# Convolution Engine:
## Balancing Efficiency & Flexibility in Specialized Computing

Did the heavy lifting but could not come today

**Wajahat Qadeer, Rehan Hameed,**

**Ofer Shacham,** That's me ☺

**Preethi Venkatesan, Christos Kozyrakis, Mark Horowitz**

**Stanford University**

- **By show of hands, who here has an (HD) camera on them?**
- **How many CPU's/GPU's in the room?**
- **How many of those xPU's are used for the image processing?**

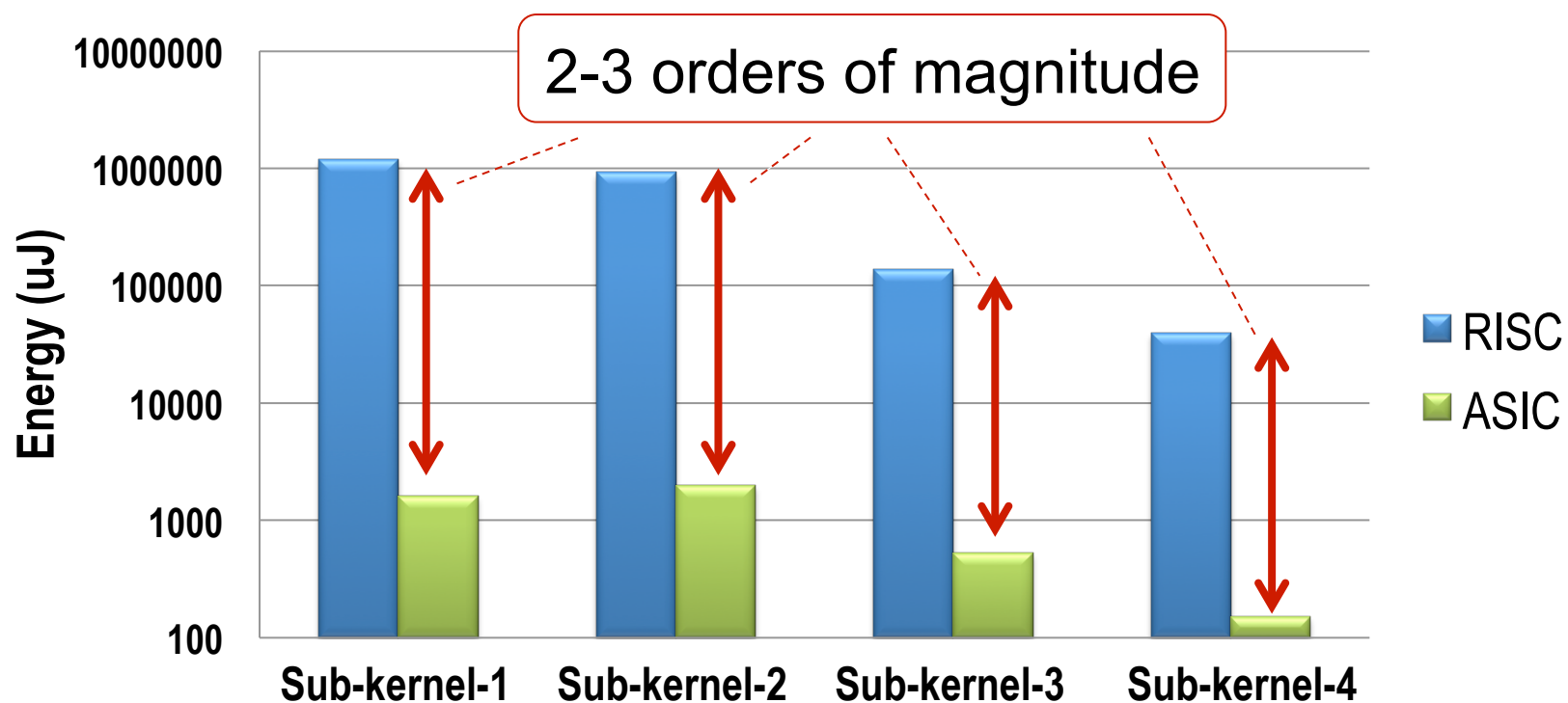shacham@alumni.sta

# Imaging and video systems

- **High computational requirements, low power budget**
  - Stills: ~10M pixels  x  10 frames per second
  - Video: ~2M pixels  x  30 frames per second
  - ~400 math operations per pixel (just for the image acquisition)

- **On CPU… not enough horse power**

- **On GPU… too much power**

- **Typically use special purpose custom HW**
  - About 500X better performance, 500X lower energy than CPU

- **By coupling compute and storage closely together, ASIC's are orders of magnitude performance and energy more efficient**



2-3 orders of magnitude

* R. Hameed et. al., Understanding Sources of Inefficiency in General-Purpose Chips. ISCA '10
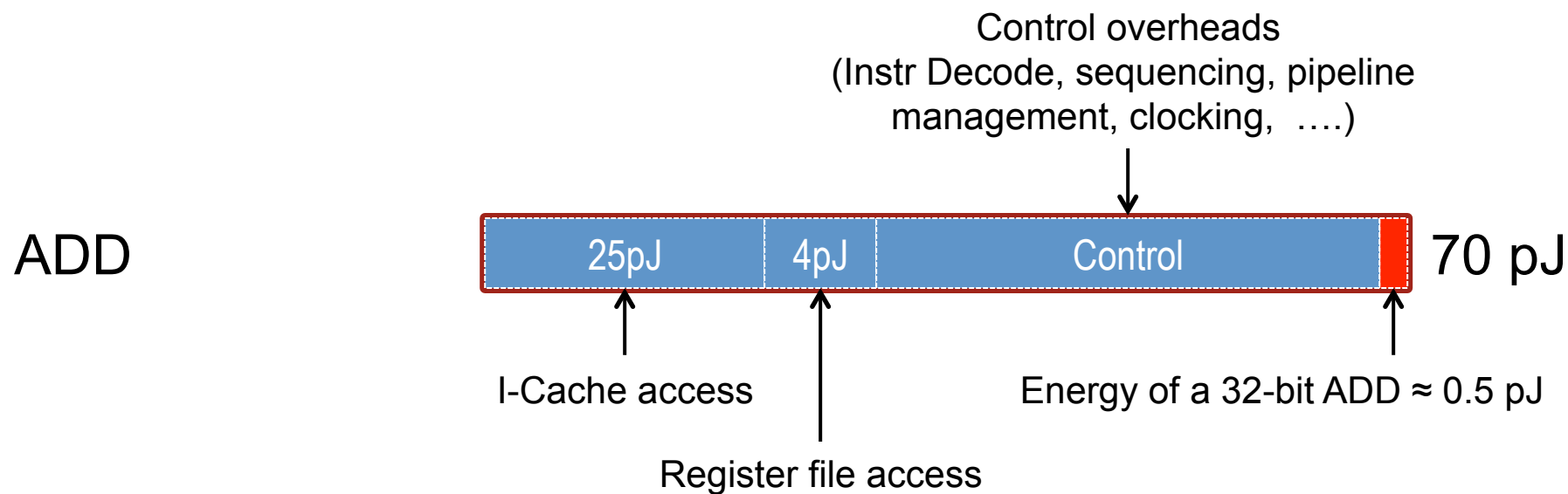
# We are solving the wrong problem!

- **Yes, ASIC is 1000X more efficient than general purpose**
- **Yes, general purpose is more programmable than ASIC**
- **Yes, we can make each one marginally better**

- **But those are good answers to all the wrong questions!**

- **The right questions:**

    **Why is the RISC energy so high?**

    **What type of computation can we make efficient?**

    **Can we make it just 100X better but keep it programmable?**
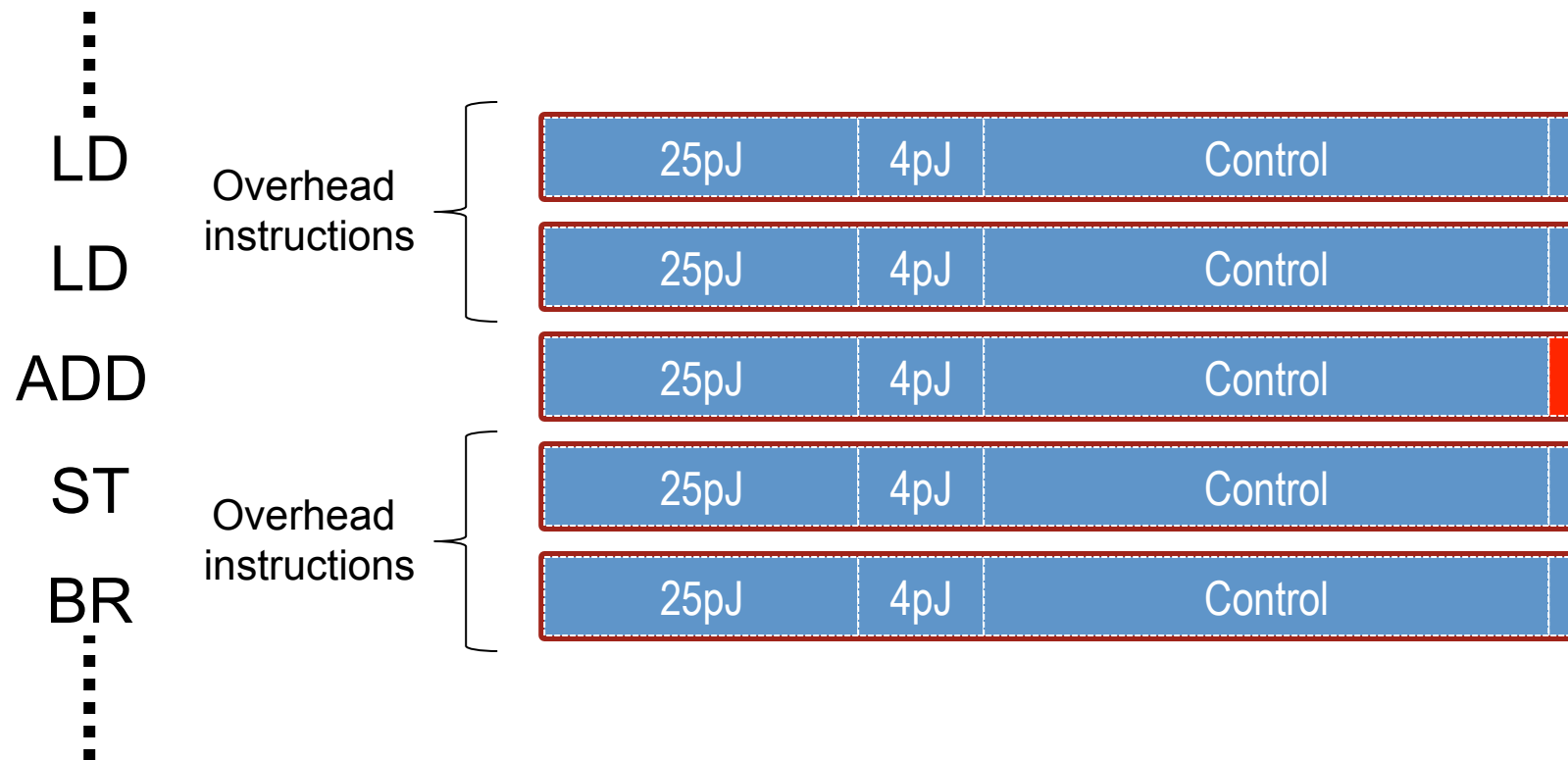
# Anatomy of a RISC Instruction

Control overheads
(Instr Decode, sequencing, pipeline
management, clocking,  ….)

ADD

| 25pJ | 4pJ | Control | |
|------|-----|---------|--|

70 pJ

I-Cache access

Register file access

Energy of a 32-bit ADD ≈ 0.5 pJ

\* Assuming a typical 32-bit embedded RISC in 45mn @ 0.9V technology

# Other instructions overhead

LD

LD

ADD

ST

BR

Overhead instructions

Overhead instructions

| 25pJ | 4pJ | Control |
|------|-----|---------|

| 25pJ | 4pJ | Control |
|------|-----|---------|

| 25pJ | 4pJ | Control |
|------|-----|---------|

| 25pJ | 4pJ | Control |
|------|-----|---------|

| 25pJ | 4pJ | Control |
|------|-----|---------|

\* Assuming a typical 32-bit embedded RISC in 45mn @ 0.9V technology

shacham@alumni.stanford.edu

# D-Cache accesses overhead

| LD | 25pJ | 25pJ | 4pJ | Control |
| LD | 25pJ | 25pJ | 4pJ | Control |
| ADD | | 25pJ | 4pJ | Control |
| ST | 25pJ | 25pJ | 4pJ | Control |
| BR | | 25pJ | 4pJ | Control |

↑
D-Cache access
overheads

\* Assuming a typical 32-bit embedded RISC in 45mn @ 0.9V technology

# SIMD machines give some improvement

- **SIMD units amortize overhead <u>and</u> improve performance**

ADD

| I-Cache | RF | Control | |

SIMD ADD

| I-Cache | RF | Control | |

- **Achieves 10X better energy and performance AND is programmable**

- **Can we do 100X and keep it programmable?**

shacham@alumni.stanford.edu

**Each memory and instruction fetch must be amortized by hundreds of operations**

# What we want to see

D-Cache accesses much
narrower than functional path

| LD | D-Cache | I-Cache | Reg File | Control |

| OP | I-Cache | Reg File | Control | |
| OP | I-Cache | Reg File | Control | |
| OP | I-Cache | Reg File | Control | |

Many ALU instructions
per LD/ST instruction

Many ops per instruction

| OP | I-Cache | Reg File | Control | |
| OP | I-Cache | Reg File | Control | |
| OP | I-Cache | Reg File | Control | |

| ST | D-Cache | I-Cache | Reg File | Control |

- **Most of the computation is performed over (overlapping) stencils**

- **Looks like convolution:** $\left(Img \otimes f\right)_{[n,m]} = \sum_{l=-c}^{c} \sum_{k=-c}^{c} Img_{[k,l]} \cdot f_{[n-k,m-l]}$
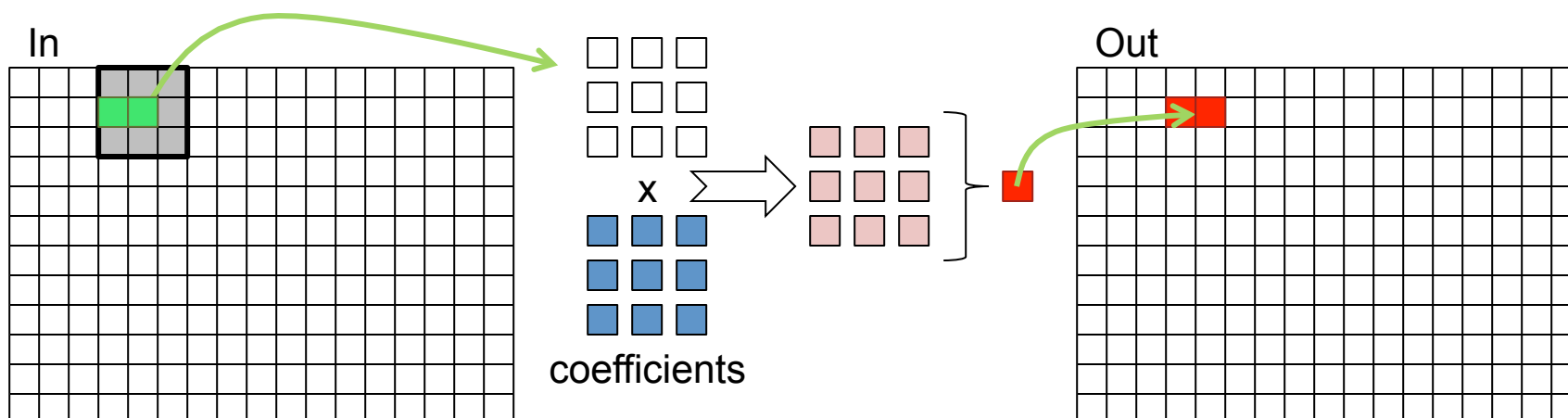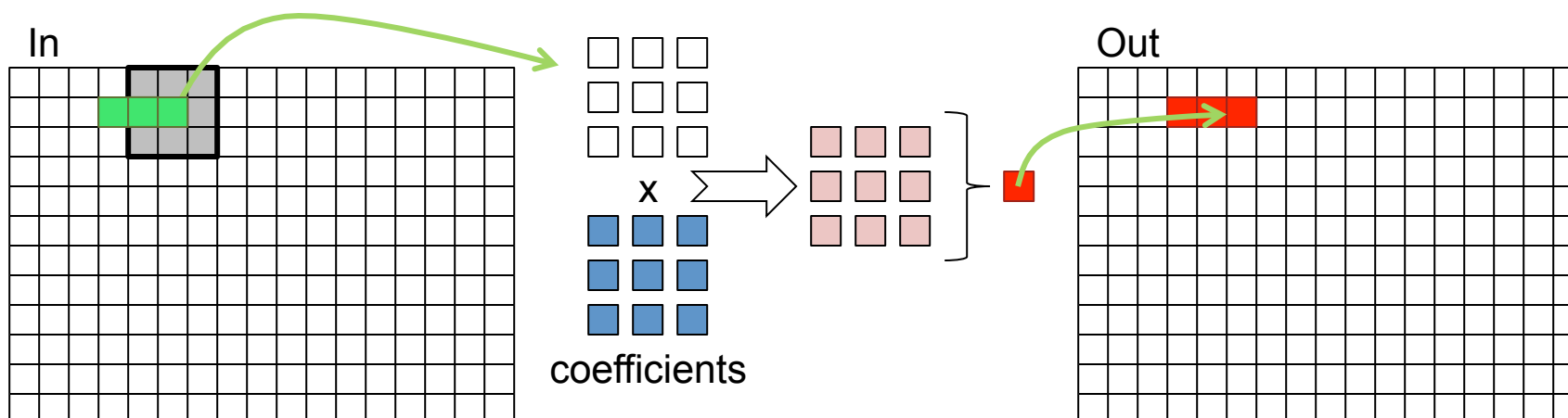
In

x

coefficients

Out

- **Most of the computation is performed over (overlapping) stencils**

- **Looks like convolution:** $\left(Img \otimes f\right)_{[n,m]} = \sum_{l=-c}^{c} \sum_{k=-c}^{c} Img_{[k,l]} \cdot f_{[n-k,m-l]}$
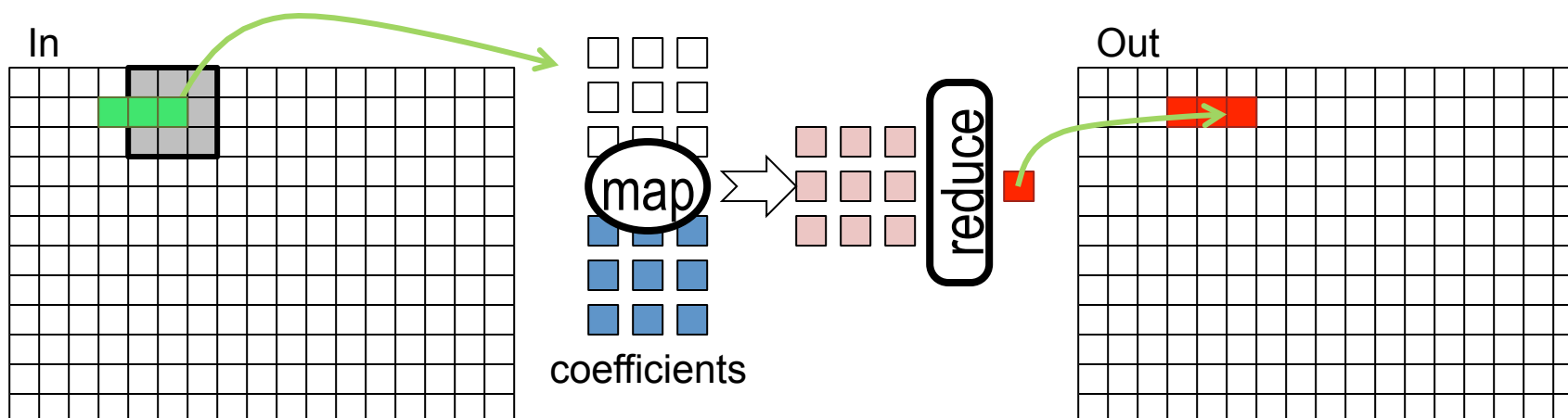
In

x

coefficients

Out

- **Most of the computation is performed over (overlapping) stencils**

- **Looks like convolution:** $\left(Img \otimes f\right)_{[n,m]} = \sum_{l=-c}^{c} \sum_{k=-c}^{c} Img_{[k,l]} \cdot f_{[n-k,m-l]}$

In

coefficients

Out

- **It only looks like convolution:**

$$\left( Img \overset{CE}{\otimes} f \right)_{[n,m]} = Reduce_{l=-c}^{c}\left[ Reduce_{k=-c}^{c}\left[ map\left( Img_{[k,l]}, f_{[n-k,m-l]} \right) \right] \right]$$
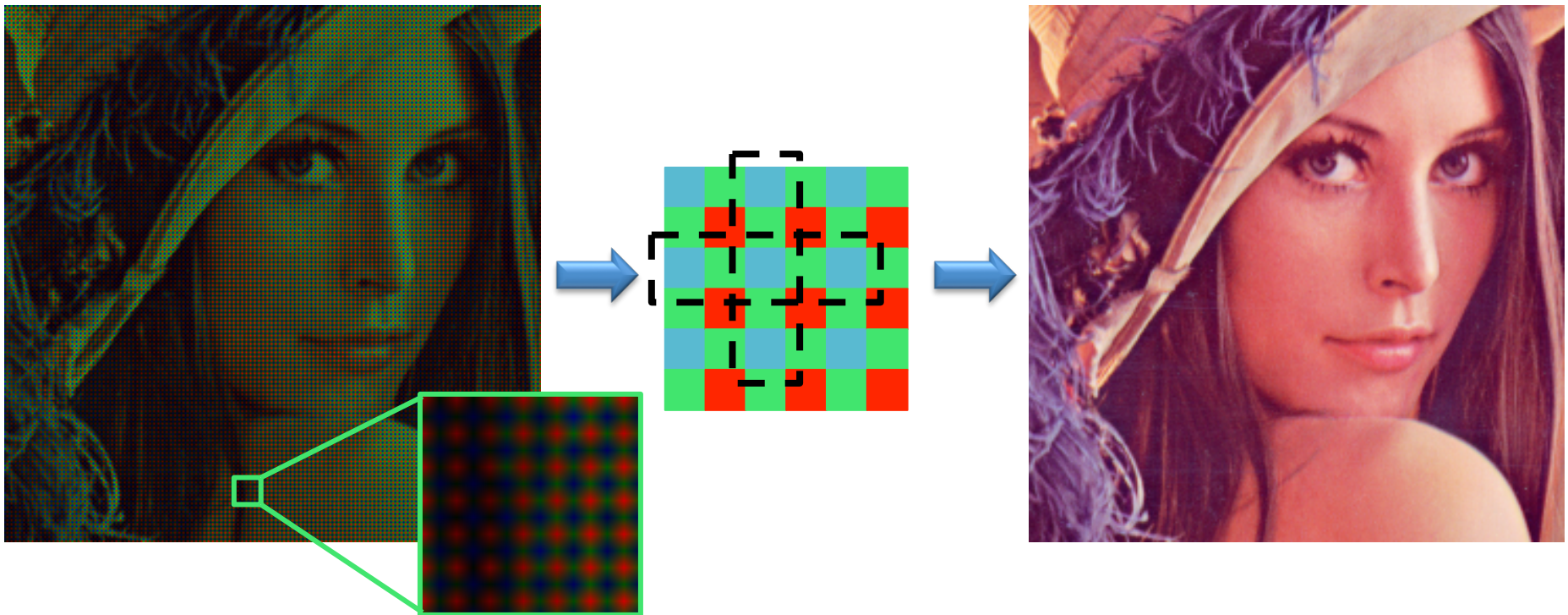


In

map

coefficients

reduce

Out

- **De-mosaic:**
  - Adaptive color plane interpolation (ACPI)*: image gradients followed by a three-tap filter in the direction of smallest gradient.
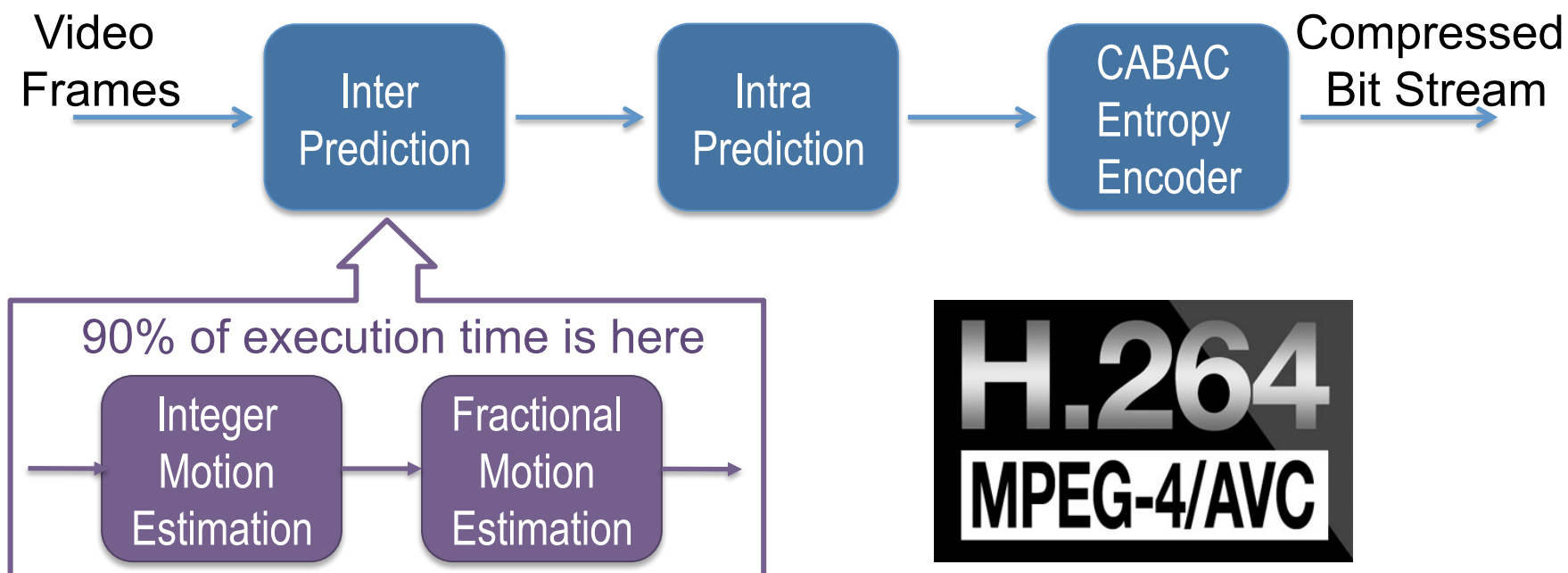


* Y. Cheng et. al. An adaptive color plane interpolation method based on edge detection.  Journal of Electronics (China), 2007.
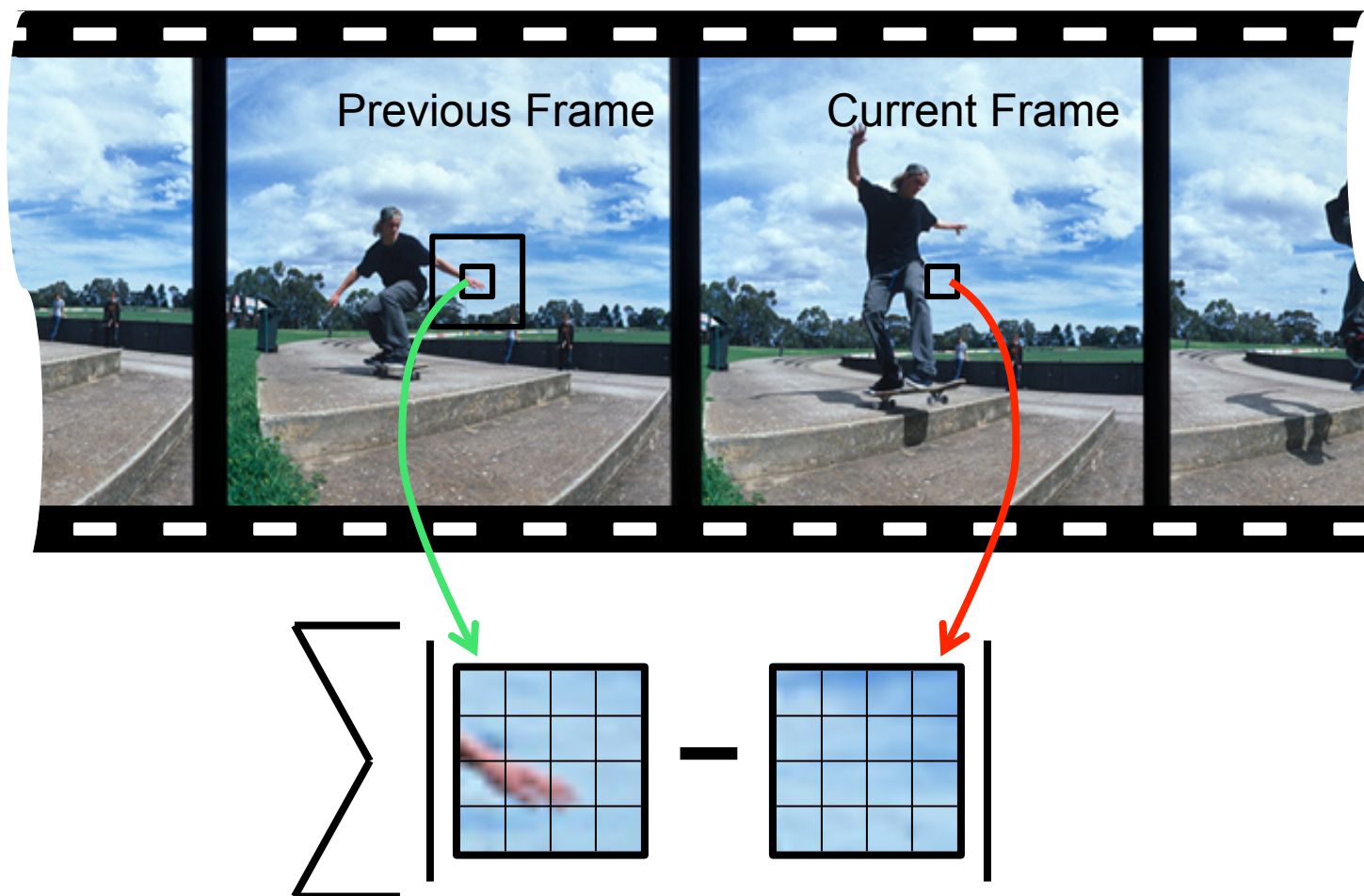
# Let's look at more convolution-like workloads

- **H.264 (high definition) video encoder:**
  - IME: 2D-Sum of absolute differences
  - FME: Half pixel interpolation, quarter pixel interpolation, 2D SAD

Video Frames → Inter Prediction → Intra Prediction → CABAC Entropy Encoder → Compressed Bit Stream

90% of execution time is here

Integer Motion Estimation → Fractional Motion Estimation

H.264 MPEG-4/AVC

- **Trying to find best match for a stencil within a small neighborhood**

shacham@alumni.stanford.edu

# The convolution engine must support different ops

| | Map | Reduce | Stencil Size | Data Flow |
|---|---|---|---|---|
| IME SAD | Abs Diff | Add | 4x4 | 2D Convolution |
| FMW ½ pixel up-sample | Multiply | Add | 6 | 1D Horizontal & vertical conv. |
| FME ¼ pixel up-sample | Average | None | -- | 2D Matrix operation |
| SIFT Gaussian blur | Multiply | Add | 9, 13, 15 | 1D Horizontal & vertical conv. |
| SIFT DoG | Subtract | None | -- | 2D Matrix operation |
| SIFT Extreme | Compare | Logic AND | 3 | 1D Horizontal & vertical conv. |
| Demosaic interpolation | Multiply | Complex | 3 | 1D Horizontal & vertical conv. |

# Convolution Engine:
# An architecture for convolution-like kernels

Coefficients

Stencil neighborhood

2D Regfile

| 0 | 0 | 1 | | 15 |
| 1 | 0 | 1 | | 15 |
| | | | | |
| 15 | 0 | 1 | | 15 |

2D **shift** Regfile

| 0 | 0 | 1 | | 15 | 16 | 17 | | 31 |
| 1 | 0 | 1 | | 15 | 16 | 17 | | 31 |
| | | | | | | | | |
| 15 | 0 | 1 | | 15 | 16 | 17 | | 31 |

Wide 64-lane SIMD "map" unit

ALU  ALU  ALU   . . . . . . . . . . . . . . .   ALU

Flexible "reduce" step

Arithmetic / Logical reduction

Current
frame pixels

Reference
frame pixels

2D Regfile

2D shift
Regfile

Wide 64-
lane SIMD
"map" unit

ALU  ALU  ALU  . . . . . . . . . . . . . . .  ALU

Flexible
"reduce" step

Arithmetic / Logical reduction

# Example from H.264's motion estimation: Mapping Sum of Absolute Differences (SAD)
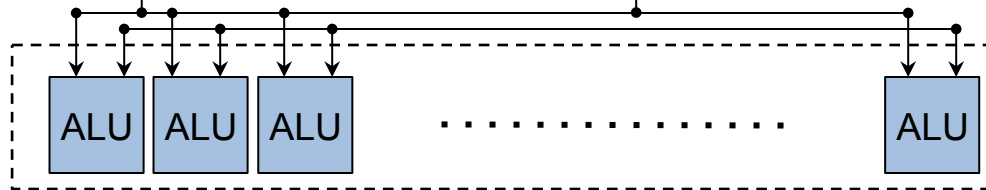


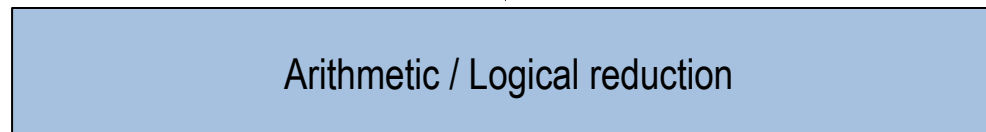Current frame pixels

Reference frame pixels

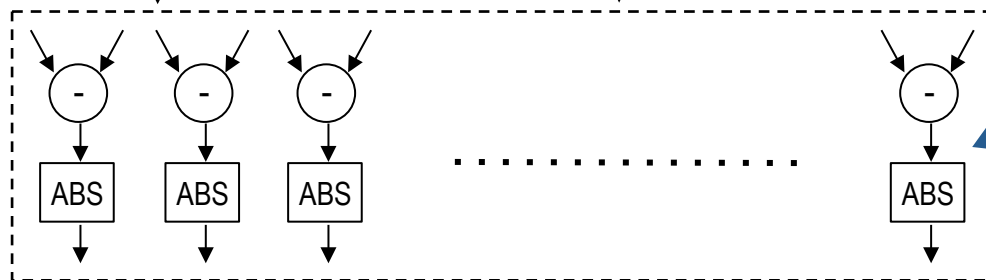2D Regfile

2D shift Regfile

Wide 64-lane SIMD "map" unit

ALU's instruction set to |a-b|

ABS   ABS   ABS   . . . . . . . . . . . . .   ABS

Flexible "reduce" step

Arithmetic / Logical reduction

# Example from H.264's motion estimation: Mapping Sum of Absolute Differences (SAD)



Current frame pixels

Reference frame pixels

pixels shift left

2D Regfile

2D shift Regfile

Wide 64-lane SIMD "map" unit

ALU's instruction set to |a-b|

Flexible "reduce" step

Sum (Reduction)

Summation tree

# Example from H.264's motion estimation: Mapping Sum of Absolute Differences (SAD)



Reference frame pixels

pixels shift left

2D Regfile

|   |   |   |   |     |
|---|---|---|---|-----|
| 0 | 0 | 1 |   | 15  |
| 1 | 0 | 1 |   | 15  |
|   |   |   |   |     |
| 15| 0 | 1 |   | 15  |

2D shift Regfile

|    |   |   |     |    |    |    |     |    |   |
|----|---|---|-----|----|----|----|-----|----|---|
| 0  | 1 | 2 |     | 16 | 17 | 18 |     | 31 | 0 |
| 1  | 1 | 2 |     | 16 | 17 | 18 |     | 31 | 0 |
| 15 | 1 | 2 |     | 16 | 17 | 18 |     | 31 | 0 |

Wide 64-lane SIMD "map" unit

- ABS    - ABS    - ABS    . . . . . . . . . . . . .    - ABS

Flexible "reduce" step

Sum (Reduction)

# Example from H.264's motion estimation: Mapping Sum of Absolute Differences (SAD)

Reference frame pixels

pixels shift left

2D Regfile

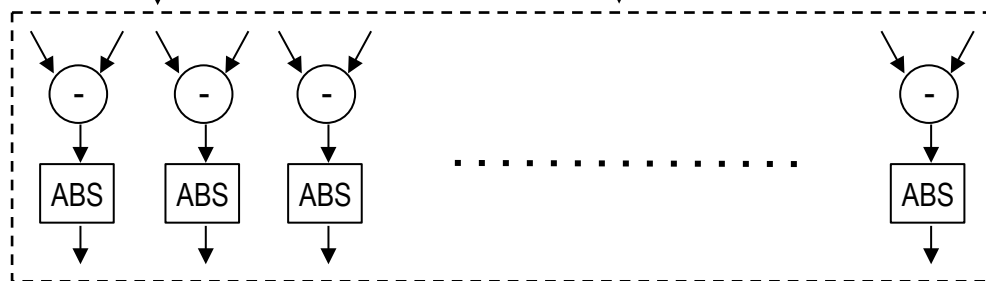2D shift Regfile

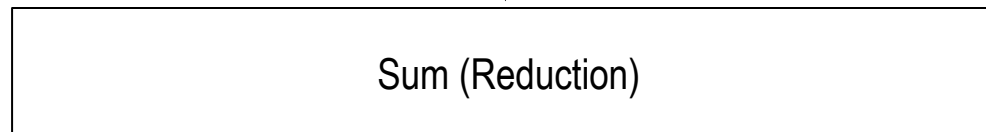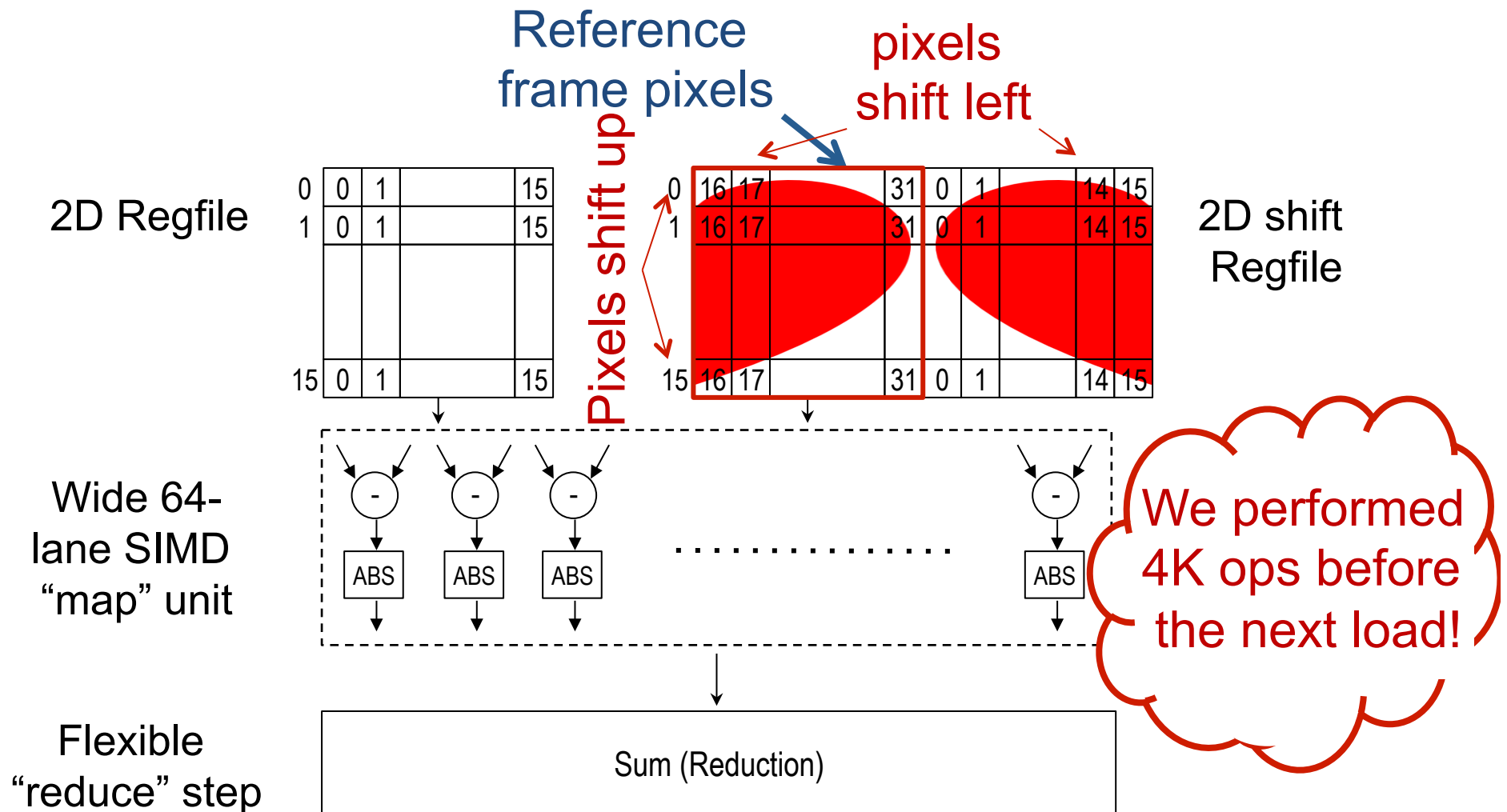Wide 64-lane SIMD "map" unit

Flexible "reduce" step

Sum (Reduction)

# Example from H.264's motion estimation: Mapping Sum of Absolute Differences (SAD)

Reference frame pixels

Pixels shift up

2D Regfile

2D shift Regfile

Wide 64-lane SIMD "map" unit

Flexible "reduce" step

Sum (Reduction)

# Example from H.264's motion estimation: Mapping Sum of Absolute Differences (SAD)

# Our *Convolution Engine* as implemented

2D Register

16x16

bit

2D / Co
Acce  s IF

ALU

**Get full implementation details in the paper:**

- **How we accomplished complex reduce steps using a "fused instructions graph"**

- **How we work on BIG stencils by combining multiple convolution slices**

- **The details of the ISA for the engine**

- **And so on, and so forth…**

16-way
SIMD

16 x
10bit lane

# Result #1:
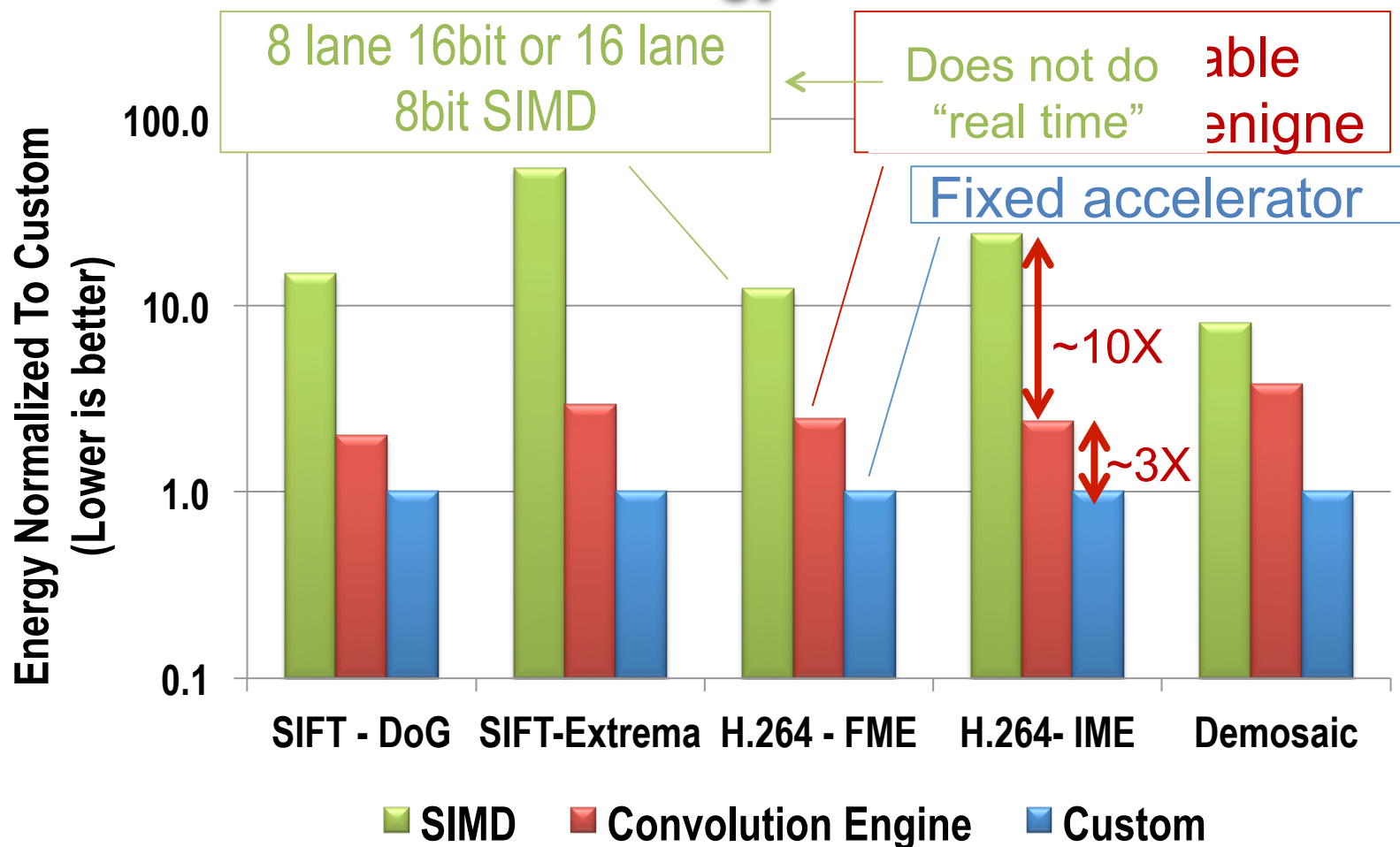## CE is user programmable in C!

```
SET_CE_OPS (CE_ABSDIFF, CE_ADD);   // Set map & reduce funcs to abs-diff and add
SET_CE_OPSIZE(16);                 // Set convolution size 16x16

// Load the 16x16 current macroblock into 2D coefficients register
for (int i=0; i<16; i++ {
    LD_COEFF_REG_128(curMBPtr, i);    // Load 16 pixels to row i of coefficient register
    curMBPtr += imgWidth;
}
// Load the first 32x16 current reference window into 2D input register
for (int i=0; i<16; i++ {
    LD_2D_REG_128(refPtr, 0, SHIFT_ENABLED);        // Load & shift-up 16 pixels to 2D Reg
    LD_2D_REG_128(refPtr+16, 1, SHIFT_DISABLED);    // Load next 16 pixels
    refPtr += imgWidth;
}
// Calculate one row of SAD output
for (int x = 0; x < 16; x++) {
    CONVOLVE_2D(ROTATE_LEFT, x);    // 16x16 2D convolution step and shift left
}
// Store 16 output SAD results
ST_OUT_REG_128(outPtr);
```

*IME code snippet*

## CE is 100X more energy efficient than RISC



8 lane 16bit or 16 lane 8bit SIMD

Does not do "real time"

Fixed accelerator

~10X

~3X

Energy Normalized To Custom (Lower is better)

100.0

10.0

1.0

0.1

SIFT - DoG    SIFT-Extrema    H.264 - FME    H.264- IME    Demosaic

■ SIMD    ■ Convolution Engine    ■ Custom

■ All variations were implemented as Tensilica extensions (TIE)

# Conclusions

- **There are classes of computations for which we can build efficient hardware, and we typically build them in ASIC**

- **Image and video are ubiquitous and represents one of those classes as their computation is convolution-like**

- **But when we restrict the domain, two orders of magnitude better programmable engines are also possible!**

- **Flexible specialized engines are <u>not</u> an oxymoron**
    - Flexible convolution engine improves power & performance by ~100X
    - Only 2-3X worse off than a dedicated (not flexible) accelerator

# THANK YOU FOR LISTENING!
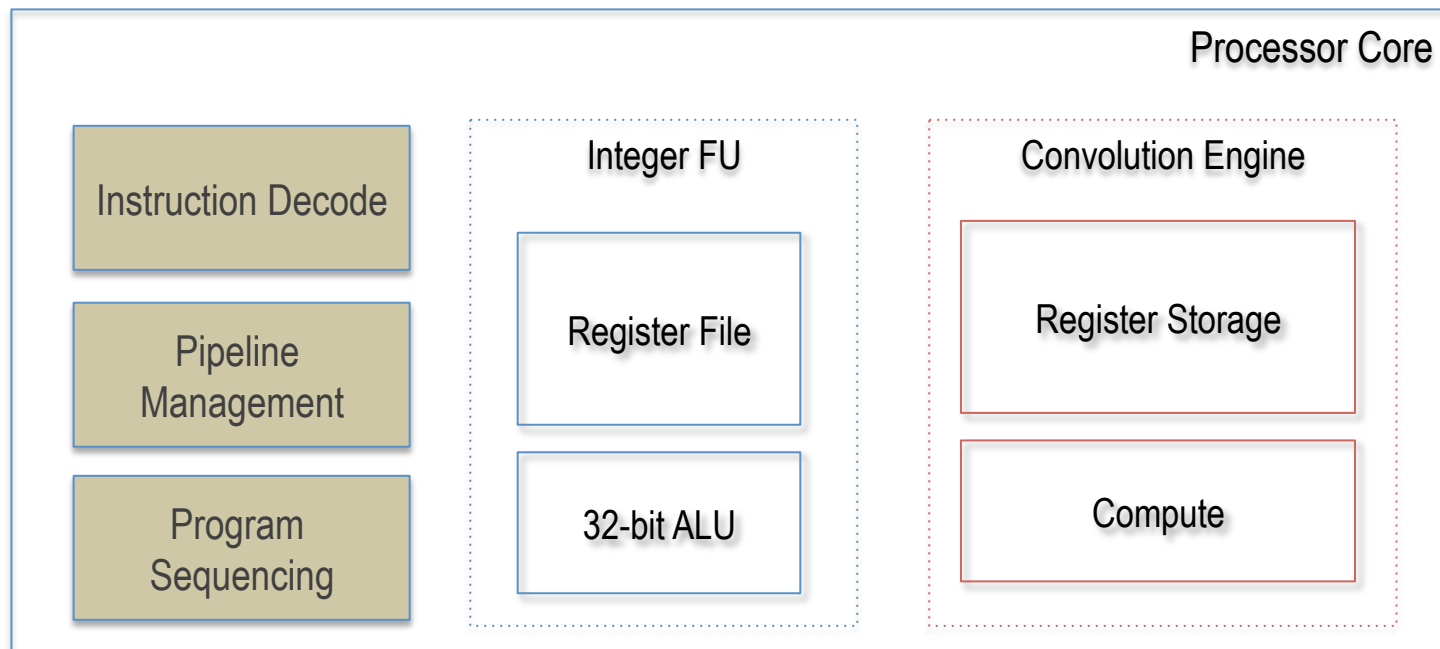
# BACKUP SLIDES…

# Energy dissipation in RISC machines

- **Let's do a breakdown of a typical RISC Instruction**

- **Keep in mind (at 45nm):**
  - Addition is ~0.1pJ for 8bits (ASIC) or ~0.5pJ for 32bits (RISC)
  - Multiplication is ~0.2pJ for 8bits (ASIC) or ~3.1pJ for 32bits (RISC)
  - But a single RISC instruction is 70pJ

- **Need to see where the overhead is, and how we can mitigate it**

# Processor Integration

- **Specialized Functional Unit**
  - Adds about 30 instructions to the processor ISA
  - The execution flow is controlled by the processor



Processor Core

Instruction Decode

Pipeline Management

Program Sequencing

Integer FU

Register File

32-bit ALU

Convolution Engine

Register Storage

Compute

# Evaluating the Convolution Engine
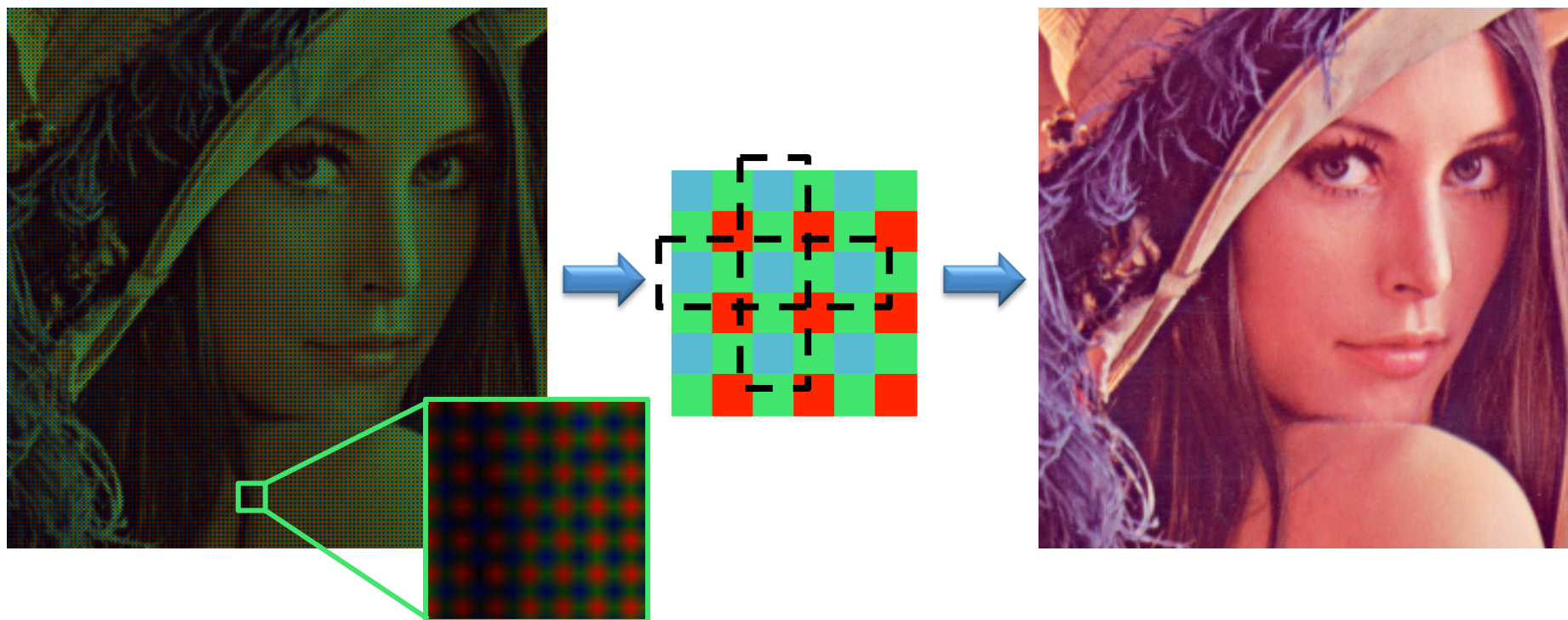
- **Applications**
- **SIFT Feature extraction**
  - Often a basic step for computational photography algorithms
    - HDR Imaging
    - Panorama stitching
    - Smart zoom / Super resolution
    - Multi-frame noise reduction
    - Synthetic aperture
    - Augmented reality
    - Flash – No-Flash photography
    - Video de-shake
    - ……

- **H.264 encoder**
  - Every video system has one

- **De-mosaic:**
  - Adaptive color plane interpolation (ACPI)*: image gradients followed by a three-tap filter in the direction of smallest gradient.



\* Y. Cheng et. al. An adaptive color plane interpolation method based on edge detection. Journal of Electronics (China), 2007.