# Lightweight Arithmetic for Mobile Multimedia Devices

Tsuhan Chen 陳祖翰
Carnegie Mellon University
tsuhan@cmu.edu

Thanks to Fang Fang and Rob Rutenbar

---

## Multimedia Applications on Mobile Devices

- **Multimedia Processing**
  - More and more applications are ported from PCs to mobile devices
  - **Floating-point** computational intensive

- **Multimedia System Development**
  - Media designers use 32-64bit floats in C++ for algorithms
  - ASIC designers use 10-20bit fixed-point units for hardware
  - Serious design **disconnect**

## Multimedia Applications on Mobile Devices

- Multimedia Processing
  - More and more applications are ported from PCs to mobile devices
  - **Floating-point** computational intensive

- **Multimedia System Development**
  - Media designers use 32-64bit floats in C++ for algorithms
  - ASIC designers use 10-20bit fixed-point units in hardware
  - Serious design **disconnect**

---

# Fixed-Point vs. Floating-Point

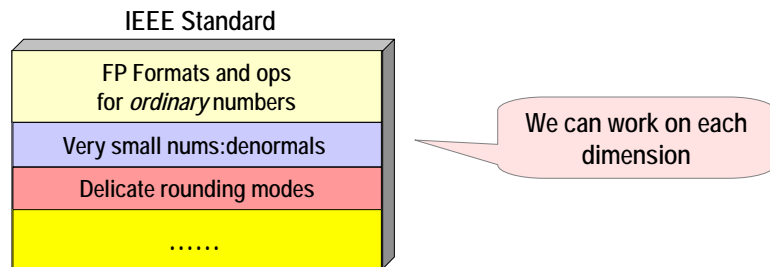| Fixed point | Floating-point |
|---|---|
| s  **Integer**   **fraction** | s   **exp**    **fraction** |
| ✗ Limited dynamic range and precision | ✓ Wide dynamic range & high precision |
| ✓ Small, less power consumption | ✗ Big, power intensive |
| ✗ From SW to HW: time-consuming and error-prone | ✓ Easy translation from SW to HW |

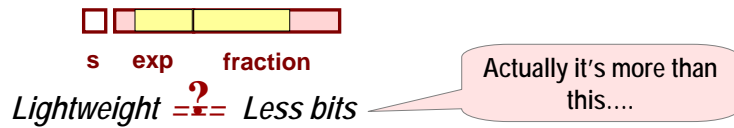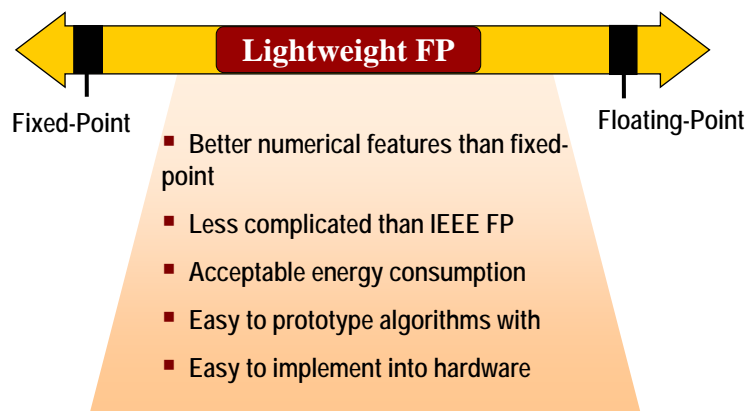*How about make this lightweight?*
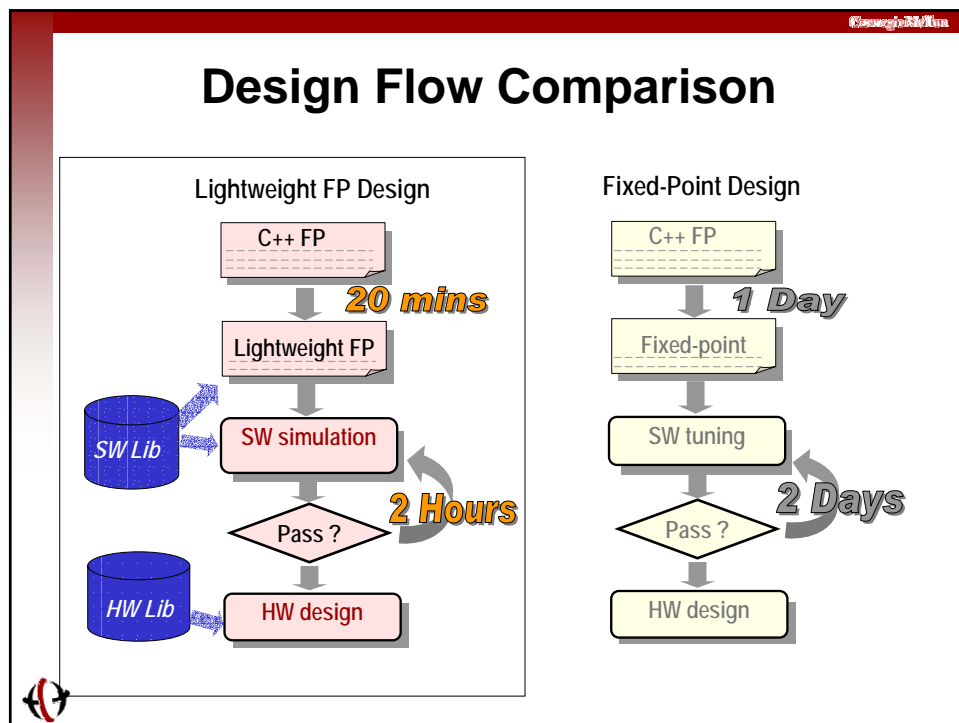
*Don't use more than necessary.*

# What Does "Lightweight" Mean

**s**    **exp**     **fraction**

*Lightweight* $\overset{?}{=}$ *Less bits*

Actually it's more than this....

IEEE Standard

| FP Formats and ops for *ordinary* numbers |
| Very small nums:denormals |
| Delicate rounding modes |
| ...... |

We can work on each dimension
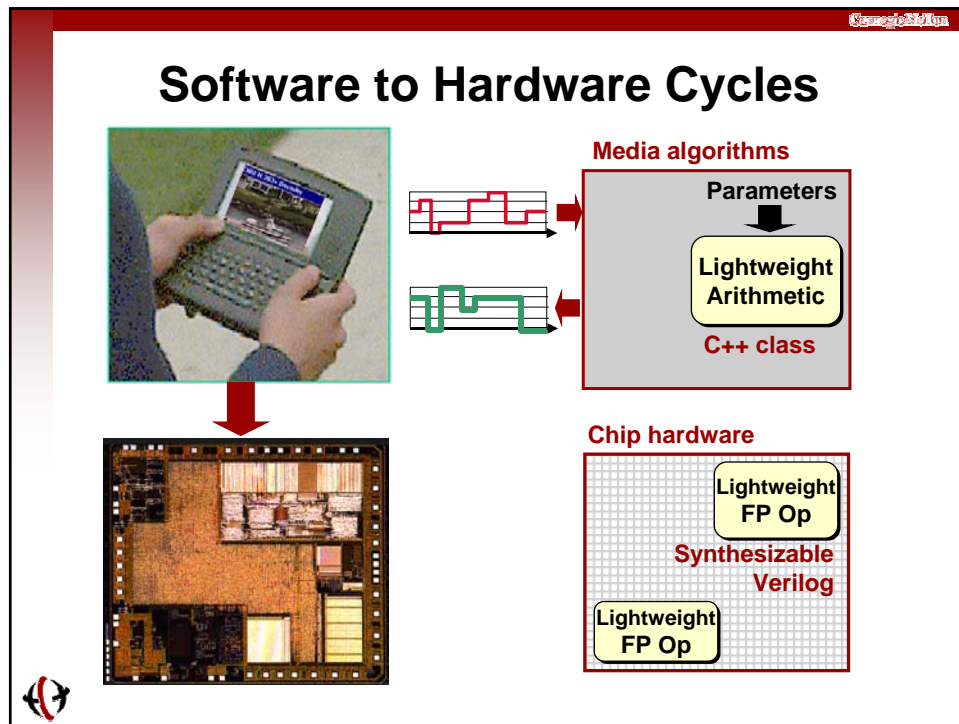
---

# Lightweight Floating-Point Arithmetic

- Lightweight FP arithmetic is a *middle-ground* solution

**Lightweight FP**

Fixed-Point                Floating-Point

- Better numerical features than fixed-point

- Less complicated than IEEE FP

- Acceptable energy consumption

- Easy to prototype algorithms with

- Easy to implement into hardware

# Software to Hardware Cycles

**Media algorithms**

Parameters

Lightweight
Arithmetic

**C++ class**

**Chip hardware**

Lightweight
FP Op

**Synthesizable
Verilog**

Lightweight
FP Op

---

# Design Flow Comparison

Lightweight FP Design

C++ FP

*20 mins*

Lightweight FP

*SW Lib*

SW simulation

*2 Hours*

Pass ?

*HW Lib*

HW design

Fixed-Point Design

C++ FP

*1 Day*

Fixed-point

SW tuning

*2 Days*

Pass ?

HW design

**4**

# Design Flow Comparison

**Lightweight FP Design**

```
C++ FP
  │  20 mins
  ▼
Lightweight FP
  │
  ▼
SW simulation  ◄── SW Lib
  │        ⟳ 2 Hours
  ▼
Pass ?
  │
  ▼
HW design  ◄── HW Lib
```

**Fixed-Point Design**

```
C++ FP
  │  1 Day
  ▼
Fixed-point
  │
  ▼
SW tuning
  │        ⟳ 2 Days
  ▼
Pass ?
  │
  ▼
HW design
```

---

# IEEE Standard vs. Lightweight IP

## IEEE FP Standard

- **32 / 64 bits**
  - 8 / 11 bits exponent
  - 23 / 52 bits mantissa
  - 1 sign bit

- **Specs**
  normal numbers *as well as* special values (infinity), edge cases (INF - INF), etc.

## Lightweight Arithmetic IP

- *Fewer bits*
  - Fewer bits of fraction
    → less numerical precision
  - Fewer bits of exponent
    → less dynamic range

- **Which of the special cases/numbers should be supported?**

# IEEE Floats vs. CMUfloats

CMUfloats

IEEE Floats

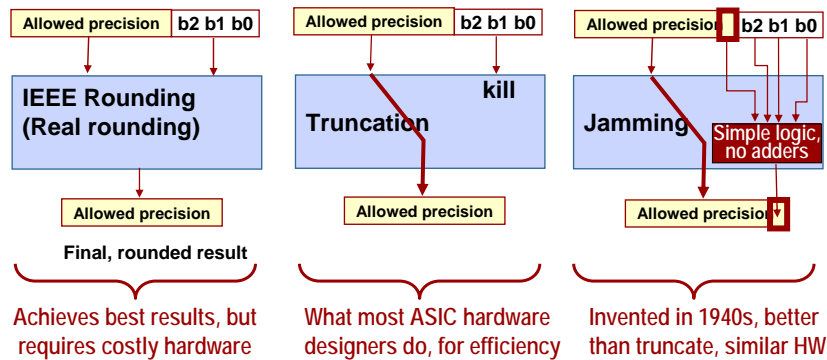| FP Formats and ops for *ordinary* numbers |
| Very small nums:denormals |
| Delicate rounding modes |
| ...... |

- Customizable format providing variable dynamic range and precision
  Fraction [1, 23],    exponent width [1,8]

- On-off switch for denormalization

- Multiple choices for rounding mode
  Real-rounding / Jamming / Truncation

---

# Rounding in CMUfloat

- We support not only IEEE rounding, but also two "quick & dirty" modes

| Allowed precision | b2 b1 b0 |

**IEEE Rounding (Real rounding)**

| Allowed precision |

**Final, rounded result**

Achieves best results, but requires costly hardware

| Allowed precision | b2 b1 b0 |

kill

**Truncation**

| Allowed precision |

What most ASIC hardware designers do, for efficiency

| Allowed precision | b2 b1 b0 |

**Jamming**   Simple logic, no adders

| Allowed precision |

Invented in 1940s, better than truncate, similar HW

# C++ CMUfloat library

- Supported operators

| Cmufloat double float int short | = | Cmufloat | + == - >=, > * <=, < / != | Cmufloat double float int short |
|---|---|---|---|---|

- Other supported C++ features
  - Pointer *Cmufloat * a;*
  - Reference *Cmufloat & a ;*
  - Array *Cmufloat a[10][10] ;*
  - Argument passing *func ( Cmufloat a )*
  - I/O stream *cout << a;*

# C++ Cmufloat Library

- Supported operators

| Cmufloat double float int short | = | Cmufloat | + == - >=, > * <=, < / != | Cmufloat double float int short |
|---|---|---|---|---|

- Other supported C++ features
  - Pointer *Cmufloat * a;*
  - Reference *Cmufloat & a ;*
  - Array *Cmufloat a[10][10] ;*
  - Argument passing *func ( Cmufloat a )*
  - I/O stream *cout << a;*

# Software Library: Advantages

- Transparent mechanism to embed 'Cmufloat' in the algorithm
  - The overall structure of the source code can be preserved
  - Minimal effort in translating standard FP to lightweight FP

```
Cmufloat <14,5> a = 0.5;  // 14 bit fraction and 5 bit exponent
Cmufloat <> b= 1.5;       // Default Cmufloat is IEEE float
Cmufloat <18,6> c[2];     //  Define an array
float fa;

c[1] = a + b;
fa   = a * b;             // Assign the result to float
c[2] = fa + b;            // Operation between float and Cmfloat
```

---

# Software Library: Advantages (Cont.)

- Arithmetic operators are implemented by bit-level manipulation: more precise

Our approach:
Emulates the hardware implementation exactly

Previous approach

```
Add( b,  c) {
a' = b + c;
a = round (a');
}
```
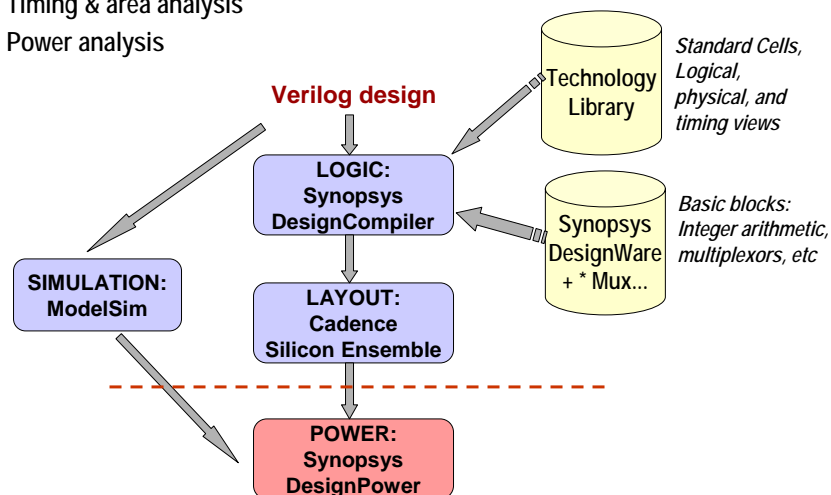
Built-in FP operator

Round to limited bit-width

# Summary: Features Supported

- Bit widths
  - Variable from 2 bits (1 sign + 1 exp + 0 man) to 32 bits ( IEEE std)
- Rounding
  - Use jamming (1.00011 rounds to 1.01)
  - Experiments show jamming is nearly as good as full IEEE rounding, always superior to truncation, yet same complexity as truncation
- Denormalized numbers
  - Not supported--our experiments on video/audio codecs suggest that denormal numbers do not improve the performance
- Exceptions
  - Support only the exceptional values for infinity, zero and NAN
  - Helps make the smaller FP sizes more robust
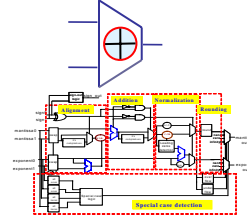
---

# Hardware Library:  ASIC Design Flow

- Verilog to layout flow
- Timing & area analysis
- Power analysis

**Verilog design**

Technology Library
*Standard Cells, Logical, physical, and timing views*

**LOGIC: Synopsys DesignCompiler**

Synopsys DesignWare + * Mux...
*Basic blocks: Integer arithmetic, multiplexors, etc*

**SIMULATION: ModelSim**

**LAYOUT: Cadence Silicon Ensemble**

**POWER: Synopsys DesignPower**
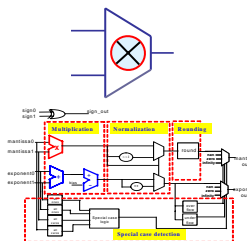
# Lightweight FP Adders/Multipliers

- Feature Supported
  - Bit widths:
    Variable from 3 bits (1 sign + 1 exp + 1 frac) to 32 bits ( IEEE std)
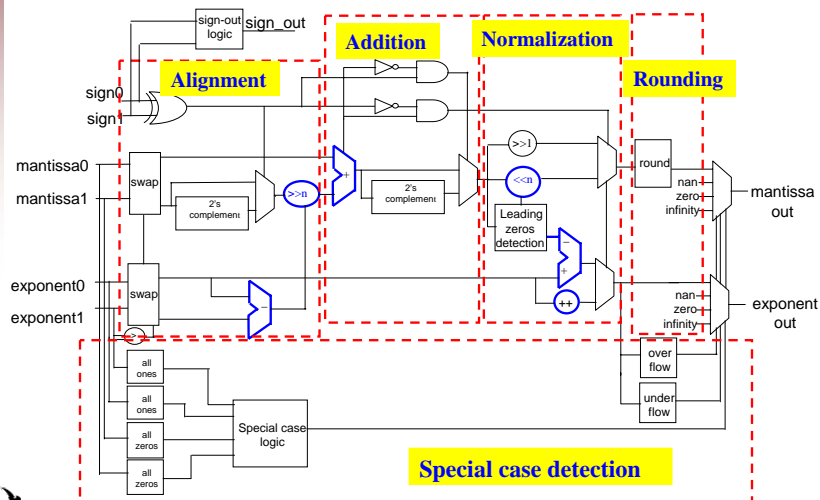  - Rounding:
    Jamming / Truncation

- Design Issues
  - Design method
  - Subcomponent structures
    - Core integer adder structure?
    - Core shifter structure?
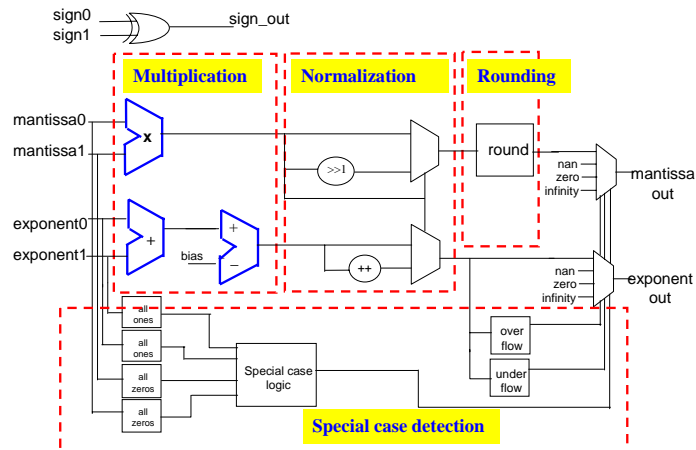    - Core integer multiplier structure?



---

# Floating Pt Adder

Blue modules have large area and / or delay
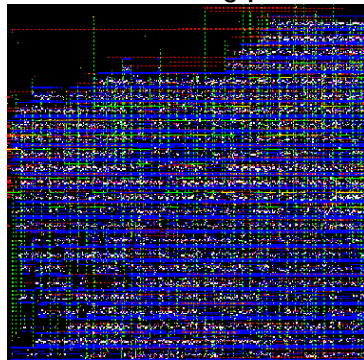
# Floating Pt Multiplier
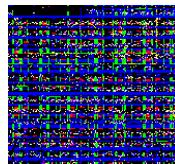
Blue modules have large area and / or delay



We see that the multiplier has less 'over-head' than the adder

---

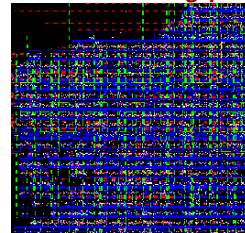# Design Examples: Adders

32-bit Floating-point



14-bit Floating-pt

20-bit Fixed-pt

| | 32-bit FP | 20-bit FIX | 14-bit FP |
|---|---|---|---|
| Area( $um^2$ ) - post layout | 26634 | 4866 | 10096 |
| Delay(ns) - post synthesys | 48.95 | 2.44 | 25.77 |

# Adder Analysis

Fixed-point adder is *very* simple;  floating point adder is complex

### Fixed Point

```
aaaa.aaaa
bbbb.bbbb
```
$$0pppp.pppp$$

If not 0, then overflow

No rounding

```
pppp.pppp
```
Final answer

### Floating Point

$$aa \times 2^{aa}$$
$$bb \times 2^{bb}$$
$$1.aa \quad \times 2^{MAX(aa,bb)}$$
$$0.00bb \times 2^{MAX(aa,bb)}$$
$$0.0ppp \times 2^{MAX(aa,bb)}$$
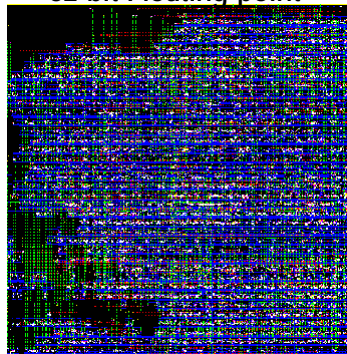
Alignments are expensive

Extra add & normalize steps

$$pp \times 2^{pp}$$
Final answer

---
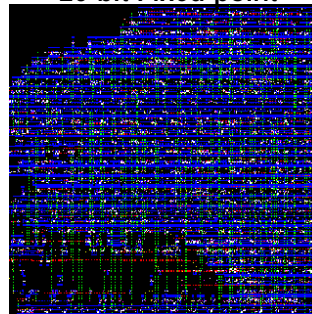
# Design Examples: Multipliers

32-bit Floating-point
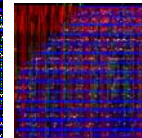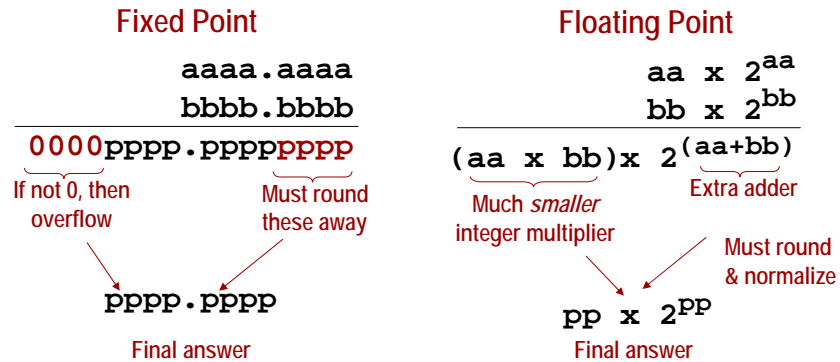
20-bit Fixed-point

*14-bit Floating-pt*



|  | 32-bit FP | 20-bit FIX | 14-bit FP |
|---|---|---|---|
| Area( um$^2$) - post layout | 60713 | 40738 | 8851 |
| Delay(ns) - post synthesis | 24.14 | 22.82 | 15.89 |

**12**

# Multiplier Analysis

Fixed-pt needs *more* bits to get the *same* dynamic range,
*increasing* the size of the multiplier unit

### Fixed Point

```
    aaaa.aaaa
    bbbb.bbbb
0000pppp.pppppppp
```

If not 0, then overflow

Must round these away

```
pppp.pppp
```

Final answer

### Floating Point

$$aa \times 2^{aa}$$
$$bb \times 2^{bb}$$
$$(aa \times bb) \times 2^{(aa+bb)}$$

Much *smaller* integer multiplier

Extra adder

Must round & normalize

$$pp \times 2^{pp}$$

Final answer

---

# Radix-2 vs. Radix-16

- Higher radix has less complexity in the shifter

**Radix-2 FP ( 8-bit fraction + 1 leading bit )**

**Radix-16 FP ( 11-bit fraction )**

- More fraction width increases the complexity of multiplier
- If the number of adders is much more than the number of multipliers, then radix-16 is preferred

| Radix | Area (um$^2$) | Delay (ns) | Radix | Area (um$^2$) | Delay (ns) |
|-------|---------------|------------|-------|---------------|------------|
| 2 | 8401 | 10.21 | 2 | 7893 | 5.8 |
| 16 | 7389 (-12%) | 8.48 (-17%) | 16 | 11284 (+43%) | 6.62 (+14%) |
| | **Adder** | | | **Multiplier** | |

**13**

# Power Analysis

- IDCT in
  - 32-bit IEEE FP
  - 15-bit radix-16 lightweight FP
  - Fixed-point implementation
    - 12-bit accuracy for constants
    - Widest bit-width is 24 in the whole algorithm (not fine tuned)

| Implementation | Area(um$^2$) | Delay(ns) | Power(mw) |
|---|---|---|---|
| IEEE FP | 926810 | 111 | 1360 |
| Lightweight FP | 216236 | 46.75 | 143 |
| Fixed-point | 106598 | 36.11 | 110 |

---

# Multimedia Encoding/Decoding

Encoding of Media

Video camera

or **Encoder** 101110010...

Uncompressed multimedia file

Encoded bitstream

At the decoding end, we have a choice as to how accurately we wish to decode the data for playback

Playback of Media

**Decoder**

# Video Codec

- H.261/263, MPEG-1/2/4, and even JPEG

**8-bit**  **9-bit**  **Floating point**  **DCT** → **Q** **12-bit** *Transmit* **12-bit**

IDCT requires floating point, and has an IEEE quality spec (1180-1990) that requires comparison against a 64-bit IEEE double implementation

**IQ**

**IQ**

**Floating point**

**Floating point**

**8-bit**

**IDCT**

**IDCT**

**9-bit**

**9-bit**

+

+

**Motion Comp**

**D**

**D**

**Motion Comp**

**8-bit**

---

# Video Quality vs. Bit-width

- Use PSNR (Peak-Signal-to-Noise) to measure perceptual video quality

**Test video**  →  **Proposed Codec CMUfloat**

Qi

Pi

**Measure PSNR**
**(Noise = |Qi – Pi| )**

- CMUfloat can go *very small*, ~14bits
( 5 exponent + 8 fraction + 1 sign bits = 14 total bits )

Yellow pts show where PSNR decreases by 0.2dB from asymptotic value

**PSNR**
**(dB)**

39
37
35
33
31

5  7  9  11  13  15  17  19  21  23 **Fraction-width (bit)**

**15**

# Rounding Modes

- Compare 3 rounding modes using IDCT video streams

**Comparison of Rounding Methods**

Jamming is nearly as good as real rounding in precision, but as simple as truncation in hardware.

Legend:
- Real
- Jamming
- Truncation

PSNR (dB) axis: 33, 34, 35, 36, 37, 38, 39

Fraction-width(bit) axis: 6 7 8 9 10 11 12 13 14 15 16 17

---

# Video Demo

– IEEE double vs. variable-precision CMUfloats

Decoded with 14-bit "lightweight" IDCT

Decoded with 64-bit "double" IDCT

Decoded with 11-bit "lightweight" IDCT

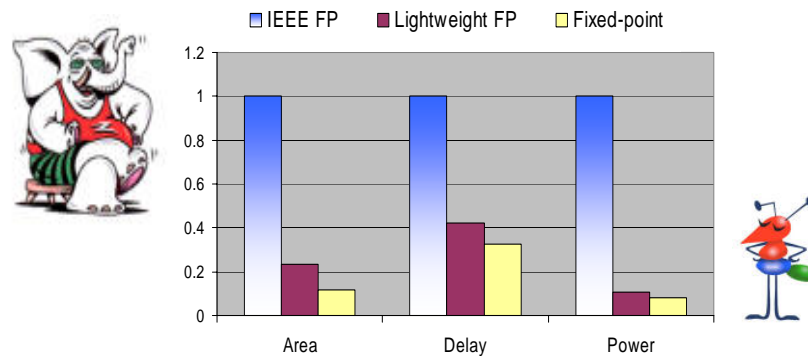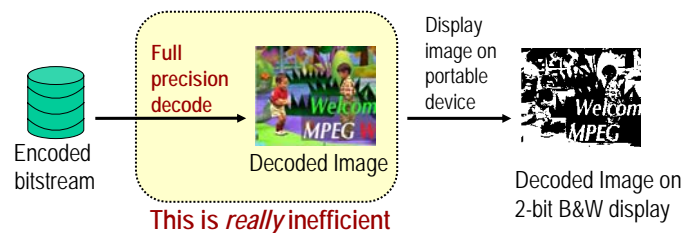## Hardware Reduction Using Lightweight FP

- Comparison in Area/Delay/Power
  - 32-bit IEEE FP IDCT / 14-bit lightweight FP IDCT with Jamming rounding / 20-bit fixed point IDCT
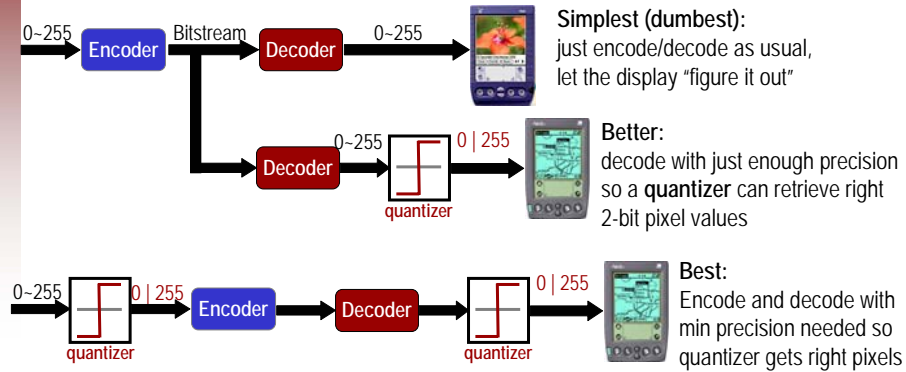


---

# Low-Resolution Display

- Media software commonly done in full precision (32–64 bits)
  - *Why do this if the display cannot handle it?*
  - On a portable video player:



Encoded bitstream → Full precision decode → Decoded Image → Display image on portable device → Decoded Image on 2-bit B&W display

This is *really* inefficient

- Can't we do better than this, with smarter operators?

# Low-Resolution Display (cont.)

0~255 → **Encoder** → Bitstream → **Decoder** → 0~255 →

**Simplest (dumbest):**
just encode/decode as usual, let the display "figure it out"

→ **Decoder** → 0~255 → [quantizer] → 0 | 255 →

**quantizer**

**Better:**
decode with just enough precision so a **quantizer** can retrieve right 2-bit pixel values

0~255 → [quantizer] → 0 | 255 → **Encoder** → **Decoder** → [quantizer] → 0 | 255 →

**quantizer**              **quantizer**

**Best:**
Encode and decode with min precision needed so quantizer gets right pixels

- Results
  - Simplest: needs ~**20-bit** lightweight floats to work
  - Better: needs **16-bit** lightweight floats; even just **11-bits** looks decent
  - Best: needs just **9-bit** floats (4 fraction bits) to work just fine.

---

# Video Demo

- Full Precision (64 bit)
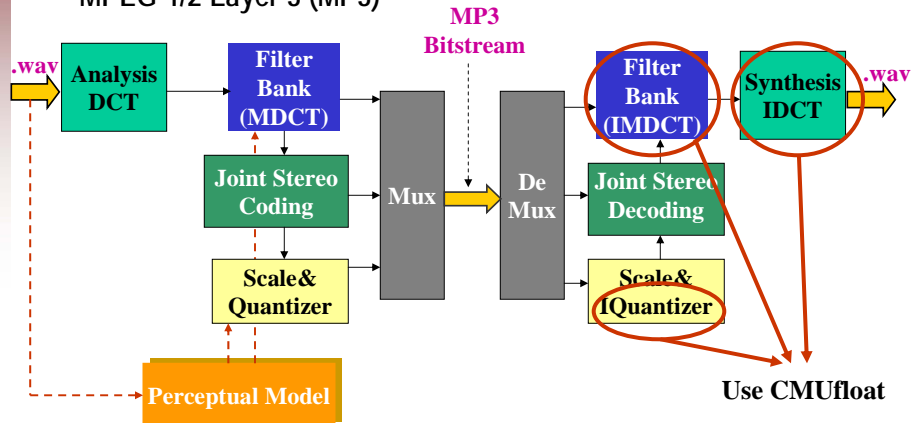
- Using 23 bits (IEEE 1180 passed)

- Using 11 bits (IEEE 1180 failed)
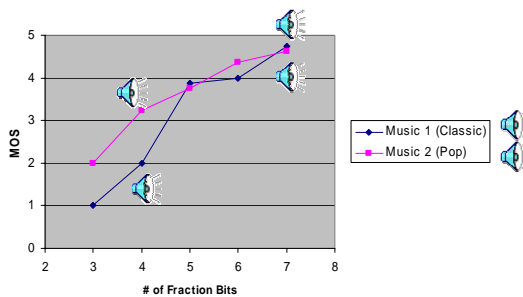
**18**

# How About Audio?

- MPEG-1/2 Layer 3 (MP3)

**MP3 Bitstream**

.wav → **Analysis DCT** → **Filter Bank (MDCT)** → **Mux** → **De Mux** → **Filter Bank (IMDCT)** → **Synthesis IDCT** → .wav

**Joint Stereo Coding** → **Joint Stereo Decoding**

**Scale& Quantizer** **Scale& IQuantizer**

**Perceptual Model**

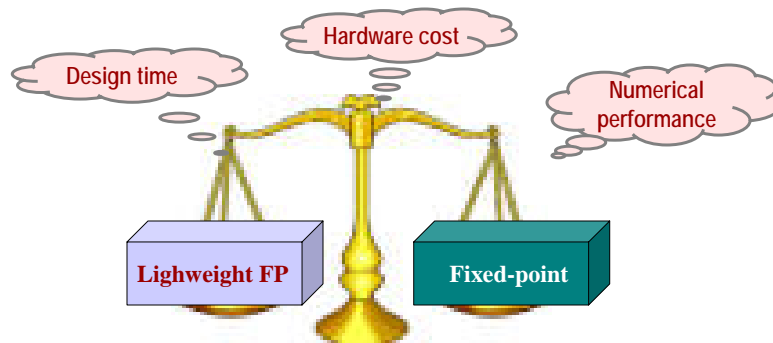**Use CMUfloat**

- No standard tests for quality

---

# Audio Quality

- Need to rely on subjective testing on perceptual quality
  - Mean Opinion Score (MOS)
    - From 5 "imperceptible difference" to 1 "really annoying"
- Results
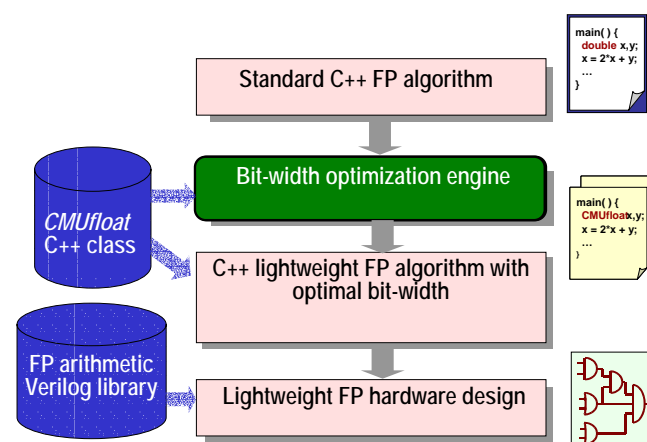  - 8 subjects. 6-bit exponent and 3~7 bit fraction



Legend: Music 1 (Classic), Music 2 (Pop)

Y-axis: MOS (0 to 5)
X-axis: # of Fraction Bits (2 to 8)

---

# Recap…

- **Accomplishments**
  - C++ lightweight FP arithmetic library
  - Verilog lightweight FP arithmetic library
  - Extensive experiments on video/audio/speech

- **Is the lightweight FP solution universal?**
  - No, tradeoff between fixed-point solution and lightweight FP solution

- **Ongoing work**
  - Automatic design flow

- **Important for low-power mobile devices**

---

## Advanced Multimedia Processing Lab

Please visit us at:

http://amp.ece.cmu.edu