the Metropolis method with different perturbation sizes, and for quasi-Monte Carlo samples are shown by figure 8.7.

Considering the performance of the Metropolis method for our algorithms, we have to conclude that for homogeneous scenes, it cannot provide significant noise reduction compared to quasi-Monte Carlo walks. This is due to the fact that the variation of the integrand in equation (8.5) is modest. The combined and bi-directional walking techniques cause even further smoothing which is good for the quasi-Monte Carlo but bad for the Metropolis sampling. On the other hand, the number of samples was quite low (we used a few thousand samples). For so few samples the Metropolis method suffers from the problems of initial bias and correlated samples. Due to the smooth integrand, the drawbacks are not compensated by the advantages of importance sampling.



Figure 8.7: Error measurements for the "difficult scene" (left) and the image rendered by Metropolis walks (right)



Figure 8.8: Convergence of the first-bounce as computed by 100, 500, 1000 and 5000 Metropolis samples

Evaluation of the start-up bias

In order to theoretically evaluate the start-up bias, let us examine a simplified, 1-dimensional case when the importance is constant, thus the transition proposed by the tentative transition function is always accepted. In this case, the probability density in the equilibrium is constant. The question is how quickly the Metropolis method approaches to this constant density (figure 8.9).



Figure 8.9: Convergence of Metropolis method to the uniform distribution

Metropolis method can generate samples following a given probability density in a closed interval. Since random mutations may result in points that are outside the closed interval, the boundaries should be handled in a special way. If the variable of an integrand denotes "angle of direction", then the interval can be assumed to be "circular", that is, the external points close to one boundary are equivalent to the internal points of the other boundary. Using this assumption, let us suppose that the domain of the integration is $[-\pi, \pi]$ and the integrand is periodic with 2π . Let the probability distribution at step n be p_n . The Metropolis method is initiated from a single seed at 0, thus $p_0 = \delta(x)$. Assume that transition probability $P(y \to x)$, which is equal to the tentative transition probability for constant importance, is homogeneous, that is $P(y \to x) = P(x - y)$. Using the total probability theorem, the following recursion can be established for the sequence of p_n :

$$p_{n+1}(x) = \int_{-\infty}^{\infty} p_n(y) \cdot P(x \to y) \, dy = \int_{-\infty}^{\infty} p_n(y) \cdot P(x-y) \, dy = p_n * P, \tag{8.23}$$

where * denotes the convolution operation. Applying Fourier transformation to this iteration formula, we can obtain:

$$p_{n+1}^* = p_n^* \cdot P^*, \tag{8.24}$$

where $p_{n+1}^* = \mathcal{F}p_{n+1}$, $p_n^* = \mathcal{F}p_n$ and $P^* = \mathcal{F}P$. Since the domain is "circular", i.e. x denotes the same sample point as $x + 2k\pi$ for any integer k, the probability density is periodic, thus it can be obtained as a Fourier series:

$$p_n(x) = \sum_{k=-\infty}^{\infty} a_k^{(n)} e^{jkx},$$
(8.25)

where $j = \sqrt{-1}$. The Fourier transform is thus a discrete spectrum:

$$p_n^*(f) = \sum_{k=-\infty}^{\infty} a_k^{(n)} \cdot \delta(f-k).$$
(8.26)

Substituting this into equation (8.24), we get

$$p_{n+1}^{*}(f) = \left(\sum_{k=-\infty}^{\infty} a_{k}^{(n)} \cdot \delta(f-k)\right) \cdot P^{*}(f) = \sum_{k=-\infty}^{\infty} a_{k}^{(n)} \cdot P^{*}(k) \cdot \delta(f-k),$$
(8.27)

thus $a_k^{(n+1)} = a_k^{(n)} \cdot P^*(k)$. Using the same concept *n* times, and taking into account that the initial distribution is $\delta(x)$, we can obtain:

$$p_n^*(f) = \sum_{k=-\infty}^{\infty} \left(P^*(k)\right)^n \cdot \delta(f-k)$$
(8.28)

thus in the original domain

$$p_n(x) = \sum_{k=-\infty}^{k=\infty} (P^*(k))^n \cdot e^{jkx}.$$
(8.29)

The L_2 error between p_n and the stationary distribution is then

$$|p_n - p_{\infty}||_2 = \sqrt{\int_0^1 |p_n(x) - a_0^{(\infty)}|^2} \, dx.$$
(8.30)

Note that according to the definition of the Fourier series

$$a_0^{(n)} = \frac{1}{2\pi} \cdot \int_{-\pi}^{\pi} p_n(x) \, dx = 1 \tag{8.31}$$

independently of *n*, thus $a_0^{(\infty)}$ is also 1.

Using this and substituting equation (8.29) in equation (8.30), we get the following error for the distribution:

$$||p_n - p_{\infty}||_2 = \sqrt{\sum_{k=-\infty, k \neq 0}^{k=\infty} |P^*(k)|^{2n}}.$$
(8.32)

Starting from multiple seeds

So far we have assumed that the integrand is estimated from a single random walk governed by the Markov process. One way of reducing the startup bias is to use several walks initiated from different starting points, called seeds, and combine their results. If the initial point is generated from seedpoints x_1, x_2, \ldots, x_N randomly selecting x_i with probability α_i ($\sum_{i=1}^N \alpha_i = 1$), then the initial probability distribution is the following

$$p_0(x) = \sum_{i=1}^{N} \alpha_i \cdot \delta(x - x_i).$$
(8.33)

Using the same concept as before, the probability density after n steps can be obtained in the following form

$$p_n(x) = \sum_{k=-\infty}^{k=\infty} \sum_{i=1}^{N} \alpha_i \cdot (P^*(k))^n \cdot e^{jk(x-x_i)}$$
(8.34)

The error of the probability distribution after step n is then

$$||p_n - p_{\infty}||_2 = \sqrt{\sum_{k=-\infty, k\neq 0}^{k=\infty} \left|\sum_{i=1}^{N} \alpha_i \cdot (P^*(k))^n \cdot e^{jk(x-x_i)}\right|^2}.$$
(8.35)

Analysis of uniform random perturbations

Let the perturbation be the selection of a point following uniform distribution from an interval of size Δ centered around the current point. Formally the transition probability is

$$P(x \to y) = \begin{cases} 1/\Delta & \text{if } |x - y| < \Delta, \\ 0 & \text{otherwise.} \end{cases}$$
(8.36)

The Fourier transform of this function is

$$P^*(k) = \frac{\sin k\pi\Delta}{k\pi\Delta}$$
(8.37)

which can be rather big even for large k values. This formula, together with equation (8.32) allows to generate the graph of the startup errors for different sample numbers and for different perturbation sizes (figure 8.10). The presented analysis can be generalized to higher dimensions thus we can obtain a method to estimate how many sample points must be ignored to get rid of the start-up bias. Note that the probability density is not accurate for many iterations if the perturbation size is small compared to the size of the domain. This situation gets just worse for higher dimensions.



Figure 8.10: Startup error for different perturbation size Δ

8.4 Stochastic iteration using ray-bundles

Let us recall that finite-element formulation needed by ray-bundles modifies the rendering equation to the following form (equation (8.7)):

$$\mathbf{L}(\omega) = \mathbf{L}^{e}(\omega) + \int_{\Omega} \mathbf{F}(\omega', \omega) \cdot \mathbf{A}(\omega') \cdot \mathbf{L}(\omega') \, d\omega' = \mathbf{L}^{e}(\omega) + \mathcal{T}_{F}\mathbf{L}(\omega).$$
(8.38)

This section proposes a stochastic iteration type method for the solution of this equation. Let us define a random operator \mathcal{T}_F^* that behaves like the projected transport operator in the average case in the following way:

A random direction is selected using a uniform distribution and the radiance of all patches is transferred into this direction.

Formally, the definition, which denotes the random direction by ω' , is

$$\mathcal{T}^*(\omega')\mathbf{L}(\omega) = 4\pi \cdot \mathbf{T}(\omega',\omega) \cdot \mathbf{L}(\omega') = 4\pi \cdot \mathbf{F}(\omega',\omega) \cdot \mathbf{A}(\omega') \cdot \mathbf{L}(\omega').$$
(8.39)

If direction ω' is sampled from a uniform distribution — i.e. the probability density of the samples is $1/(4\pi)$ — then the expected value of the application of this operator is

$$E[\mathcal{T}^*(\omega')\mathbf{L}(\omega)] = \int_{\Omega} 4\pi \cdot \mathbf{T}(\omega',\omega) \cdot \mathbf{L}(\omega') \frac{d\omega'}{4\pi} = \mathcal{T}_F \mathbf{L}(\omega).$$
(8.40)

In the definition of the random operator ω is the actually generated and ω' is the previously generated directions. Thus a "randomization point" is a global direction in this method.

The resulting algorithm is quite simple. In a step of the stochastic iteration an image estimate is computed by reflecting the previously computed radiance estimate towards the eye, and a new direction is found and this direction together with the previous direction are used to evaluate the random transport operator. Note that the algorithm requires just one variable for each patch i, which stores the previous radiance L[i].

The ray-bundle based stochastic iteration algorithm is:

Generate the first random global direction ω_1 for each patch *i* do $L[i] = L_i^e(\omega_1)$ for m = 1 to *M* do // iteration cycles Calculate the image estimate reflecting the incoming radiance $L[1], \ldots, L[n]$ from ω_m towards the eye Average the estimate with the Image Generate random global direction ω_{m+1} for each patch *i* do $L^{\text{new}}[i] = L_i^e(\omega_{m+1}) + 4\pi \cdot \sum_{j=1}^n \tilde{f}_i(\omega_m, \omega_{m+1}) \cdot \mathbf{A}(\omega_m)_{ij} \cdot L[j]$ endfor Display Image

Due to the fact that $\mathbf{A}(\omega_m)_{ij}$, $\mathbf{A}(-\omega_m)_{ij}$ can simultaneously be computed, the algorithm can easily be generalized to simultaneously handle bi-directional transfers.

Note that this algorithm is quite similar to the global walk algorithm, but it does not reinitialize the irradiance vector after each *D*th step. In fact, it generates a single infinite walk, and adds the effect of the lightsources to the reflected light field and computes the image accumulation after each step.



Figure 8.11: Stochastic iteration versus bi-directional walking techniques of length 5 and of length 10. The test scene was the homogeneous Cornell box (section 3.4)

The convergence rates are compared for bi-directional walking techniques and for stochastic iteration in figure 8.11. Note that walking methods have bias which depends on the length of the walk. Stochastic iteration, on the other hand, not only free from this bias problem, but its convergence speed is significantly better than that of the walking techniques.

8.4.1 Can we use quasi-Monte Carlo techniques in iteration?

Stochastic iteration can also be viewed as a single walk which uses a single sequence of usually 4dimensional randomization points (for ray-bundle tracing 2-dimensional randomization points), and the $\mathcal{T}_{i+k}^*\mathcal{T}_{i+k-1}^*\ldots\mathcal{T}_i^*L^e$ terms are included in integral quadratures simultaneously for all k.

It means that the randomization points should support not only 4-dimensional integration, but using subsequent pairs also 8-dimensional integration, using the subsequent triplets 12-dimensional integration, etc. Sequences that support k-dimensional integrals when subsequent k-tuples are selected are called k-uniform sequences [Knu81]. The widely used Halton or Hammersley sequences are only



Figure 8.12: Ray-bundle based stochastic iteration with random and quasi-random numbers. The test scene was the Cornell box shown in figure 8.21

1-uniform, thus theoretically they should provide false results. This is obvious for the Hammersley sequence, in which the first coordinate is increasing. Such a sequence would search for only those multiple reflections where the angle corresponding to the first coordinate always increases in subsequent reflections. It is less obvious, but is also true for the Halton sequence. Due to its construction using radical inversion, the subsequent points in the sequence are rather far, thus only those reflections are considered, where the respective angle changes drastically.

In order to avoid this problem without getting rid of the quasi-Monte Carlo sequences, in [SKFNC97] we proposed the random scrambling of the sample points. The same problem arises, for example, when generating uniform distributions on a sphere, for which [CMS98] proposed to increase the dimension of the low-discrepancy sequence.

Note that this problem is specific to quasi-Monte Carlo integration and does not occur when classical Monte-Carlo method is used to select the sample points (a random sequence is ∞ -uniform [Knu81]).

In order to demonstrate these problems, we tested the ray-bundle based iteration for different random (i.e. pseudo-random) and low-discrepancy sequences. The test scene was the Cornell box. In figure 8.12 we can see that the Hammersley sequence gives completely wrong result and the Halton sequence also deteriorates from the real solution. The two random generators (rand and drand48), however, performed well.

The figure also included a modification of the $q_n = \{\pi^n\}$ quasi-Monte Carlo sequence (operator $\{\}$ selects the fractional part of a number). This is believed to be (but has not been proven to be) ∞ -uniform [Deá89]. However, this sequence is very unstable numerically, therefore we used the following scheme started at $\kappa_0 = 1$:

 $\kappa_n = \kappa_{n-1} \cdot (\pi - 2)$ if $\kappa_n > 100000$ then $\kappa_n = \kappa_n - 100000$ $q_n = \{\kappa_n\}$

8.5 Calculation of the radiance transport in a single direction

The key of the calculation of the radiance transported by a ray-bundle is the determination of geometry matrix **A**. Examining the elements of the geometry matrix

$$\mathbf{A}(\omega')|_{ij} = \langle b_j(h(\vec{x}, -\omega')) \cdot \cos \theta', \tilde{b}_i(\vec{x}) \rangle = \int\limits_S b_j(h(\vec{x}, -\omega')) \cdot \tilde{b}_i(\vec{x}) \cdot \cos \theta' \ d\vec{x},$$

we can realize that its computation requires the determination of whether a point of the support of basis function *i* is visible from a point of the support of basis function *j* in a given direction $-\omega'$ (note that $b_j(h(\vec{x}, -\omega')) \cdot \tilde{b}_i(\vec{x})$ is non zero only if \vec{x} is in the support of \tilde{b}_i and the point $h(\vec{x}, -\omega')$ is in the support of b_j). For the algorithms to be introduced, the support of a basis function is a planar triangle, while the support of an adjoint basis function is either a planar triangle or a single vertex. Direction ω' is called the **transillumination direction** [Neu95, SKFNC97].

Note that an element of the geometry matrix is non zero only if $\cos \theta' \ge 0$, that is if patch *i* is facing towards the transillumination direction. Faces meeting this requirement are called **front faces**, while those faces which cannot meet this are called **back faces**. Obviously, only front faces can get radiance contribution from a transillumination direction (this can be lifted easily to allow transparent materials).



Figure 8.13: Global visibility algorithms

The determination of the geometry matrix is a **global visibility problem**, since only the viewing direction is fixed but the eye position is not. In fact, the eye position should visit all surface points or all vertices depending on the selected adjoint base. Looking at figure 8.13, it is easy to see that the global visibility problem can be solved in an incremental way if the patches are visited in the order of their position in the transillumination direction. In fact, what is visible from a patch differs just in a single patch from what is visible from the next patch. This single patch may appear as a new and may hide other patches. The required sorting is not obvious if the patches overlap in the transillumination direction, but this can be solved in a way as proposed in the **painter's algorithm** [NNS72]. On the other hand, in our case the patches are usually small, thus simply sorting them by their center introduces just a negligible error. Although sorting seems worthwhile, it is not obligatory. Thus the proposed visibility algorithms will be classified according to whether or not an initial sorting is required.

At a given point of all global visibility algorithms the objects visible from the points of a patch must be known. This information is stored in a data structure called the **visibility map**. The visibility map can also be regarded as an image on the plane perpendicular to the transillumination direction. This plane is called the **transillumination plane** (figure 8.13).

The algorithms that generate the visibility map can be either discrete or continuous. **Discrete algorithms** decompose the transillumination plane to small pixels of size δP . Their visibility map is simply a rasterized image where each pixel can store either the index of the visible patch or the radiance of the visible point. For **continuous algorithms**, the visibility map identifies those regions in which the visibility information is homogeneous (either the same patch is seen or no patch is seen). Discrete algorithms are faster and the rendering hardware and the z-buffer of workstations can also be exploited [SKe95, SKM94] but in order to handle all patches simultaneously, the "window" of the algorithm should include all patches. The large window, however, should be decomposed into sufficiently small pixels to provide the required precision, which might result in high resolution requirements for sparse scenes. Continuous algorithms are free from these resolution problems, but are usually more difficult to implement and are much slower.

The computation of the geometry matrix also depends on the selected basis functions and the finiteelement algorithm (adjoint basis). We consider two different sets of basis functions and two finite element approaches. In the first case the Galjerkin method is applied for piece-wise constant basis functions. Secondly, piece-wise linear basis functions are used in a point-collocation algorithm.

8.5.1 Galerkin's method with piece-wise constant basis functions

Let us use the following basis functions

$$b_j(\vec{x}) = \begin{cases} 1 \text{ if } \vec{x} \in A_j, \\ 0 \text{ otherwise.} \end{cases}$$
(8.41)

In **Galerkin's method**, the unknown directional functions $L_i(\omega)$ are found to ensure that approximation (8.1) is the real solution of the radiance equation in the subspace induced by the basis functions b_i . To satisfy normalization criteria, the adjoint base is selected as follows:

$$\tilde{b}_j(\vec{x}) = \begin{cases} 1/A_j \text{ if } \vec{x} \in A_j, \\ 0 \text{ otherwise.} \end{cases}$$
(8.42)

Since $\langle b_i(\vec{x}), \tilde{b}_i(\vec{x}) \rangle$ is 1 if i = j and zero otherwise, the element i, j of the geometry matrix is

$$\mathbf{A}(\omega')|_{ij} = \langle b_j(h(\vec{x}, -\omega')) \cdot \cos \theta', \tilde{b}_i(\vec{x}) \rangle = \frac{1}{A_i} \cdot \int_{A_i} b_j(h(\vec{x}, -\omega')) \cdot \cos \theta' \, d\vec{x}.$$
(8.43)

Since the integrand of this equation is piece-wise constant, the integral can also be evaluated analytically:

$$\frac{1}{A_i} \cdot \int_{A_i} b_j(h(\vec{x}, -\omega')) \cdot \cos \theta' \, d\vec{x} = \frac{A(i, j, \omega')}{A_i},\tag{8.44}$$

where $A(i, j, \omega')$ expresses the projected area of patch *j* that is visible from patch *i* at direction $-\omega'$. In the unoccluded case this is the intersection of the projections of patch *i* and patch *j* onto the transillumination plane. If occlusion occurs, the projected areas of other patches that are in between patch *i* and patch *j* should be subtracted as shown in figure 8.14.



Figure 8.14: Interpretation of $A(i, j, \omega')$

Having the visibility map of patches visible from A_i , the computation of $A(i, j, \omega')$ requires to determine which regions are inside the projection of A_i and to sum these areas. Note that the geometry matrix is symmetric in the following sense

$$A(i, j, \omega') = A(j, i, -\omega')$$

thus bi-directional transfers really do not have visibility computation overhead.

In the following sections different continuous and discrete visibility algorithms are presented to determine the necessary visibility information.

Continuous algorithm with initial sorting: Local visibility map

This algorithm processes the patches in the order defined by the transillumination direction and maintains the visibility graph dynamically [SKFNC97].



transillumination direction



When the processing of patch *i* is started, the visibility map shows which patches are visible from patch *i*. To calculate the $A(i, j, \omega')$ values, those patches which have projections either entirely or partly in the projection of patch *i* are selected from the visibility map, and are clipped onto patch *i* to clearly separate inner regions. The process of clipping of the patches onto each other is quite similar to the **Weiler-Atherthon algorithm** [WA77]. Then the projected areas of those patch parts which are inside the projection of patch *i* are summed to find the $A(i, j, \omega')$ values.

When we step onto the patch next to patch i, a new visibility map is created by replacing those region parts that are inside the projection of A_i by the projection of A_i . Then, if patch i can reflect energy onto the next patch (it is a back facing patch with respect to the transillumination direction), then patch ishould be added to the visibility map, otherwise, the place of the projection of patch i will be empty.

The algorithm, that maintains a list L for the sorted patches, V for the projected patches that are currently in the visibility map, O for those clipped, not-hidden patches whose projections are outside patch i, and I for those clipped, not-hidden patches whose projections are inside patch i, is as follows:

```
list L = \text{Sort patches in direction } \omega'
visibility map V = \{ \}
for each patch i in L do
Clip patches in V onto patch i and generate: O = \text{outside list}, I = \text{inside list}
if patch i is front facing then
A(i, j, \omega') = \sum_{j \in I} A(j)
V = O
else V = O + \text{patch } i
endfor
```

If the number of patches is n, then the size of the visibility map is $\mathcal{O}(n)$ in the average case but $\mathcal{O}(n^2)$ in the worst-case. Thus the resulting algorithm that compares each patch with the actual visibility map will have $\mathcal{O}(n^2)$ average-case and $\mathcal{O}(n^3)$ worst-case time-complexity. This is not acceptable when complex scenes are processed. In order to reduce the complexity, we can use the wide selection of computational geometry methods that usually apply spatial decomposition on the plane to reduce the number of unnecessary comparisons [SKFP98b].

Another alternative is to introduce randomization (or quasi-randomization) in the computation of the visible areas, as suggested by **Russian-roulette** [SKP99b]. When the outside list O is generated, the patches that have small area times radiance are randomly (or quasi-randomly) deleted from the list. Suppose that those patches are considered for deletion for which $A(j) \cdot L(j) < \epsilon$ where A(j) is the projected area of the patch j and L(j) is the radiance of this patch, and the probability of not deleting this patch is $A(j) \cdot L(j)/\epsilon$. Whenever a small patch is kept, its radiance is multiplied by $\epsilon/A(j) \cdot L(j)$ to compensate for occasions when it is deleted.

Thus the following instructions should be inserted in the algorithm:

```
for each patch j in O do

if A(j) \cdot L(j) < \epsilon then

Generate uniformly distributed random number r in [0, 1]

if A(j) \cdot L(j) < r \cdot \epsilon then Delete patch i from O

else L(j) = \epsilon/A(j)

endif

endfor
```

In quasi-Monte Carlo integration r is an additional integration variable that should be generated by an independent low-discrepancy sequence.

Continuous algorithm without initial sorting: Global visibility map

This algorithm first projects all the polygon vertices and edges onto the transillumination plane, then determines all the intersection points between the projected edges, and form a planar graph that is a superset of the set of projected edges of the polygons [SKFP98b, SKFNC97]. In the resulting planar graph, each territory represents a list of patches that can be projected onto the territory. Furthermore, if the patches do not intersect, the order of patches is also unique in each territory.

Thus, to compute $A(i, j, \omega')$ for some *i*, the lists of the territories should be visited to check whether *i* is included. If patch *i* is found, then the patch next to it on the list should be retained to find index *j*, and the area of the territory should be added to $A(i, j, \omega')$.



transillumination direction

Figure 8.16: Global visibility map
The draft of the algorithm to generate the data structure is the following:

Clear geometry matrix $A(i, j, \omega')$ Project vertices and edges onto the transillumination plane Calculate all intersection points between projected edges Compute the graph of the induced planar subdivision for each region of this graph do Sort patches visible in this region for each patch *i* in the list of patches visible in this region j = Next patch in the list if patch *i* is front facing and patch *j* is back facing then $A(i, j, \omega') +=$ area of the region endfor

The speed of the algorithm is considerably affected by how well its steps are implemented. A simplistic implementation of the intersection calculation, for example, would test each pair of edges for possible intersection. If the total number of edges is n, then the time complexity of this calculation would be $\mathcal{O}(n^2)$. Having calculated the intersection points, the structure of the subdivision graph has to be built, that is, incident nodes and arcs have to be assigned to each other somehow. The number of intersection points is $\mathcal{O}(n^2)$, hence both the number of nodes and the number of arcs fall into this order. A simplistic implementation of the graph computation would search for the possible incident arcs for each node, giving a time complexity of $\mathcal{O}(n^4)$. This itself is inadmissible in practice, not to mention the possible time complexity of the further steps. However, applying the results of computational geometry, we can do it much better. Algorithms are available that can do it in $\mathcal{O}((n+i)\log n)$ [Dév93, SKe95] time where n is the number of patches (or edges) and i is the number of edge intersections, or even in $\mathcal{O}(n^{1+\varepsilon}\sqrt{k})$ time [dB92] where k is the number of edges in the visibility map. The number of intersections i and the number of edges k are in $\mathcal{O}(n^2)$ in the worst-case, but are in $\mathcal{O}(n)$ in practical scenes.

Discrete algorithm with initial sorting: Global painter's algorithm

Discrete algorithms determine the visible patches for each front facing patch through a discretized window. This is a visibility problem, and the result is an "image" of the patches, assuming the eye to be on patch i, the window to be on the transillumination plane and the color of patch j to be j if the patch is facing to patch i and to be 0 otherwise.



Figure 8.17: Application of painter's algorithm

If the patches are sorted in the transillumination direction and processed in this order, the computation of $A(i, j, \omega')$ requires the determination of the pixel values inside the projection of patch *i*. Then, to proceed with the next patch in the given order, the pixels covered by patch *i* are filled with *i* if patch

i is not front facing and 0 otherwise. The two steps can be done simultaneously by a modified scanconversion algorithm that reads the value of the image buffer before modifying it.

This is summarized in the following algorithm [SKF97]:

```
Sort patches in direction \omega' (painter's algorithm)
Clear image
for each patch i in the sorted order do
if patch i is front facing then
for each pixel in the projection of patch i
j = \text{Read pixel}
A(i, j, \omega') += \delta P
Write 0 to the pixel
endfor
else Render patch i with color i
endfor
```

// δP is the area of a pixel

Sorting a data set is known to have $\mathcal{O}(n \log n)$ time complexity, so does the painter's algorithm in the average case. A single cycle of the second **for** loop contains only instructions that work with a single patch and an "image", thus the time required for a single cycle is independent of the number of patches. Since the **for** loop is executed *n* number of times, the time complexity of the **for** loop is $\mathcal{O}(n)$. Consequently the algorithm requires $\mathcal{O}(n \log n)$ time.

Discrete algorithm without initial sorting: software z-buffer

For the sake of completeness, we mention that the **global z-buffer algorithm** [Neu95] can also be used for our purposes. This method stores not just the closest patch index and its z value in the buffer, but the whole list of those patches which can be projected onto this pixel. The patches are scan-converted by the z-buffer algorithm, and are inserted into the lists associated with the covered pixels. The lists of pixels can be used to compute $A(i, j, \omega')$ similarly to the continuous global visibility map algorithm.

Discrete algorithm without initial sorting: exploitation of the hardware z-buffer

In this section another method is proposed that traces back the visibility problem to a series of z-buffer steps to allow the utilization of the z-buffer hardware of workstations [SKP98a].



Figure 8.18: Calculating the power transfer using the z-buffer

In this algorithm the radiance is transferred by dynamically maintaining two groups of patches, an **emitter group** and a **receiver group**, in a way that no patch in the receiver group is allowed to hide a patch in the emitter group looking from the given transillumination direction.

Let the two classes of patches be rendered into two image buffers — called the emitter and receiver images, respectively — setting the color of patch j to j and letting the selected direction be the viewing direction for the receiver set and its inverse for the emitter set.

Looking at figure 8.18, it is obvious that a pair of such images can be used to calculate the radiance transfer of all those patches which are fully visible in the receiver image. The two images must be scanned parallely and when *i* (that is the index of patch *i*) is found in the receiver image, the corresponding pixel in the emitter image is read and its value (*j*) is used to decide which $A(i, j, \omega')$ should be increased by the area of the pixel.

In order to find out which patches are fully visible in the receiver image, the number of pixels they cover is also computed during scanning and then compared to the size of their projected area. For those patches whose projected area is approximately equal to the total size of the covered pixels, we can assume that they are not hidden and their accumulated irradiances are valid, thus these patches can be removed from the receiver set and rendered into the emitter image to calculate the radiance transfer for other patches (this is the strategy to maintain the emitter and receiver sets automatically).

This leads to an incremental algorithm that initially places all patches in the receiver set. Having calculated the receiver image by the z-buffer algorithm, the radiance transfer for the fully visible patches are evaluated, and then they are moved from the receiver set to the emitter set. The algorithm keeps doing this until no patch remains in the receiver set (cyclic overlapping would not allow the algorithm to stop, but this can be handled by a clipping as in the painter's algorithm [NNS72]). The number of z-buffer steps required by the algorithm is quite small even for complex practical scenes [Sbe96]. Exploiting the built-in z-buffer hardware of advanced workstations, the computation can be fast.

When checking whether or not the visible size is approximately equal to the projected patch size, the allowed tolerance is the total area of the pixels belonging to the edge of the patch, which in turn equals to the sum of the horizontal and vertical sizes of a triangular patch.

This algorithm "peels" the scene by removing the layers one by one. The sequences of evolving receiver and emitter images are shown in figure 8.19 and in figure 8.20, respectively. Note that the first receiver image contains all patches, thus its pair is an empty image that is not included figure 8.20. The pair of the second receiver image is the first shown emitter image, etc. The last emitter image includes all patches thus its receiver pair is empty, and therefore is not shown in figure 8.19.

In the following algorithm R denotes the collection of the receiver patches.

```
R = all patches
Clear emitter_image
while R is not empty
       Clear receiver_image
       for each patch r in R
           Render patch r into receiver_image via z-buffer
           Clear row r of geometry matrix A(r, j, \omega')
           patch[r].visible_size = 0
      endfor
      for each pixel P
           r = \text{receiver\_image}[P]
           e = \text{emitter\_image}[P]
           patch[r].visible_size += \delta P
           A(r, e, \omega') \models \delta P
       endfor
                                                                               // Move patches to the emitter set
       for each patch p
           if patch[p].visible_size \approx projected size of patch p then
             Remove patch p from R and render it to emitter_image
           endif
       endfor
endwhile
```



Figure 8.19: Steps of the evolution of receiver images



Figure 8.20: Steps of the evolution of emitter images

8.5.2 Analysis of the finite resolution problem of discrete methods

In order to find out how important the resolution of the visibility map, a Cornell box scene (figure 8.21) consisting of 3705 triangular patches has been rendered with the global painter's algorithm having set the resolution to different values from 50×50 to 1000×1000 . Since the resolution can only be interpreted when compared to the size of the patches, table 8.2 summarizes the average projected patch sizes in pixels and also the residual errors of the iteration.



Figure 8.21: A Cornell box as rendered using 100×100 (left) and 1000×1000 (right) pixels for the visibility map



Figure 8.22: Comparison of the error curves using visibility maps of different resolutions as a function of iterations (left) and of computation time (right)

Note that the resolution plays negligible role until the average patch per pixel ratio is significantly greater than one (left of figure 8.22). For instance, the error curves of resolutions 1000×1000 and 500×500 can hardly be distinguished. This can be explained by the stochastic nature of the algorithm. Each radiance transfer uses a different direction, thus a different discrete approximation of the size of the patch. Although this approximation can be quite inaccurate in a single step, the expected value of these approximations will still be correct. As the algorithm generates the result as the average of the estimates,

resolution	pixel-per-patch	residual error
50×50	0.5	0.25
100×100	2	0.05
200×200	8	0.02
500×500	55	≤ 0.01
1000×1000	220	≤ 0.01

Table 8.2: Average number of pixels per patch and the discretization errors

these approximation errors will be eliminated. The effect of the low resolution is just an "additional noise". However, when the pixel size becomes comparable to the projected size of the patches, then the iteration will deteriorate from the real solution as we can clearly see it in figure 8.22 when the resolution is lower than 200×200 . The core of the problem is that discrete filling algorithms always assume that both the width and the height of the patch is at least 1. Thus, for patches having smaller height or width, even the expected value will be incorrect. The computation time is roughly proportional to the number of pixels in the visibility map thus it is desirable to keep the resolution low (right of figure 8.22). For example, the 500 stochastic iterations that generated the left and the right of figure 8.21 needed 3 and 8 minutes respectively. The optimal selection of the resolution is the minimal number which guarantees that even the smallest patches are projected onto a few pixels.

Even patches of subpixel size can be correctly handled by a Russian-roulette like method introduced in section 8.5.1. The basic idea is a random rendering of small patches to sustain the correct expectation of the projected size time the radiance. Suppose that those patches are considered for random rendering for which $A(j) \cdot L(j) < \epsilon$ where A(j) is the projected area of patch j and L(j) is the radiance of the patch, and the probability of rendering pixels is $A(j) \cdot L(j)/\epsilon$. When a patch is not rendered, the pixels are filled up with a reference value that corresponds to zero radiance. Whenever a small patch is rendered, its radiance is multiplied by $\epsilon/A(j) \cdot L(j)$ to compensate for occasions when it is not rendered.

8.5.3 Point collocation method with piece-wise linear basis functions

In this method [SKFP98a] the radiance variation is assumed to be linear on the triangles. Thus, each vertex *i* of the triangle mesh will correspond to a "tent shaped" basis function b_i that is 1 at this vertex and linearly decreases to 0 on the triangles incident to this vertex (figure 8.23). Assume that the shading normals are available at the vertices.

In the point-collocation method, the unknown directional functions $L_j(\omega)$ are determined to ensure that the residual of the approximation is zero at the vertices of the triangle mesh. This corresponds to Dirac-delta type adjoint basis functions, where $\tilde{b}_j(\vec{x})$ is non-zero at vertex *i* only. Thus the geometry matrix is

$$\mathbf{A}(\omega')|_{ij} = \langle b_j(h(\vec{x}, -\omega')) \cdot \cos\theta', \delta(\vec{x} - \vec{x}_i) \rangle = b_j(h(\vec{x}_i, -\omega')) \cdot \cos\theta'(\vec{x}_i).$$
(8.45)

Calculation of the irradiance at vertices

Since now the irradiance is not piece-wise constant but piece-wise linear, it is better to evaluate $\mathbf{A}(\omega') \cdot \mathbf{I}$ directly than evaluating the geometry matrix and the irradiance separately. Thus we have to find

$$(\mathbf{A}(\omega') \cdot \mathbf{I})[i] = \sum_{j=1}^{n} b_j(h(\vec{x}_i, -\omega')) \cdot \cos \theta'(\vec{x}_i) \cdot \mathbf{I}[j].$$
(8.46)

In this formula the $b_j(h(\vec{x}_i, -\omega'))$ factor is non-zero for those *j* indices which represent a vertex of the patch visible from \vec{x}_i at direction ω' . The exact value can be derived from the calculation of the height of the linear "tent" function at point $h(\vec{x}_i, -\omega')$. This means that having identified the patch visible from



Figure 8.23: Linear basis function in 3-dimension (left) and global visibility algorithm for the vertices (right)

 \vec{x}_i at direction ω' , the required value is calculated as a linear interpolation of the irrandiances of the vertices of this patch.

To solve it for all patches, the triangular patches are sorted in direction ω' , then painted one after the other into an image buffer. For vertex *i* of each triangle, the "color" is set to

$$\mathbf{L}_{i}^{e}(\omega') + 4\pi \cdot \mathbf{F}_{ii}(\omega'_{\text{previous}},\omega') \cdot \mathbf{I}[i]$$

at step d and the linear interpolation hardware (Gouraud shading) is used to generate the color (or irradiance) inside the triangle. For back-facing patches this step clears the place of the triangle in the "image". If the triangles that are in front of the given triangle in direction ω' are rendered into the image buffer, then the radiance illuminating the vertices of the given triangle is readily available in the current image buffer. Assuming that patches are processed in the order of the transillumination direction, every patch should be rendered only once into the image buffer.

Thus the calculation of the irradiances at a given transillumination direction is:

```
Sort patches in direction \omega' (painter's algorithm)

Clear image-buffer

for each patch i in sorted order do

if patch i is front facing then

for each vertex v[i] of the patch i do \operatorname{color}[v[i]] = \mathbf{L}_{v[i]}^{e}(\omega') + 4\pi \cdot \mathbf{F}_{v[i],v[i]}(\omega'_{\operatorname{previous}}, \omega') \cdot \mathbf{I}[v[i]]

Render patch i into the image-buffer with Gouraud shading

else

for each vertex v[i] of the patch i do \mathbf{I}[v[i]] = (\text{image buffer at projection of } v[i]) \cdot \cos \theta'[v[i]]

Render patch i with color 0

endif

endfor
```

The processing of a single direction for all patches requires a sorting step and the rendering of each triangle into a temporary buffer. This can be done in $O(n \log n)$ time.

8.6 Handling sky-light illumination

The visibility methods introduced so far can easily be extended for sky-light illumination by initializing the image on the transillumination plane by a special value if the direction points downwards (sky is usually above the horizon). When the radiance is transferred and this value is found in a given pixel, then the irradiance of the receiver is updated according to the intensity of the sky-light (figures 9.3 and 9.9).

8.7 Improving the efficiency

The error of Monte-Carlo and quasi-Monte Carlo techniques depends on the number of samples and on how flat the integrand is. For Monte-Carlo integration this is expressed by equation (5.27) while for quasi-Monte Carlo quadrature this is shown by the Koksma-Hlawka inequality. Thus the efficiency of the methods can be improved by the proper formulation of the global illumination problem as a numerical integration where the integrand is flatter. If iteration approach is used, then another source of the error is the energy defect. This defect depends on the initial radiance function. The closer the initial function is to the solution, the less the defect is.

8.7.1 Self-correcting iteration

Self-correcting iteration (section 7.2.1) means that averaging is applied not only in the image estimation but also for the radiance to continue the iteration. This means that a new random transport operator is combined not only with the last operator but all preceding operators. This seems advantageous, but it also has a drawback. Namely, the results of the previous iteration steps must be stored to allow such a combination. In order to avoid this additional storage, self-correcting iteration is applied only to the diffuse reflection, the non-diffuse reflection is estimated by normal stochastic iteration. Let us decompose the radiance and the random light transport operator to separately model diffuse and non-diffuse reflections,

$$L = L^e + L_d + L_{nd}, \qquad \mathcal{T}^* = \mathcal{T}^*_d + \mathcal{T}^*_{nd},$$

and rewrite the iteration in the following way

$$L'_{d}(m) = \mathcal{T}^{*}_{d}(m)(L^{e} + L_{nd}(m-1) + L_{d}(m-1)),$$

$$L_{d}(m) = \tau_{m} \cdot L'_{d}(m) + (1 - \tau_{m}) \cdot L_{d}(m-1),$$

$$L_{nd}(m) = \mathcal{T}^{*}_{nd}(m)(L^{e} + L_{nd}(m-1) + L_{d}(m-1)),$$

$$P_{m} = \mathcal{M}(L^{e} + L_{d}(m) + L_{nd}(m)),$$
(8.47)

where τ_m is an appropriate sequence that converges to 0. In practice, the $\tau_m = 1/m$ or $\tau_m = 1/\sqrt{m}$ have been found to be appropriate.



Figure 8.24: Efficiency improvements of self-correcting iteration in the Cornell box scene (figure 8.21)

The efficiency gained from proper initialization of the radiance function and by self-correcting iteration is shown in figure 8.24. The sequence of images are included by figure 9.13.

8.7.2 Preprocessing the small lightsources

The primary causes of the high variation of the integrand are the smaller, bright lightsources. Especially, if the lightsource is infinitely small but its power is not zero — that is, its radiance is infinite then the error bounds of the Monte-Carlo and quasi-Monte Carlo quadratures will be infinite as well. Consequently, these small, bright lightsources should be handled by special techniques.

The problem of small lightsources can be solved by a modified version of the "**first-shot**" that shoots the power of the point lightsources onto other surfaces, then removes them from the scene. This method [CMS98] works well in the radiosity setting, since in this case, the representation of the reflected radiance requires a diffuse "emission" in each patch, thus the memory overhead of the first-shot is just one variable per patch. However, in non-diffuse scenes the classical first-shot has prohibitive memory requirements, since even if the original light-sources are diffuse, their reflection may have general directional function, which requires the representation of the complete reflected, non-diffuse radiance function. If the directional variation of the radiance is represented by n basis functions (i.e. n is the number of small solid angles in which the radiance can be supposed to be constant) in each patch, then the method requires nnew variables for each patch. To solve this problem, a new first-shot technique is proposed that is called the **incoming first-shot**. Incoming first-shot precomputes and stores the incoming radiance received by the patches from each point sample on the lightsources) [SKP98a]. The secondary, non-diffuse emission to a direction is computed from these incoming radiances on the fly. The method is feasible if l is small, thus it can be used for a few point lightsources and smaller area lightsources.



Figure 8.25: Incoming first shot technique

Formally, the unknown radiance L is decomposed into two terms:

$$L = L^{ep} + L^{np} \tag{8.48}$$

where L^{ep} is the emission of the small area and point-like lightsources, L^{np} is the emission of the large area lightsources and the reflected radiance. Substituting this into the rendering equation we have:

$$L^{ep} + L^{np} = L^e + \mathcal{T}(L^{ep} + L^{np}).$$
(8.49)

Expressing L^{np} we obtain:

$$L^{np} = (L^e - L^{ep} + \mathcal{T}L^{ep}) + \mathcal{T}L^{np}.$$
(8.50)

Introducing the new lightsource term

$$L^{e*} = L^e - L^{ep} + \mathcal{T}L^{ep}$$
(8.51)

which just replaces the small lightsources (L^{ep}) by their single reflection ($\mathcal{T}L^{ep}$), the equation for L^{np} is similar to the original rendering equation:

$$L^{np} = L^{e*} + \mathcal{T}L^{np}. \tag{8.52}$$

It means that first the direct illumination caused by the small lightsources must be computed, then they can be removed from the scene and added again at the end of the computation.

Incoming first-shot of point lightsources

Suppose that the scene contains l point lightsources at locations $\vec{y}_1, \ldots, \vec{y}_l$ with powers Φ_1, \ldots, Φ_l , respectively, then their reflection at point \vec{x} is:

$$(\mathcal{T}L^{ep})(\vec{x},\omega) = \sum_{i=1}^{l} \frac{\Phi_l \cdot v(\vec{y}_i, \vec{x})}{4\pi |\vec{y}_i - \vec{x}|^2} \cdot f_r(\omega'_i, \vec{x}, \omega) \cdot \cos \theta'_i,$$
(8.53)

where ω'_i is the direction from lightsource *i* to point \vec{x} , θ'_i is the angle between ω'_i and the surface normal, and $v(\vec{y}_i, \vec{x})$ indicates the mutual visibility of \vec{x} and \vec{y}_i . Suppose that the patch under consideration is patch *j* and its area is A_j . Having projected the reflected radiance into adjoint base \tilde{b}_i , we get:

$$L_{j}^{e*}(\omega) = \langle \mathcal{T}L^{ep}(\vec{x},\omega), \tilde{b}_{j}(\vec{x}) \rangle = \sum_{i=1}^{l} \frac{1}{A_{j}} \cdot \int_{A_{j}} \frac{\Phi_{i} \cdot v(\vec{y}_{i},\vec{x})}{4\pi |\vec{y}_{i}-\vec{x}|^{2}} \cdot f_{r}(\omega_{i}',\vec{x},\omega) \cdot \cos\theta_{i}' \, d\vec{x}.$$
(8.54)

To compute the reflection of a lightsource at a point, the visibility of the lightsource from the point must be determined. We can use **shadow rays** evaluated by ray-shooting, but this is rather slow. Another alternative is to exploit the image synthesis hardware in the following way. The eye is put at the lightsource and the window is defined as one of the faces of a cube placed around the eye. Rendering the images for each faces with constant shading and using the index of the patches as color values, the visible areas of the patches from the lightsource can be determined.



Figure 8.26: Computation of the lightsource visibility by hardware

The integral in equation (8.54) can also be evaluated on the six window surfaces (W) that form a cube around the lightsource (figure 8.26). To find formal expressions, let us express the solid angle $d\Omega_p$, in which a differential surface area $d\vec{x}$ is seen through pixel area $d\vec{p}$, both from the surface area and from the pixel area:

$$d\Omega_p = \frac{d\vec{x} \cdot \cos\theta'_i}{|\vec{y}_i - \vec{x}|^2} = \frac{d\vec{p} \cdot \cos\theta_p}{|\vec{y}_i - \vec{p}|^2},\tag{8.55}$$

where θ_p is the angle between direction pointing to \vec{x} from \vec{y}_i and the normal of the window. The distance $|\vec{y}_i - \vec{p}|$ between pixel point \vec{p} and the lightsource \vec{y}_i equals to $f/\cos\theta_p$ where f is the distance from \vec{y}_i to the window plane, that is also called the **focal distance**. Using this and equation (8.55), differential area $d\vec{x}$ can be expressed and subsituted into equation (8.54), thus we can obtain:

$$L_j^{e*}(\omega) = \sum_{i=1}^l \frac{1}{A_j} \cdot \int\limits_W \frac{\Phi_i \cdot v(\vec{y}_i, \vec{x})}{4\pi} \cdot f_r(\omega'_i, \vec{x}, \omega) \cdot \frac{\cos \theta_p^3}{f^2} d\vec{p}.$$

Let P_j be the set of pixels in which patch j is visible from the lightsource. P_j is computed by running a z-buffer/constant shading rendering step for each sides of the window surface, assuming that the color

of patch j is j, then reading back the "images" (figure 8.28). The reflected radiance on patch j is approximated by a discrete sum as follows:

$$L_j^{e*}(\omega) \approx \sum_{i=1}^l \frac{\Phi_i}{4\pi f^2 A_j} \cdot \sum_{p \in P_j} f_r(\omega_i', \vec{x}, \omega) \cdot \cos \theta_p^3 \cdot \delta P,$$
(8.56)

where δP is the area of a single pixel in the image. If R is the resolution of the image — i.e. the top of the hemicube contains $R \times R$ pixels, while the side faces contain $R \times R/2$ pixels – then $\delta P = 4f^2/R^2$. If the BRDF can be assumed to be $\tilde{f}_j(\omega'_i, \omega)$ in patch j, then the reflected radiance can be decomposed into 3 factors: the power spectrum of the lightsource Φ_i , the BRDF $\tilde{f}_j(\omega'_i, \omega)$ which is also a spectrum and is the only factor which depends on viewing direction ω , and a scalar factor defined as:

$$r_{ij} = \frac{1}{\pi R^2 A_j} \cdot \sum_{p \in P_j} \cos \theta_p^3.$$

These scalar factors are computed and stored at each patch, which requires just one float variable per each patch and each point lightsource.

If variables r_{ij} are available, then the incoming first-shot phase is complete. During global illumination when the reflected radiance $L_j^{e*}(\omega)$ is needed at point \vec{x} of patch j, this is computed on the fly from the stored scalar parameters r_{ij} , from the directions pointing from \vec{x} to the lightsources and from the power of the lightsources according to the following formula:

$$L_j^{e*}(\omega) = \sum_{i=1}^l \Phi_i \cdot r_{ij} \cdot \tilde{f}_j(\omega'_i, \omega).$$
(8.57)

Smaller area lightsources

Now let us discuss the computation of a single reflection of the light coming from a small area lightsource S of emission $L^e(\vec{y}, \omega)$ to a point \vec{x} . The reflection at point \vec{x} is

$$(\mathcal{T}L^{ep})(\vec{x},\omega) = \int_{\Omega_S} L^e(h(\vec{x},-\omega'),\omega') \cdot f_r(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\omega' = \int_S \frac{L^e(\vec{y},\omega') \cdot \cos\theta \cdot v(\vec{y},\vec{x})}{|\vec{y}-\vec{x}|^2} \cdot f_r(\omega',\vec{x},\omega) \cdot \cos\theta' \, d\vec{y}, \tag{8.58}$$

where Ω_S is the solid angle in which lightsource S is visible, \vec{y} is a running point on the lightsource and θ is the angle between ω' and the surface normal of the lightsource at \vec{y} . The projected reflected radiance of patch j is

$$L_{j}^{e*}(\omega) = \frac{1}{A_{j}} \cdot \int_{A_{j}} (\mathcal{T}L^{ep})(\vec{x},\omega) d\vec{x} =$$

$$\int_{S} \frac{1}{A_{j}} \cdot \int_{A_{j}} \frac{L^{e}(\vec{y},\omega') \cdot \cos\theta \cdot v(\vec{y},\vec{x})}{|\vec{y}-\vec{x}|^{2}} \cdot \cos\theta' \cdot f_{r}(\omega',\vec{x},\omega) d\vec{x} d\vec{y}, \qquad (8.59)$$

The outer integral is estimated by a trapezoidal-like rule. It means that the lightsource area is tessellated to triangles (or quadrilaterals). The integrand is evaluated at the common vertices and is assumed to be linear between the vertices. If the number of vertices is l, then the quadrature rule is:

$$L_j^{e*}(\omega) \approx \sum_{i=1}^l \frac{1}{A_j} \cdot \int\limits_{A_j} \frac{L^e(\vec{y}_i, \omega_i') \cdot \cos \theta_i \cdot S_{ti} \cdot v(\vec{y}_i, \vec{x})}{3|\vec{y}_i - \vec{x}|^2} \cdot \cos \theta_i' \cdot f_r(\omega_i', \vec{x}, \omega) \ d\vec{x},$$

where S_{ti} is the total area of the lightsource triangles that share vertex *i* and factor 1/3 comes from the fact that a triangle has 3 vertices.

Note that the inner integral is the same as the integral in equation (8.54), with the substitution

$$\frac{\Phi_i}{4\pi} \leftarrow L^e(\vec{y}_i, \omega_i') \cdot \cos \theta_i \cdot \frac{S_{ti}}{3}.$$

There is another slight difference in the window surface. A one-sided area lightsource can emit light into that halfspace which is "above" the plane of lightsource. Thus the window surface becomes a **hemicube** (figure 8.28). An even better window surface is the **cubic tetrahedron** [BKP91], since it has just 3 faces while the hemicube has 5.

Summarizing the incoming first-shot from a smaller area lightsource consists of the following steps. First the lightsource is decomposed into a triangle mesh. A hemicube or a cubic tetrahedron is placed at each vertex \vec{y}_i of the mesh and the visibility of the other surfaces are determined. Scalar factors

$$r_{ij} = \frac{4S_{ti} \cdot \cos \theta_i}{3R^2 A_j} \cdot \sum_{p \in P_j} \cos \theta_p^3$$

are stored in each patch *j*. The reflected radiance can be obtained from this scalar factor during the global illumination computation in the following way:

$$L_j^{e*}(\omega) = \sum_{i=1}^l L^e(\vec{y}_i, \omega_i') \cdot r_{ij} \cdot f_r^j(\omega_i', \omega).$$
(8.60)



Figure 8.27: Error of ray-bundle stochastic iteration with and without first-shot for the Sierpiensky set scene (figure 9.1). The left image has been plotted as a function of time to demonstrate that the overhead of the first-shot step amortizes quickly

The error curves of figure 8.27 compare the stochastic iteration with and without the incoming firstshot. We can conclude that the incoming first-shot step has its overhead, but it is worth doing. We can see that the stochastic iteration is about 20 times faster with the incoming first-shot than without it.

Diffuse shot

Getting rid of the point and small area lightsources reduces the variation of integrand, but medium size lightsources might still pose problems. The incoming first-shot method proposed in the previous section is of limited use for these lightsources since the incoming first-shot requires l additional variables per patch, where l is the number of point samples on the lightsources, which may become very memory demanding. Thus a different approach is needed, which decomposes the transport operator instead of subdividing the lightsources.



looking to the right

looking forward

looking backward

Figure 8.28: Placement of the hemicube around a lightsource point and the images on the 5 hemicube faces



Figure 8.29: Ray-bundle stochastic iteration with incoming first-shot

Let us express the BRDF of the surfaces as a sum of the diffuse f_d and non-diffuse (specular) f_{nd} terms (note that the available BRDF representations do exactly this),

$$f_r(\omega', \vec{x}, \omega) = f_d(\vec{x}) + f_{nd}(\omega', \vec{x}, \omega)$$

and let us express the transport operator as the sum of diffuse and non-diffuse reflections:

$$\mathcal{T} = \mathcal{T}_d + \mathcal{T}_{nd},$$
$$\mathcal{T}_d L = \int_{\Omega_H} L(h(\vec{x}, -\omega'), \omega') \cdot \cos \theta' \cdot f_r(\vec{x}) \, d\omega', \quad \mathcal{T}_{nd} L = \int_{\Omega} L(h(\vec{x}, -\omega'), \omega') \cdot \cos \theta' \cdot f_{nd}(\omega', \vec{x}, \omega) \, d\omega'.$$

The basic idea of the "diffuse shot" technique is that $\mathcal{T}_d L^e$ can be calculated in a preprocessing phase. The storage of the found $\mathcal{T}_d L^e$ requires just one variable per patch (this is why we handled the diffuse reflection separately).



Figure 8.30: First step of the diffuse shot technique

Then, during the global walks, the first step should only be responsible for the non-diffuse reflection. The diffuse part is added to the result of this first step. The method can also be explained as a restructuring of the Neumann-series expansion of the solution of the rendering equation in the following way:

$$L = L^e + \mathcal{T}L^e + \mathcal{T}^2L^e + \mathcal{T}^3L^e + \dots = (L^e + \mathcal{T}_dL^e) + \mathcal{T}_{nd}L^e + \mathcal{T}(\mathcal{T}_dL^e + \mathcal{T}_{nd}L^e) + \dots$$
(8.61)

where $\mathcal{T}_d L^e$ is known after the preprocessing phase. Note that this method handles the first step in a special way, thus it requires the different bounces to be stored separately, as it is done by the self-correcting iteration or by combined and bi-directional walk methods. In self-correcting iteration, the formulae are reorganized in the following way to include ($\mathcal{T}_d L^e$) that is already known:

$$L'_{d}(m) = (\mathcal{T}_{d}L^{e}) + \mathcal{T}_{d}^{*}(m)(L_{nd}(m-1) + L_{d}(m-1)),$$

$$L_{d}(m) = \tau_{m} \cdot L'_{d}(m) + (1 - \tau_{m}) \cdot L_{d}(m-1),$$

$$L_{nd}(m) = \mathcal{T}_{nd}^{*}(m)(L^{e} + L_{nd}(m-1) + L_{d}(m-1)),$$

$$P_{m} = \mathcal{M}(L^{e} + L_{d}(m) + L_{nd}(m)),$$
(8.62)

Note that this restructuring replaced random variable $\mathcal{T}_d^* L^e$ by its mean $(\mathcal{T}_d L^e)$ which is a common variance reduction technique.

In order to present the idea for the walk methods, let us denote the result of the diffuse shot by L_{dif} . The calculation of the *d*-bounce irradiance J_d for d = 1, 2, ... is modified as follows:

$$\begin{aligned} \mathbf{J}_0 &= \mathbf{A}(\omega'_D) \cdot \mathbf{L}^e(\omega'_D), \\ \mathbf{J}_1 &= 4\pi \cdot \mathbf{A}(\omega'_{D-d}) \cdot \left(\mathbf{F}_{nd}(\omega'_{D-d+1}, \omega'_{D-d}) \cdot \mathbf{J}_0 + \mathbf{L}_{dif}\right), \\ \mathbf{J}_d &= 4\pi \cdot \mathbf{A}(\omega'_{D-d}) \cdot \mathbf{F}(\omega'_{D-d+1}, \omega'_{D-d}) \cdot \mathbf{J}_{d-1}, \end{aligned}$$

where \mathbf{F}_{nd} is the diagonal matrix of non-diffuse reflectance functions.

Now let us discuss the computation of a single diffuse reflection of the light coming from patch S of emission $L^e(\vec{y}, \omega)$ to a point \vec{x} . We can follow the hardware supported approach of the previous section, which places hemicubes at the lightsource, determines the visible patches and adds up the diffuse reflections of all sample points. On the other hand, another technique can also be used that works well for larger area lightsources and small receiver patches. If the patches are small, then L^e can be assumed to be constant for all \vec{y} points and directions pointing form \vec{y} to \vec{x} . Thus the diffuse reflection is

$$L_d(\vec{x}) = \mathcal{T}_d L^e = \int_{\Omega_H} L^e(h(\vec{x}, -\omega'), \omega') \cdot \cos \theta' \cdot f_r(\vec{x}) \, d\omega' \approx L^e(\vec{y}, \omega_{\vec{y} \to \vec{x}}) \cdot f_r(\vec{x}) \cdot \int_{\Omega_S} \cos \theta' \, d\omega',$$
(8.63)

where Ω_S is the solid angle in which patch S is visible, that is the projection of the patch onto the unit hemisphere. Note also that $\cos \theta'$ means projecting onto the base plane of the hemisphere, thus the computation of the diffuse reflection requires the area obtained by projecting the patch first onto the unit hemisphere then to the base circle.



Figure 8.31: Hemispherical projection of a planar polygon ($\vec{R_l}$ and $\vec{R_{l\oplus 1}}$ are two consecutive vertices)

To simplify the problem, consider only one edge line of the polygon first, and two consecutive vertices, \vec{R}_l and $\vec{R}_{l\oplus 1}$, on it (figure 8.31). Operator \oplus stands for modulo addition which can handle the problem that the next of vertex l is usually vertex l + 1, except for the last vertex which is followed by vertex 0. The hemispherical projection of this line is a half great circle. Since the radius of this great circle is 1, the area of the sector formed by the projections of \vec{R}_l and $\vec{R}_{l\oplus 1}$ and the center of the hemisphere is simply half the angle of \vec{R}_l and $\vec{R}_{l\oplus 1}$. Projecting this sector orthographically onto the equatorial plane, an ellipse sector is generated, having the area of the great circle sector multiplied by the cosine of the angle of the surface normal \vec{N} and the normal of the segment ($\vec{R}_l \times \vec{R}_{l\oplus 1}$). The area of the doubly projected polygon can be obtained by adding and subtracting the areas of the ellipse sectors of the different edges, as is demonstrated in figure 8.31, depending on whether the projections of vectors \vec{R}_l and $\vec{R}_{l\oplus 1}$ follow each other clockwise. This sign value can also be represented by a signed angle of the two vectors, expressing the area of the double projected polygon as a summation:

$$\sum_{l} \frac{1}{2} \cdot \operatorname{angle}(\vec{R}_{l}, \vec{R}_{l\oplus 1}) \frac{(\vec{R}_{l} \times \vec{R}_{l\oplus 1})}{|\vec{R}_{l} \times \vec{R}_{l\oplus 1}|} \cdot \vec{N}.$$
(8.64)

This method has supposed that the lightsource patch is above the plane of \vec{x} and is totally visible. Surfaces below the equatorial plane do not pose any problems, since we can get rid of them by the application of a clipping algorithm. When partial occlusion occurs, then either a continuous (object precision) visibility algorithm is used to select the visible parts of the surfaces, or the visibility term is estimated by firing several rays to surface element j and averaging their 0/1 associated visibilities. Since this means extra visibility computation, we prefer this method to the hemicube approach only if the lightsources are great and are seen from almost everywhere.

8.7.3 Adaptive importance sampling and resolution control in iteration

The light transfer might be significantly different in particular directions. According to importance sampling, important directions should be selected more often. If the probability density of the direction is $p(\omega)$, then the definition of the random transport operator is then:

$$\mathcal{T}^{*}(\omega')\mathbf{L}(\omega) = \frac{\mathbf{T}(\omega',\omega)\cdot\mathbf{L}(\omega')}{p(\omega)}.$$
(8.65)

Similarly, when the light is transferred into a single direction, particular regions might provide different contributions. To emphasize important regions, we can use higher resolution discrete visibility maps where the transfer is significant. To establish the required probability density $p(\omega)$ and the resolution plan, an adaptive strategy is proposed. Let us decompose the set of directions into finite subsets, using, for example, a spacing of the ϕ , θ angles, and associate with each subset an importance value which defines a discrete probability density $p(\omega)$ and a low-resolution visibility map storing also importance values that define where higher resolution is needed. Initially all importance values are equal. Running the algorithm, the transferred radiance is obtained and registered into the two importance maps. After a predefined number of steps, the importance is set proportionally to the average radiance values.

8.7.4 Reducing the power defect of the iteration

The power defect of the iteration can be reduced by the proper selection of the initial radiance function. In the ideal case the initial radiance is close to the real solution, thus the initial radiance is selected as a rough estimate of the final solution. To get this rough estimate, either the iteration can be run for a few samples without taking into account the measured values [Neu95, Sbe96], or the initial function is made equal to a single bounce of the lightsources which are computed by the methods of the previous sections.

Let us consider another method. Assume that the incoming radiance of each patch is constant and the initial radiance is obtained as the reflection of this constant radiance.

We shall suppose that the environment is closed. This is not a real limitation since any scene can be closed by adding closing surfaces and setting their BRDF to zero. Let the constant radiance be \tilde{L} . The initial radiance is then its reflection

$$L_0(\vec{x},\omega) = \mathcal{T}\tilde{L} = \int_{\Omega_H} \tilde{L} \cdot f_r(\omega', \vec{x}, \omega) \cdot \cos\theta' \, d\omega' = \tilde{L} \cdot a(\vec{x}, \omega),$$

where $a(\vec{x}, \omega)$ is the albedo. To find the constant \tilde{L} value, the total energy in the scene is set to be equal to the expected one. Using the rendering equation we have:

$$\int_{\Omega_H} \int_{S} \tilde{L} \cos \theta \, d\vec{x} \, d\omega = \int_{\Omega_H} \int_{S} \tilde{L}^e(\vec{x}, \omega) \cos \theta \, d\vec{x} \, d\omega + \int_{\Omega_H} \int_{S} \int_{\Omega_H} \tilde{L} \cdot f_r(\omega', \vec{x}, \omega) \cdot \cos \theta' \, d\omega' \cos \theta \, d\vec{x} \, d\omega.$$

Simplifying this, we obtain:

$$\tilde{L} = \frac{\int \int \tilde{L}^e(\vec{x},\omega)\cos\theta \, d\vec{x} \, d\omega}{S\pi \cdot (1-a_{\text{mean}})}, \quad \text{where} \quad a_{\text{mean}} = \frac{1}{S\pi} \cdot \int \int \int a(\vec{x},\omega)\cos\theta \, d\vec{x} \, d\omega.$$

In figure 9.13, the image sequence shows how the iteration enhances this initial estimate. The error measurements of figure 8.32 indicate that setting the initial radiance really improves the initial behaviour, but the benefits vanish later. The improvements take longer for self-correcting iteration than for normal stochastic iteration.



Figure 8.32: Efficiency improvements of the stochastic iteration (left) and for the self-correcting iteration (right) provided by initializing the radiance in the Cornell box scene (figure 8.21)

8.7.5 Constant radiance step

The average radiance \tilde{L} can be used not only for initializing the radiance, but also simplifying the global illumination problem similarly to [Neu96]. Let us decompose the unknown radiance function into this average \tilde{L} and a distance from this average ΔL . Substituting $L = \tilde{L} + \Delta L$ into and rendering equation we can obtain:

$$\Delta L(\vec{x},\omega) = L^e(\vec{x},\omega) + (a(\vec{x},\omega) - 1) \cdot L + \mathcal{T}\Delta L.$$
(8.66)

This is an easier — i.e. lower variation — problem than the original rendering equation if the new lightsource term $L^{e*} = L^e(\vec{x}, \omega) + (a(\vec{x}, \omega) - 1) \cdot \tilde{L}$ is flatter than $L^e(\vec{x}, \omega)$. This is the case if the lightsources do not reflect and the non-lightsource surfaces have high albedo.



Figure 8.33: Efficiency improvements of the constant radiance step in the Cornell box scene (figure 8.21)

As shown in figure 8.33, for the Cornell box-scene the improvements are negligible.

Chapter 9

Simulation results

The presented algorithms have been implemented in C++ in OpenGL environment and the program has been executed on Silicon Graphics Indigo2, Silicon Graphics O2 workstations and on a PC with 400 MHz Pentium II processor. The running times given in the following sections are measured on the PC.

The images have been generated by the ray-bundle walk method and by the different versions of stochastic iteration. The resolutions of the image and the visibility map are 800×800 and 600×600 , respectively.

Figure 9.1 compares the images obtained by the diffuse shot, first shot, diffuse radiosity and nondiffuse global illumination type ray-bundle tracing. The scene that contains a 3D Sierpiensky set has 22768 patches. The diffuse albedo of the patches in this set is (0.18, 0.06, 0.12) on the wavelengths 400 nm, 552 nm and on 700 nm, respectively. The specular albedo is wavelength independent and is between 0.8 and 0.4 depending on the viewing angle. The specular reflection has been modeled by the stretched-Phong model [NNSK98b]. The "shine" parameter is 3. These images demonstrate that simplifications of local illumination algorithms and ignoring of non-diffuse reflections might result in very inaccurate images.

9.1 Testing the walk method

Figure 9.2 shows a scene as rendered after the incoming first-shot applied only to the point lightsource and after 500 walks of length 5. The scene contains specular, metallic objects tessellated to 9605 patches, and is illuminated by both area (ceiling) and point (right-bottom corner) lightsources. The specular reflection has been modeled by a physically plausible modification of the Phong model that is particularly suitable for metals [NNSK98b]. A global radiance transfer took about 0.5 second on a Pentium II/400MHz computer. Since the radiance information of a single patch is stored on 3 wavelengths in 3×18 float variables (3 for the emission, 3 for the incoming radiance generated by the point lightsource, 3D(D+1)/2 = 45 for the incoming radiances of the steps of the walk and 3 for the accumulating radiance perceived from the eye), the extra memory used in addition to storing the scene is only 2.1 Mbytes. In figure 9.3 a fractal terrain can be seen that was rendered by global ray-bundle walks without any special preprocessing. The mountain is a realization of the 2D Brownian motion, which has been generated by a random midpoint perturbation algorithm setting the recursion level to 4 [SK99d]. The scene is illuminated by a spherical lightsource (the moon) and also by sky-light illumination. Note that the moon has been placed in front of the mountain to enable it to simultaneously illuminate the mountain and to have reflection in the lake. This is a rather difficult scene for the ray-bundle tracing to render since the lake is highly specular (the shine parameter is 30) and the snow has 0.98 albedo.



diffuse first shot

general first shot



radiosity solution

global illumination solution

Figure 9.1: Sierpiensky set (22.7k patches) after the diffuse first shot, the incoming first-shot (local illumination solution), the diffuse radiosity solution (diffuse global illumination) and the ray-bundle stochastic iteration (non-diffuse global illumination)



Figure 9.2: A golden Beethoven with metallic spheres (9.6k patches) after the incoming first-shot of the point lightsource at the right-bottom corner (left) and after 500 global ray-bundle walks (right)

9.2 Testing stochastic iteration and self-correcting iteration

Figure 9.7 shows a scene as rendered after 500 stochastic iterations and after 3000 iterations when the algorithm is fully converged. This pair of images demonstrates that this algorithm can provide good image quality even after relatively few number of iterations. The scene contains specular, metallic objects tessellated to 9519 patches. The calculation of the left image took 6 minutes. Figure 9.8 shows a scene containing 56745 patches after 300 stochastic iterations, which provide an accuracy within 5 percents (20 minutes computation time).



Figure 9.4: Convergence of stochastic iteration for the sphere-flake (figure 9.8) and for the golden Beethoven with teapot (figure 9.7)

The speed of convergence of stochastic iteration has been measured for the sphere-flake (figure 9.8) and for the room containing a Beethoven and a teapot (figure 9.7). The measurement results are shown in figure 9.4. Note that the algorithm converges faster for sphere-flake scene, which is due to the larger lightsources.

9.2.1 Self-correcting stochastic iteration

Figure 9.9 shows a fractal terrain containing 59614 patches rendered by self-correcting stochastic iteration (45 minutes computation time). The scene is similar to that of figure 9.3, but here the Brownian surface has been tessellated to more patches increasing the recursion level of the random midpoint perturbation algorithm to 7.

Figure 9.12 contains again the teapot and the Beethoven head, but now the BRDFs are not metallic. Figure 9.13 shows the evolution of the images for normal iteration, and also when the initial radiance estimation step has been applied. The measurement results for the self-correcting, stochastic iteration are in figure 9.5. The poorer asymptotic convergence for the mountain scene is the effect of the highly specular lake.

9.2.2 Self-correcting stochastic iteration with incoming first-shot

Figure 8.27 compares the speed of the convergence of stochastic iteration with and without the proposed incoming first-shot step for the Sierpiensky set. In figures 9.10, 9.11 and 9.13 the timing and the image quality of the two methods can also be compared. In the mountain scene of figure 9.10, 15 point samples are selected from the "moon". The incoming first-shot phase took 112 seconds, which were needed by the 15×5 z-buffer/constant-shading rendering steps.



Figure 9.5: Convergence of self-correcting stochastic iteration for the Sierpiensky set (figure 9.1), for the mountain with lake and moon (figure 9.10) and for the ceramic Beethoven with teapot (figure 9.12)

When rendering Sierpiensky set of figure 9.11 the area lightsource has been subdivided into a mesh of 8 triangles and 9 vertices. The incoming first-shot phase took 55 seconds. A single radiance transfer by a ray-bundle took 1.5 seconds without the first-shot results and 2 seconds when the incoming first-shot was also used. The 0.5 second overhead is due to the reflection of the result stored by the incoming first-shot both towards the eye and towards to next global direction. Despite the overhead, we can conclude that incoming first-shot is worth for this small extra time, since the resulting algorithm converges very quickly, and the image is almost fully converged after 2.5 minutes. Comparing the error curves, we can see that the stochastic iteration is about 20 times faster with the incoming first-shot than without it.



Figure 9.6: Convergence of self-correcting stochastic iteration having applied the incoming first-shot step for the Sierpiensky set (figure 9.1), for the mountain with lake and moon (figure 9.10) and for the ceramic Beethoven with teapot (figure 9.12)

The error curves of the renderings with self-correcting stochastic iteration combined with the incoming first-shot are shown in figure 9.6. These curves and also the image sequences in figures 9.10, 9.11 and 9.13 demonstrate that this alternative is far superior to the others.



Figure 9.7: A room with golden Beethoven and a teapot (9.5k patches) rendered by stochastic iteration after 500 steps i.e. 6 minutes (left) and when fully converged (right)



Figure 9.8: A golden sphere-flake (56k patches) illuminated by area lightsources rendered by stochastic iteration (300 iterations, 20 minutes)



Figure 9.9: Highly tessellated mountain with moon (60k patches) rendered with stochastic, self-correcting iteration



0 iteration, 0 sec

without incoming first-shot







200 iterations, 424 secs



first-shot + 0 iteration, 116 secs

with incoming first-shot







first-shot + 50 iterations, 272 secs first-shot + 100 iterations, 432 secs

Figure 9.10: Mountain with moon rendered with self-correcting iteration without and with the incoming radiance step



0 iteration, 0 secs

50 iterations, 68 secs

with incoming first-shot



100 iterations, 137 secs



first-shot + 0 iteration, 50 secs

first-shot + 10 iterations, 70 secs first-shot + 50 iterations, 150 secs

Figure 9.11: Comparison of stochastic iteration without (upper-row) and with incoming first-shot (lower-row)



Figure 9.12: Ceramic teapot and Beethoven rendered with stochastic, self-correcting iteration



0 iteration

stochastic iteration



100 iterations, 65 secs

self-correction





0 iteration



100 iterations, 66 secs





0 iteration



100 iterations, 67 secs



first-shot + 0 iteration, 48 secs



first-shot + 10 iterations, 59 secs



200 iterations, 132 secs



200 iterations, 134 secs



first-shot + 70 iterations, 125 secs

Figure 9.13: Ceramic teapot and Beethoven (12.7k patches) rendered by stochastic iteration (first row), by self-correction (second row), by self-correction with initial radiance estimation (third row), and by self-correction and incoming first-shot (fourth row)

Chapter 10

Conclusions

This thesis has proposed methods aiming at the efficient solution of the global illumination problem.

The thesis has reviewed the state-of-the-art of global illumination methods, emphasizing those that incorporate Monte-Carlo techniques. These reviews and comparisons have been published in [SKe95, SK99d, SK99a, SK99b]. The state-of-the-art has been improved by several techniques using coherence, quasi-Monte Carlo methods and importance sampling. The particular results and the references where they were published are as follows:

1 Stochastic iteration for the solution of the non-diffuse global illumination problem

In chapter 7 a general framework, called **stochastic iteration** has been proposed, to solve the rendering equation [SK98c, SK99c]. Stochastic iteration replaces the light-transport operator by a random transport operator generated from a finite-dimensional distribution. The expected value of the application of the random transport operator should be equivalent to the application of the light-transport operator. In the iteration the random transport operator is used and the measured value (i.e. pixel color) is obtained as the average of the measured radiances of the iteration steps. Compared to previous techniques involving Monte-Carlo techniques and iteration, this method is unique in the sense that it can attack the general form of the rendering equation, that is the non-diffuse global illumination problem. The proposed methods provide asymptotically correct results if the random operators are contractions and are not strongly correlated. It was also shown that stochastic iteration can eliminate the error accumulation problem and the prohibitive memory demand of classical iteration. A methodology of the elaboration of practical algorithms is also presented. An analogous approach of stochastic iteration could also be successfully used in very different application areas [SK94, SKMFF96, SKMFH00].

1.1 Single-ray stochastic iteration algorithm

In section 7.4.1 a single-ray based iteration algorithm is proposed [SK99a, SK99c], which uses the value of the radiance function in a single point and direction, thus it requires no tessellation and finite-element representation. Unlike other random-walk algorithms that make decisions based on local characteristics, this method also uses the information gathered at previous steps. Since this method requires BRDF sampling and albedo computation, BRDF models have also been developed that allow fast importance sampling [NNSK98b, NNSK99a, NNSK98a].

1.2 Quasi-Monte Carlo methods in iteration

In section 8.4.1, the application of low-discrepancy sequences in iteration has been formally analyzed and examined by simulation [SK99a, SK99c]. It turned out that only ∞ -distribution sequences can be

expected to provide correct results, Halton and Hammersley sequences fail. Simulation results have indicated that the $\{\pi^n\}$ deterministic series is appropriate.

2 Quasi-Monte Carlo methods for the solution of the rendering equation

In section 6.2 it has formally proven that quasi-Monte Carlo methods can effectively be used for the solution of the rendering equation [SKF97], and they are better than Monte-Carlo techniques for practical scenes, in spite of the fact that the integrand of the rendering equation is not of finite-variation in the sense of Hardy and Krause. We also concluded that quasi-Monte Carlo techniques are primarily worth applying for the estimation of lower bounces [SKP98b, SKP99a].

2.1 Importance sampling and Russian-roulette for quasi-Monte Carlo quadrature

The classical Monte-Carlo importance sampling and Russian-roulette have been generalized for quasi-Monte Carlo methods [SKCP98, SKCP99, SK98d].

3 Global illumination algorithm using ray-bundles

Applying finite-element representation for the positional variation of the radiance function, a projected form of the rendering equation has been established, which is the theoretical basis of the light-transport using complete wavefronts, called ray-bundles. Ray-bundles allow us to exploit the coherence of the radiance function. Two finite-element strategies, the Galerkin method with piece-wise constant basis functions and the point-collocation method with piece-wise linear basis functions have been examined and the related formulae have been developed.

3.1 Global visibility algorithms

Global visibility algorithms have been developed to allow fast tracing of ray-bundles in section 8.5 [SKF97, SKFNC97, SKFP98b, SKFP98a, SKP98a, SKP99b, SKM94], and the complexity of visibility algorithms have been analyzed [SKe95, HMSK92, MSK95, FMSK96, SKM96, SKF97, SKM97, SKM98a, SKM98b]. These algorithms can be classified as continuous and discrete. Continuous algorithms determine the visibility for all points in the scene, while discrete algorithms consider only sample points organized into a regular, raster grid. Continuous algorithms have mainly theoretical interest. Discrete algorithms, that are better than the continuous algorithms both in speed and in storage complexities, on the other hand, result in very fast solutions and also allow the utilization of the built-in rendering hardware. The resolution limit of discrete methods has been examined and a Russian-roulette-like random visibility algorithm has been proposed to reduce the resolution requirement of discrete methods [SK98a]. The computation time depends on the total surface area and the resolution of the visibility map. Since only the expected value of the visible surface area should be accurately computed, and an actual rasterization can be very coarse, lower resolution visibility maps can also be used. The limit is when the projected patch size becomes comparable to the pixel size, since classical filling algorithms always generate an approximation whose height and width are at least 1. This problem can be solved by modifying the filling algorithm to handle patches or spans randomly if their width or height is less than 1. These low resolution maps increase the variance, thus slow down the convergence a little bit, but still provide unbiased results and significantly reduce the computation time of a single transfer. The random visibility approach can also be applied to speed-up continuous algorithms [SK98a].

3.2 Ray-bundle tracing

In section 8.2 a new, combined finite-element and random-walk algorithm, called the **ray-bundle tracing**, has been presented to solve the rendering problem of complex scenes including also glossy surfaces [SKFP98b, SKP98a, SK98a, SKP99b]. The method is fundamentally different from other global radiosity algorithms, such as the method of global lines or the transillumination method in its ability to solve the general, non-diffuse rendering problem. It applies a formulation that makes the integrand have finite variation, thus the efficiency of low-discrepancy sequences can be fully exploited [SKF97]. The basic idea of the method is to form bundles of parallel rays that can be traced efficiently, taking advantage of image-coherence and the built-in rendering hardware. The basic algorithm has also been extended to combine the results of the steps of a walk and also to take bi-directional steps.

Unlike other random walk methods using importance sampling, this approach does not emphasize the locally important directions, but handles a large number (1 million or even infinite) parallel rays simultaneously instead, thus it is more efficient than those methods when the surfaces are not very specular.

The time complexity of the algorithm depends on the used global visibility algorithm. For example, the global painter's algorithm has $O(n \log n)$ complexity (*n* is the number of patches) [SKF97], which is superior to the $O(n^2)$ complexity of classical, non-hierarchical radiosity algorithms [SKM95].

The memory requirement is comparable to that of the diffuse radiosity algorithms, although the new algorithm is also capable to handle non-diffuse reflections or refractions. Thus this method does not suffer from the prohibitive memory hunger of other non-diffuse finite-element approaches. Since global ray-bundle walks are computed independently, the algorithm is very well suited for parallelization.

The application of importance sampling for the ray-bundle tracing has been studied in section 8.3.2. The Metropolis and VEGAS methods have been considered, but only the Metropolis method was examined in details, using both theoretical considerations and simulation study [SKP98a, SKDP99]. For homogeneous scenes, Metropolis sampling could not provide significant noise reduction compared to quasi-Monte Carlo walks. This is due to the fact that the integrand of equation (8.5) is continuous and is of finite variation unlike the integrand of the original rendering equation, thus if its variation is modest then quasi-Monte quadrature is almost unbeatable. If the radiance distribution has high variation (difficult lighting conditions), then the Metropolis method becomes more and more superior. On the other hand, in section 8.7.3 an adaptive — i.e. sequential Monte-Carlo type — importance sampling approach was proposed for ray-bundle based iteration.

3.3 Stochastic iteration with ray-bundles

In section 8.4 a stochastic iteration algorithm using ray-bundles is also proposed that can efficiently render scenes of moderate specularity illuminated by larger area lightsources [SK98c, SK99c]. This algorithm provides unbiased estimates and seems to be significantly better than the finite-length approaches, in terms of both speed and storage space. The performance of stochastic iteration has been further improved by finding a good initial radiance distribution and by incorporating self-correcting iteration estimates for the diffuse part of the radiance function. This allows more lightpaths to be included into the integral quadrature provided by a certain number of iterations.

3.4 Incoming first-shot

To handle small lightsources, the **incoming first-shot** method has been proposed for non-diffuse algorithms in section 8.7.2 [SKP98a, SKSMT00]. On the other hand, a similar method, called **diffuse first-shot** has been introduced to handle medium size lightsources. These preprocessing techniques speed-up stochastic iteration by considerable amount, and allow the global illumination rendering of complex scenes in a few minutes.

Since the proposed methods are able to render complex scenes with physical correctness in reasonable time, the primary application areas of the methods include CAD, lighting design, terrain visualization, computer aided motion picture development, etc.

10.1 Future improvements

The single-ray based stochastic iteration is being implemented with heuristic self-correction. If selfcorrecting iteration were used directly, then each iteration step would introduce a new point to compete for random selection. To limit the hunger for memory, after a certain steps the population should be decimized randomly. By incorporating visual importance into the survival probabilities, bi-directional algorithms can be obtained. We concluded that the single-ray based algorithm automatically estimates the average contraction of the light transport and uses this estimate to randomly terminate the walk.

We intend to incorporate an **adaptive tessellation** method into the ray-bundle tracing algorithm. Note that when computing the radiance transport for a given patch, the variation of the incoming radiance can also be easily estimated. If this variation exceeds a given limit, then we cannot assume that the outgoing radiance of the patch is homogeneous, thus the patch has to be subdivided. This tessellation scheme is much more robust than those methods which examine the radiosity gradient. This tessellation scheme also detects highlights that are completely inside a patch and can even provide information where the patch should be subdivided. Considering this, the method is also able to do **discontinuity meshing**.

The algorithm seems to be particularly efficient to handle **participating media** since it can handle very many parallel lines simultaneously. Participating media can be modeled as a set of partly transparent "points" that always project onto a single pixel in the visibility map. Since this point field can be rendered very quickly, the radiance transfer in a single direction can be computed very efficiently.

The Metropolis method seems to be worth using only for very difficult lighting conditions. On the other hand, the Metropolis method is sensitive to its parameters such as the extent of perturbation. Future research should concentrate on the automatic and "optimal" determination of these control parameters.

Global walks are computed completely independently. Stochastic iteration, on the other hand, depends on the previous sample, but it can also be decomposed into several parallel runs using different randomization sequences. Thus the algorithm can easily be ported to **parallel computing** systems.

The presented methods are particularly efficient if the surfaces are not highly specular and the scene is illuminated by larger area lightsources. This is the application domain where other random walks becomes inefficient. Thus it is worth joining the two approaches using **multiple importance sampling** in a way that the strengths of the two approaches can be preserved.

Since the iteration is basically view-independent, just the result of each iteration step is projected to the eye, images for many cameras can be computed simultaneously, which can be used to produce **animation**. Since the radiance is stored in the object space, if the surfaces are not highly specular, then the same radiance information remains valid for a wider range of viewing directions. It means that the images must be computed by stochastic iteration just for a few camera orientations (viewing directions), and for the inbetweening frames, the radiance of the patches can be interpolated. In order to save space, the radiance is stored just for a few directions at less specular surfaces.

Due to the high-speed of self-correcting stochastic iteration combined with incoming first-shot, the method seems appropriate for interactive walk-through in scenes of glossy objects of moderate complexity. Note that in this approach the incoming radiance and the diffuse reflection are view independent thus are valid for all view positions. When either the view position of the viewing direction change, the only factor which should be recomputed by stochastic iteration is the indirect specular (i.e. glossy) reflection. In order to speed-up this calculation, each patch is associated with a single variable I that represents the average indirect incoming radiance from the whole hemisphere. Variable I is also view independent. Recall that the response to the homogeneous illumination is the incoming homogeneous radiance times the albedo, thus when the viewing direction changes from ω to ω_{new} , then the reflected glossy illumination should be corrected approximately by $(a(\omega_{new}) - a(\omega)) \cdot I$. Due to the assumption of homogeneous indirect illumination, this is just an approximation, but is quite accurate and can be accepted as a good starting point of subsequent iteration steps. An even more precise approach would associate not one but several average incoming radiance variables with each patch, that can represent different incoming solid angles. Since when the actual radiance is close to the solution, the stochastic iteration converges after a few steps, and a single iteration for a scene of $10^3 - 10^4$ patches requires about 0.05 - 0.5 seconds, which results in a new frame in about a second.

Acknowledgements

This work has been supported by the *National Scientific Research Fund* (OTKA), ref.No.: F 015884 and ref.No.: T029135, by the *Austrian-Hungarian Action Fund*, ref.No.: 29p4, 32öu9 and 34öu28, and by the *Spanish-Hungarian Research Fund*, ref.No.: E9. The research work has been carried out at the Department of Control Engineering and Information Technology of the Technical University of Budapest, and partly at the Institute of Computer Graphics of the Vienna University of Technology and at the Department of Applied Mathematics and Informatics of the University of Girona. Prof. Péter Arató department head, Prof. Béla Lantos and the whole staff of the Department of Control Engineering and Information Support during the work.

Chapters of this thesis have also been presented as Ph.D. courses at the Technical University of Budapest, University of Girona and at the Vienna University of Technology [SK99b]. Special thanks go to Prof. Werner Purgathofer for the valuable discussions and for making the resources of the Institute of Computer Graphics of the Vienna University of Technology available for this research and to Prof. Mateu Sbert from the University of Girona whose ideas inspired many problems considered in this research. The discussions with professors and the most active Ph.D. students of the related courses, including among others Robert F. Tobler, Francesc Castro, Roel Martinez, Jan Prikryl, Anna Vilanova Bartoli, Jiri Hladuvka and Helmut Mastal, also helped to put the ideas in a simpler form.

I am also grateful to the Hungarian Telematics Co. where some results of this work have been built into a commercial system visualizing industrial processes and where some technical parts of the work have been realized, and to its employees for their technical help and logistical support, including Péter Risztics, István Jankovits, Elefteria Szertaridisz and Mária Keszthelyi. When I implemented the algorithms, I re-used classes and modules from the software developed by the members of the computer graphics group over the years, including the sphere-flake generator and parts from the radiosity program written by Tibor Fóris, a ray-tracing program by Balázs Csébfalvi and I also utilized model conversion and text, figure and image processing tools made by Péter Dornbach and Tamás Horváth. Thanks also go to mathematicians Sarolta Inczédy, László Neumann, Attila Neumann and István Deák for drawing my attention to the mathematical and computer graphics techniques on which the presented methods are built, including, for example, the transillumination method, and for urging me to pursue this research. I appreciate the much-needed moral support and tolerance provided by my friends and family members.

And at last, we cannot forget the tremendous compiling, computation, rendering and text processing work carried out by Jenny (Silicon Graphics Indigo 2), Kata (PC 400 Mhz Pentium II), Susan (HP-Apollo 720), Thais, Hengest (Sun), Ringlotte, Marille, Ribisl, Zwetschke, Knieriem (Silicon Graphics O2) and Schnoferl (Silicon Graphics Indy).

Related publications

Books and scripts

- [SK99b] L. Szirmay-Kalos. *Monte-Carlo Methods in Global Illumination*. Institute of Computer Graphics, Vienna University of Technology, Vienna, 1999.
- [SK99d] L. Szirmay-Kalos. Számítógépes grafika. ComputerBooks, Budapest, 1999.
- [SKe95] L. Szirmay-Kalos (editor). Theory of Three Dimensional Computer Graphics. Akadémia Kiadó, Budapest, 1995.

Journals

- [SK99c] L. Szirmay-Kalos. Stochastic iteration for non-diffuse global illumination. Computer Graphics Forum (Eurographics'99), 18(3):233–244, 1999.
- [SKP99b] L. Szirmay-Kalos and W. Purgathofer. Global ray-bundle tracing with infinite number of rays. *Computers and Graphics*, 23(2):193–202, 1999.
- [SKCP99] L. Szirmay-Kalos, B. Csébfalvi, and W. Purgathofer. Importance driven quasi-random walk solution of the rendering equation. *Computers and Graphics*, 23(2):203–212, 1999.
- [NNSK99a] L. Neumann, A. Neumann, and L. Szirmay-Kalos. Compact metallic reflectance models. Computer Graphics Forum (Eurographics'99), 18(3):161–172, 1999.
- [NNSK99b] L. Neumann, A. Neumann, and L. Szirmay-Kalos. Reflectance models by pumping up the albedo function. *Machine Graphics and Vision*, 8(1):3–18, 1999.
- [NNSK99c] L. Neumann, A. Neumann, and L. Szirmay-Kalos. Reflectance models with fast importance sampling. *Computer Graphics Forum*, 18(4):249–265, 1999.
- [SKMFH00] L. Szirmay-Kalos, G. Márton, T. Fóris, and T. Horváth. Development of process visualization systems - an object oriented approach. *Journal of System Architecture*, 46:275–296, 2000.
- [SKFP98a] L. Szirmay-Kalos, T. Fóris, and W. Purgathofer. Non-diffuse, random-walk radiosity algorithm with linear basis functions. *Machine Graphics and Vision*, 7(1):475–484, 1998.
- [SKM98b] L. Szirmay-Kalos and G. Márton. Worst-case versus average-case complexity of ray-shooting. Journal of Computing, 61(2):103–133, 1998.
- [SKM98a] L. Szirmay-Kalos and G. Márton. Construction and analysis of worst-case optimal ray-shooting algorithms. *Computers and Graphics*, 22(2):793–806, 1998.
- [CSK98] B. Csébfalvi and L. Szirmay-Kalos. Interactive volume rotation. Machine Graphics and Vision, 7(4):793–806, 1998.
- [SKFNC97] L. Szirmay-Kalos, T. Fóris, L. Neumann, and B. Csébfalvi. An analysis to quasi-Monte Carlo integration applied to the transillumination radiosity method. *Computer Graphics Forum (Eurographics'97)*, 16(3):271–281, 1997.
- [SKHR97] L. Szirmay-Kalos, T. Horváth, and P. Risztics. Generalization and solution of reward models. *IASTED Int. Journal of Modelling and Simulation*, June 1997.
- [SK96] L. Szirmay-Kalos. Application of variational calculus in the radiosity method. *Periodica Polytechnica (Electrical Engineering)*, 40(2):123–138, 1996.
- [HMSK92] T. Horváth, P. Márton, G. Risztics, and L. Szirmay-Kalos. Ray coherence between sphere and a convex polyhedron. *Computer Graphics Forum*, 2(2):163–172, 1992.

Chapters in books

- [SKP98a] L. Szirmay-Kalos and W. Purgathofer. Global ray-bundle tracing with hardware acceleration. In *Rendering Techniques '98*, Springer, pages 247–258, 1998.
- [SK94] L. Szirmay-Kalos. Dynamic layout algorithm to display general graphs. In Paul Heckbert, editor, *Graphics Gems IV*. Academic Press, Boston, 1994.

Conference proceedings

- [SK99a] L. Szirmay-Kalos. Monte-Carlo methods for global illumination. In *Spring Conference of Computer Graphics '99*, pages 1–28, Budmerice, 1999. invited talk.
- [SKP99a] L. Szirmay-Kalos and W. Purgathofer. Analysis of the quasi-monte carlo integration of the rendering equation. In *Winter School of Computer Graphics '99*, pages 281–288, Plzen, 1999.
- [SKDP99] L. Szirmay-Kalos, P. Dornbach, and W. Purgathofer. On the start-up bias problem of metropolis sampling. In *Winter School of Computer Graphics* '99, pages 273–280, Plzen, 1999.
- [SKFP98b] L. Szirmay-Kalos, T. Fóris, and W. Purgathofer. Quasi-Monte Carlo global ray-bundle tracing with infinite number of rays. In *Winter School of Computer Graphics* '98, pages 386–393, Plzen, 1998.
- [SKCP98] L. Szirmay-Kalos, B. Csébfalvi, and W. Purgathofer. Importance-driven quasi-Monte Carlo solution of the rendering equation. In *Winter School of Computer Graphics* '98, pages 377–386, Plzen, 1998.
- [SKF97] L. Szirmay-Kalos and T. Fóris. Radiosity algorithms running in sub-quadratic time. In Winter School of Computer Graphics '97, pages 562–571, Plzen, 1997.
- [SKM97] L. Szirmay-Kalos and G. Márton. On the limitations of worst-case optimal ray-shooting algorithms. In Winter School of Computer Graphics '97, pages 552–561, Plzen, 1997.
- [CMSK97] B. Csébfalvi, G. Márton, and L. Szirmay-Kalos. Fast opacity control of volumetric ct data. In Winter School of Computer Graphics '97, pages 79–87, Plzen, 1997.
- [SKMFF96] L. Szirmay-Kalos, G. Márton, T. Fóris, and J. Fábián. Application of object-oriented methods in process visualisation. In *Winter School of Computer Graphics* '96, pages 349–358, Plzen, 1996.
- [FMSK96] T. Fóris, G. Márton, and L. Szirmay-Kalos. Ray-shooting in logarithmic time. In Winter School of Computer Graphics '96, pages 84–90, Plzen, 1996.
- [SKM96] L. Szirmay-Kalos and G. Márton. On the complexity of ray-tracing. In *Dagstuhl Seminar on Rendering*, Dagstuhl, Germany, June 1996.
- [SK95] L. Szirmay-Kalos. Stochastic sampling of two-dimensional images. In *COMPUGRAPHICS* '95, Alvor, 1995.
- [SKM95] L. Szirmay-Kalos and G. Márton. On convergence and complexity of radiosity algorithms. In Winter School of Computer Graphics '95, pages 313–322, Plzen, 1995.
- [MSK95] G. Márton and L. Szirmay-Kalos. On average-case complexity of ray tracing algorithms. In Winter School of Computer Graphics '95, pages 187–196, Plzen, 1995.
- [SKM94] L. Szirmay-Kalos and G. Márton. On hardware implementation of scan-conversion algorithms. In 8th Symp. on Microcomputer Appl., Budapest, 1994.
- [SK93] L. Szirmay-Kalos. Global element method in radiosity calculation. In *COMPUGRAPHICS '93*, Alvor, 1993.
- [DRSK92] B. Dobos, P. Risztics, and L. Szirmay-Kalos. Fine-grained parallel processing of scan conversion with i860 microprocessor. In 7th Symp. on Microcomputer Appl., Budapest, 1992.
- [MRSK92] G. Márton, P. Risztics, and L. Szirmay-Kalos. Quick ray-tracing exploiting ray coherence theorems. In 7th Symp. on Microcomputer Appl., Budapest, 1992.
- [HKRSK91] T. Horváth, E. Kovács, P. Risztics, and L. Szirmay-Kalos. Hardware-software-firmware decomposition of high-performace 3d graphics systems. In 6th Symp. on Microcomputer Appl., Budapest, 1991.

Technical reports

- [NNSK98b] L. Neumann, A. Neumann, and L. Szirmay-Kalos. New simple reflectance models for metals and other specular materials. Technical Report TR-186-2-98-17, Institute of Computer Graphics, Vienna University of Technology, 1998. www.cg.tuwien.ac.at/.
- [SK98a] L. Szirmay-Kalos. Global ray-bundle tracing. Technical Report TR-186-2-98-18, Institute of Computer Graphics, Vienna University of Technology, 1998. www.cg.tuwien.ac.at/.
- [SK98b] L. Szirmay-Kalos. Object-oriented framework and methodology to process visualization system development. Technical Report TR-186-2-98-19, Institute of Computer Graphics, Vienna University of Technology, 1998. www.cg.tuwien.ac.at/.
- [NNSK98a] L. Neumann, A. Neumann, and L. Szirmay-Kalos. Analysis and pumping up the albedo function. Technical Report TR-186-2-98-20, Institute of Computer Graphics, Vienna University of Technology, 1998. www.cg.tuwien.ac.at/.
- [SK98c] L. Szirmay-Kalos. Stochastic iteration for non-diffuse global illumination. Technical Report TR-186-2-98-21, Institute of Computer Graphics, Vienna University of Technology, 1998. www.cg.tuwien.ac.at/.
- [SK98d] L. Szirmay-Kalos. Stochastic methods in global illumination state of the art report. Technical Report TR-186-2-98-23, Institute of Computer Graphics, Vienna University of Technology, 1998. www.cg.tuwien.ac.at/.
- [SKP98b] L. Szirmay-Kalos and W. Purgathofer. Quasi-Monte Carlo solution of the rendering equation. Technical Report TR-186-2-98-24, Institute of Computer Graphics, Vienna University of Technology, 1998. www.cg.tuwien.ac.at/.
- [SKSMT00] L. Szirmay-Kalos, M. Sbert, R. Martinez, and R. F. Tobler. Incoming first-shot for non-diffuse global illumination. Technical Report TR-186-2-00-04, Institute of Computer Graphics, Vienna University of Technology, 2000. www.cg.tuwien.ac.at/.
- [SK88] L. Szirmay-Kalos. Árnyalási modellek a háromdimenziós raszter grafikában (Szakszemináriumi Füzetek 30). BME, Folyamatszabályozási Tanszék, 1988.

BIBLIOGRAPHY

[Ábr97]	Gy. Ábrahám. Optika. Panem-McGraw-Hill, Budapest, 1997.
[AH93]	L. Aupperle and P. Hanrahan. A hierarchical illumination algorithms for surfaces with glossy reflection. <i>Computer Graphics (SIGGRAPH '93 Proceedings)</i> , pages 155–162, 1993.
[AK90]	J. Arvo and D. Kirk. Particle transport and image synthesis. In <i>Computer Graphics (SIGGRAPH '90 Proceedings)</i> , pages 63–66, 1990.
[Ant80]	J. Antal. Fizikai kézikönyv műszakiaknak. Műszaki Könyvkiadó, Budapest, 1980.
[Arv95]	J. Arvo. Stratified sampling of spherical triangles. In <i>Computer Graphics (SIGGRAPH '95 Proceedings)</i> , pages 437–438, 1995.
[BBS96]	G. Baranoski, R. Bramley, and P Shirley. Fast radiosity solutions for environments with high average reflectance. In <i>Rendering Techniques '96</i> , pages 345–355, 1996.
[Bek97]	P. Bekaert. Error control for radiosity. In Rendering Techniques '97, Porto, Portugal, 1997.
[Bek99]	Ph. Bekaert. <i>Hierarchical and stochastic algorithms for radiosity</i> . PhD thesis, University of Leuven, 1999. http://www.cs.leuven.ac.be/cwis/research/graphics/CGRG.PUBLICATIONS/PHBPPHD.
[BF89]	C. Buckalew and D. Fussell. Illumination networks: Fast realistic rendering with general reflectance functions. <i>Computer Graphics (SIGGRAPH '89 Proceedings)</i> , 23(3):89–98, July 1989.
[BKP91]	J. C. Beran-Koehn and M. J. Pavicic. A cubic tetrahedral adaptation of the hemicube algorithm. In James Arvo, editor, <i>Graphics Gems II</i> , pages 299–302. Academic Press, Boston, 1991.
[Bli77]	J. F. Blinn. Models of light reflection for computer synthesized pictures. In <i>Computer Graphics</i> (SIGGRAPH '77 Proceedings), pages 192–198, 1977.
[BNN ⁺ 98]	P. Bekaert, L. Neumann, A. Neumann, M. Sbert, and Y. Willems. Hierarchical Monte-Carlo radios- ity. In <i>Rendering Techniques '98</i> , pages 259–268, 1998.
[BS95]	P. Bodrogi and J. Schanda. Testing the calibration model of colour crt monitors. <i>Displays</i> , 16(3):123–133, 1995.
[BS96]	G. Besuievsky and M. Sbert. The multi-frame lighting method - a monte-carlo based solution for radiosity in dynamic environments. In <i>Rendering Techniques '96</i> , pages pp 185–194, 1996.
[CG85]	M. Cohen and D. Greenberg. The hemi-cube, a radiosity solution for complex environments. In <i>Computer Graphics (SIGGRAPH '85 Proceedings)</i> , pages 31–40, 1985.
[Chi88]	H. Chiyokura. Solid Modelling with DESIGNBASE. Addision Wesley, 1988.
[CLSS97]	P. H. Christensen, D. Lischinski, E. J. Stollnitz, and D. H. Salesin. Clustering for glossy global illumination. <i>ACM Transactions on Graphics</i> , 16(1):3–33, 1997.
[CMS98]	F. Castro, R. Martinez, and M. Sbert. Quasi Monte-Carlo and extended first-shot improvements to the multi-path method. In L. Szirmay-Kalos, editor, <i>Spring Conference on Computer Graphics '98</i> , pages 91–102, 1998.
[CPC84]	R. Cook, T. Porter, and L. Carpenter. Distributed ray tracing. In <i>Computer Graphics (SIGGRAPH '84 Proceedings)</i> , pages 137–145, 1984.
[Cse97]	B. Csébfalvi. A review of Monte-Carlo ray tracing algorithms. In <i>CESCG '97</i> , Central European Seminar on Computer Graphics, Vienna, pages 87–103, 1997.
[CSSD96]	P. H. Christensen, E. J. Stollnitz, D. H. Salesin, and T. D. DeRose. Global illumination of glossy

environments using wavelets and importance. ACM Transactions on Graphics, 15(1):37-71, 1996.

[CT81]	R. Cook and K. Torrance. A reflectance model for computer graphics. <i>Computer Graphics</i> , 15(3), 1981.
[dB92]	M. de Berg. <i>Efficient Algorithms for Ray Shooting and Hidden Surface Removal</i> . PhD thesis, Rijksuniversiteit te Utrecht, The Nederlands, 1992.
[Do98]	P. Dornbach. Implementation of bidirectional ray-tracing algorithm. In <i>CESCG '98</i> , Central European Seminar on Computer Graphics, Budmerice, 1998.
[DBW97]	Ph. Dutre, Ph. Bekaert, and Y. D. Willems. Bidirectional radiosity. In <i>Rendering Techniques '97</i> , pages 205–216, 1997.
[Deá89]	I. Deák. Random Number Generators and Simulation. Akadémia Kiadó, Budapest, 1989.
[Deá97]	I. Deák. Monte Carlo módszerek a többdimenziós térben elhelyezkedő halmazok valószínűségének meghatározására normális eloszlás esetén. PhD thesis, MTA, Hungary, 1980.
[Dév93]	F. Dévai. Computational Geometry and Image Synthesis. PhD thesis, MTA, Hungary, 1993.
[DLW93]	Ph. Dutre, E. Lafortune, and Y. D. Willems. Monte Carlo light tracing with direct computation of pixel intensities. In <i>Compugraphics '93</i> , pages 128–137, Alvor, 1993.
[DW96]	Ph. Dutre and Y. D. Willems. Potential-driven Monte Carlo particle tracing for diffuse environments with adaptive probability functions. In <i>Rendering Techniques</i> '96, pages 306–315, 1996.
[EH94]	Gröeller E. and Löffelmann H. Extended camera specification for image synthesis. <i>Machine Graphics and Vision</i> , 3:513–530, 1994.
[Erm75]	S.M. Ermakow. <i>Die Monte-Carlo-Methode und verwandte Fragen</i> . R. Oldenbourg Verlag, Wien, 1975.
[Fed95]	M. Feda. <i>Radiosity</i> . Institute of Computer Graphics, Vienna University of Technology, Vienna, 1995.
[FP94]	M. Feda and W. Purgathofer. A median cut algorithm for efficient sampling of radiosity functions. In <i>Eurographics'94</i> , 1994.
[Ga99]	A. Gascón and C. Coll. Virtual Newton Telescope. In <i>CESCG '99</i> Central European Seminar on Computer Graphics, Budmerice, 1999.
[Gla95]	A. Glassner. <i>Principles of Digital Image Synthesis</i> . Morgan Kaufmann Publishers, Inc., San Francisco, 1995.
[Hec91]	P. S. Heckbert. <i>Simulating Global Illumination Using Adaptive Meshing</i> . PhD thesis, University of California, Berkeley, 1991.
[Her91]	Ivan Herman. The Use of Projective Geometry in Computer Graphics. Springer-Verlag, Berlin, 1991.
[HMF98]	M. Hyben, I. Martisovits, and A. Ferko. Scene complexity for rendering in flatland. In L. Szirmay-Kalos, editor, <i>Spring Conference on Computer Graphics '98</i> , pages 112–120, 1998.
[HSA91]	P. Hanrahan, D. Salzman, and L. Aupperle. Rapid hierachical radiosity algorithm. <i>Computer Graphics (SIGGRAPH '91 Proceedings)</i> , 1991.
[HTSG91]	X. He, K. Torrance, F. Sillion, and D. Greenberg. A comprehensive physical model for light reflection. <i>Computer Graphics</i> , 25(4):175–186, 1991.
[ICG86]	D. S. Immel, M. F. Cohen, and D. P. Greenberg. A radiosity method for non-diffuse environments. In <i>Computer Graphics (SIGGRAPH '86 Proceedings)</i> , pages 133–142, 1986.
[JC95]	H. W. Jensen and N. J. Christensen. Photon maps in bidirectional Monte Carlo ray tracing of complex objects. <i>Computers and Graphics</i> , 19(2):215–224, 1995.
[JC98]	H. W. Jensen and P. H. Christensen. Efficient simulation of light transport in scenes with participating media using photon maps. <i>Computers and Graphics (SIGGRAPH '98 Proceedings)</i> , pages 311–320, 1998.
[Jen95]	H. W. Jensen. Importance driven path tracing using the photon maps. In <i>Rendering Techniques '95</i> , pages 326–335, 1995.
[Jen96]	H. W. Jensen. Global illumination using photon maps. In <i>Rendering Techniques '96</i> , pages 21–30, 1996.
-----------------------	--
[Ko73]	A. Kósa. Variációszámítás. Tankönyvkiadó, Budapest, 1973.
[Kaj86]	J. T. Kajiya. The rendering equation. In <i>Computer Graphics (SIGGRAPH '86 Proceedings)</i> , pages 143–150, 1986.
[Kel95]	A. Keller. A quasi-Monte Carlo algorithm for the global illumination in the radiosity setting. In H. Niederreiter and P. Shiue, editors, <i>Monte-Carlo and Quasi-Monte Carlo Methods in Scientific Computing</i> , pages 239–251. Springer, 1995.
[Kel96a]	A. Keller. Quasi-Monte Carlo Radiosity. In X. Pueyo and P. Schröder, editors, <i>Rendering Techniques</i> '96, pages 101–110. Springer, 1996.
[Kel96b]	A. Keller. The fast Calculation of Form Factors using Low Discrepancy Sequences. In <i>Proc. Spring Conference on Computer Graphics (SCCG '96)</i> , pages 195–204, Bratislava, Slovakia, 1996. Comenius University Press.
[Kel97]	A. Keller. Instant radiosity. Computer Graphics (SIGGRAPH '97 Proceedings), pages 49–55, 1997.
[Kel98]	A. Keller. Quasi-Monte Carlo Methods for Photorealistic Image Synthesis. Shaker-Verlag, 1998.
[Knu81]	D.E. Knuth. <i>The art of computer programming. Volume 2 (Seminumerical algorithms).</i> Addison-Wesley, Reading, USA, 1981.
[Kra89]	G. Krammer. Notes on the mathematics of the PHIGS output pipeline. <i>Computer Graphics Forum</i> , 8(8):219–226, 1989.
[Kra99]	G. Krammer. <i>Bevezetés a számítógépi grafikába</i> . Jegyzet, ELTE, 1999. http://valerie.inf.elte.hu/ ~ krammer/eltettk/grafika/jegyzet/index.html.
[LA98]	G. Lukács and L. Andor. Photometric ray tracing. In <i>The Mathematics of Surfaces VIII, Editors: R. J. Cripps and R. R. Martin</i> , pages 325–338, The Institute of Mathematics and its Applications, Information Geometers Ltd, 1998.
[Lan91]	B. Lantos. Robotok Irányítása. Akadémiai Kiadó, Budapest, Hungary, 1991.
[LB94]	B. Lange and B. Beyer. Rayvolution: An evolutionary ray tracing algorithm. In <i>Photorealistic Rendering Techniques</i> , pages 136–144, 1994.
[Lep80]	G. P. Lepage. An adaptive multidimensional integration program. Technical Report CLNS-80/447, Cornell University, 1980.
[Lew93]	R. Lewis. Making shaders more physically plausible. In <i>Rendering Techniques '93</i> , pages 47–62, 1993.
[LW93]	E. Lafortune and Y. D. Willems. Bi-directional path-tracing. In <i>Compugraphics '93</i> , pages 145–153, Alvor, 1993.
[LW96]	E. Lafortune and Y. D. Willems. A 5D tree to reduce the variance of Monte Carlo ray tracing. In <i>Rendering Techniques '96</i> , pages 11–19, 1996.
[Már95a]	G. Márton. Acceleration of ray tracing via Voronoi-diagrams. In Alan W. Paeth, editor, <i>Graphics Gems V</i> , pages 268–284. Academic Press, Boston, 1995.
[Már95b]	G. Márton. Sugárkövető algoritmusok átlagos bonyolultságának vizsgálata. PhD thesis, MTA, Hungary, 1995.
[Mát81]	L. Máté. Funkcionálanalízis műszakiaknak. Műszaki Könyvkiadó, Budapest, 1981.
[Min41]	M. Minnaert. The reciprocity principle in lunar photometry. <i>Astrophysical Journal</i> , 93:403–410, 1941.
[Mit92]	D. Mitchell. Ray Tracing and Irregularities of Distribution. In <i>Rendering Techniques '92</i> , pages 61–69, Bristol, UK, 1992.
[Mit96]	D. P. Mitchell. Consequences of stratified sampling in graphics. <i>Computer Graphics (SIGGRAPH '96 Proceedings)</i> , pages 277–280, 1996.
[MRR ⁺ 53]	N. Metropolis, A. Rosenbluth, M. Rosenbluth, A. Teller, and E. Teller. Equations of state calculations by fast computing machines. <i>Journal of Chemical Physics</i> , 21:1087–1091, 1953.

- [Neu96] L. Neumann. Constant radiance term. Technical Report TR-186-2-96-12, Institute of Computer Graphics, Vienna University of Technology, 1996. www.cg.tuwien.ac.at/.
- [Neu95] L. Neumann. Monte Carlo radiosity. *Computing*, 55:23–42, 1995.
- [NFKP94] L. Neumann, M. Feda, M. Kopp, and W. Purgathofer. A new stochastic radiosity method for highly complex scenes. In *Proc. of the 5th. EG Workshop on Rendering*, 1994.
- [Nie92] H. Niederreiter. *Random number generation and quasi-Monte Carlo methods*. SIAM, Pennsilvania, 1992.
- [NNB97] L. Neumann, A. Neumann, and P. Bekaert. Radiosity with well distributed ray sets. *Computer Graphics Forum (Eurographics*'97), 16(3):261–270, 1997.
- [NNS72] M. E. Newell, R. G. Newell, and T. L. Sancha. A new approach to the shaded picture problem. In Proceedings of the ACM National Conference, pages 443–450, 1972.
- [NPT⁺95] L. Neumann, W. Purgathofer, R. F. Tobler, A. Neumann, P. Elias, M. Feda, and X. Pueyo. The stochastic ray method for radiosity. In *Rendering Techniques* '95, pages 206–218, 1995.
- [ON94] M. Oren and S. Nayar. Generalization of lambert's reflectance model. *Computer Graphics (SIG-GRAPH '94 Proceedings)*, pages 239–246, 1994.
- [Pel97] M. Pellegrini. Monte Carlo approximation of form factors with error bounded a priori. Discrete and Comp. Geometry, 1997.
- [PFTV92] W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. T. Vetterling. Numerical Recipes in C (Second Edition). Cambridge University Press, Cambridge, USA, 1992.
- [Pho75] B. T. Phong. Illumination for computer generated images. *Communications of the ACM*, 18:311–317, 1975.
- [PM95] S. N. Pattanaik and S. P. Mudur. Adjoint equations and random walks for illumination computation. ACM Transactions on Graphics, 14(1):77–102, 1995.
- [Pop87] Gy. Popper. *Bevezetés a végeselem-módszer matematikai elméletébe*. BME Mérnöktovábbképző Intézet, Budapest, 1987.
- [PP98] J. Prikryl and W. Purgathofer. Perceptually based radiosity. In *Eurographics '98, STAR State of the Art Report*, 1998.
- [Ro76] P. Rózsa. Lineáris algebra és alkalmazásai. Műszaki Könyvkiadó, Budapest, 1976.
- [Rén62] A. Rényi. *Wahrscheinlichkeitsrechnung*. VEB Deutscher Verlag der Wissenschaften, Berlin, 1962.
- [RVW98] G. Renner, T. Várady, and V. Weiss. Reverse engineering of free-form features. In PROLAMAT 98, CD proc., Trento, 1998.
- [SAWG91] F. X. Sillion, J. R. Arvo, S. H. Westin, and D. P. Greenberg. A global illumination solution for general reflectance distributions. *Computer Graphics (SIGGRAPH '91 Proceedings)*, 25(4):187–198, 1991.
- [Sbe96] M. Sbert. *The Use of Global Directions to Compute Radiosity*. PhD thesis, Catalan Technical University, Barcelona, 1996.
- [Sbe97] M. Sbert. Error and complexity of random walk Monte-Carlo radiosity. *IEEE Transactions on Visualization and Computer Graphics*, 3(1), 1997.
- [Sbe99] M. Sbert. Optimal absorption probabilities for random walk radiosity. *to be published*, 1999.
- [SC94] F. Sillion and Puech C. *Radiosity and Global Illumination*. Morgan Kaufmann Publishers, Inc., San Francisco, 1994.
- [Sch93] Ch. Schlick. A customizable reflectance model for everyday rendering. In Fourth Eurographics Workshop on Rendering, pages 73–83, France, 1993.
- [Sch96] J. Schanda. CIE colorimetry and colour displays. In IS&T/SID Conf. Scottsdale, 1996.
- [SDS95] F. Sillion, G. Drettakis, and C. Soler. Clustering algorithm for radiance calculation in general environments. In *Rendering Techniques* '95, pages 197–205, 1995.

[Se66]	J. Shrider (editor). <i>The Monte-Carlo Method</i> . Pergamon Press, Oxford, 1966. Also in <i>The method of statistical trials (The Monte Carlo Method)</i> , Fizmatgiz, Moscow, 1965.
[SGCH94]	P. Schröder, S.J. Gortler, M.F. Cohen, and P. Hanrahan. Wavelet projections for radiosity. <i>Computer Graphics Forum</i> , 13(2):141–151, 1994.
[SH81]	R. Siegel and J. R. Howell. <i>Thermal Radiation Heat Transfer</i> . Hemisphere Publishing Corp., Washington, D.C., 1981.
[Shi90]	P. Shirley. A ray-tracing method for illumination calculation in diffuse-specular scenes. In <i>Proc. Graphics Interface</i> , pages 205–212, 1990.
[Shi91a]	P. Shirley. Discrepancy as a quality measure for sampling distributions. In <i>Eurographics '91</i> , pages 183–194. Elsevier Science Publishers, 1991.
[Shi91b]	P. Shirley. Time complexity of Monte-Carlo radiosity. In <i>Eurographics '91</i> , pages 459–466. Elsevier Science Publishers, 1991.
[SMP98]	M. Sbert, R. Martinez, and X. Pueyo. Gathering multi-path: a new Monte-Carlo algorithm for radiosity. In <i>Winter School of Computer Graphics '98</i> , pages 331–338, Plzen, 1998.
[Sob91]	I. Sobol. Die Monte-Carlo Methode. Deutscher Verlag der Wissenschaften, 1991.
[SP89]	F. Sillion and C. Puech. A general two-pass method integrating specular and diffuse reflection. In <i>Computer Graphics (SIGGRAPH '89 Proceedings)</i> , pages 335–344, 1989.
[SP92]	V. Székely and A. Poppe. <i>Számítógépes grafika alapjai IBM PC-n</i> . ComputerBooks, Budapest, 1992.
[SPS98]	M. Stamminger, Slussalek P., and H-P. Seidel. Three point clustering for radiance computations. In <i>Rendering Techniques '98</i> , pages 211–222, 1998.
[ST93]	G. Stoyan and G. Takó. Numerikus mószerek. ELTE, Budapest, 1993.
[SW93]	P. Shirley and C. Wang. Luminaire sampling in distribution ray tracing. <i>SIGGRAPH '93 Course Notes</i> , 1993.
[SWZ96]	P. Shirley, C. Wang, and K. Zimmerman. Monte Carlo techniques for direct lighting calculations. <i>ACM Transactions on Graphics</i> , 15(1):1–36, 1996.
[Tob98]	R. F. Tobler. ART – Advanced Rendering Toolkit, 1998. http://www.cg.tuwien.ac.at/research /rendering/ART.
[Vea97]	E. Veach. <i>Robust Monte Carlo Methods for Light Transport Simulation</i> . PhD thesis, Stanford University, http://graphics.stanford.edu/papers/veach_thesis, 1997.
[VG95]	E. Veach and L. Guibas. Bidirectional estimators for light transport. In <i>Computer Graphics (SIG-GRAPH '95 Proceedings)</i> , pages 419–428, 1995.
[VG97]	E. Veach and L. Guibas. Metropolis light transport. <i>Computer Graphics (SIGGRAPH '97 Proceedings)</i> , pages 65–76, 1997.
[VMC97]	T. Várady, R. R. Martin, and J. Cox. Reverse engineering of geometric models - an introduction. <i>Computer-Aided Design</i> , 29(4):255–269, 1997.
[WA77]	K. Weiler and P. Atherton. Hidden surface removal using polygon area sorting. In <i>Computer Graphics (SIGGRAPH '77 Proceedings)</i> , pages 214–222, 1977.
[War92]	G. Ward. Measuring and modeling anisotropic reflection. Computer Graphics, 26(2):265–272, 1992.
[War95]	T. Warnock. Computational investigations of low-discrepancy point sets. In H. Niederreiter and P. Shiue, editors, <i>Monte-Carlo and Quasi-Monte Carlo Methods in Scientific Computing</i> , pages 354–361. Springer, 1995.
[WCG87]	J. R. Wallace, M. F. Cohen, and D. P. Greenberg. A two-pass solution to the rendering equation: A synthesis of ray tracing and radiosity methods. In <i>Computer Graphics (SIGGRAPH '87 Proceedings)</i> , pages 311–324, 1987.
[WS82]	G. Wyszecki and W. Stiles. <i>Color Science: Concepts and Methods, Quantitative Data and Formulae.</i> Wiley, New York, 1982.
[ZS95]	K. Zimmerman and P. Shirley. A two-pass solution to the rendering equation with a source visibility preprocess. In <i>Rendering Techniques '95</i> , pages 284–295, 1995.

SUBJECT INDEX

1

1-equidistribution 36

5D adaptive tree 52

A

abstract lightsource model 12 acceptance probability 45, 71, 93 adaptive importance sampling 52 adaptive tessellation 134 adjoint basis 29, 84 adjoint operator 11 albedo 12, 53, 82, 86 ambient light 19 animation 134

B

back faces 101 base 51 Bi-directional algorithms 69 bi-directional path-tracing 47, 69 Bi-directional Reflection Distribution Function 12 bi-directional reflection/refraction function 8 bi-directional transfer 103 bi-directional transport matrix 84 Bi-directional walking 89 BRDF 8, 12 BRDF sampling 52 brick-rule 34

С

camera 2, 13 camera parameter 14 caustics 64 central limit theorem 35 CIE XYZ 4 closed environment 26 clustering 24 coherent component 7 color 3 color matching functions 3 combined walking 89 continuous algorithms 101 contraction 16 contribution indicator function 58 cubic tetrahedron 32, 116

D

d-bounce irradiance 86, 118 density 51 detailed balance 46 differential solid angle 5 diffuse first-shot 133 diffuse radiosity 24 dimensional core 35, 47 dimensional explosion 35, 47 direct contribution 9 direct lightsource calculations 55 directional distributions 24 directional-lightsource 12 discontinuity meshing 134 discrepancy 37 Discrete algorithms 101 distributed ray-tracing 47, 63 domain of discontinuity 49

Е

emission 8 emitter group 106 Energy balance 12 equidistribution 36 error measure 15 expansion 20 explicit equation 11 eye 2

F

filtering 15 finite-element method 20, 24, 28 first-shot 113 flux 6 focal distance 14, 114 form factor 31 front faces 101 fundamental law of photometry 7

G

Galerkin's method 30, 102 gamma-correction 13 gathering 61 gathers 22 Gaussian quadrature 34 genetic algorithms 52 geometry matrix 84 geometry of the virtual world 1 global directions 83 global illumination 17 Global illumination solution 17 global importance sampling 52 global method 48 global pass 1, 2 global visibility problem 101 global z-buffer algorithm 106

Gouraud shading 87 in radiosity method 88 Grassmann laws 3

Η

Halton-sequence 41 Hardy-Krause variation 38 hemicube 32, 116 hemisphere algorithm 32 hidden surface problem 16 hierarchical methods 24 hierarchical radiosity 25 human eye 1, 13

I

illumination hemisphere 5 illumination networks 24 illumination sphere 5 image synthesis 1 image-buffer 13 image-space error measure 15 implicit equation 11 importance function 56, 91 importance sampling 40, 47 incoming first-shot 113, 133 Instant radiosity 73 intensity 6 interactive walk-through 134 inversion 20 iteration 20 Iteration techniques 24

K

k-uniform sequences 99kd-tree 72Koksma-Hlawka inequality 38

L

Light 3 light power 1 light-tracing 47, 66 lighting 1, 2 lightsource 12 lightsource sampling 52 links 52 Local illumination methods 17 local importance sampling 52 local method 48 local pass 2 look-up table 13 low-discrepancy 41 low-discrepancy series 40 luminance 56, 91

M

Mach-banding 88 material properties 2 max *d*-bounce irradiance 86 measurement device 2 metamers 4 Metropolis light-transport 72 Metropolis sampling 52 monitor 13 Monte-Carlo 35 Monte-Carlo radiosity 47 multi-path method 48 multiple importance sampling 134 multiresolution methods 24

Ν

norm 16

0

object-space error measure 15 optical material properties 1

Р

painter's algorithm 101 parallel computing 24, 134 participating media 2, 134 partitioned hemisphere 24 path-tracing 47, 64 photon 6 Photon tracing 65 photon-map 52, 72 physically plausible 12 pinhole camera model 13 point collocation 30 point-lightsource 12 potential 9 potential equation 9 potential measuring operator 10, 11 power 6 power equation 32 pupil 13

Q

quasi-Monte Carlo 35 quasi-Monte Carlo quadrature 36

R

radiance 1, 2, 6 radiance measuring operator 9, 11 radiosity method 31 randomization point 75 ray-bundle 83 ray-bundle tracing 132 Ray-casting 61 receiver group 106 Reciprocity 12 Recursive ray-tracing 17 reflected/refracted component 8 rendering 1 rendering equation 2, 8 short form 8 rendering problem 10 RGB 4 Russian-roulette 57, 104

S

scalar product 11

scaling value 9 Self-correcting iteration 77 sensitivity function 9 shadow ray 114 shooting 24, 61, 65 Simpson-rule 34 sky-light illumination 12 solid angle 5 spherical coordinate system 5 splitting 22 standard NTSC phosphors 4 star-discrepancy 37 stochastic iteration 25, 74, 131 stochastic radiosity 78 Stochastic ray-radiosity 79 stratification 40 surface 11 surface dependent camera parameter 15

Т

tentative path 71 tentative sample 45, 93 tentative transition function 45, 71, 93 tessellation 11 theorem of iterated logarithm 41 theorems of large numbers 35 tone mapping 2, 4 total expected value theorem 75 transfer probability density 7 transillumination direction 101 transillumination directions 78 transillumination plane 101 transillumination radiosity method 78 trapezoidal-rule 34 triangle mesh 11 tristimulus 3 two-pass methods 47

U

uniform 36

V

Van der Corput sequence 41 variation 37 variation in the sense of Vitali 37 VEGAS algorithm 91 VEGAS method 52 virtual world 1 visibility calculation 16 visibility function 8 visibility indicator 31 visibility map 101 visibility ray-tracing 17, 62

W

walk 83 wavelet 24 Weiler-Atherthon algorithm 103 well-distributed ray-sets 25 white point 4 ∞ -uniform 100