# Chapter 4

## Data-Level Parallelism (DLP) in Vector, SIMD, and GPU Architectures

## Part 3: SIMD Extensions and GPU Intro

"We call these algorithms *data parallel* algorithms because their parallelism comes from simultaneous operations across large sets of data, rather than from multiple thread of control."

- W. Daniel Hillis and Guy L. Steele
"Data Parallel Algorithms," *Comm. ACM* (1986)

"If you were plowing a field, which would you rather use, two strong oxen or 1024 chickens?"

- Seymour Cray, Father of the Supercomputer
(arguing for two powerful vector processors versus many simple processors)

# Acknowledgements

- Thanks to many sources for slide material

# SIMD Extensions

- Media applications operate on data types narrower than the native word size
  - Examples: graphics and sound

- Major limitations (compared to vector instructions)
  - Number of data operands encoded into op code
    - Consequence?
  - No sophisticated addressing modes (strided, scatter-gather)
    - Consequence?
  - No mask registers
    - Consequence?

# SIMD Implementations

- Implementations
  - Intel MMX (1996)
    - Eight 8-bit integer ops or four 16-bit integer ops
  - Streaming SIMD Extensions (SSE) (1999)
    - Eight 16-bit integer ops
    - Four 32-bit integer/fp ops or two 64-bit integer/fp ops
  - Advanced Vector Extensions (2010)
    - Four 64-bit integer/fp ops

  - Operands must be consecutive and aligned memory locations

# Example SIMD Code: DAXPY

```
    L.D          F0,a            ;load scalar a
    MOV          F1, F0          ;copy a into F1 for SIMD MUL
    MOV          F2, F0          ;copy a into F2 for SIMD MUL
    MOV          F3, F0          ;copy a into F3 for SIMD MUL
    DADDIU       R4,Rx,#512      ;last address to load
Loop:
    L.4D         F4,0[Rx]        ;load X[i], X[i+1], X[i+2],
    X[i+3]
    MUL.4D       F4,F4,F0
        ;a×X[i],a×X[i+1],a×X[i+2],a×X[i+3]
    L.4D         F8,0[Ry]        ;load Y[i], Y[i+1], Y[i+2],
    Y[i+3]
    ADD.4D       F8,F8,F4        ;a×X[i]+Y[i], ...,
    a×X[i+3]+Y[i+3]
    S.4D         0[Ry],F8        ;store into Y[i], ..., Y[i+3]
    DADDIU       Rx,Rx,#32       ;increment index to X
    DADDIU       Ry,Ry,#32       ;increment index to Y
    DSUBU        R20,R4,Rx       ;compute bound
    BNEZ         R20,Loop        ;check if done
```

# SSE (Streaming SIMD)

- Intel 4-wide floating point
- Pentium III: 2 fully-pipelined, SIMD, single-precision FP units
- Eight 128-bit registers added to ISA
- In HW, (historically) broke 4-wide instructions into two 2-wide instructions
  – Interrupts are a problem
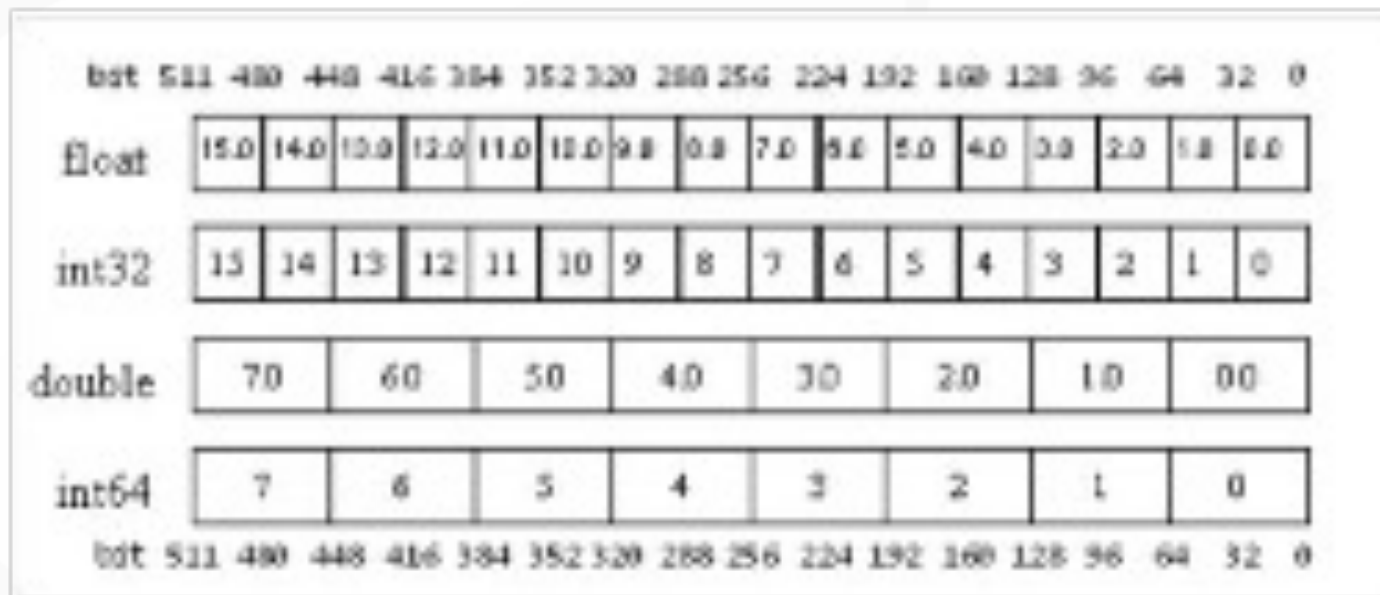- MMX + SSE: added 10% to PIII die

# Multimedia Extensions

- Very short vectors added to existing ISAs for micros
- Usually 64-bit registers split into 2x32b or 4x16b or 8x8b
  - Newer designs have 128-bit registers (Altivec, SSE2)
- Limited instruction set:
  - no vector length control
  - no strided load/store or scatter/gather
  - unit-stride loads must be aligned to 64/128-bit boundary
- Limited vector register length:
  - Requires superscalar dispatch to keep multiply/add/load units busy
  - Requires loop unrolling to hide latencies but increases register pressure
  - Trend towards fuller vector support in microprocessors
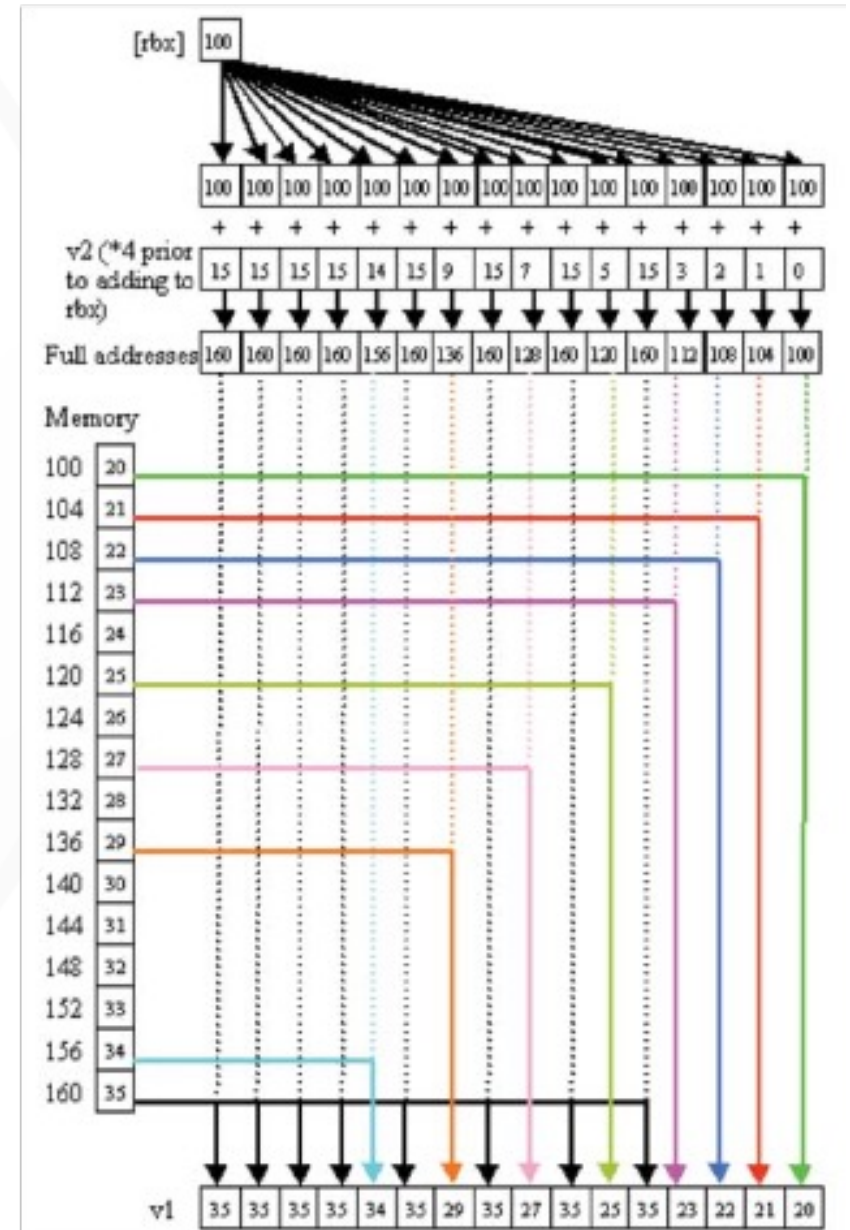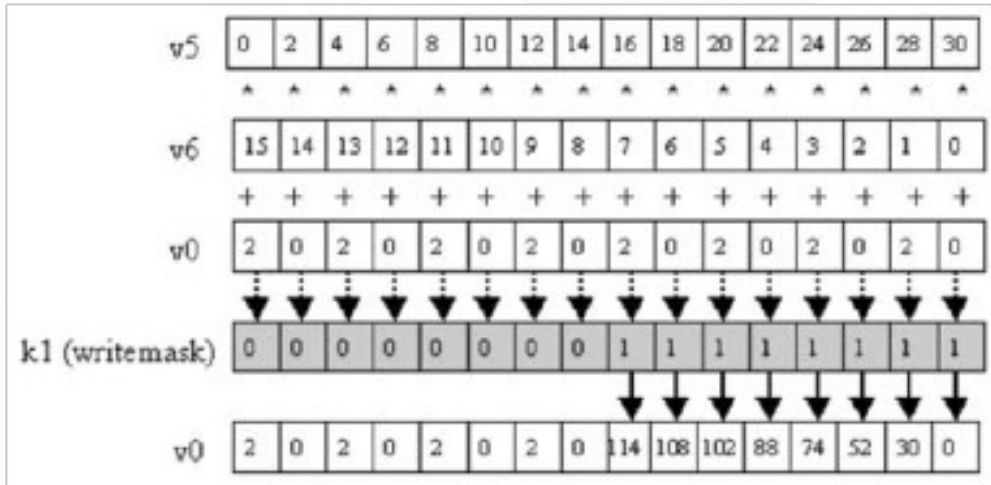
# Larrabee Vectors
(Intel's Initial "GPU" Offering)

- 32 512b-wide new registers
- 8 16b-wide vector mask registers
- Register ops generally ternary (a=b op c)
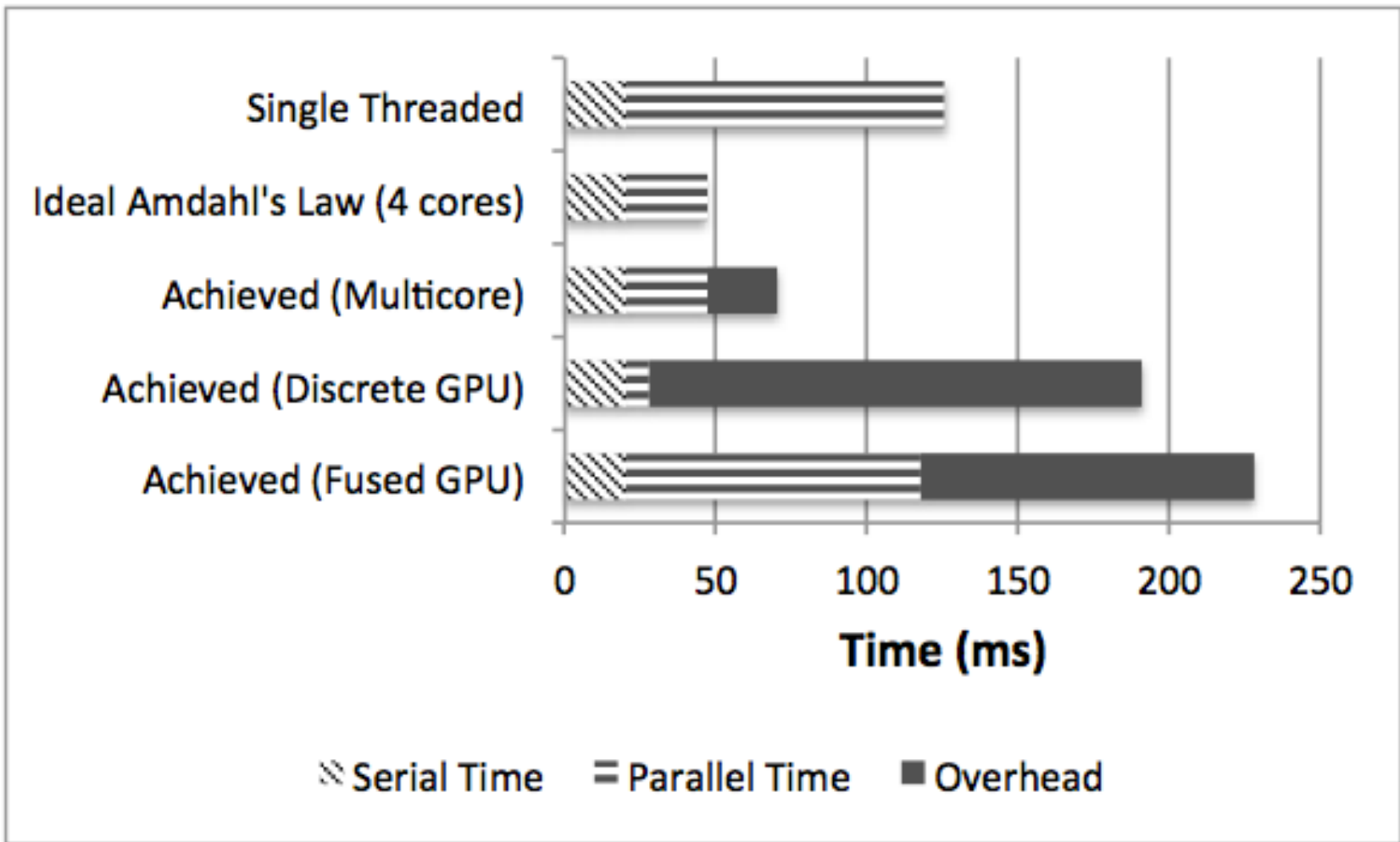- Extensive use of masks

# Larrabee Example
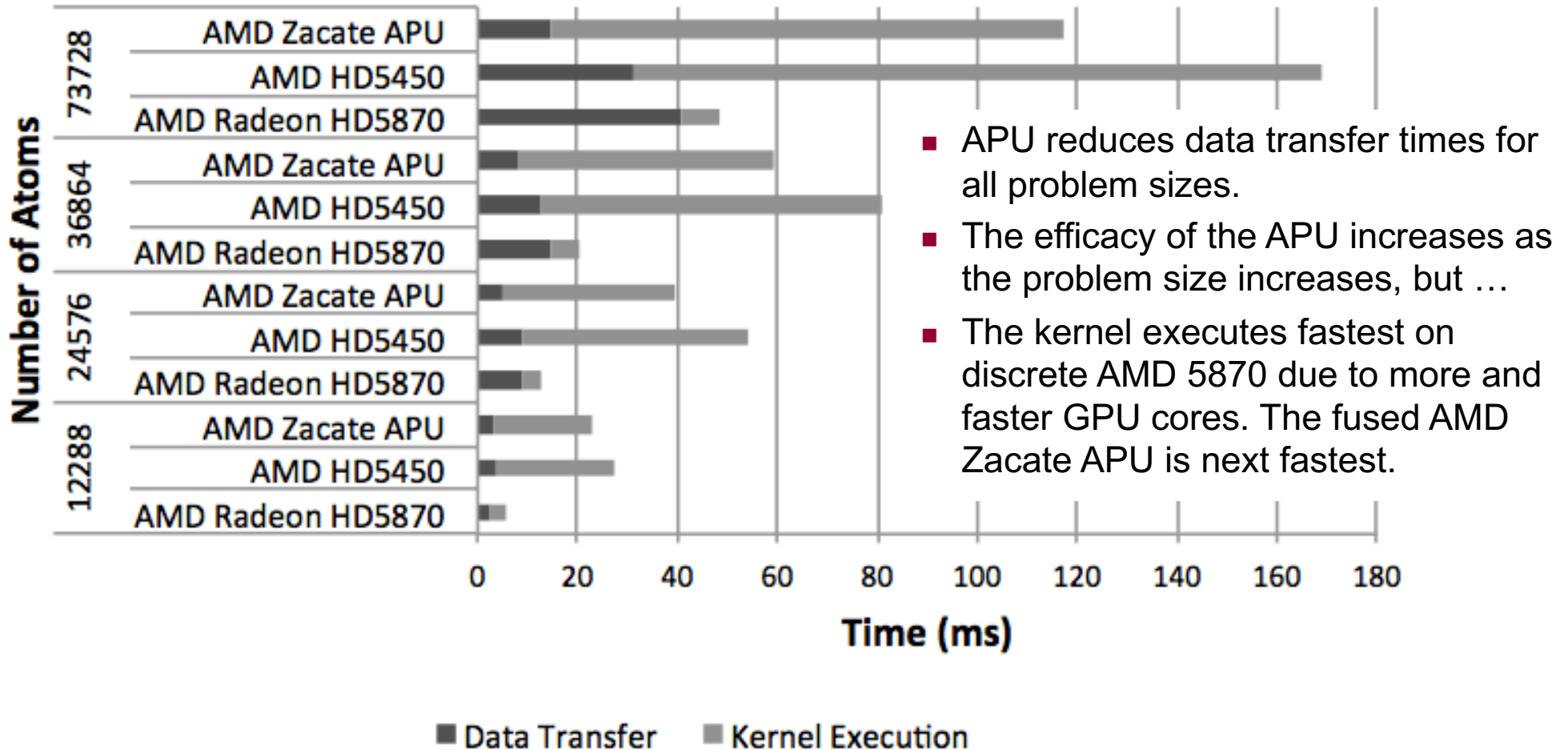
```
vmadd231ps v0 {k1}, v5, v6
```

# Graphical Processing Units

- Given the hardware invested to do graphics well, how can be supplement it to improve performance of a wider range of applications?

- Basic Idea
  - Heterogeneous execution model
    - CPU is the *host*, GPU is the *device*
    - What makes this model challenging?
  - Develop a C/C++-like programming language for GPU
  - Unify all forms of GPU parallelism as *CUDA thread* (NVIDIA)
    - OpenCL: Emerging vendor-independent language for multiple platforms
  - Programming model is "Single Instruction Multiple Thread"
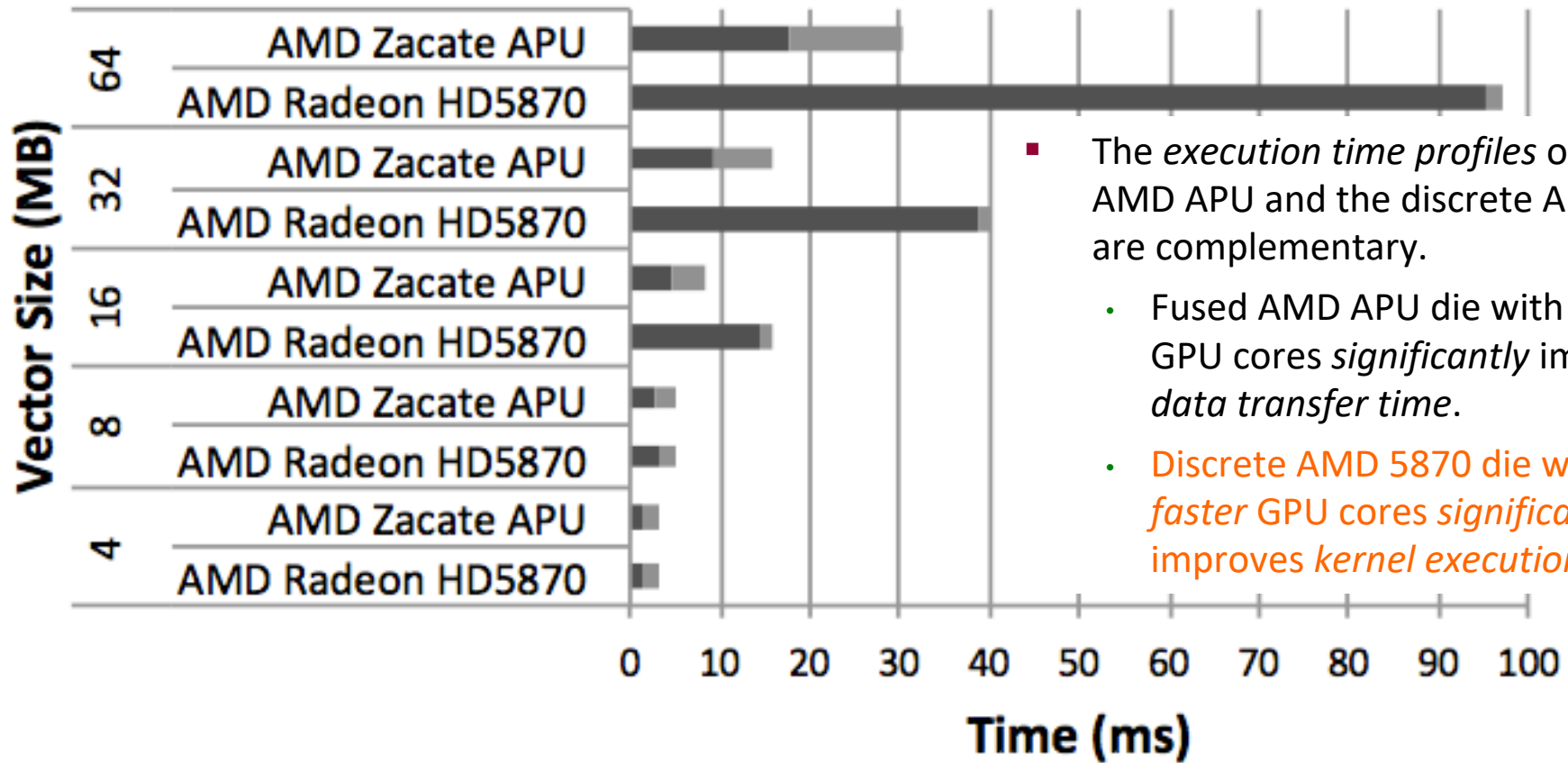
# Amdahl's Law:
# Different Parallel Devices



Chart: Time (ms) for Single Threaded, Ideal Amdahl's Law (4 cores), Achieved (Multicore), Achieved (Discrete GPU), Achieved (Fused GPU) — Serial Time, Parallel Time, Overhead

# Performance:
# Molecular Modeling (N-Body)



- APU reduces data transfer times for all problem sizes.
- The efficacy of the APU increases as the problem size increases, but …
- The kernel executes fastest on discrete AMD 5870 due to more and faster GPU cores. The fused AMD Zacate APU is next fastest.

# Performance: Reduction (Dense Linear Algebra)



- The *execution time profiles* of the fused AMD APU and the discrete AMD 5870 are complementary.
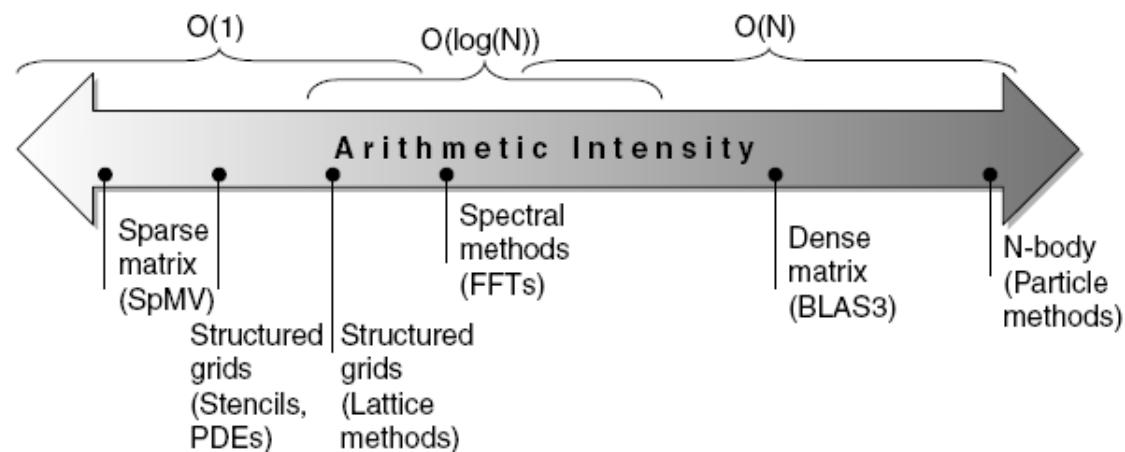  - Fused AMD APU die with only 80 GPU cores *significantly* improves *data transfer time*.
  - Discrete AMD 5870 die with 1600 *faster* GPU cores *significantly* improves *kernel execution time.*

# Roofline Performance Model

[Williams et al., 2009]

- Basic Idea
  - Plot peak floating-point throughput as a function of arithmetic intensity
  - Ties together floating-point performance, memory performance for a target machine, and arithmetic intensity
- Arithmetic Intensity
  - Floating-point operations per byte read

# Examples

- Attainable GFLOPs/sec = min (Peak Memory BW × Arithmetic Intensity, Peak Floating Point Perf.)