# The Goal:

## Sustained ExaFLOPS on Problems of Interest

## …

## at reasonable cost



NVIDIA.

# The End of Historic Scaling



Transistors (thousands)

Single-thread Performance (SpecINT)

Frequency (MHz)

Typical Power (Watts)

Number of Cores

Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten
Dotted line extrapolations by C. Moore

CORAL
150-300PF (5-10x)
11MW (1.1x)
14-27 GFLOPs/W (7-14x)
Lots of Threads

**2023**

1,000PF (50x)
72,000HCNs (4x)
20MW (2x)
50 GFLOPs/W (25x)
~$10^{10}$ Threads (1000x)

**2017**

**2013**

**You Are Here**

20PF
18,000 GPUs
10MW
2 GFLOPs/W
~$10^{7}$ Threads
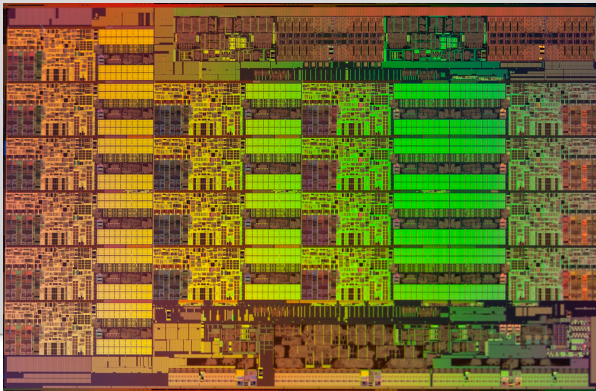
# HETEROGENEOUS NODE



5 NVIDIA.

How do we get to 50GFlops/Watt?

NVIDIA.

# Start with an energy-efficient architecture

# CPU
## 130 pJ/flop (Vector SP)

Optimized for Latency
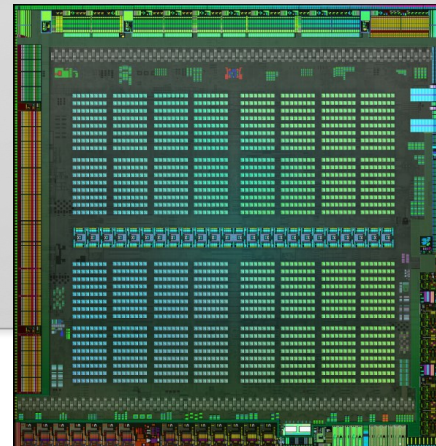
Deep Cache Hierarchy



Haswell
22 nm

# GPU
## 30 pJ/flop (SP)

Optimized for Throughput

Explicit Management
of On-chip Memory



Maxwell
28 nm

NVIDIA.

## CPU
## 2 nJ/flop (Scalar SP)

Optimized for Latency

Deep Cache Hierarchy



Haswell
22 nm

## GPU
## 30 pJ/flop (SP)

Optimized for Throughput

Explicit Management
of On-chip Memory



Maxwell
28 nm

NVIDIA.

# HOW IS POWER SPENT IN A CPU?

## In-order Embedded

Data Supply 17%
Clock + Control Logic 24%
ALU 6%
Register File 11%
Instruction Supply 42%

Dally [2008] (Embedded in-order CPU)

## OOO Hi-perf

ALU 4%
Supply 5%
RF 14%
Issue 11%
Clock + Pins 45%
Rename 10%
Fetch 11%

Natarajan [2003] (Alpha 21264)

NVIDIA.

# Latency-Optimized Core (LOC)



# Throughput-Optimized Core (TOC)

How do we continue to scale energy efficiency

…in a world where technology scaling is diminished?

**NVIDIA.**

Do Less Work

Eliminate waste and redundancy

Move fewer bits

Move data more efficiently

NVIDIA.

# DO LESS WORK
## Mixed Precision Arithmetic

double-precision

| 1 | 11 | 52 |
|---|----|----|

5x precision bits
60x range

single-precision

| 1 | 8 | 23 |
|---|---|----|

Only use as much precision as you need

Exploit mix of representations

Scaled arithmetic

4x throughput
4x bandwidth
4x capacity
< ¼ energy/op

half-precision

| 1 | 5 | 10 |
|---|---|----|

# ELIMINATE WASTE
## Temporal SIMT

**Spatial SIMT (current GPUs)**

**Pure Temporal SIMT**

32-wide datapath

1-wide

thread 0 ... thread 31

time — 1 cyc

ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld ld
ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml ml
ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad ad
st st st st st st st st st st st st st st st st st st st st st st st st st st st st st st st st

*1 warp instruction = 32 threads*

time — 1 cyc

(threads)

ld 0
ld 1
ld 2
ld 3
ld 4
ld 5
ld 6
ld 7
ld 8
ld 9
ld 10

15  NVIDIA.

# ELIMINATE WASTE

## Temporal SIMT

**32-wide** **(41%)**

**4-wide** **(65%)**

**1-wide** **(100%)**

Increase efficiency on divergent code

NVIDIA.

# ELIMINATE WASTE

## Variable Warp Sizing

scheduler

Schedule several warps, *different PCs*

scheduler

Gang schedule warps with *same PC*

Time

Small warps

+ Improved perf for divergent code

+ Better SIMD utilization

Emulate wide warp HW

+ Wider converged execution

+ Memory locality/convergence

+ Reduced power (frontend)

*Rogers [ISCA 2015]*

17 NVIDIA.

# ELIMINATE REDUNDANCY

## Scalarization

SIMT Execution

| scalar op | LD R2←<A> | LD R2←<A> | LD R2←<A> | LD R2←<A> |
| vector load | LD R3←R2, 1 | LD R3←R2, 2 | LD R3←R2,3 | LD R3←R2, 4 |
| vector op | ADD R4←R3, 2 | ADD R4←R3, 2 | ADD R4←R3, 2 | ADD R4←R3, 2 |

Scalarized SIMT Execution

| LD SR2←<A> | | | | |
| VLD SR3←SR2, 1 | | | | |
| | ADD R4←SR3, 2 | ADD R4←SR3, 2 | ADD R4←SR3, 2 | ADD R4←SR3, 2 |

*Lee [CGO 2013]*

NVIDIA.

# MOVE FEWER BITS

## Register File Cache (RFC)

Small multi-ported register file

Capture locality of commonly used operands
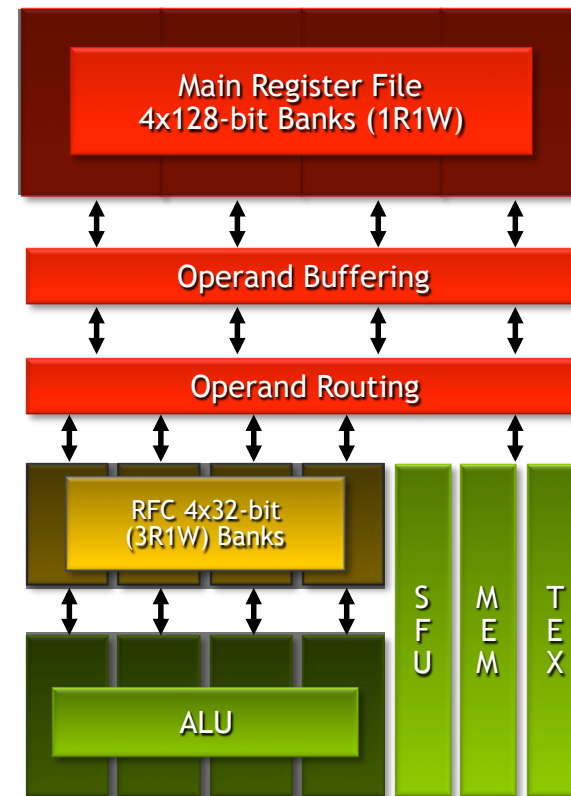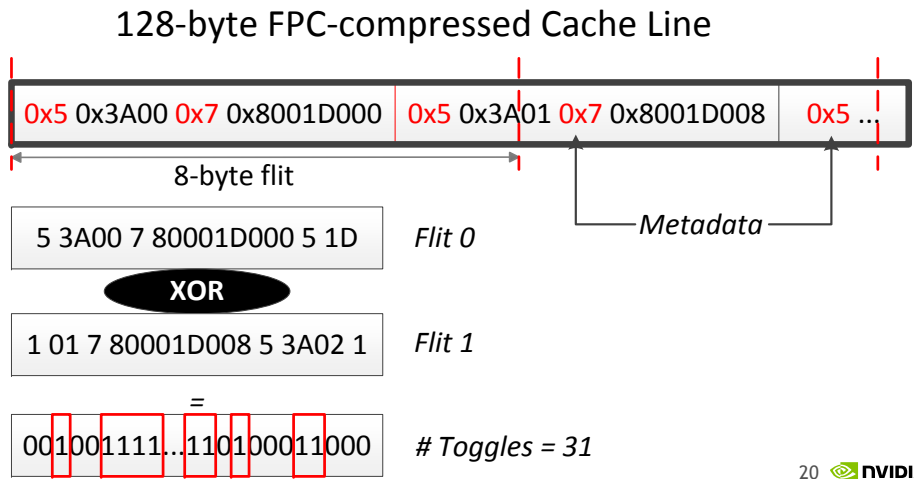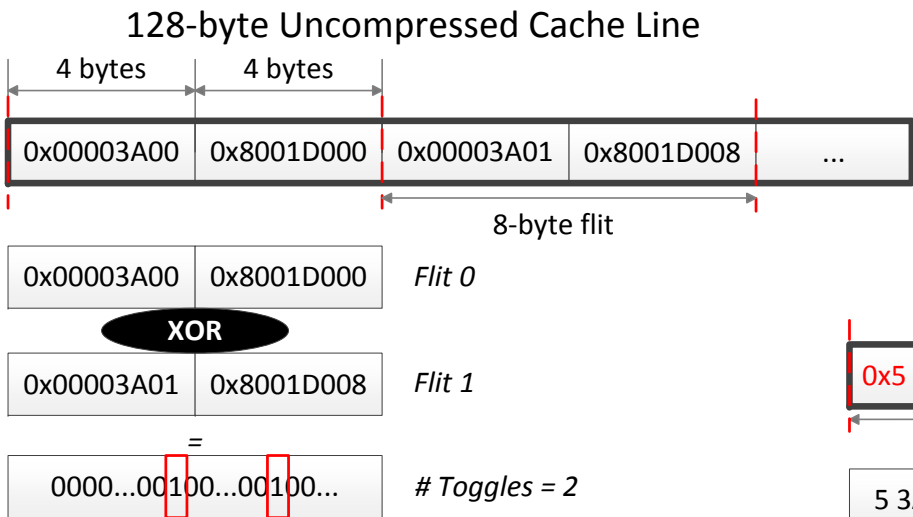
Can reduce RF energy by 50%

*Gebhart [ISCA 2011]*



Main Register File
4x128-bit Banks (1R1W)

Operand Buffering

Operand Routing

RFC 4x32-bit
(3R1W) Banks

SFU  MEM  TEX

ALU

# MOVE DATA MORE EFFICIENTLY

## Toggle-aware Compression

128-byte Uncompressed Cache Line

| 4 bytes | 4 bytes | | | |
|---|---|---|---|---|
| 0x00003A00 | 0x8001D000 | 0x00003A01 | 0x8001D008 | ... |

8-byte flit

| 0x00003A00 | 0x8001D000 | *Flit 0* |
|---|---|---|

**XOR**

| 0x00003A01 | 0x8001D008 | *Flit 1* |
|---|---|---|

=

| 0000...00100...00100... | *# Toggles = 2* |
|---|---|

*Pekhimenko [HPCA 2016]*

**Compression can increase power consumption**

**Goal: reduce bus toggling**

128-byte FPC-compressed Cache Line

| 0x5 0x3A00 0x7 0x8001D000 | 0x5 0x3A01 0x7 0x8001D008 | 0x5 ... |
|---|---|---|

8-byte flit

Metadata

| 5 3A00 7 80001D000 5 1D | *Flit 0* |
|---|---|

**XOR**

| 1 01 7 80001D008 5 3A02 1 | *Flit 1* |
|---|---|

=

| 001001111...110100011000 | *# Toggles = 31* |
|---|---|

# MINIMIZE DATA MOVEMENT
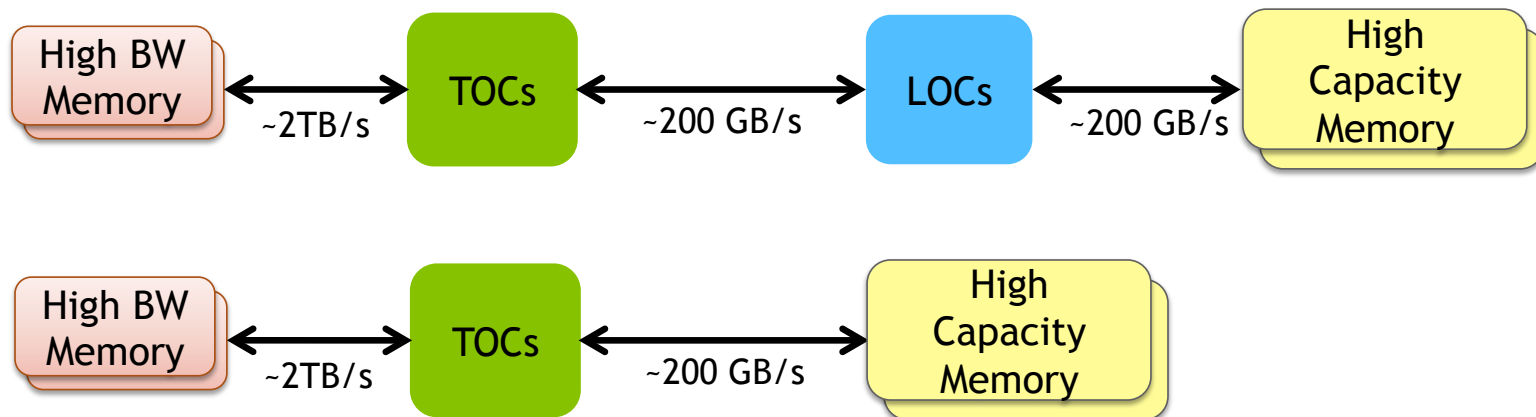Packaging

High-bandwidth on-package memory

Reduces distance

Increases bandwidth

Offers opportunity to optimize signaling circuits

# MINIMIZE DATA MOVEMENT

## Heterogeneous DRAM Architectures

| High BW Memory | ←→ | TOCs | ←→ | LOCs | ←→ | High Capacity Memory |
|---|---|---|---|---|---|---|
| | ~2TB/s | | ~200 GB/s | | ~200 GB/s | |

| High BW Memory | ←→ | TOCs | ←→ | High Capacity Memory |
|---|---|---|---|---|
| | ~2TB/s | | ~200 GB/s | |

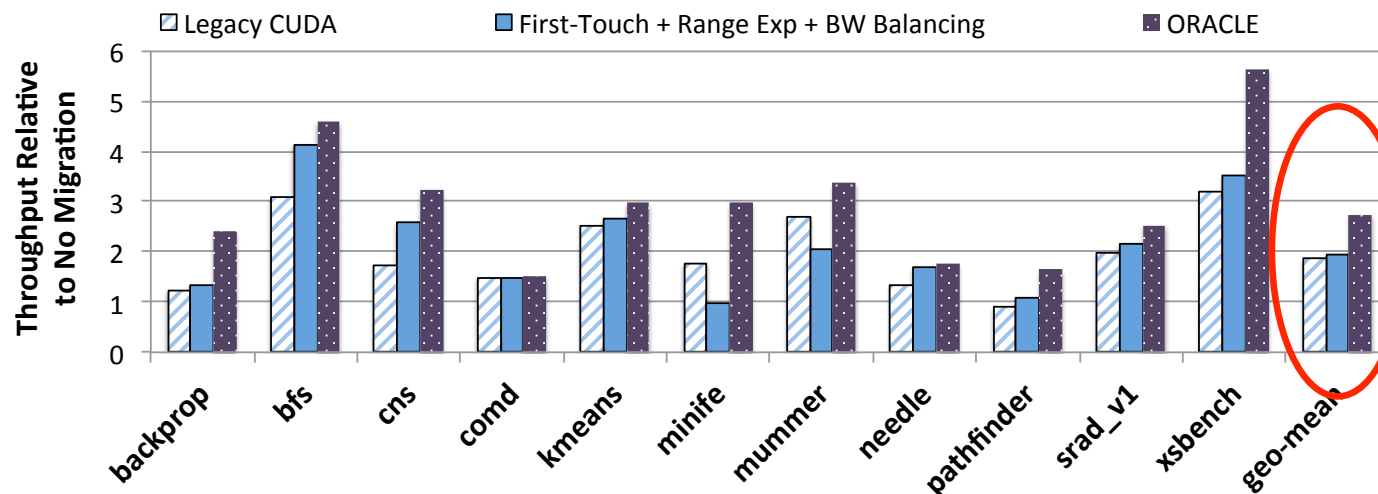Challenges

    Exploiting all available bandwidth

    Maximizing locality for frequently accessed data

# MINIMIZE DATA MOVEMENT

## Software-managed Caching with On-Package Memory

Strategies

Aggressively migrate pages upon First-Touch to GDDR memory

Pre-fetch neighbors of touched pages to reduce TLB shootdowns

Throttle page migrations when nearing peak BW



Competitive with manual memory copy

Close to "perfect" prefetch

*Agarwal [HPCA 2015]*

23

# MINIMIZE DATA MOVEMENT

## Hardware Managed DRAM Cache

| Tag | Set Index | Offset |
|---|---|---|

(512M lines) 29-bits     5-bits

Tag overhead: hundreds of MB

    Alloy tag and data in same DRAM row (Micro12)

tag* (2-byte)    data (64-byte)

Cache organization: optimize for bandwidth

    Direct mapped, consecutive sets in same row

16-byte DRAM bus

ROW BUFFER

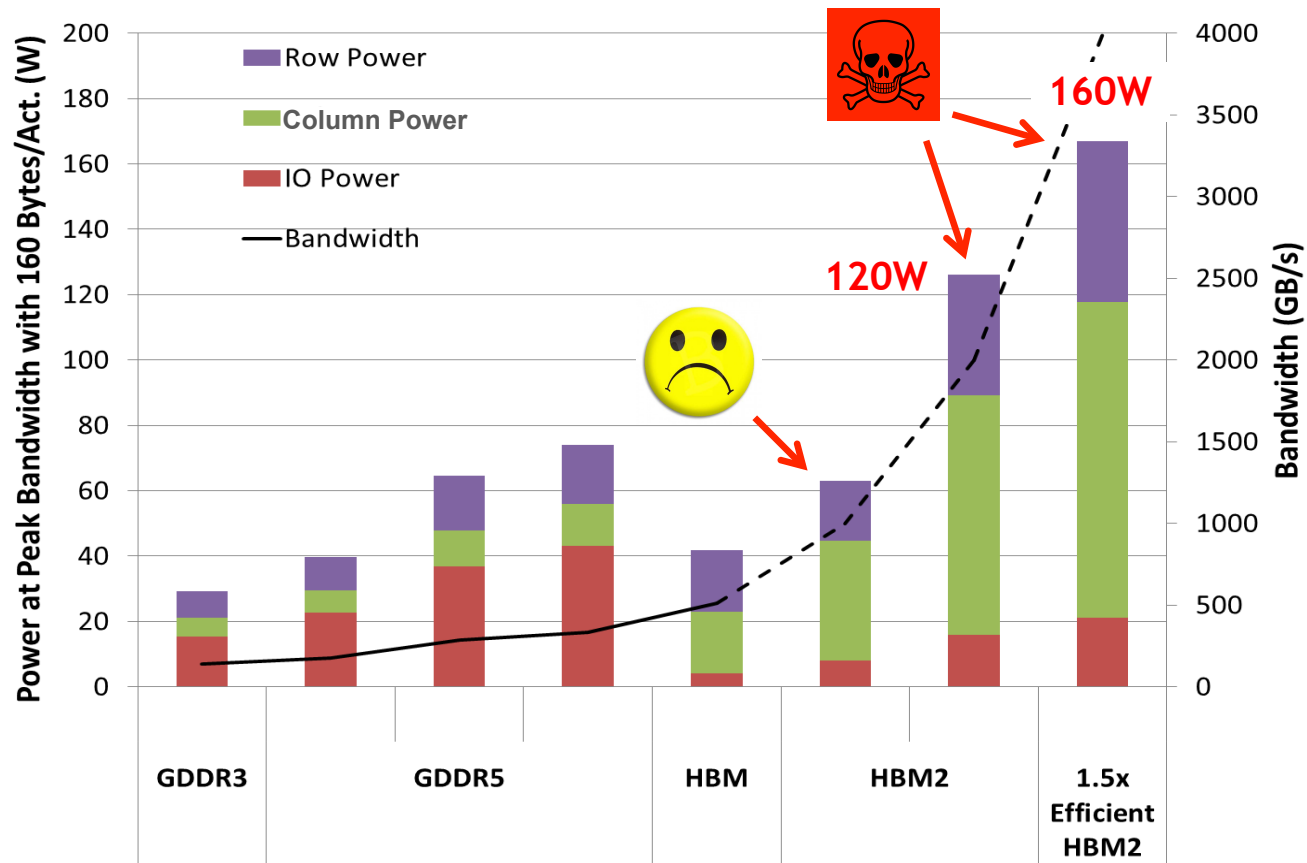Only medium of accessing DRAM array

1GB DRAM

Results

    Fine-grained transfers good for lower locality apps

    Can eliminate some page migration overheads

# LOOMING MEMORY POWER CRISIS

# SUMMARY

Do Less Work

Eliminate waste and redundancy

Move fewer bits

Move data more efficiently

NVIDIA.