

```
ics
& (depth < MAXDEPTH)
{
    inside / inside
    nt = nt / nc;
    cos2t = 1.0f - nt;
    D, N );
    )
    at a = nt - nc, b = nt;
    at Tr = 1 - (R0 + (1 - R0) *
    Tr) R = (D * nnt - N * (a
    E * diffuse;
    = true;
    refl + refr)) && (depth < MAXDEPTH)
    D, N );
    refl * E * diffuse;
    = true;
    MAXDEPTH)
    survive = SurvivalProbability( diffuse );
    estimation - doing it properly, clearly
    if;
    radiance = SampleLight( &rand, I, &t, &align,
    e.x + radiance.y + radiance.z) > 0) && (depth <
    v = true;
    at brdfPdf = EvaluateDiffuse( L, N ) * Survive;
    at3 factor = diffuse * INVPI;
    at weight = Mis2( directPdf, brdfPdf );
    at cosThetaOut = dot( N, L );
    E * ((weight * cosThetaOut) / directPdf) * (radiance
    random walk - done properly, closely following
    survive)
    ;
    at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf
    survive;
    pdf;
    r1 = E * brdf * (dot( N, R ) / pdf);
    sion = true;
```

/INFOMOV/

Optimization & Vectorization

J. Bikker - Sep-Nov 2016 - Lecture 9: "Fixed Point Math"

Welcome!



Today's Agenda:

- Introduction
- Float to Fixed Point and Back
- Operations
- Fixed Point & Accuracy
- Demonstration



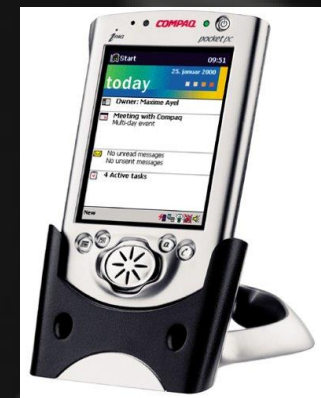
Introduction

The Concept of Fixed Point Math

Basic idea: *emulating floating point math using integers.*

Why?

- Not every CPU has a floating point unit.
- Specifically: cheap DSPs do not support floating point.
- Mixing floating point and integer is Good for the Pipes.
- Some floating point ops have long latencies (div).
- Data conversion can be a significant part of a task.
- Fixed point can be more accurate.



```

...
    & (depth < MAXDEPTH)
...
    inside / ...
    nt = nt / nc;
    os2t = 1.0f - nnt;
    D, N );
    );
...
    at a = nt - nc; b = nt;
    at Tr = 1 - (R0 + (1 - R0)
    Tr) R = (D * nnt - N * (a
...
    E * diffuse;
    = true;
...
    refl + refr)) && (depth < MAXDEPTH)
    D, N );
    refl * E * diffuse;
    = true;
...
MAXDEPTH)
survive = SurvivalProbability( diffuse,
estimation - doing it properly, clear
if;
    radiance = SampleLight( @rand, I, R, R);
    e.x + radiance.y + radiance.z) > 0) && (survive)
...
    w = true;
    at brdfPdf = EvaluateDiffuse( L, N ) * Pdf;
    at3 factor = diffuse * INVPI;
    at weight = Mix2( directPdf, brdfPdf );
    at cosThetaOut = dot( N, L );
    E * ((weight * cosThetaOut) / directPdf) * (radiance
...
random walk - done properly, closely following
survive)
...
    at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf
    survive;
    pdf;
    n = E * brdf * (dot( N, R ) / pdf);
    sion = true;

```



Introduction

The Concept of Fixed Point Math

Basic idea: we have π : 3.1415926536.

- Multiplying that by 10^{10} yields 31415926536.
- Adding 1 to π yields 4.1415926536.
- But, we scale up 1 by 10^{10} as well:
adding $1 \cdot 10^{10}$ to the scaled up version of π yields 41415926536.

➔ In base 10, we get N digits of fractional precision if we multiply our numbers by 10^N (and remember where we put that dot).

```

...
    & (depth < MAXDEPTH)
...
    c = inside / inside;
    nt = nt / nc;
    cos2t = 1.0f - nnt;
    D, N );
    )
...
    at a = nt - nc, b = nt;
    at Tr = 1 - (R0 + (1 - R0));
    Tr) R = (D * nnt - N * (a
...
    E * diffuse;
    = true;
...
    refl + refr) && (depth < MAXDEPTH)
...
    D, N );
    -refl * E * diffuse;
    = true;
...
MAXDEPTH)
...
survive = SurvivalProbability( diffuse );
estimation - doing it properly, clearly
if;
    radiance = SampleLight( @rand, I, R, Alignment
    e.x + radiance.y + radiance.z) > 0) && (survive
...
    v = true;
    at brdfPdf = EvaluateDiffuse( L, N ) * Survive;
    at3 factor = diffuse * INVPI;
    at weight = Mis2( directPdf, brdfPdf );
    at cosThetaOut = dot( N, L );
    E * ((weight * cosThetaOut) / directPdf) * (radiance
...
random walk - done properly, closely following
ive)
...
    at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf
    survive;
    pdf;
    n = E * brdf * (dot( N, R ) / pdf);
    sion = true;

```



Introduction

The Concept of Fixed Point Math

Addition and subtraction are straight-forward with fixed point math.

We can also use it for interpolation:

```
void line( int x1, int y1, int x2, int y2 )
{
    int dx = (x2 - x1) * 10000;
    int dy = (y2 - y1) * 10000;
    int pixels = max( abs( x2 - x1 ), abs( y2 - y1 ) );
    dx /= pixels;
    dy /= pixels;
    int x = x1 * 10000, y = y1 * 10000;
    for( int i = 0; i < pixels; i++, x += dx, y += dy )
        plot( x / 10000, y / 10000 );
}
```



Introduction

The Concept of Fixed Point Math

For multiplication and division things get a bit more complex.

- $\pi \cdot 2 \equiv 31415926536 * 2000000000 = 62831853072000000000$
- $\pi / 2 \equiv 31415926536 / 2000000000 = 1$ (or 2, if we use proper rounding).

Multiplying two fixed point numbers yields a result that is 10^{10} too large (in this case).

Dividing two fixed point numbers yields a result that is 10^{10} too small.

```

ics
& (depth < MAXDEPTH)
{
    // Inside of the sphere
    int nt = nt / nc;
    double cos2t = 1.0f - nt;
    double D, N;
    // ...
    int a = nt - nc, b = nt;
    double Tr = 1 - (R0 + (1 - R0) * cos(a));
    double R = (D * nnt - N * (a < b));
    // ...
    E * diffuse;
    // ...
    refl + refr)) && (depth < MAXDEPTH)
    {
        D, N;
        refl * E * diffuse;
        // ...
    }
}
MAXDEPTH)
survive = SurvivalProbability( diffuse );
// estimation - doing it properly, close to
if;
radiance = SampleLight( &rand, I, &t, &align,
e.x + radiance.y + radiance.z) && (rand <
v = true;
at brdfPdf = EvaluateDiffuse( L, N ) * survive;
at3 factor = diffuse * INVPI;
at weight = Mis2( directPdf, brdfPdf );
at cosThetaOut = dot( N, L );
E * ((weight * cosThetaOut) / directPdf) * (radiance
random walk - done properly, closely following
ive)
;
at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf );
survive;
pdf;
n = E * brdf * (dot( N, R ) / pdf);
sion = true;

```



Introduction

The Concept of Fixed Point Math

On a computer, we obviously do not use base 10, but base 2. Starting with π again:

- Multiplying by 2^{16} yields 205887.
- Adding $1 \cdot 2^{16}$ to the scaled up version of π yields 271423.

In binary:

- $205887 = 00000000\ 00000011\ 00100100\ 00111111$
- $271423 = 00000000\ 00000100\ 00100100\ 00111111$

Looking at the first number (205887), and splitting in two sets of 16 bit, we get:

- 00000000000011 (base 2) = 3 (base 10);
- 10010000111111 (base 2) = 9279 (base 10); $\frac{9279}{2^{16}} = 0.141586304$.



Introduction

The Concept of Fixed Point Math

Interpolation, using base 2:

```
void line( int x1, int y1, int x2, int y2 )
{
    int dx = (x2 - x1) << 16;
    int dy = (y2 - y1) << 16;
    int pixels = max( abs( x2 - x1 ), abs( y2 - y1 ) );
    dx /= pixels;
    dy /= pixels;
    int x = x1 << 16, y = y1 << 16;
    for( int i = 0; i < pixels; i++, x += dx, y += dy )
        plot( x >> 16, y >> 16 );
}
```



Introduction

Practical example

Texture mapping in Quake 1: Perspective Correction

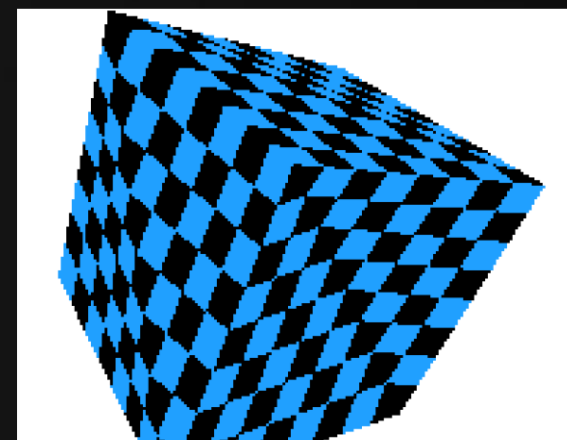
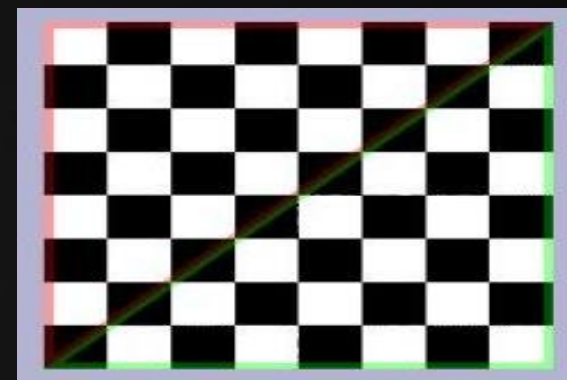
- Affine texture mapping: interpolate u/v linearly over polygon
- Perspective correct texture mapping: interpolate 1/z, u/z and v/z.
- Reconstruct u and v per pixel using the reciprocal of 1/z.

Quake’s solution:

- Divide a horizontal line of pixels in segments of 8 pixels;
- Calculate u and v for the start and end of the segment;
- Interpolate linearly (fixed point!) over the 8 pixels.

And:

Start the floating point division (21 cycles) for the next segment, so it can complete while we execute integer code for the linear interpolation.



```

...
    & (depth < MAXDEPTH)
...
    inside / 1.0;
    nt = nt / nc;
    cos2t = 1.0f - nnt;
    D, N );
    );
...
    at a = nt - nc, b = nt;
    at Tr = 1 - (R0 + (1 - R0)
    Tr) R = (D * nnt - N * (a
...
    E * diffuse;
    = true;
...
    refl + refr)) && (depth < MAXDEPTH)
...
    D, N );
    refl * E * diffuse;
    = true;
...
MAXDEPTH)
...
survive = SurvivalProbability( diffuse,
estimation - doing it properly, close
if;
radiance = SampleLight( @rand, I, R, Align
e.x + radiance.y + radiance.z) > 0) && (survive)
...
v = true;
at brdfPdf = EvaluateDiffuse( L, N ) * Pear
at3 factor = diffuse * INVPI;
at weight = Mis2( directPdf, brdfPdf );
at cosThetaOut = dot( N, L );
E * ((weight * cosThetaOut) / directPdf) *
...
random walk - done properly, closely follow
ive)
...
at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf
survive;
pdf;
n = E * brdf * (dot( N, R ) / pdf);
sion = true;

```



Conversions

Practical Things

Converting a floating point number to fixed point:

Multiply the float by a power of 2 represented by a floating point value, and cast the result to an integer. E.g.:

```
int fp_pi = (int)(3.141593f * 65536.0f); // 16 bits fractional
```

After calculations, cast the result to int by discarding the fractional bits. E.g.:

```
int result = fp_pi >> 16; // divide by 65536
```

Or, get the original float back by casting to float and dividing by 2^{fractionalbits}:

```
float result = (float)fp_pi / 65536.0f;
```

Note that this last option has significant overhead, which should be outweighed by the gains.



Conversions

Practical Things - Considerations

Example: precomputed sin/cos table

```
#define FP_SCALE 65536.0f 1073741824.0f
int sintab[256], costab[256];
for( int i = 0; i < 256; i++ )
    sintab[i] = (int)(FP_SCALE * sinf( (float)i / 128.0f * PI )),
    costab[i] = (int)(FP_SCALE * cosf( (float)i / 128.0f * PI ));
```

What is the best value for FP_SCALE in this case? And should we use int or unsigned int for the table?

Sine/cosine: range is [-1, 1]. In this case, we need 1 sign bit, and 1 bit for the whole part of the number. So:

- We use 30 bits for fractional precision, 1 for sign, 1 for range.
- In base 10, the fractional precision is ~10 digits (float has 7).



Conversions

Practical Things - Considerations

Example: values in a z-buffer

A 3D engine needs to keep track of the depth of pixels on the screen for depth sorting. For this, it uses a z-buffer.

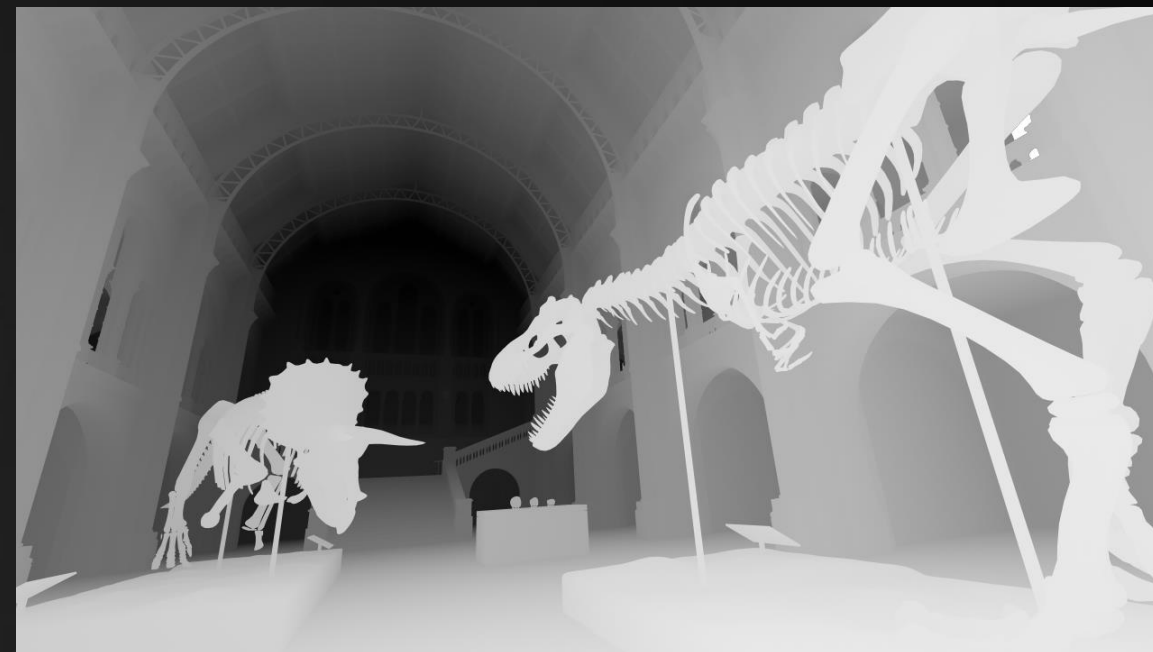
We can make two observations:

1. All values are positive (no objects behind the camera are drawn);
2. Further away we need less precision.

By adding 1 to z, we guarantee that z is in the range [1..infinity].

The reciprocal of z is then in the range [0..1].

We store $1/(z+1)$ as a 0:32 unsigned fixed point number for maximum precision.



```

...
    & (depth < MAXDEPTH)
...
    c = inside / (1.0f - c);
    nt = nt / nc;
    cos2t = 1.0f - nt;
    D, N );
    )
...
    at a = nt - nc; b = nt;
    at Tr = 1 - (R0 + (1 - R0) *
    Tr) R = (D * nnt - N * (a0
...
    E * diffuse;
    = true;
...
    refl + refr)) && (depth < MAXDEPTH)
...
    D, N );
    refl * E * diffuse;
    = true;
...
MAXDEPTH)
...
survive = SurvivalProbability( diffuse,
estimation - doing it properly, close
if;
radiance = SampleLight( @rand, I, R, Align
e.x + radiance.y + radiance.z) > 0) && (surv
...
v = true;
at brdfPdf = EvaluateDiffuse( L, N ) * Pear
at3 factor = diffuse * INVPI;
at weight = Mis2( directPdf, brdfPdf );
at cosThetaOut = dot( N, L );
E * ((weight * cosThetaOut) / directPdf) * (radi
...
random walk - done properly, closely follow
ive)
...
at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf
survive;
pdf;
n = E * brdf * (dot( N, R ) / pdf);
sion = true;

```



Conversions

Practical Things - Considerations

We pick the right precision based on the problem at hand.

Sin/cos: original values [-1..1];

→ sign bit + 31 fractional bits;

→ 0:31 signed fixed point.

Storing 1/(z+1): original values [0..1];

→ 32 fractional bits;

→ 0:32 unsigned fixed point.

Particles: original values [-50..50];

→ sign bit + 6 integer bits, 32-7=25 fractional bits;

→ 6:25 signed fixed point.

In general:

- first determine if we need a sign;
- then, determine how many bits are need to represent the integer range;
- use the remainder as fractional bits.



Today's Agenda:

- Introduction
- Float to Fixed Point and Back
- Operations
- Fixed Point & Accuracy
- Demonstration



Operations

Basic Operations on Fixed Point Numbers

Operations on mixed fixed point formats:

- $A+B$ ($I_A:F_A + I_B:F_B$)

To be able to add the numbers, they need to be in the same format.

Example: $I_A:F_A=4:28$, $I_B:F_B=16:16$

Option 1: $A \gg= 12$ (to make it 16:16)

Option 2: $B \ll= 12$ (to make it 4:28)

Problem with option 2: we do not get 4:28, we get 16:28!

Problem with option 1: we drop 12 bits from A.



Operations

Basic Operations on Fixed Point Numbers

Operations on mixed fixed point formats:

- $A * B \ (I_A : F_A * I_B : F_B)$

We can freely mix fixed point formats for multiplication.

Example: $I_A : F_A = 18 : 14$, $I_B : F_B = 14 : 18$

Result: 32:32, shift to the right by 18 to get a ..:14 number, or by 14 to get a ..:18 number.

Problem: the intermediate result doesn't fit in a 32-bit register.

```

...
    & (depth < MAXDEPTH)
...
    inside / inside;
    nt = nt / nc;
    cos2t = 1.0f - nnt;
    D, N );
    )
...
    at a = nt - nc, b = nt;
    at Tr = 1 - (R0 + (1 - R0));
    Tr) R = (D * nnt - N * (a
...
    E * diffuse;
    = true;
...
    refl + refr)) && (depth < MAXDEPTH)
...
    D, N );
    refl * E * diffuse;
    = true;
...
MAXDEPTH)
...
    survive = SurvivalProbability( diffuse,
    estimation - doing it properly, clean
    if;
    radiance = SampleLight( &rand, I, &t, &align;
    e.x + radiance.y + radiance.z) > 0) && (survive
...
    v = true;
    at brdfPdf = EvaluateDiffuse( L, N ) * Psurvive;
    at3 factor = diffuse * INVPI;
    at weight = Mis2( directPdf, brdfPdf );
    at cosThetaOut = dot( N, L );
    E * ((weight * cosThetaOut) / directPdf) * (radiance
...
    random walk - done properly, closely following
    survive)
...
    at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf;
    survive;
    pdf;
    n = E * brdf * (dot( N, R ) / pdf);
    sion = true;

```



Operations

Multiplication

- *“Ensure that intermediate results never exceed 32 bits.”*

Suppose we want to multiply two 20:12 unsigned fixed point numbers:

1. `(fp_a * fp_b) >> 12;` // good if fp_a and fp_b are very small
2. `(fp_a >> 12) * fp_b;` // good if fp_a is a whole number
3. `(fp_a >> 6) * (fp_b >> 6);` // good if fp_a and fp_b are large
4. `((fp_a >> 3) * (fp_b >> 3)) >> 6;`

Which option we chose depends on the parameters:

```
fp_a = PI;
fp_b = 0.5f * 2^12;
int fp_prod = fp_a >> 1; // ☺
```



Operations

Division

- *“Ensure that intermediate results never exceed 32 bits.”*

Dividing two 20:12 fixed point numbers:

1. `(fp_a << 12) / fp_b;` // good if fp_a and fp_b are very small
2. `fp_a / (fp_b >> 12);` // good if fp_b is a whole number
3. `(fp_a << 6) / (fp_b >> 6);` // good if fp_a and fp_b are large
4. `((fp_a << 3) / (fp_b >> 3)) << 6;`

Note that a division by a constant can be replaced by a multiplication by its reciprocal:

- `fp_reciprocal = (1 << 12) / fp_b;`
- `fp_prod = (fp_a * fp_reciprocal) >> 12;` // or one of the alternatives



Operations

Multiplication, Take 2

- *“Use a 64-bit intermediate result.”*

$$A * B \quad (I_A : F_A * I_B : F_B)$$

Example: $I_A : F_A = 16:16, I_B : F_B = 16:16$

Result: 32:32

*Calculate a 64-bit result (with enough room for 32:32),
throw out 32 bits afterwards.*

x86 MUL instruction:

MUL EDX

Functionality:

multiplies EDX by EAX, stores the result in EDX:EAX.

➔ Tossing 32 bits: ignore EAX.

➔ x86 is designed for 16:16.



Operations

Multiplication

Special case: multiply by a 32:0 number.

```
int fp_pi = (int)(3.141593f * 65536.0f); // 16 bits fractional
int fp_2pi = fp_pi * 2; // 16 bits fractional
```

We did this in the line function:

```
dx /= pixels; // dx is 16:16, pixels is 32:0
dy /= pixels;
```



Operations

Square Root

For square roots of fixed point numbers, optimal performance is achieved via `_mm_rsqrt_ps` (via float). If precision is of little concern, use a lookup table, optionally combined with interpolation and / or a Newton-Raphson iteration.

Sine / Cosine / Log / Pow / etc.

Almost always a LUT is the best option.

```

...
    & (depth < MAXDEPTH)
...
    inside / 1.0f;
    nt = nt / nc;
    cos2t = 1.0f - nnt;
    D, N );
}

at a = nt - nc, b = nt;
at Tr = 1 - (R0 + (1 - R0) *
Tr) R = (D * nnt - N * (a
E * diffuse;
= true;

efl + refr)) && (depth < MAXDEPTH)
D, N );
efl * E * diffuse;
= true;

MAXDEPTH)
survive = SurvivalProbability( diffuse,
estimation - doing it properly, close
if;
radiance = SampleLight( &rand, I, &t, &light
e.x + radiance.y + radiance.z) > 0) && (survive
v = true;
at brdfPdf = EvaluateDiffuse( L, N ) * Survive;
at3 factor = diffuse * INVPI;
at weight = Mis2( directPdf, brdfPdf );
at cosThetaOut = dot( N, L );
E * ((weight * cosThetaOut) / directPdf) * (radiant
random walk - done properly, closely following
ive)

at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf
urvive;
pdf;
n = E * brdf * (dot( N, R ) / pdf);
sion = true;

```



Operations

Fixed Point & SIMD

For a world of hurt, combine SIMD and fixed point:

```

_mm_mul_epu32
_mm_mullo_epi16
_mm_mulhi_epi16
_mm_srl_epi32
_mm_srai_epi32

```

See MSDN for more details.

```

...
    & (depth < MAXDEPTH)
...
    inside / inside;
    nt = nt / nc;
    cos2t = 1.0f - nt * nt;
    D, N );
}
...
    at a = nt - nc, b = nt;
    at Tr = 1 - (R0 + (1 - R0) *
    Tr) R = (D * nnt - N * (a
...
    E * diffuse;
    = true;
...
    refl + refr)) && (depth < MAXDEPTH)
...
    D, N );
    refl * E * diffuse;
    = true;
...
MAXDEPTH)
...
    survive = SurvivalProbability( diffuse,
    estimation - doing it properly, close
    if;
    radiance = SampleLight( &rand, I, &t, &light
    e.x + radiance.y + radiance.z) > 0) && (survive
...
    v = true;
    at brdfPdf = EvaluateDiffuse( L, N ) * survive;
    at3 factor = diffuse * INVPI;
    at weight = Mis2( directPdf, brdfPdf );
    at cosThetaOut = dot( N, L );
    E * ((weight * cosThetaOut) / directPdf) * (radiant
...
    random walk - done properly, closely following
    vive)
...
    at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf
    survive;
    pdf;
    n = E * brdf * (dot( N, R ) / pdf);
    sion = true;

```



Accuracy

Range versus Precision

Looking at the line code once more:

```
void line( int x1, int y1, int x2, int y2 )
{
    int dx = (x2 - x1) << 16;
    int dy = (y2 - y1) << 16;
    int pixels = max( abs( x2 - x1 ), abs( y2 - y1 ) );
    dx /= pixels;
    dy /= pixels;
    int x = x1 << 16, y = y1 << 16;
    for( int i = 0; i < pixels; i++, x += dx, y += dy )
        plot( x >> 16, y >> 16 );
}
```

dx=15:16, range is 32767.

precision: 16 bits,
 maximum error: $\frac{1}{2^{16}} * 0.5 = \frac{1}{2^{17}}$.
 Interpolating a 1024 pixel line,
 the maximum cumulative error
 is $2^{10} * \frac{1}{2^{17}} = \frac{1}{2^7} \approx 0.008$.



Accuracy

Error

In base 10, error is clear:

PI = 3.14 means: $3.145 > PI > 3.135$

The maximum error is thus 0.005.

In base 2, we apply the same principle:

16:16 fixed point numbers have a maximum error of $\frac{1}{2^{17}} \approx 7.6 \cdot 10^{-6}$.

➔ We get slightly more than 5 digits of decimal precision.

A 32-bit floating point number represents ~7 digits of decimal precision.

```

...
    & (depth < MAXDEPTH)
...
    inside / inside;
    nt = nt / nc;
    cos2t = 1.0f - nnt;
    D, N );
    )
...
    at a = nt - nc, b = nt;
    at Tr = 1 - (R0 + (1 - R0)
    Tr) R = (D * nnt - N * (a
...
    E * diffuse;
    = true;
...
    refl + refr)) && (depth < MAXDEPTH)
...
    D, N );
    refl * E * diffuse;
    = true;
...
MAXDEPTH)
...
survive = SurvivalProbability( diffuse,
estimation - doing it properly, close
if;
radiance = SampleLight( &rand, I, &t, &lig
e.x + radiance.y + radiance.z) > 0) && (sur
...
    v = true;
    at brdfPdf = EvaluateDiffuse( L, N ) * P
    at3 factor = diffuse * INVPI;
    at weight = Mis2( directPdf, brdfPdf );
    at cosThetaOut = dot( N, L );
    E * ((weight * cosThetaOut) / directPdf) * (radiant
...
random walk - done properly, closely following
ive)
...
    at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf
    survive;
    pdf;
    n = E * brdf * (dot( N, R ) / pdf);
    sion = true;

```



Accuracy

Error

During some operations, precision may suffer greatly:

$$x = y/z$$

$$fp_x = (fp_y \ll 8) / (fp_z \gg 8)$$

Assuming 16:16 input, fp_z briefly becomes 16:8, with a precision of only 2 decimal digits.

Similarly:

$$fp_x = (fp_y \gg 8) * (fp_z \gg 8)$$

Here, both fp_y and fp_z become 16:8, and the cumulative error may exceed $1/2^9$.

```

...
    & (depth < MAXDEPTH)
...
    inside / 1.0;
    nt = nt / nc;
    cos2t = 1.0f - nnt;
    D, N );
}

at a = nt - nc, b = nt;
at Tr = 1 - (R0 + (1 - R0) * nnt);
Tr) R = (D * nnt - N * (1 - nnt));

E * diffuse;
= true;

...
    refl + refr)) && (depth < MAXDEPTH)
    D, N );
    refl * E * diffuse;
    = true;

...
MAXDEPTH)
survive = SurvivalProbability( diffuse, r1, r2, &R, &pdf );
estimation - doing it properly, closely following the
if;
radiance = SampleLight( &rand, I, &t, &align, &N, &D );
e.x + radiance.y + radiance.z) > 0) && (depth < MAXDEPTH)

v = true;
at brdfPdf = EvaluateDiffuse( L, N );
at3 factor = diffuse * INVPI;
at weight = Mis2( directPdf, brdfPdf );
at cosThetaOut = dot( N, L );
E * ((weight * cosThetaOut) / directPdf) * (radiance.x + radiance.y + radiance.z);

...
random walk - done properly, closely following the
ive)

...
at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf );
survive;
pdf;
n = E * brdf * (dot( N, R ) / pdf);
sion = true;

```



Accuracy

Error

Careful balancing of range and precision in fixed point calculations can reduce this problem.

Note that accuracy problems also occur in float calculations; they are just exposed more clearly in fixed point. And: this time we can do something about it.

```

...
    & (depth < MAXDEPTH)
...
    t = inside / inside;
    nt = nt / nc;
    cos2t = 1.0f - nnt;
    D, N );
...
    at a = nt - nc, b = nt;
    at Tr = 1 - (R0 + (1 - R0)
    Tr) R = (D * nnt - N * (a
...
    E * diffuse;
    = true;
...
    refl + refr)) && (depth < MAXDEPTH)
...
    D, N );
    refl * E * diffuse;
    = true;
...
MAXDEPTH)
...
survive = SurvivalProbability( diffuse );
estimation - doing it properly, closely following
if;
radiance = SampleLight( &rand, I, &t, &align,
e.x + radiance.y + radiance.z) && (radiance
v = true;
at brdfPdf = EvaluateDiffuse( L, N ) * survive;
at3 factor = diffuse * INVPI;
at weight = Mis2( directPdf, brdfPdf );
at cosThetaOut = dot( N, L );
E * ((weight * cosThetaOut) / directPdf) * (radiance
...
random walk - done properly, closely following
survive)
...
at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf
survive;
pdf;
n = E * brdf * (dot( N, R ) / pdf);
sion = true;

```



Today's Agenda:

- Introduction
- Float to Fixed Point and Back
- Operations
- Fixed Point & Accuracy
- Demonstration

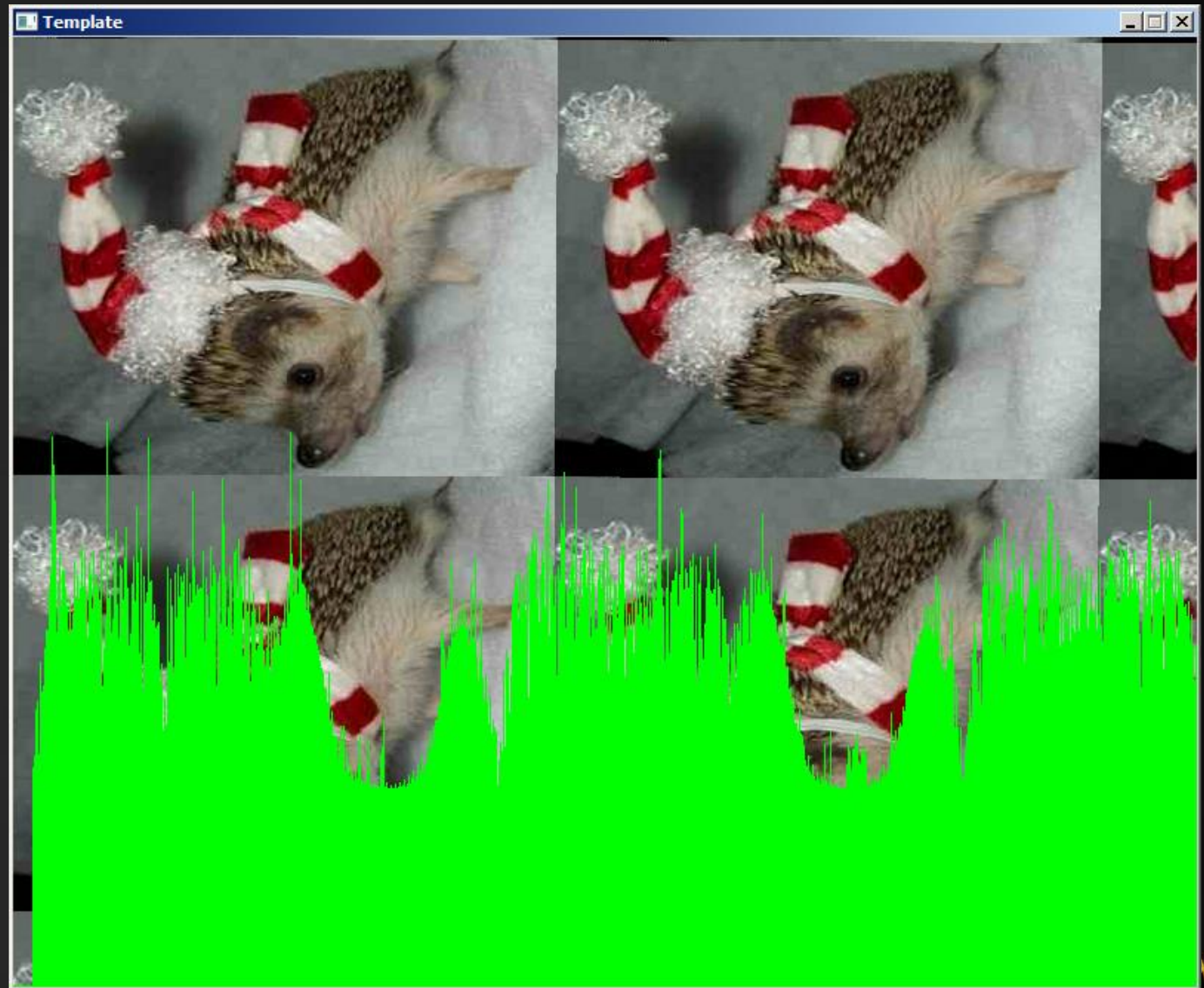


Demonstration

```

...
    & (depth < MAXDEPTH)
...
    c = inside / (1.0 - r);
    nt = nt / nc;
    cos2t = 1.0f - nt;
    D, N );
    )
...
    at a = nt - nc, b = nt;
    at Tr = 1 - (R0 + (1 - R0) * c);
    Tr) R = (D * nnt - N * (a
...
    E * diffuse;
    = true;
...
    refl + refr)) && (depth < MAXDEPTH)
...
    D, N );
    refl * E * diffuse;
    = true;
...
MAXDEPTH)
...
survive = SurvivalProbability( diffuse );
estimation - doing it properly, closely following
if;
radiance = SampleLight( &rand, I, &t, &align, &
e.x + radiance.y + radiance.z) && (rand.N <
...
w = true;
at brdfPdf = EvaluateDiffuse( L, N ) * Survive;
at3 factor = diffuse * INVPI;
at weight = Mis2( directPdf, brdfPdf );
at cosThetaOut = dot( N, L );
E * ((weight * cosThetaOut) / directPdf) * (radiance
...
random walk - done properly, closely following
ive)
...
at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf
survive;
pdf;
n = E * brdf * (dot( N, R ) / pdf);
sion = true;

```

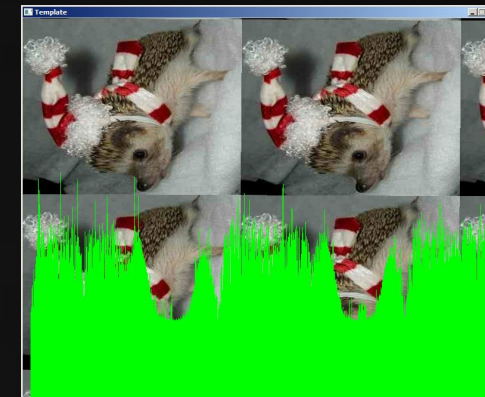


Demonstration

```

Pixel GetBilerpSample( float x, float y )
{
    float fx = x - floor( x ), fy = y - floor( y );
    int ix = (int)x & 2047, iy = (int)y & 2047;
    float w1 = (1 - fx) * (1 - fy);
    float w2 = fx * (1 - fy);
    float w3 = (1 - fx) * fy;
    float w4 = fx * fy;
    unsigned char* base = imageTest.GetBuffer();
    Pixel* pal = imageTest.GetPalette( 63 );
    int offset = ix + iy * 2048;
    Pixel p1 = ScaleColor( pal[base[offset]], (int)(w1 * 255.9f) );
    Pixel p2 = ScaleColor( pal[base[(offset + 1) & 4194303]], (int)(w2 * 255.9f) );
    Pixel p3 = ScaleColor( pal[base[(offset + 1) & 4194303]], (int)(w3 * 255.9f) );
    Pixel p4 = ScaleColor( pal[base[(offset + 1) & 4194303]], (int)(w4 * 255.9f) );
    return p1 + p2 + p3 + p4;
}

```

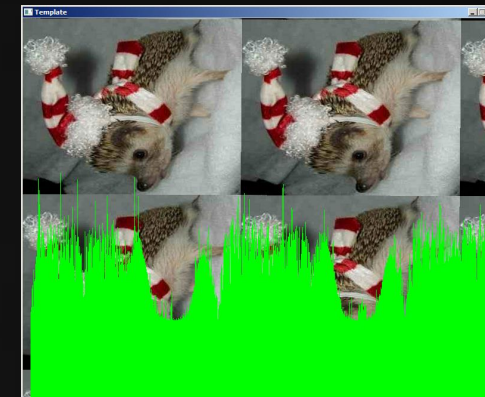


Demonstration

```

Pixel GetBilerpSample( float x, float y )
{
    int fp_x = (int)(x * 16);
    int fp_y = (int)(y * 16);
    int fp_fx = fp_x & 15;
    int fp_fy = fp_y & 15;
    int ix = (fp_x >> 4) & 2047, iy = (fp_y >> 4) & 2047;
    int w1 = (15 - fp_fx) * (15 - fp_fy);
    int w2 = fp_fx * (15 - fp_fy);
    int w3 = (15 - fp_fx) * fp_fy;
    int w4 = 255 - (w1 + w2 + w3);
    unsigned char* base = imageTest.GetBuffer();
    Pixel* pal = imageTest.GetPalette( 63 );
    int offset = ix + iy * 2048;
    Pixel p1 = ScaleColor( pal[base[offset]], w1 );
    Pixel p2 = ScaleColor( pal[base[((offset + 1) & 4194303)]], w2 );
    Pixel p3 = ScaleColor( pal[base[((offset + 1) & 4194303)]], w3 );
    Pixel p4 = ScaleColor( pal[base[((offset + 1) & 4194303)]], w4 );
    return p1 + p2 + p3 + p4;
}

```

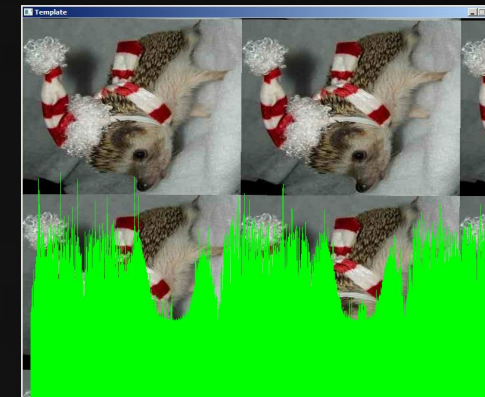


Demonstration

```

float2 dx = (r[1] - r[0]) * ((sinf( a * 2.0f ) + 1.1f) / SCRWIDTH);
float2 dy = (r[2] - r[0]) * ((sinf( a * 2.0f ) + 1.1f) / SCRHEIGHT);
int fp_dxx = (int)(dx.x * 65536);
int fp_dxy = (int)(dx.y * 65536);
int fp_dyx = (int)(dy.x * 65536);
int fp_dyy = (int)(dy.y * 65536);
for( int y = 0; y < SCRHEIGHT; y++ )
{
    int fp_x1 = fp_dyx * y, fp_y1 = fp_dyy * y;
    for( int x = 0; x < SCRWIDTH; x++, fp_x1 += fp_dxx, fp_y1 += fp_dxy )
    {
        int fp_x = ((fp_x1 + 100 * 16) * 2048) >> 12;
        int fp_y = ((fp_y1 + 100 * 16) * 2048) >> 12;
        *dst++ = GetBilerpSample( fp_x, fp_y );
    }
}

```

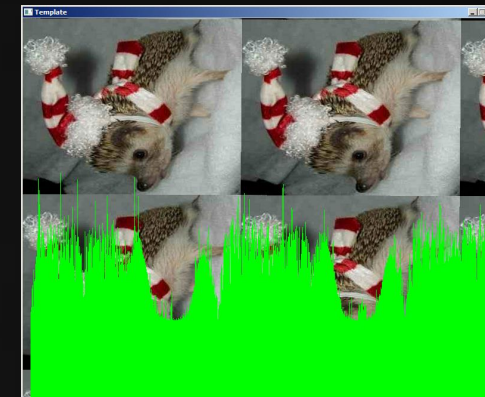


Demonstration

```

float2 dx = (r[1] - r[0]) * ((sinf( a * 2.0f ) + 1.1f) / SCRWIDTH);
float2 dy = (r[2] - r[0]) * ((sinf( a * 2.0f ) + 1.1f) / SCRHEIGHT);
int fp_dxx = (int)(dx.x * 65536 * 16384);
int fp_dxy = (int)(dx.y * 65536 * 16384);
int fp_dyx = (int)(dy.x * 65536 * 16384);
int fp_dyy = (int)(dy.y * 65536 * 16384);
for( int y = 0; y < SCRHEIGHT; y++ )
{
    int fp_x1 = fp_dyx * y, fp_y1 = fp_dyy * y;
    for( int x = 0; x < SCRWIDTH; x++, fp_x1 += fp_dxx, fp_y1 += fp_dxy )
    {
        int fp_x = (fp_x1 + 100 * 16 * 256) >> 15;
        int fp_y = (fp_y1 + 100 * 16 * 256) >> 15;
        *dst++ = GetBilerpSample( fp_x, fp_y );
    }
}

```

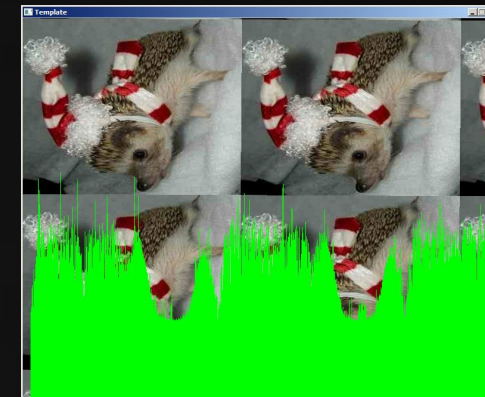


Demonstration

```

Pixel GetBilerpSample( int fp_x, int fp_y )
{
    int fp_fx = fp_x & 15;
    int fp_fy = fp_y & 15;
    int ix = (fp_x >> 4) & 2047, iy = (fp_y >> 4) & 2047;
    int w1 = (15 - fp_fx) * (15 - fp_fy);
    int w2 = fp_fx * (15 - fp_fy);
    int w3 = (15 - fp_fx) * fp_fy;
    int w4 = 255 - (w1 + w2 + w3);
    unsigned char* base = imageTest.GetBuffer();
    Pixel* pal1 = imageTest.GetPalette( w1 >> 2 );
    Pixel* pal2 = imageTest.GetPalette( w2 >> 2 );
    Pixel* pal3 = imageTest.GetPalette( w3 >> 2 );
    Pixel* pal4 = imageTest.GetPalette( w4 >> 2 );
    int offset = ix + iy * 2048;
    Pixel p1 = pal1[base[offset]];
    Pixel p2 = pal2[base[((offset + 1) & 4194303)]];
    Pixel p3 = pal3[base[((offset + 1) & 4194303)]];
    Pixel p4 = pal4[base[((offset + 1) & 4194303)]];
    return p1 + p2 + p3 + p4;
}

```



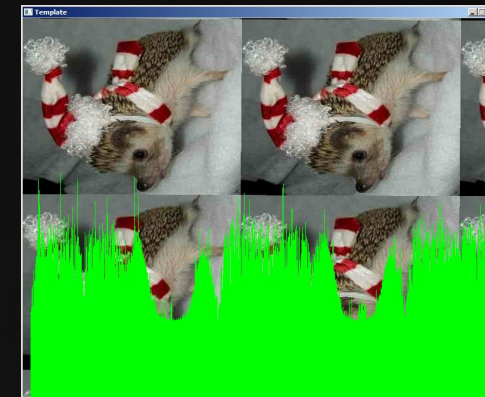
Demonstration

```

int weight[1024];

for( int i = 0; i < 256; i++ )
{
    int fp_fx = i & 15;
    int fp_fy = i >> 4;
    weight[i * 4 + 0] = ((15 - fp_fx) * (15 - fp_fy)) >> 2;
    weight[i * 4 + 1] = (fp_fx * (15 - fp_fy)) >> 2;
    weight[i * 4 + 2] = ((15 - fp_fx) * fp_fy) >> 2;
    weight[i * 4 + 3] = 63 - (weight[i * 4 + 0] + weight[i * 4 + 1] + weight[i * 4 + 2]);
}

```

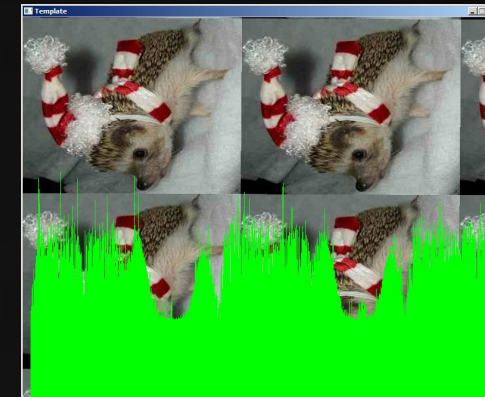


Demonstration

```

Pixel GetBilerpSample( int fp_x, int fp_y )
{
    int fp_fx = fp_x & 15;
    int fp_fy = fp_y & 15;
    int idx = (fp_fy * 16 + fp_fx) * 4;
    int ix = (fp_x >> 4) & 2047, iy = (fp_y >> 4) & 2047;
    unsigned char* base = imageTest.GetBuffer();
    Pixel* pal1 = imageTest.GetPalette( weight[idx + 0] );
    Pixel* pal2 = imageTest.GetPalette( weight[idx + 1] );
    Pixel* pal3 = imageTest.GetPalette( weight[idx + 2] );
    Pixel* pal4 = imageTest.GetPalette( weight[idx + 3] );
    int offset = ix + iy * 2048;
    Pixel p1 = pal1[base[offset]];
    Pixel p2 = pal2[base[(offset + 1) & 4194303]];
    Pixel p3 = pal3[base[(offset + 1) & 4194303]];
    Pixel p4 = pal4[base[(offset + 1) & 4194303]];
    return p1 + p2 + p3 + p4;
}

```



Demonstration

```

...
    & (depth < MAXDEPTH)
...
    c = inside / (1 + n);
    nt = nt / nc;
    cos2t = 1.0f - nt;
    D, N );
    )
...
    at a = nt - nc, b = nt;
    at Tr = 1 - (R0 + (1 - R0) * c);
    Tr) R = (D * nnt - N * (a *
...
    E * diffuse;
    = true;
...
    refl + refr) && (depth < MAXDEPTH)
    D, N );
    -refl * E * diffuse;
    = true;
...
MAXDEPTH)
...
survive = SurvivalProbability( diffuse );
estimation - doing it properly, closely following
if;
radiance = SampleLight( @rand, I, R, Align,
e.x + radiance.y + radiance.z) > 0) && (survive)
...
w = true;
at brdfPdf = EvaluateDiffuse( L, N ) * Survive;
at3 factor = diffuse * INVPI;
at weight = Mis2( directPdf, brdfPdf );
at cosThetaOut = dot( N, L );
E * ((weight * cosThetaOut) / directPdf) * (radiance
...
random walk - done properly, closely following
ive)
...
at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf );
survive;
pdf;
n = E * brdf * (dot( N, R ) / pdf);
sion = true;

```



“And that is how you rotate a Hedgehog!”



Today's Agenda:

- Introduction
- Float to Fixed Point and Back
- Operations
- Fixed Point & Accuracy
- Demonstration



/INFOMOV/

END of “Fixed Point Math”

next lecture: “GPGPU (1)”

```
ics
& (depth < MAXDEPTH)
{
    inside / 1.0f;
    nt = nt / nc;
    cos2t = 1.0f - nnt * nnt;
    D, N );
}

at a = nt - nc; b = nt + nc;
at Tr = 1 - (R0 + (1 - R0) * a);
Tr) R = (D * nnt - N * (a *

E * diffuse;
= true;

efl + refr)) && (depth < MAXDEPTH)
D, N );
efl * E * diffuse;
= true;

MAXDEPTH)

survive = SurvivalProbability( diffuse );
estimation - doing it properly, clearly
if;
radiance = SampleLight( &rand, I, &t, &light;
e.x + radiance.y + radiance.z) > 0) && (survive)
v = true;
at brdfPdf = EvaluateDiffuse( L, N ) * Survive;
at3 factor = diffuse * INVPI;
at weight = Mis2( directPdf, brdfPdf );
at cosThetaOut = dot( N, L );
E * ((weight * cosThetaOut) / directPdf) * (radiance

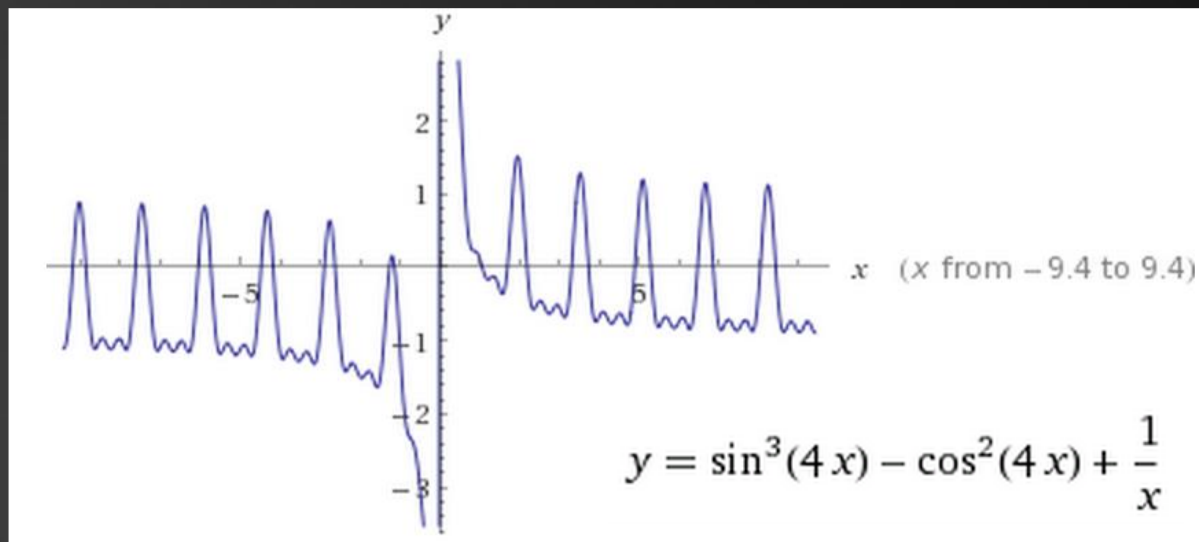
random walk - done properly, closely following
ive)

at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf;
urvive;
pdf;
n = E * brdf * (dot( N, R ) / pdf);
sion = true;
```



Accuracy

Error - Example



Backup Slides 2015



Accuracy

Improving the function.zip example

The following slides contain a step-by-step improvement of the fixed point evaluation of the function $f(x) = \sin(4x)^3 - \cos(4x)^2 + \frac{1}{x}$, which failed during the real-time session in class.

Starting point is the working, but inaccurate version available from the website.

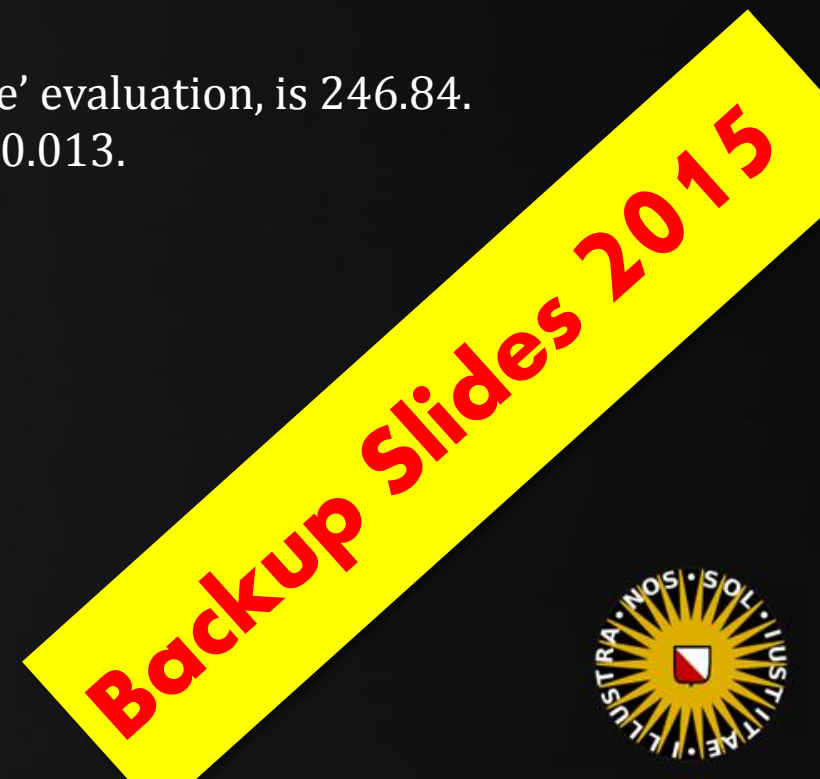
Initial accuracy, expressed as summed error relative to the ‘double’ evaluation, is 246.84.

For comparison, the summed error of the ‘float’ evaluation is just 0.013.

```

...
    & (depth < MAXDEPTH)
...
    = inside / (1 + ...);
    nt = nt / nc;
    cos2t = 1.0f - nnt;
    D, N );
...
    at a = nt - nc, b = nt;
    at Tr = 1 - (R0 + (1 - R0) * ...);
    Tr) R = (D * nnt - N * ...);
...
    E * diffuse;
    = true;
...
    refl + refr)) && (depth < MAXDEPTH)
    D, N );
    refl * E * diffuse;
    = true;
...
MAXDEPTH)
survive = SurvivalProbability( diffuse, ...);
estimation - doing it properly, clearly
if;
radiance = SampleLight( @rand, I, R, ...);
e.x + radiance.y + radiance.z > 0) && (survive)
...
    v = true;
    at brdfPdf = EvaluateDiffuse( L, N ) * ...;
    at3 factor = diffuse * INVPI;
    at weight = Mis2( directPdf, brdfPdf );
    at cosThetaOut = dot( N, L );
    E * ((weight * cosThetaOut) / directPdf) * (radiance
...
random walk - done properly, closely following ...
ive)
...
    at3 brdf = SampleDiffuse( diffuse, N, r1, r2, &R, &pdf );
    survive;
    pdf;
    n = E * brdf * (dot( N, R ) / pdf);
    sion = true;

```



Accuracy

Improving the function.zip example

```

int EvaluateFixed( double x )
{
    16:16 int fp_pi = (int)(PI * 65536.0);
    16:16 int fp_x = (int)(x * 65536.0);
    if ((fp_x >> 8) == 0) return 0; // safety net for division

    int fp_4x = fp_x * 4;
    16:16 * 3:0 = 19:16
    int a = (fp_4x << 8) / ((2 * fp_pi) >> 8); // map radians to 0..4095
    16:16 int fp_sin4x = sintab[(a >> 4) & 4095];
    16:16 int fp_sin4x3 = (((fp_sin4x >> 8) * (fp_sin4x >> 8)) >> 8) * (fp_sin4x >> 8);

    16:16 int fp_cos4x = costab[(a >> 4) & 4095];
    16:16 int fp_cos4x2 = (fp_cos4x >> 8) * (fp_cos4x >> 8);

    16:16 int fp_recix = (65536 << 8) / (fp_x >> 8);

    return fp_sin4x3 - fp_cos4x2 + fp_recix;
}
    
```

In the original code, almost everything is 16:16. This allows for a range of 0..32767 (+/-), which is a waste for most values here.

Backup Slides 20



Accuracy

Improving the function.zip example

Notice how many values do not use the full integer range: e.g, PI is 3 and needs two bits; x is -9..+9 and needs four bits, sin/cos is -1..1 and needs only one bit for range.

```

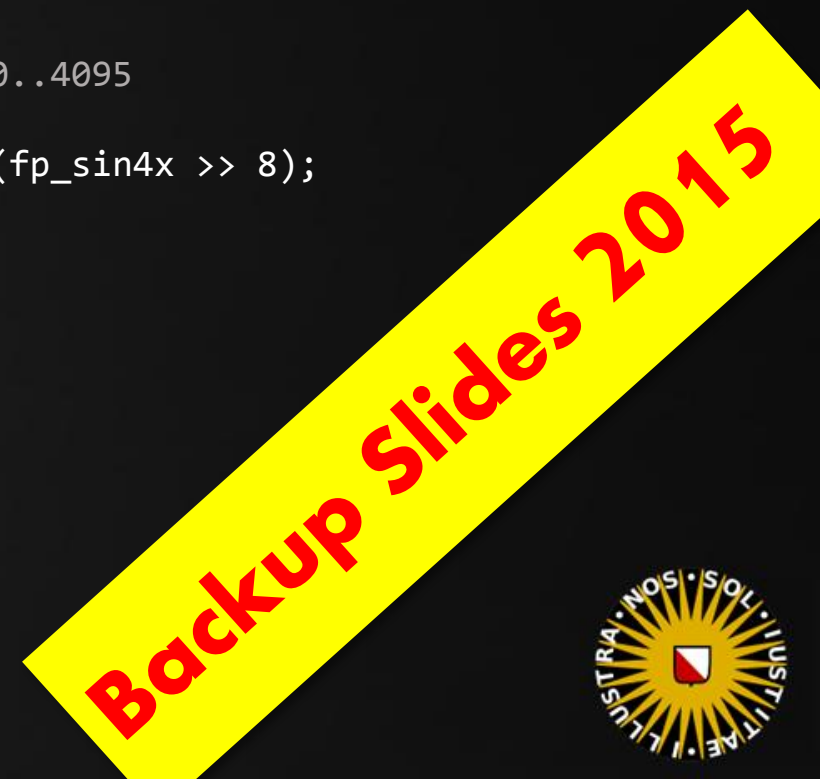
int EvaluateFixed( double x )
{
    2:16 int fp_pi = (int)(PI * 65536.0);
    4:16 int fp_x = (int)(x * 65536.0);
    if ((fp_x >> 8) == 0) return 0; // safety net for division

    int fp_4x = fp_x * 4;           16:16 * 3:0 = 19:16
    int a = (fp_4x << 8) / ((2 * fp_pi) >> 8); // map radians to 0..4095
    1:16 int fp_sin4x = sintab[(a >> 4) & 4095];
    1:16 int fp_sin4x3 = (((fp_sin4x >> 8) * (fp_sin4x >> 8)) >> 8) * (fp_sin4x >> 8);

    1:16 int fp_cos4x = costab[(a >> 4) & 4095];
    1:16 int fp_cos4x2 = (fp_cos4x >> 8) * (fp_cos4x >> 8);

    16:16 int fp_recix = (65536 << 8) / (fp_x >> 8);

    return fp_sin4x3 - fp_cos4x2 + fp_recix;  16:16
}
    
```



Accuracy

Improving the function.zip example

```
int EvaluateFixed( double x )
```

```
{
```

2:16

```
int fp_pi = (int)(PI * 65536.0);
```

4:27

```
int fp_x = (int)(x * (double)(1 << 27));
```

```
if ((fp_x >> 10) == 0) return 0; // safety net for division
```

Error is now down to 14.94.

6:25

```
int fp_4x = fp_x;
```

```
int a = fp_4x / ((2 * fp_pi) >> 3);
```

6:25 / 3:13 = 4:12

1:16

```
int fp_sin4x = sintab[a & 4095];
```

0:30

```
int fp_sin4x3 = (((fp_sin4x >> 1) * (fp_sin4x >> 1)) >> 15) * (fp_sin4x >> 1);
```

$^ 0:15 * 0:15 = 0:30; 0.15 * 0.15 = 0.30$

1:16

```
int fp_cos4x = costab[a & 4095];
```

0:39

```
int fp_cos4x2 = (fp_cos4x >> 1) * (fp_cos4x >> 1);
```

$0:15 * 0:15 = 0:30$

16:16

```
int fp_recix = (1 << 30) / (fp_x >> 13);
```

1:30 / 5:14 = 0:16

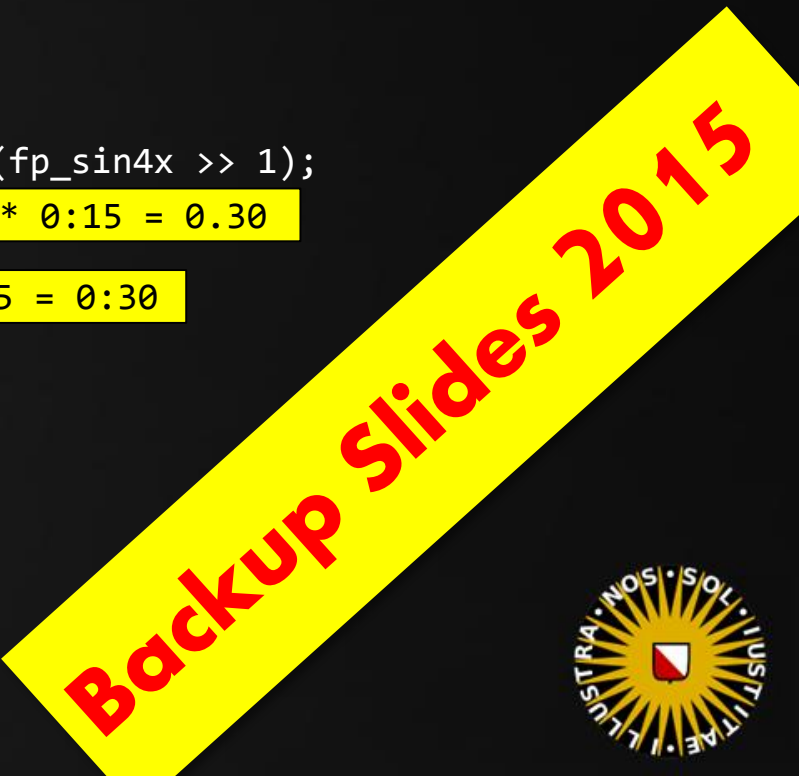
```
return ((fp_sin4x3 - fp_cos4x2) >> 14) + fp_recix;
```

16:16

```
}
```

Here, x is adjusted to use maximum precision: 4:27. 4x is then just a reinterpretation of this number, 6:25.

The calculation of sin4x3 is interesting: since sin(x) is -1..1, sin(x)^3 is also -1..1. We drop a minimal amount of bits and keep precision.



Accuracy

Improving the function.zip example

```
int EvaluateFixed( double x )
{
    int fp_pi = (int)(PI * 65536.0);
    int fp_x = (int)(x * (double)(1 << 27));
    if ((fp_x >> 10) == 0) return 0; // safety net for division

    int fp_4x = fp_x;
    int a = fp_4x / ((2 * fp_pi) >> 3);
    int fp_sin4x = sintab[a & 4095];
    int fp_sin4x3 = (((fp_sin4x >> 1) * (fp_sin4x >> 1)) >> 15) * (fp_sin4x >> 1);

    int fp_cos4x = costab[a & 4095];
    int fp_cos4x2 = (fp_cos4x >> 1) * (fp_cos4x >> 1);

    int fp_recipx = (1 << 30) / (fp_x >> 13);

    return ((fp_sin4x3 - fp_cos4x2) >> 14) + fp_recipx;
}
```

Where do we go from here?

- *The sin/cos tables still contain 1:16 data. However, the way their data is used makes that increasing precision here doesn't help.*
- *We could calculate fp_sin4x3 and fp_cos4x2 via 64-bit intermediate variables. I tried it; impact is minimal...*
- *We can return a value more precise than 16:16 (as we do currently). Problem is around $x = 0$, where the function returns large values and needs the range.*
- *Perhaps 4096 entries in the sin/cos tables is not enough?*

To be continued. ☺



Accuracy

Error – Take-away

- Fixed point code should carefully balance range and precision.
- Do not default to 16:16!
- In multiplications / divisions, carefully conserve precision.
- Use of 64-bit intermediate results is expensive in 32-bit mode. In 64-bit mode, the only disadvantage of 64-bit numbers is increased storage requirements.

