# GPU Computing

# The GPU Advantage

# The GPU Advantage

## A Tale of Two Machines

# Tianhe-1A
## at NSC Tianjin

# Tianhe-1A
## at NSC Tianjin
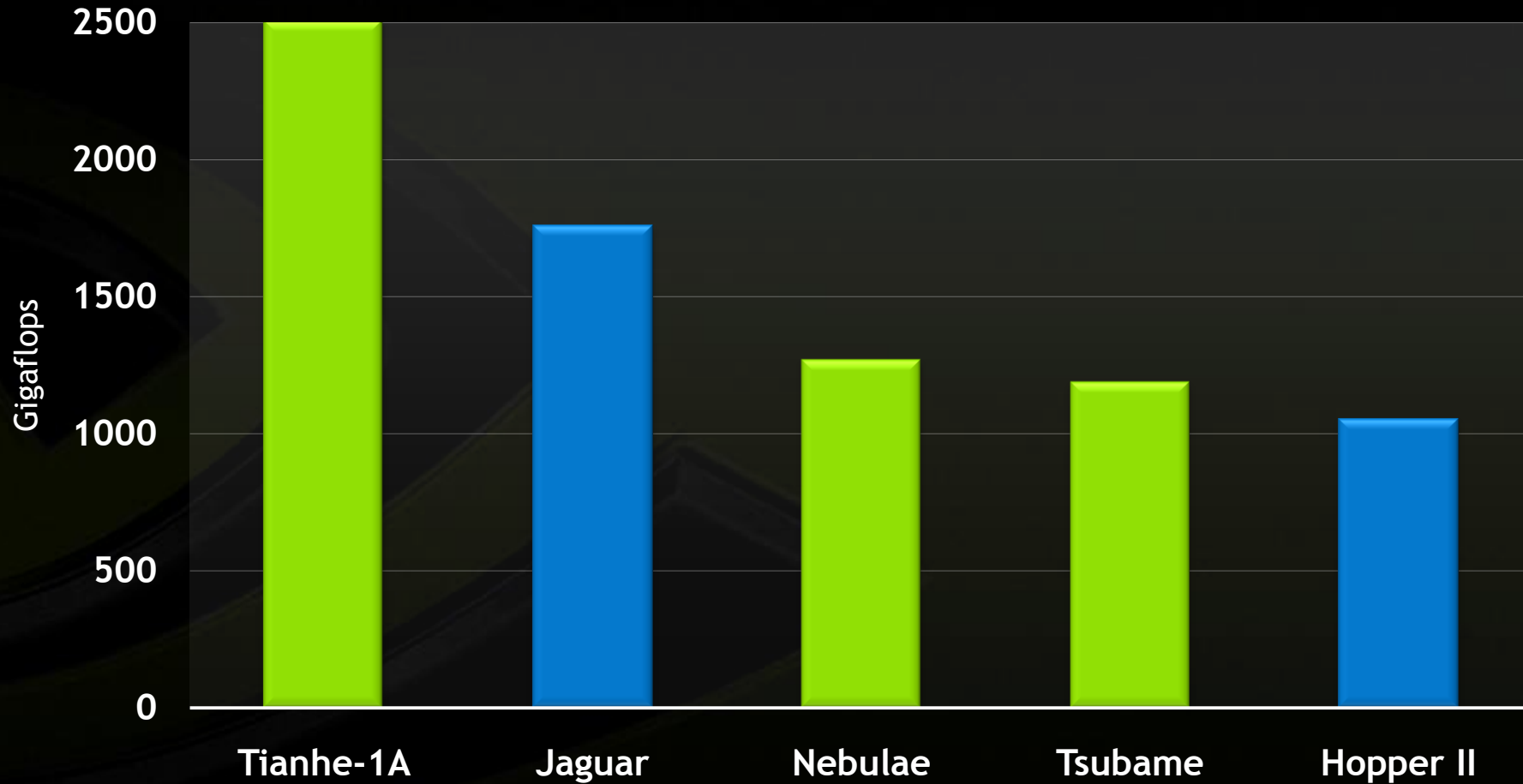
- The World's Fastest Supercomputer
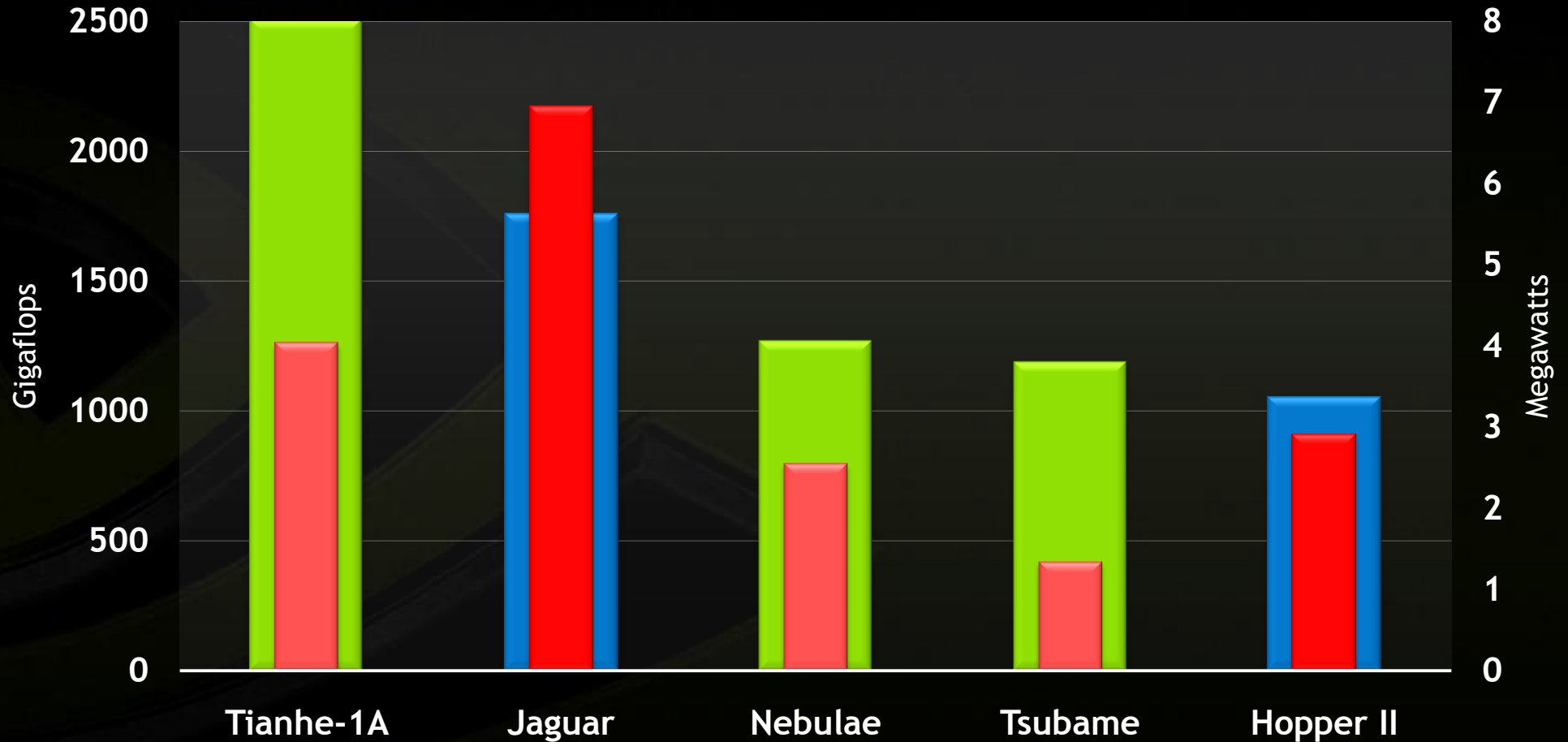- 2.507 Petaflop
- 7168 Tesla M2050 GPUs

Tesla M2050 GPUs
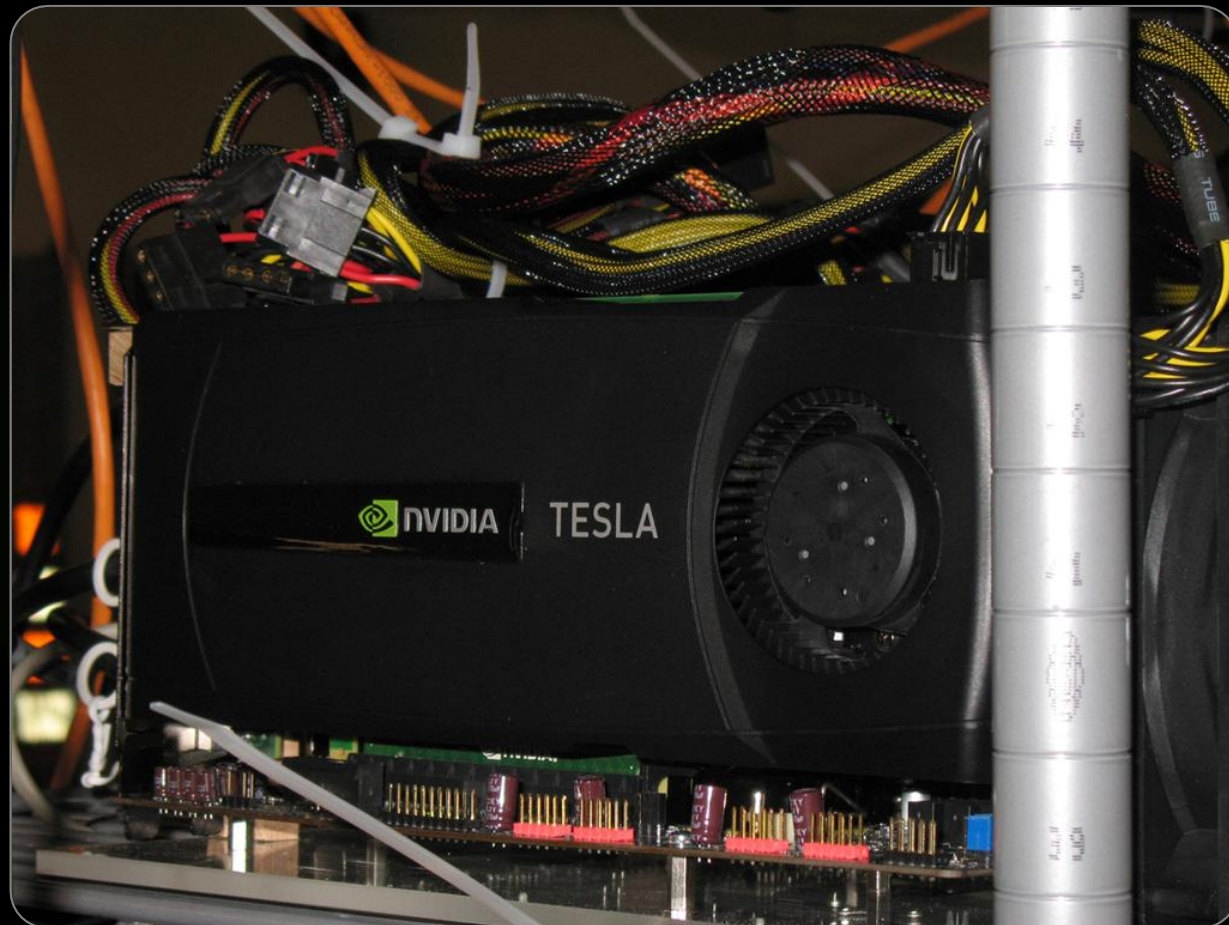
# 3 of Top5 Supercomputers

# NVIDIA/NCSA
## Green 500 Entry

# NVIDIA/NCSA

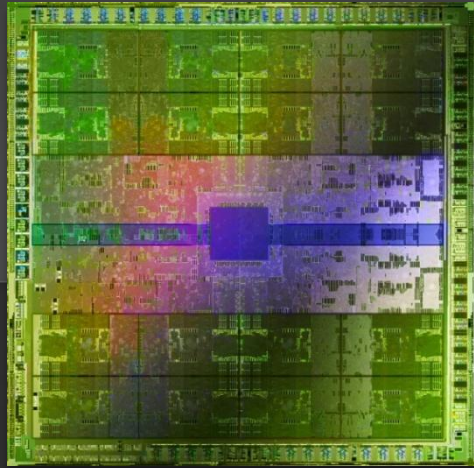Green 500 Entry

# NVIDIA/NCSA Green 500 Entry

- **128 nodes, each with:**
  - 1x Core i3 530 (2 cores, 2.93 GHz => 23.4 GFLOP peak)
  - 1x Tesla C2050 (14 cores, 1.15 GHz => 515.2 GFLOP peak)
  - 4x QDR Infiniband
  - 4 GB DRAM
- **Theoretical Peak Perf: 68.95 TF**
- **Footprint: ~20 ft^2 => 3.45 TF/ft^2**
- **Cost: $500K (street price) => 137.9 MF/$**
- **Linpack: 33.62 TF, 36.0 kW => 934 MF/W**

# The GPU Advantage
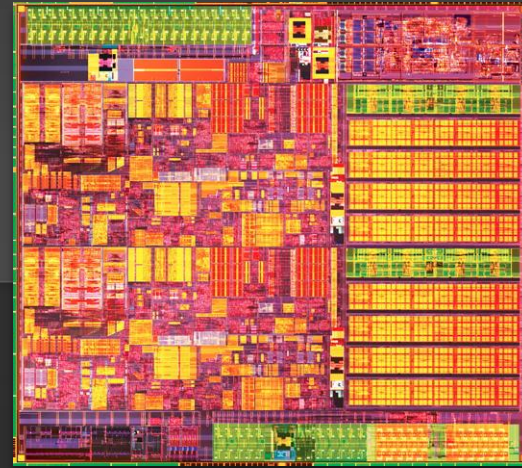
Efficiency and Programmability

# GPU
## 200pJ/Instruction
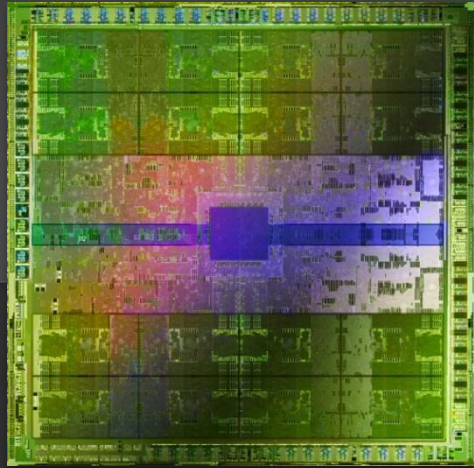
# CPU
## 2nJ/Instruction

# GPU

## 200pJ/Instruction
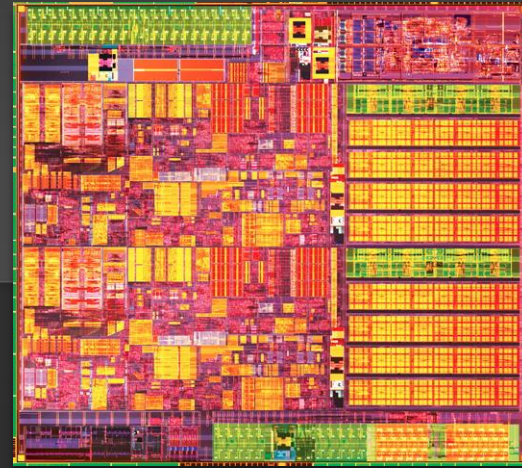
Optimized for Throughput

Explicit Management
of On-chip Memory

# CPU

## 2nJ/Instruction

Optimized for Latency

Caches

# CUDA GPU Roadmap

**Maxwell**

**Kepler**

**Fermi**

**Tesla**

DP GFLOPS per Watt

16

14

12

10

8

6

4

2

2007          2009          2011          2013

# The GPU Advantage

Efficiency and Programmability

# The GPU Advantage

CUDA Enables Programmability

# CUDA C: C with a Few Keywords

```c
void saxpy_serial(int n, float a, float *x, float *y)
{
    for (int i = 0; i < n; ++i)
        y[i] = a*x[i] + y[i];
}
// Invoke serial SAXPY kernel
saxpy_serial(n, 2.0, x, y);
```
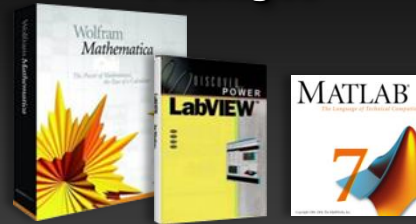
*Standard C Code*

```c
__global__ void saxpy_parallel(int n, float a, float *x, float *y)
{
    int i = blockIdx.x*blockDim.x + threadIdx.x;
    if (i < n)  y[i] = a*x[i] + y[i];
}
// Invoke parallel SAXPY kernel with 256 threads/block
int nblocks = (n + 255) / 256;
saxpy_parallel<<<nblocks, 256>>>(n, 2.0, x, y);
```

*CUDA C Code*

**Libraries**

$$\oint \mathbf{E} \cdot d\mathbf{A} = \frac{q_{enc}}{\varepsilon_0}$$

$$\oint \mathbf{B} \cdot d\mathbf{A} = 0$$

$$\oint \mathbf{E} \cdot d\mathbf{s} = -\frac{d\Phi_B}{dt}$$

$$\oint \mathbf{B} \cdot d\mathbf{s} = \mu_0 \varepsilon_0 \frac{d\Phi_E}{dt} + \mu_0 i_{enc}$$

**Mathematical Packages**

**Research & Education**

**Consultants, Training & Certification**

**DESIGNED FOR NVIDIA CUDA**

**GPU Computing Ecosystem**

**Integrated Development Environment**

Parallel Nsight for MS Visual Studio

**Tools & Partners**

**Languages & API's**

CUDA C/C++

**All Major Platforms**

# GPU Computing Today

By the Numbers:

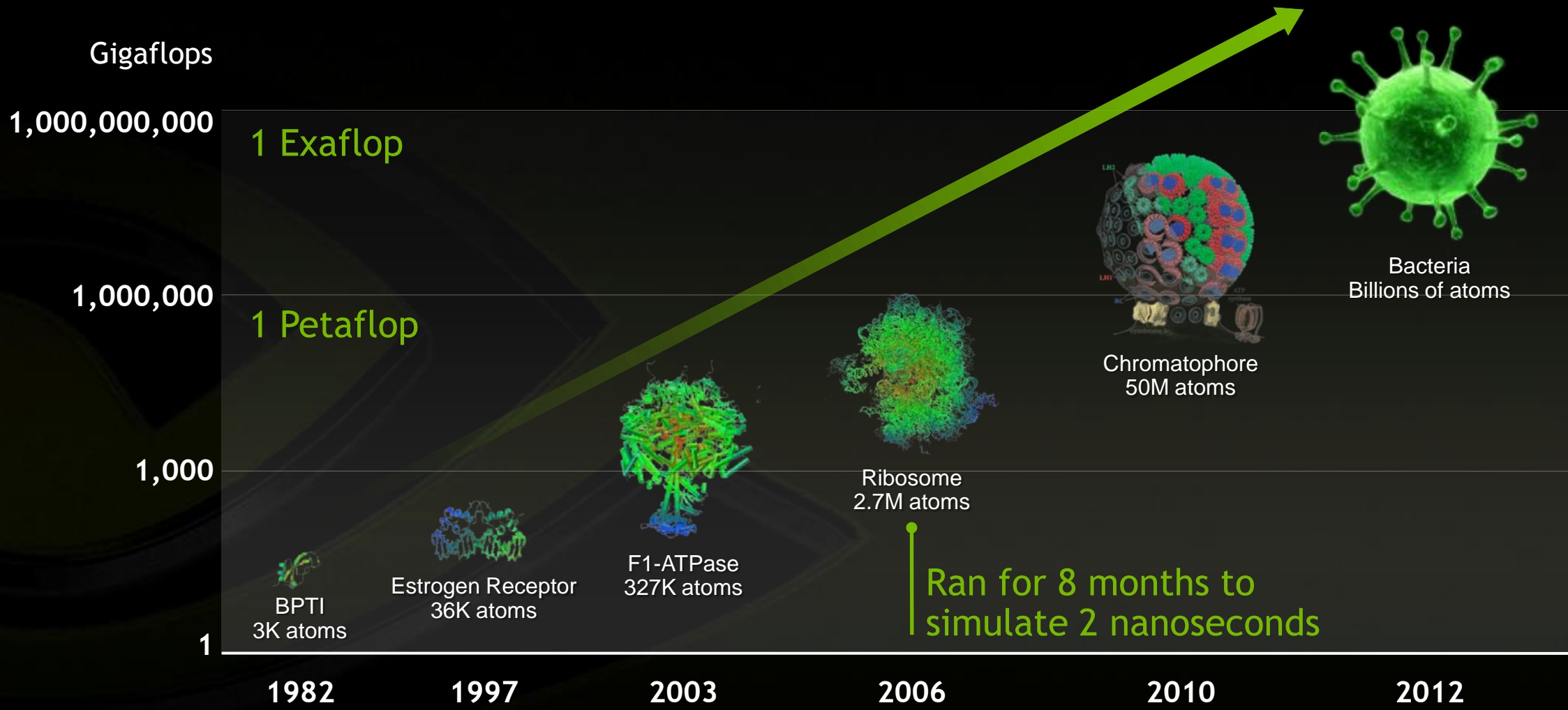| | |
|---|---|
| **200 Million** | CUDA Capable GPUs |
| **600,000** | CUDA Toolkit Downloads |
| **100,000** | Active GPU Computing Developers |
| **8,000** | Members in Parallel Nsight Developer Program |
| **362** | Universities Teaching CUDA Worldwide |
| **11** | CUDA Centers of Excellence Worldwide |

# To ExaScale and Beyond

# Science Needs 1000x More Computing



Gigaflops

1,000,000,000 — 1 Exaflop

1,000,000 — 1 Petaflop

1,000

1

BPTI
3K atoms

Estrogen Receptor
36K atoms

F1-ATPase
327K atoms

Ribosome
2.7M atoms

Ran for 8 months to simulate 2 nanoseconds

Chromatophore
50M atoms

Bacteria
Billions of atoms

1982    1997    2003    2006    2010    2012

# DARPA Study Identifies Four Challenges for ExaScale Computing



ExaScale Computing Study:
Technology Challenges in
Achieving Exascale Systems

Peter Kogge, Editor & Study Lead
Keren Bergman
Shekhar Borkar
Dan Campbell
William Carlson
William Dally
Monty Denneau
Paul Franzon
William Harrod
Kerry Hill
Jon Hiller
Sherman Karp
Stephen Keckler
Dean Klein
Robert Lucas
Mark Richards
Al Scarpelli
Steven Scott
Allan Snavely
Thomas Sterling
R. Stanley Williams
Katherine Yelick

September 28, 2008

This work was sponsored by DARPA IPTO in the ExaScale Computing Study with Dr. William Harrod as Program Manager; AFRL contract number FA8650-07-C-7724. This report is published in the interest of scientific and technical information exchange and its publication does not constitute the Government's approval or disapproval of its ideas or findings

NOTICE

Using Government drawings, specifications, or other data included in this document for any purpose other than Government procurement does not in any way obligate the U.S. Government. The fact that the Government formulated or supplied the drawings, specifications, or other data does not license the holder or any other person or corporation; or convey any rights or permission to manufacture, use, or sell any patented invention that may relate to them.

APPROVED FOR PUBLIC RELEASE, DISTRIBUTION UNLIMITED.
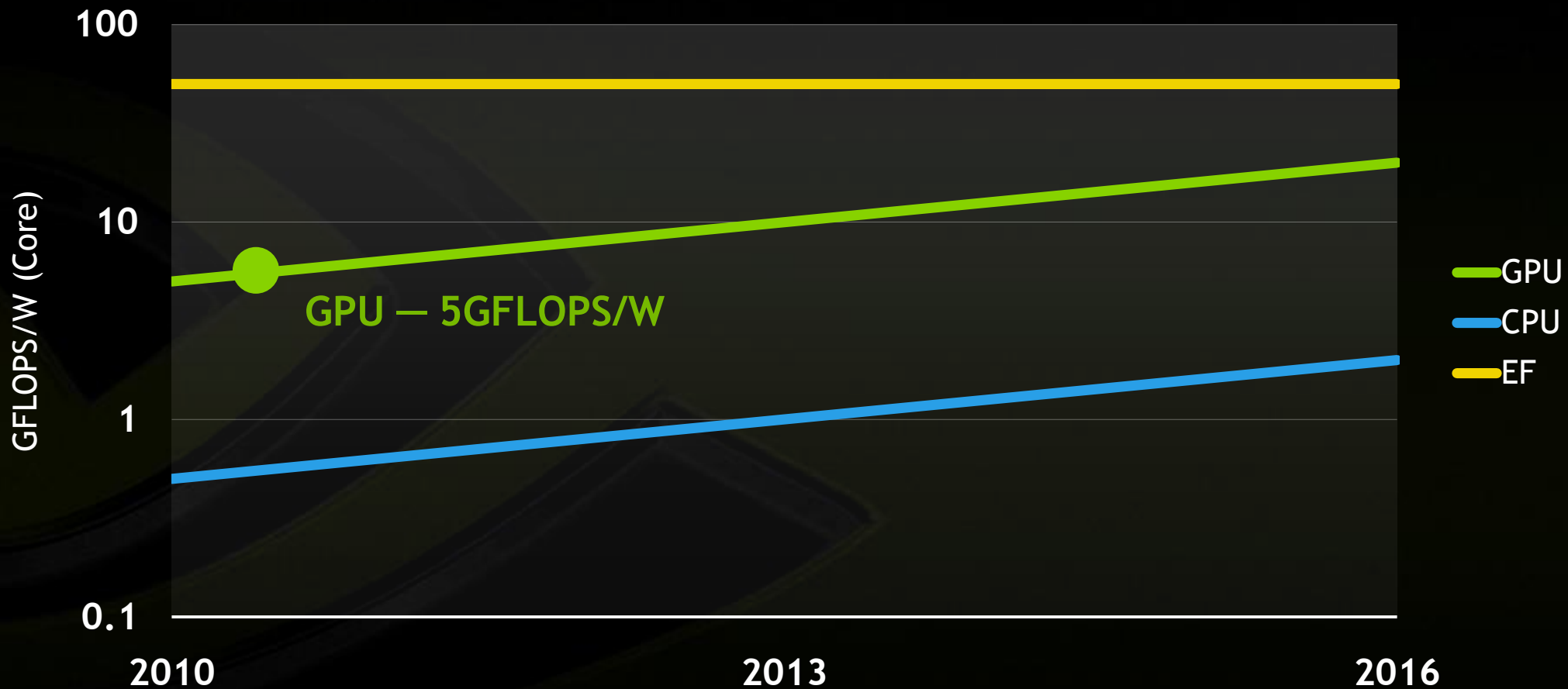
Available at
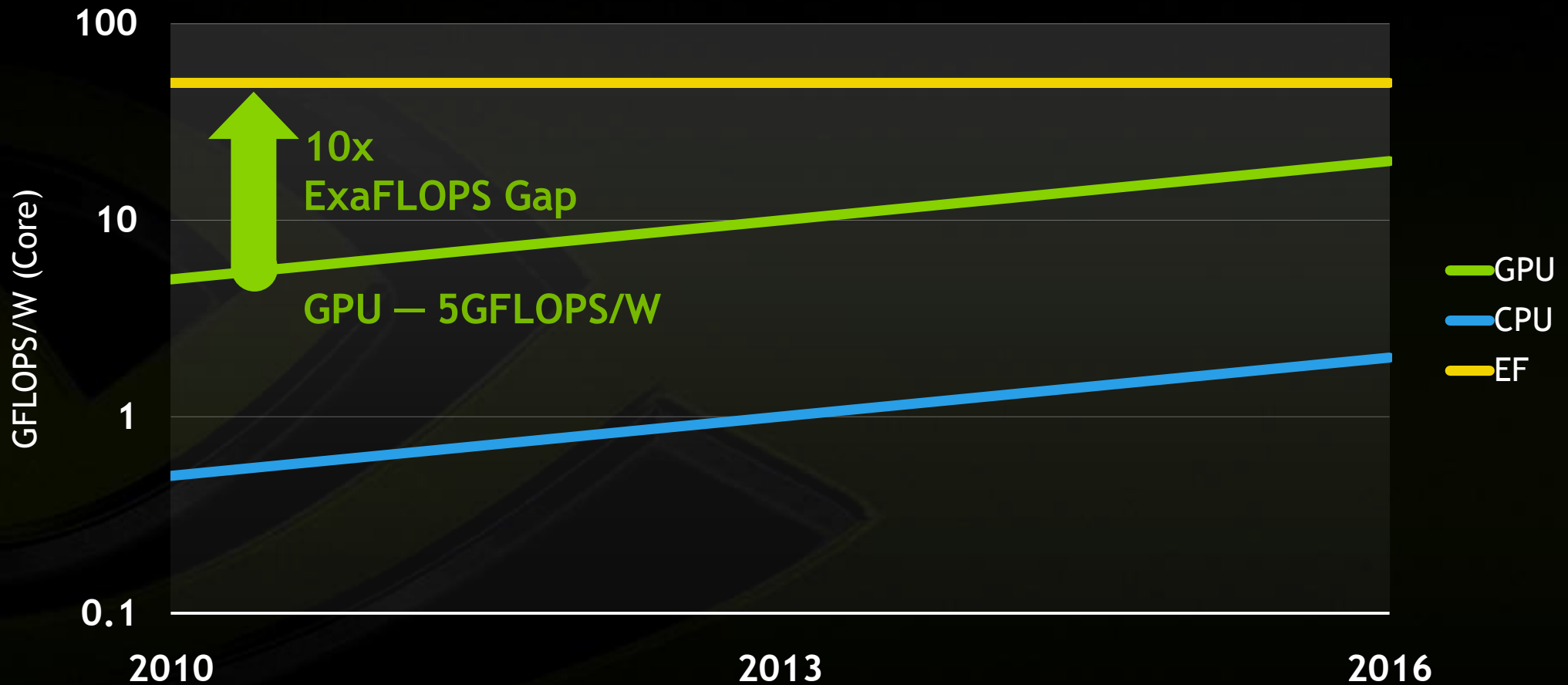www.darpa.mil/ipto/personnel/docs/ExaScale_Study_Initial.pdf

## Report published September 28, 2008:

- **Four Major Challenges**
  - Energy and Power challenge
  - Memory and Storage challenge
  - Concurrency and Locality challenge
  - Resiliency challenge
- **Number one issue is *power***
  - Extrapolations of current architectures and technology indicate over 100MW for an Exaflop!
  - Power also constrains what we can put on a chip

# Power is THE Problem

# Power is THE Problem

A GPU is the Solution
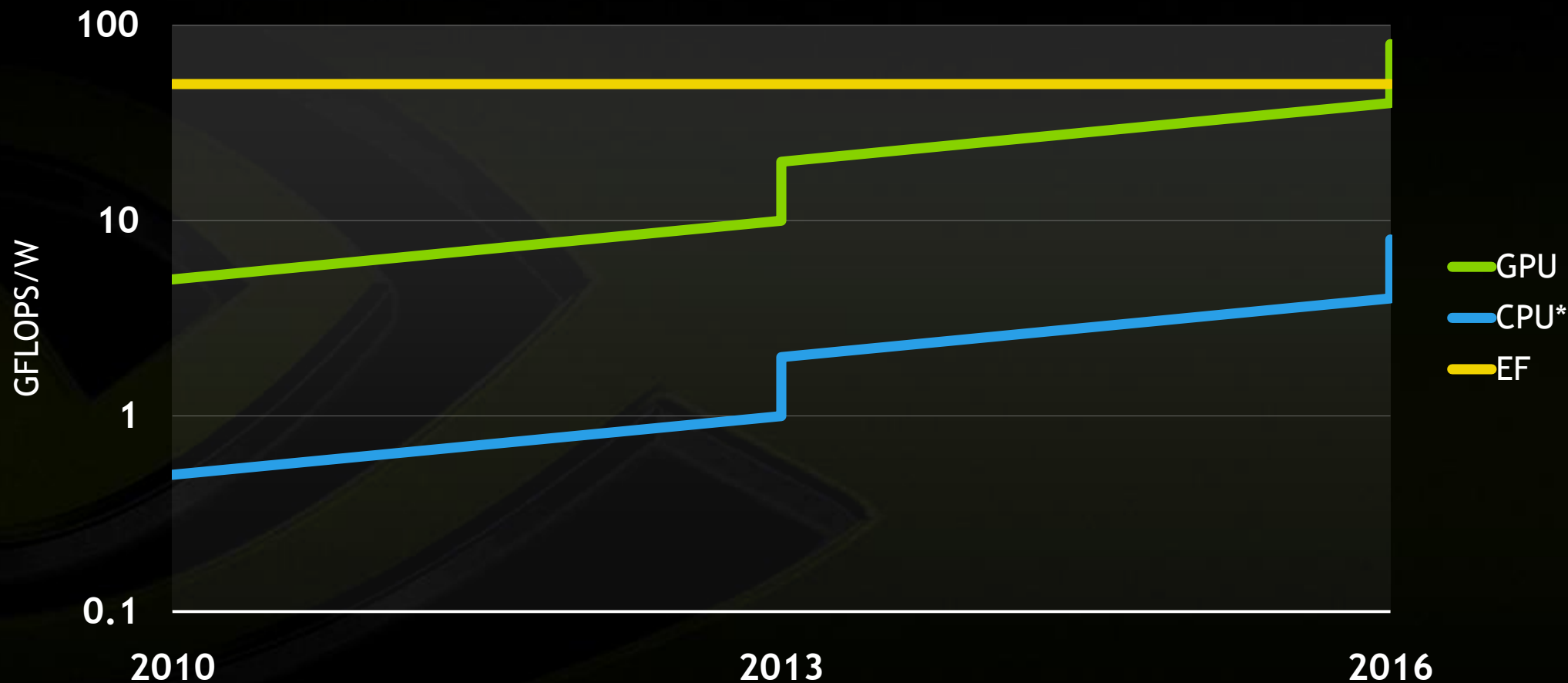
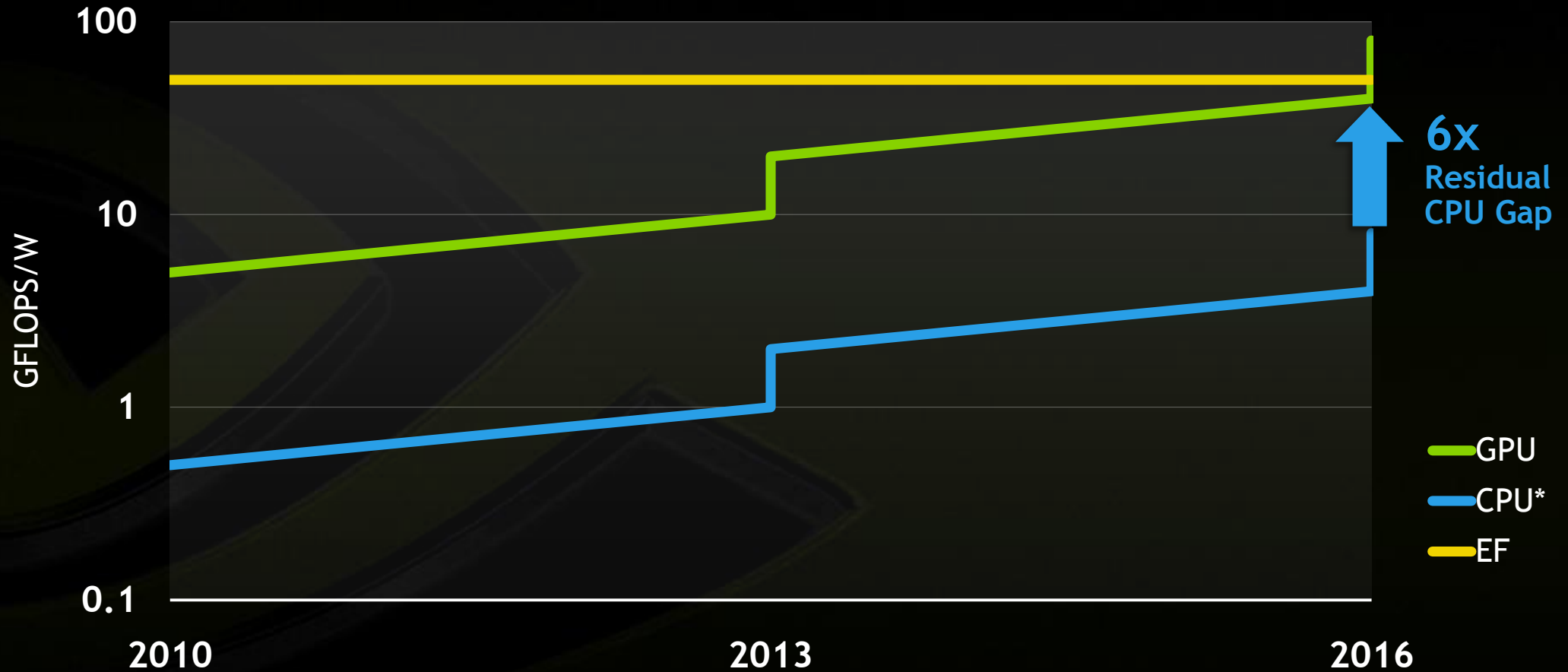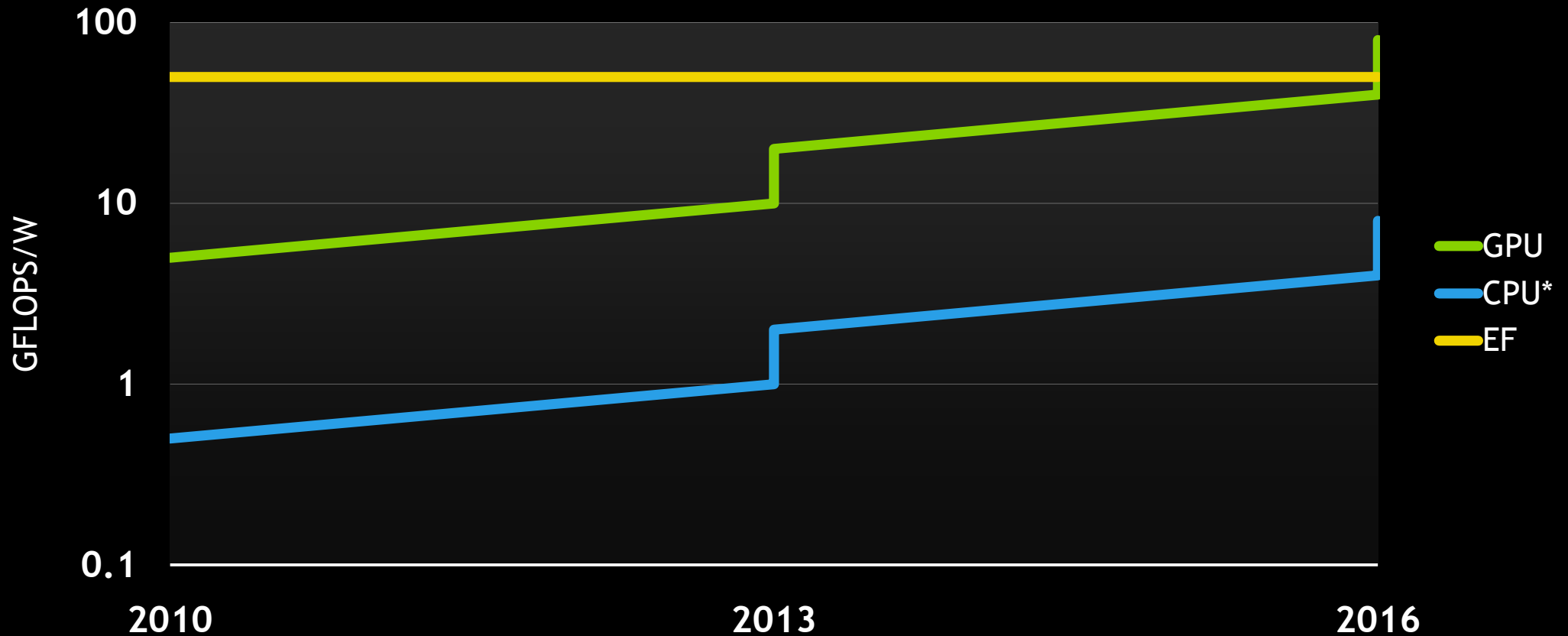# ExaFLOPS at 20MW = 50GFLOPS/W



GPU — 5GFLOPS/W

- GPU
- CPU
- EF

# GPUs Close the Gap with Process and Architecture

GPUs Close the Gap
with Process and Architecture

# GPUs Close the Gap
# With CPUs, a Gap Remains



GFLOPS/W

100

10

1

0.1

2010          2013          2016

**6x**
Residual
CPU Gap

— GPU
— CPU*
— EF

GPUs Close the Gap With CPUs, a Gap Remains

Heterogeneous Computing is Required to get to ExaScale

# Echelon Team

# System Sketch

# Execution Model

# The High Cost of Data Movement
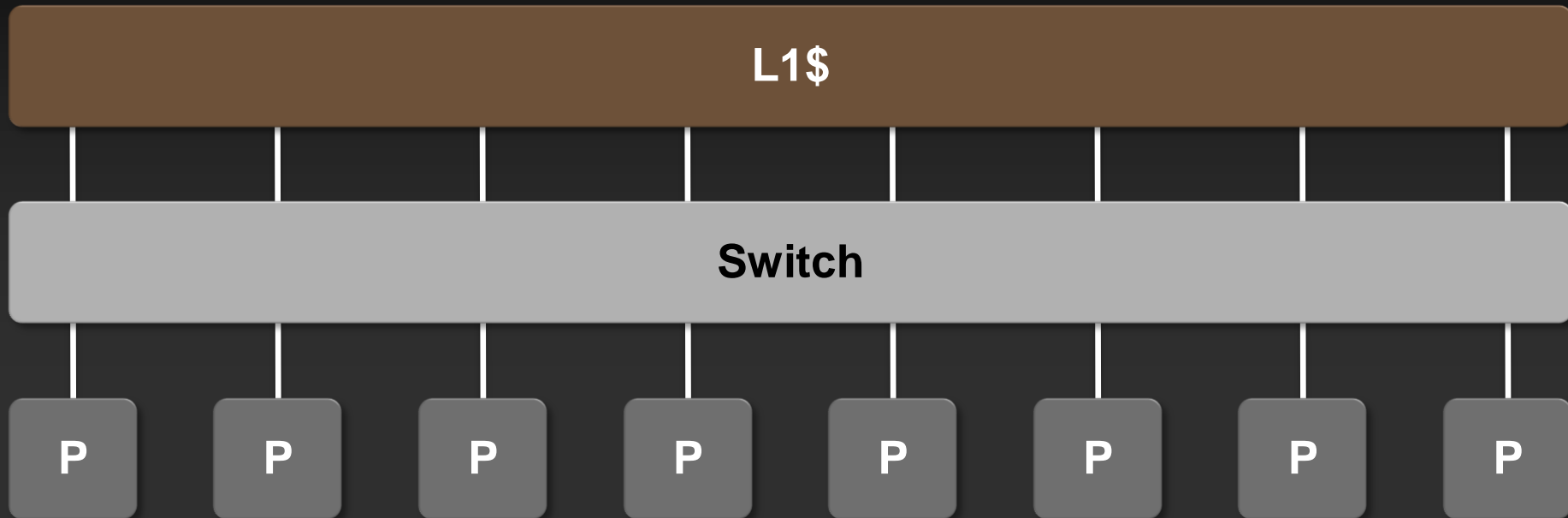
Fetching operands costs more than computing on them

20mm

64-bit DP
20pJ

26 pJ        256 pJ        16 nJ        DRAM Rd/Wr

256-bit buses

500 pJ        Efficient off-chip link

50 pJ

256-bit access
8 kB SRAM

1 nJ

28nm

# An NVIDIA ExaScale Machine

Lane – 4 DFMAs, 20GFLOPS

# SM – 8 lanes – 160GFLOPS

# Chip – 128 SMs – 20.48 TFLOPS + 8 Latency Processors

1024 SRAM Banks, 256KB each

| SRAM | SRAM | ••• | SRAM | MC | ••• | MC | NI |

**NoC**

| SM | SM | SM | SM | | SM | LP | | LP |

128 SMs 160GF each

# Node MCM – 20TF + 256GB

150GB/s
Network BW

**GPU Chip
20TF DP
256MB**

1.4TB/s
DRAM BW

**DRAM Stack**  **DRAM Stack**  ● ● ●  **DRAM Stack**  **NV Memory**

# Cabinet – 128 Nodes – 2.56PF – 38 kW



32 Modules, 4 Nodes/Module,
Central Router Module(s), Dragonfly Interconnect

# System – to ExaScale and Beyond

Dragonfly Interconnect
400 Cabinets is ~1EF and ~15MW

# CONCLUSION

# GPU Computing is the Future

**1** **GPU Computing is #1 Today**
On Top 500 AND Dominant on Green 500

**2** **GPU Computing Enables ExaScale**
At Reasonable Power

**3** **The GPU is the Computer**
A general purpose computing engine, not just an accelerator

**4** **The Real Challenge is Software**

THANK YOU

# Power is THE Problem

**1** Data Movement Dominates Power

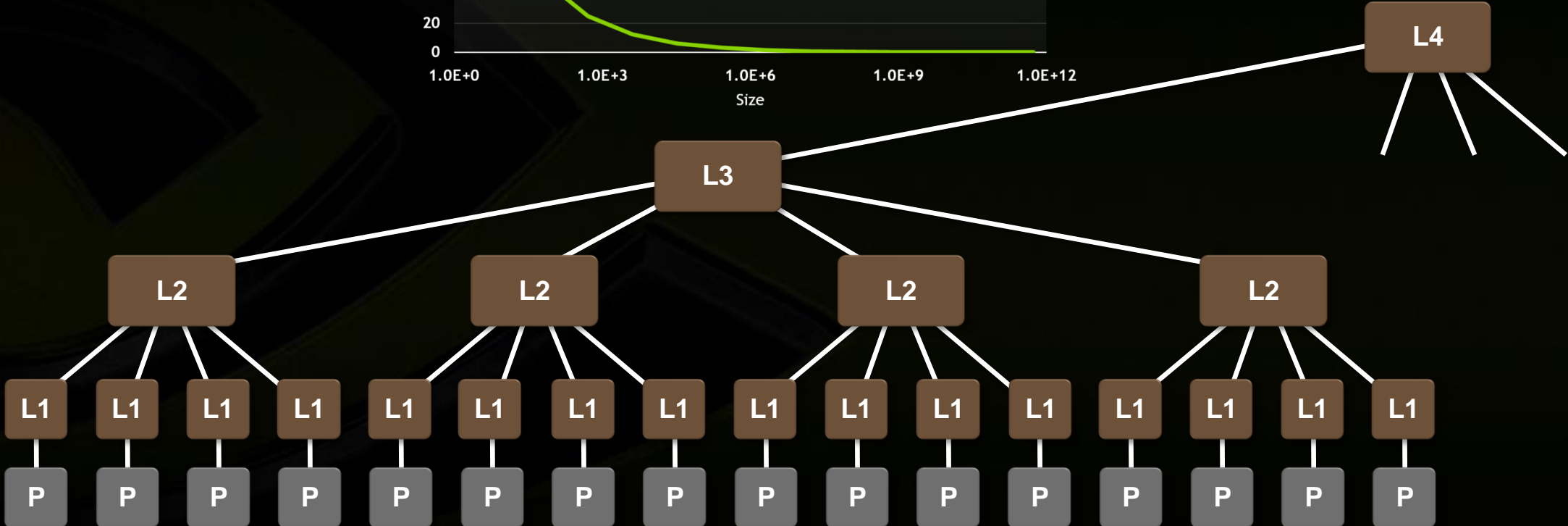**2** Optimize the Storage Hierarchy
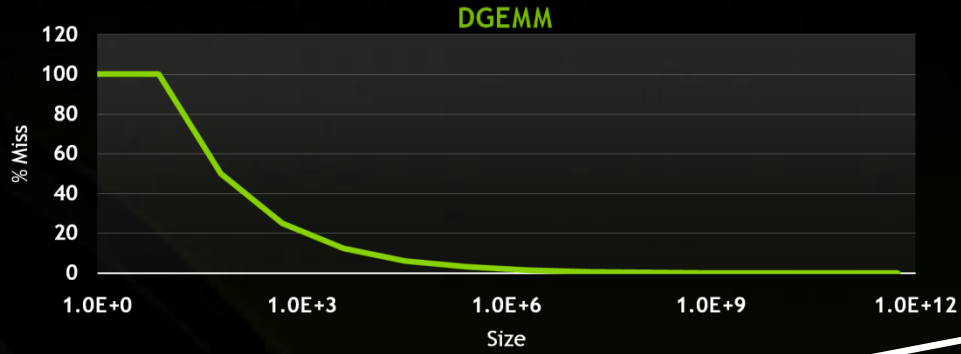
**3** Tailor Memory to the Application

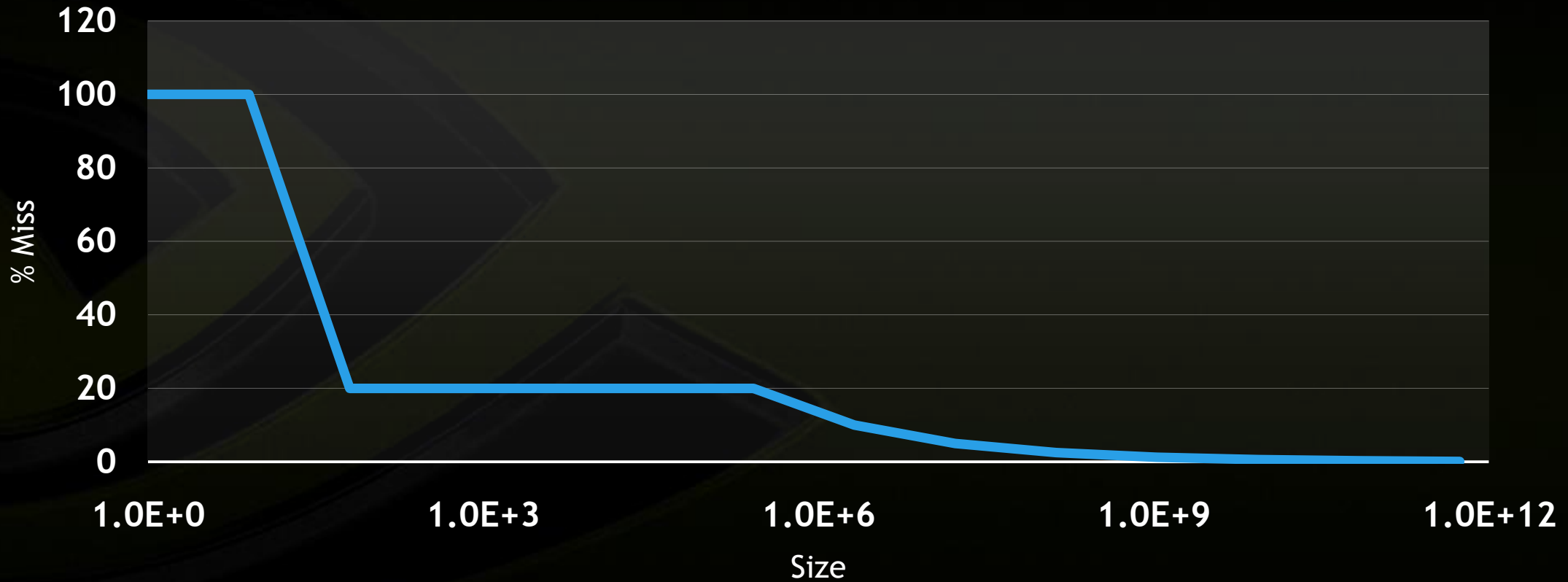# Some Applications Have Hierarchical Re-Use

## DGEMM

# Applications with Hierarchical Reuse Want a Deep Storage Hierarchy

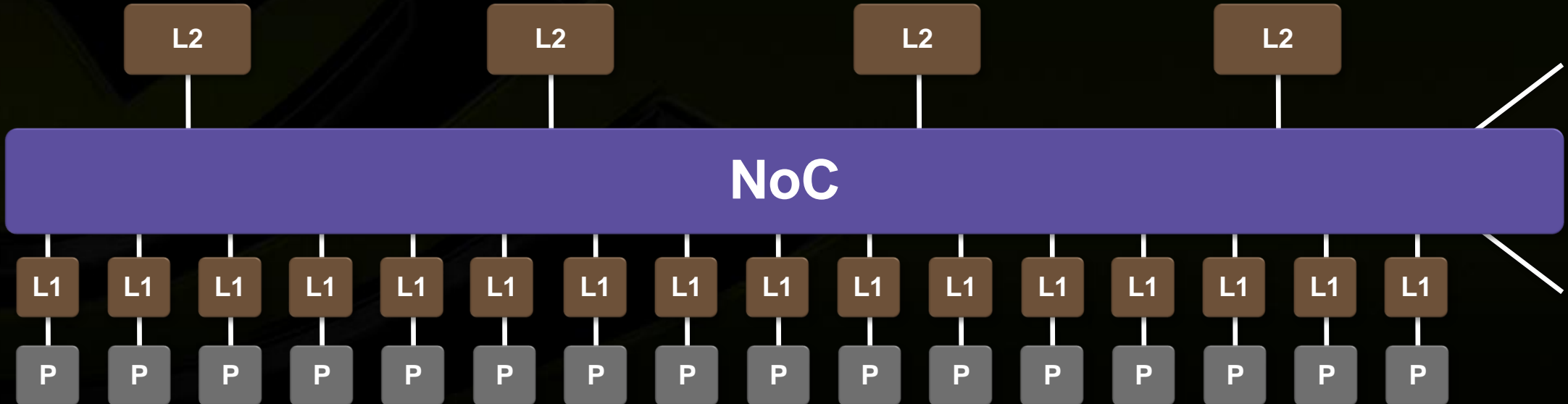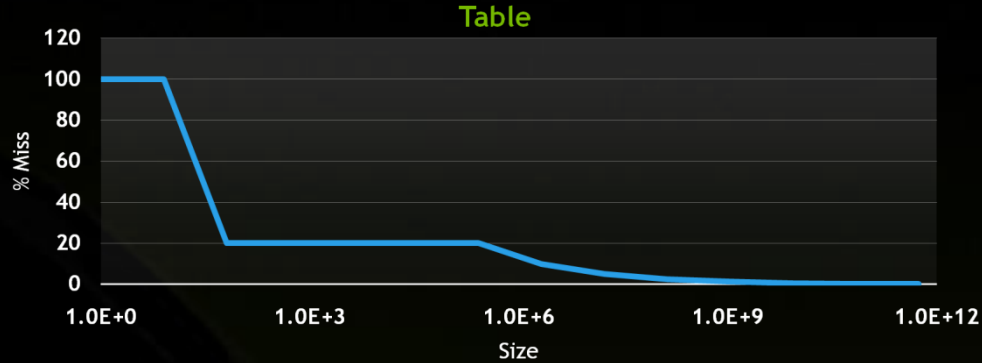# Some Applications Have Plateaus in Their Working Sets



Table

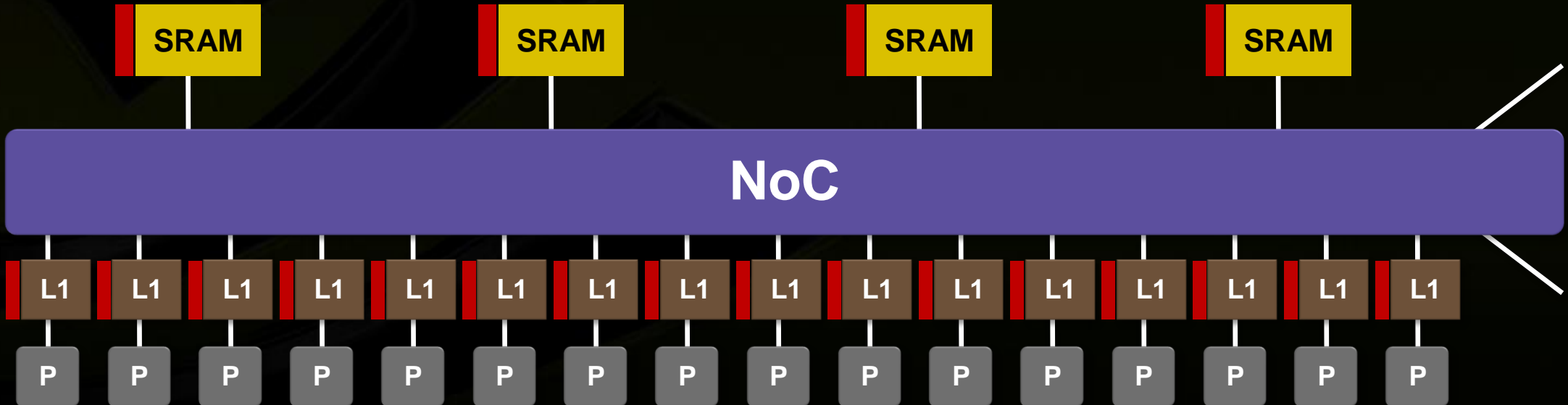# Configurable Memory Can Do Both At the Same Time

- Flat hierarchy for large working sets
- Deep hierarchy for reuse
- "Shared" memory for explicit management
- Cache memory for unpredictable sharing

# Configurable Memory
# Reduces Distance and Energy