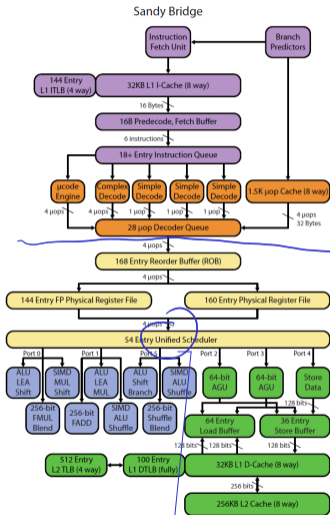February 4, 2025

## Announcements

- HW1
- Final talk topic /
  time assignment; Thu
- First talks; Feb 20.

## Goals

- Review
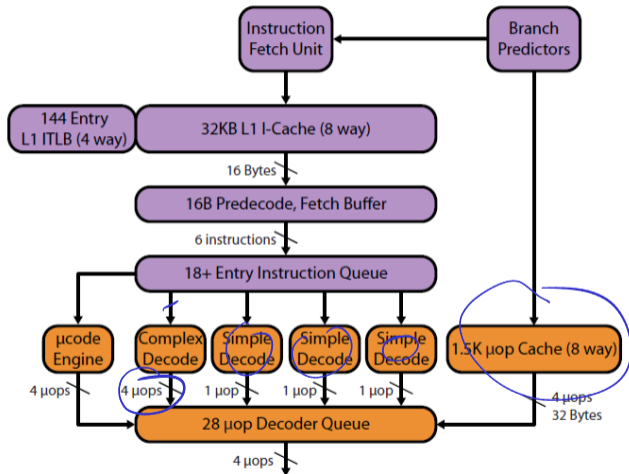- Memory subsystem

Review

# A Glimpse of a More Modern Processor



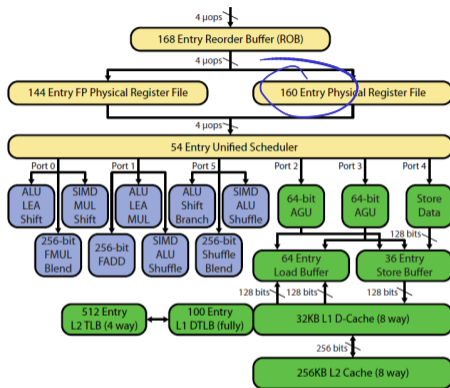[David Kanter / Realworldtech.com]

# A Glimpse of a More Modern Processor: Frontend
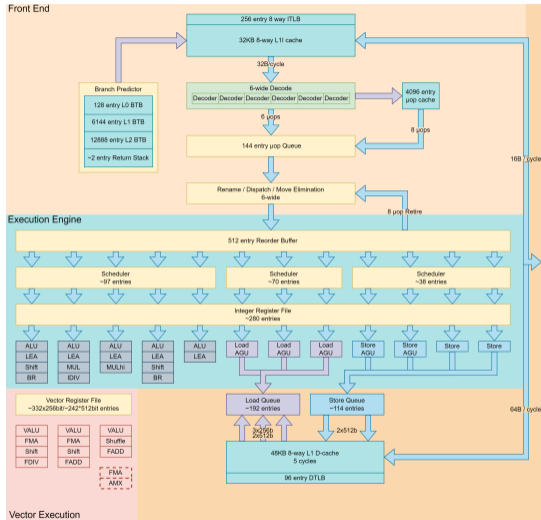


Sandy Bridge

[David Kanter / Realworldtech.com]

# A Glimpse of a More Modern Processor: Backend



- **New concept:** Instruction-level parallelism ("ILP", "superscalar")
- Where does the IPC number from earlier come from?
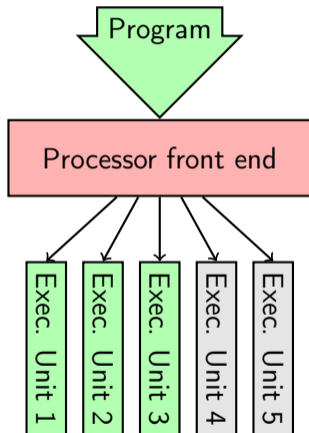
[David Kanter / Realworldtech.com]
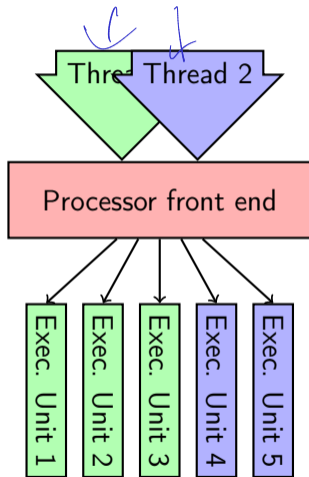
# A Glimpse of a More Modern Processor: Golden Cove



[Wikipedia ©]

# Demo

Demo: intro/More Pipeline Mysteries

# SMT/"Hyperthreading"



Q: Potential issues?

# SMT/"Hyperthreading"



Q: Potential issues?

- Contention ( not just
  cache?)

- "For losers"

- Power?
  Efficiency?

# Outline

# More Bad News from Dennard

| Parameter | Factor |
|---|---|
| Dimension | $1/\kappa$ |
| Line Resistance | $\kappa$ |
| Voltage drop | $\kappa$ |
| Response time | 1 |
| Current density | $\kappa$ |

[Dennard et al. '74, via Bohr '07]
- The above scaling law is for on-chip interconnects.
- Current $\sim$ Power *vs.* response time

Getting information from
- processor to memory
- one computer to the next

is
- slow (in *latency*)
- power-hungry

# Somewhere Behind the Interconnect: Memory

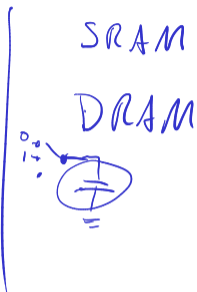Performance characteristics of memory:

- Bandwidth
- Latency

*Flops are cheap*
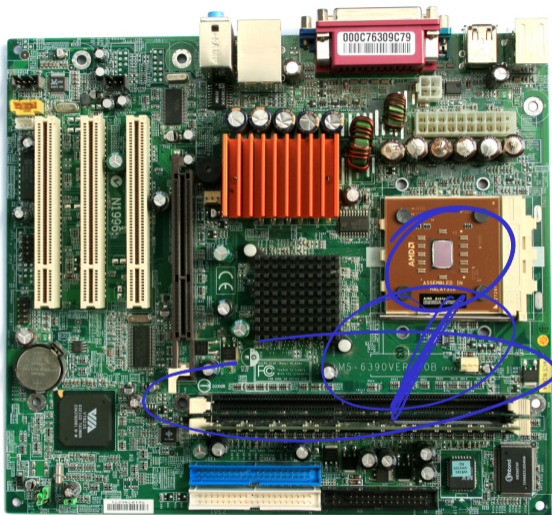*Bandwidth is money*
*Latency is physics — M. Hoemmen*

Minor addition (but important for us)?
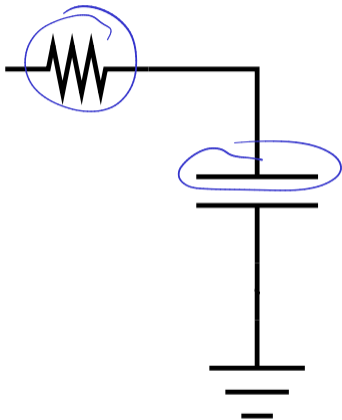
SRAM

DRAM

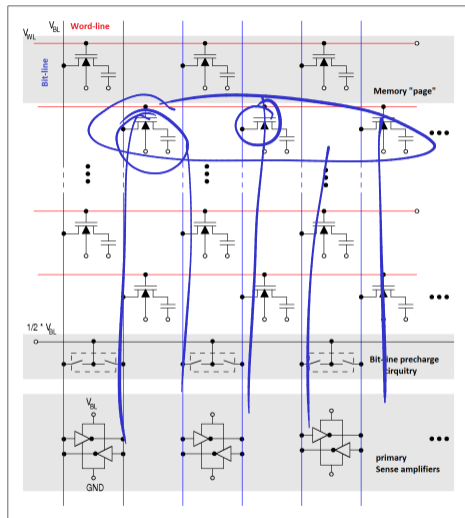> Bandwidth is money and code structure.

# Latency is Physics: Distance



[Wikipedia ©]

# Latency is Physics: Electrical Model

# Latency is Physics: DRAM

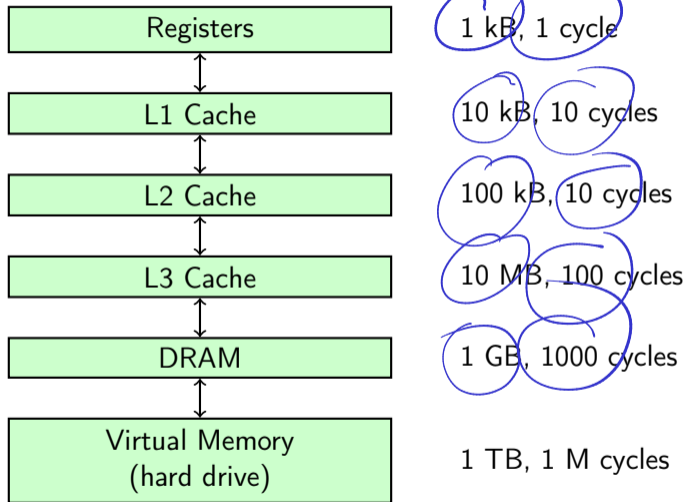# Latency is Physics: Performance Impact?

What is the performance impact of high memory latency?

> *Stall.*

Idea:

- ▶ Put a look-up table of recently-used data onto the chip.
- ▶ Cache
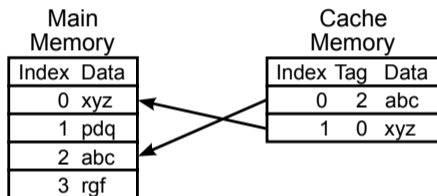
# Memory Hierarchy

| | |
|---|---|
| Registers | 1 kB, 1 cycle |
| L1 Cache | 10 kB, 10 cycles |
| L2 Cache | 100 kB, 10 cycles |
| L3 Cache | 10 MB, 100 cycles |
| DRAM | 1 GB, 1000 cycles |
| Virtual Memory (hard drive) | 1 TB, 1 M cycles |

# A Basic Cache

Demands on cache implementation:

▶ Fast, small, cheap, low power
▶ Fine-grained
▶ High "hit"-rate (few "misses")

| Main Memory | |
|---|---|
| Index | Data |
| 0 | xyz |
| 1 | pdq |
| 2 | abc |
| 3 | rgf |

| Cache Memory | | |
|---|---|---|
| Index | Tag | Data |
| 0 | 2 | abc |
| 1 | 0 | xyz |

Design Goals: at odds with each other. Why?
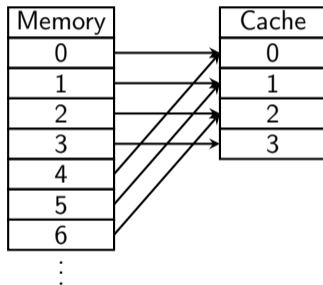
lookup logic is not free

# Caches: Engineering Trade-Offs
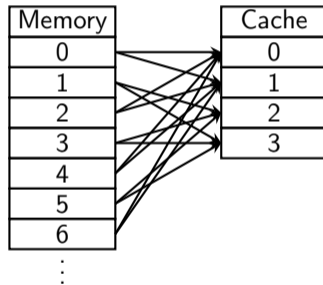
Engineering Decisions:

- ▶ More data per unit of access matching logic
  → Larger "Cache Lines"
- ▶ Simpler/less access matching logic
  → Less than full "Associativity"
- ▶ Eviction strategy
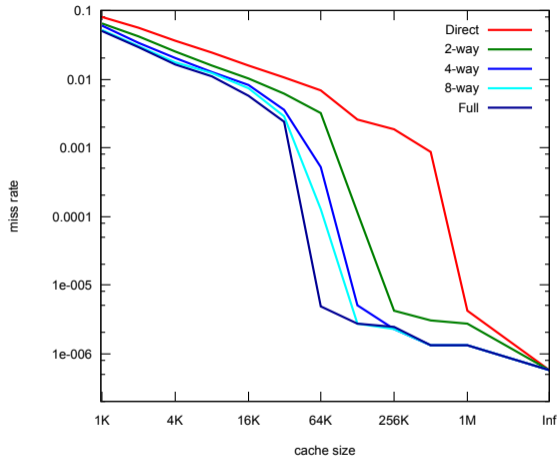- ▶ Size

# Associativity

Direct Mapped:

2-way set associative:

# Size/Associativity vs Hit Rate



Miss rate versus cache size on the Integer portion of SPEC CPU2000
[Cantin, Hill 2003]