# CS773-2022-Autumn: Computer Architecture for Performance and Security

## Lecture 3: Catch the Cache

https://www.cse.iitb.ac.in/~biswa/

Phones on silence, please

Thank You

# Microarchitecture 101: World with no caches

*North pole* ☹

**Core**

32-bit Address

Data

200 to 300 cycles

Costly DRAM accesses☹

Minimizing costly DRAM accesses is critical for performance

*4 GB DRAM*

*South pole* ☹

CASPER

3

# Remember Latency and Bandwidth
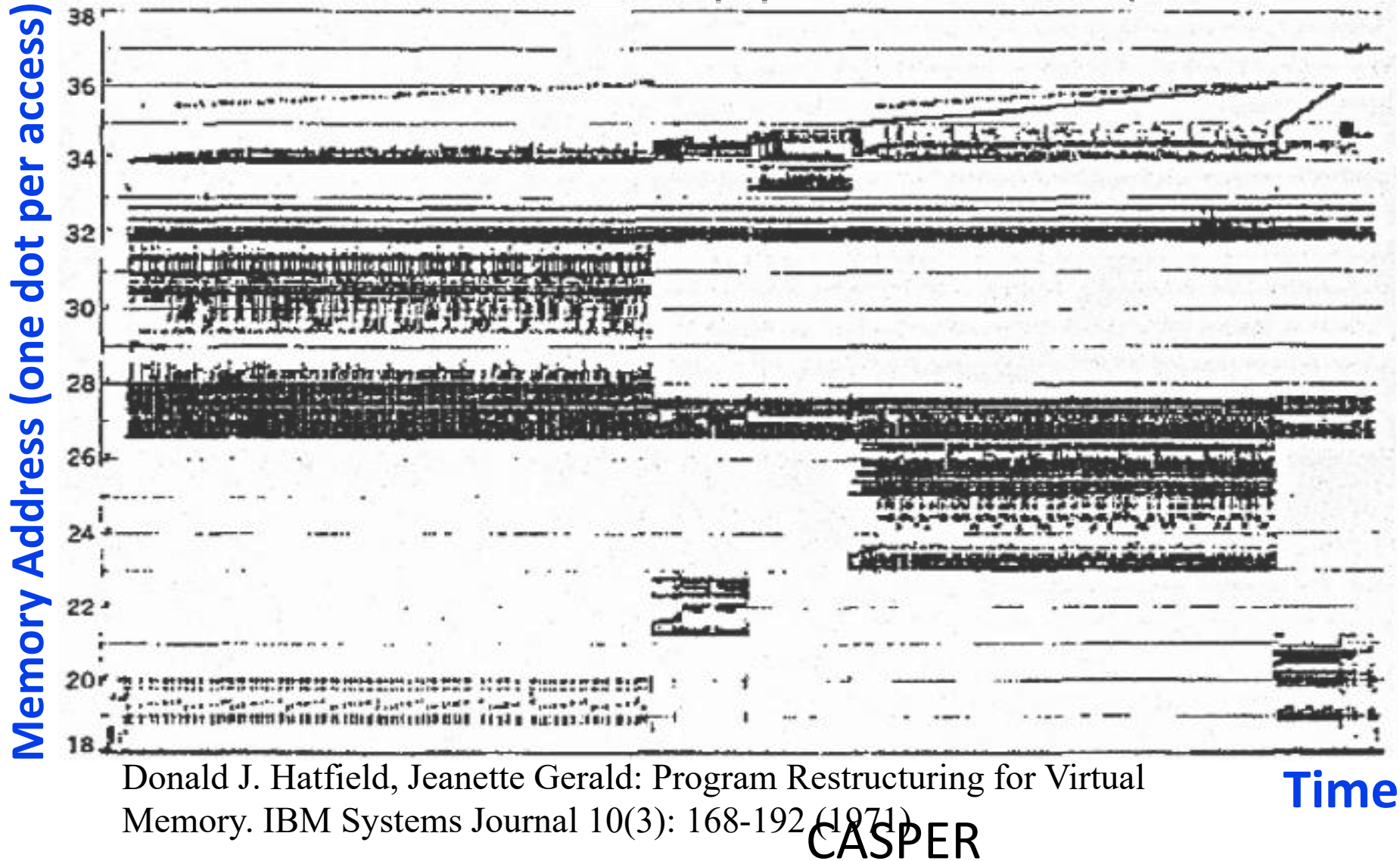
Core

Latency

Bandwidth



CASPER

# Latency ☹

*Bandwidth problems can be cured with money. Latency problems are harder because the speed of light is fixed – <span style="color:red">you can't bribe God</span>*
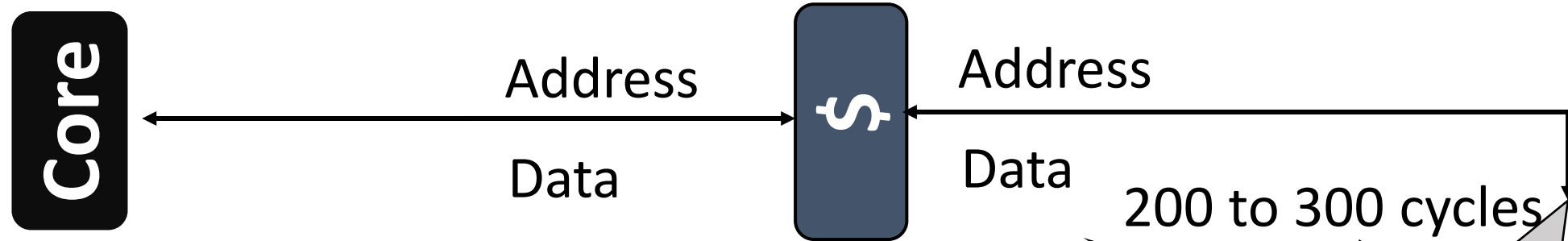
# Let's look at the Applications (benchmarks)



Memory Address (one dot per access)

Time

Donald J. Hatfield, Jeanette Gerald: Program Restructuring for Virtual Memory. IBM Systems Journal 10(3): 168-192 (1971)

# Caching: 10K Feet View

*North pole* ☺

**Core** ← Address / Data → **$** ← Address / Data

200 to 300 cycles

Caching is a *speculation* technique ☺
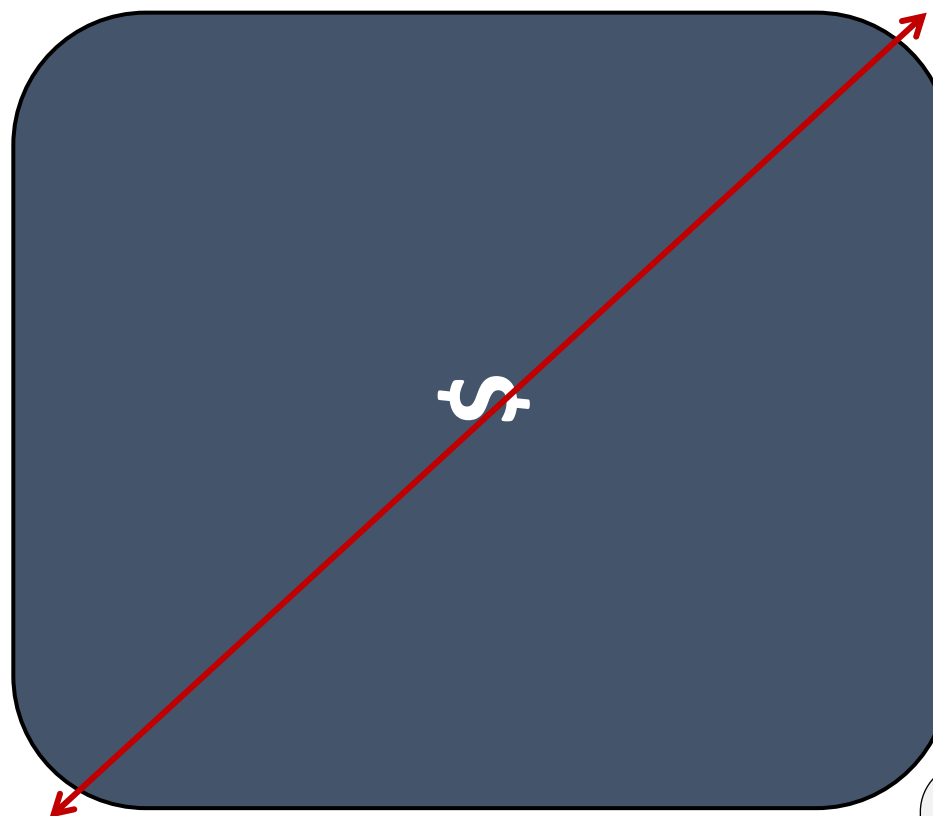Works – if locality

Costly DRAM accesses☹

# How big/small?

**Core**

Latency: low
Area: low
Capacity: low

Latency: high
Area: high
Capacity: high

CASPER

# Cache with latency

*North pole* ☺

**Core** ←— Address ——→ **$** ←— Address ——
Data: 1 cycle        Data

200 to 300 cycles

32 to 64KB $ will be available in one to four cycles ☹

Costly DRAM accesses ☹

*South pole* ☺

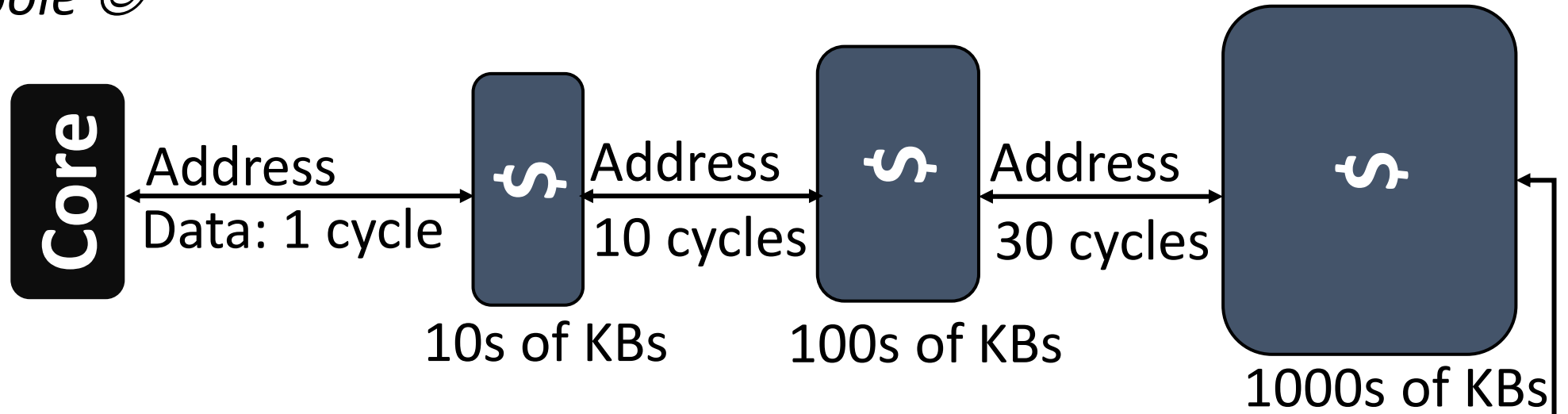# Cache hierarchy with latency

*North pole* ☺



Core — Address / Data: 1 cycle — $ — Address / 10 cycles — $ — Address / 30 cycles — $

10s of KBs          100s of KBs          1000s of KBs

Multi-level cache hierarchy

*South pole* ☺

# Cache hierarchy with latency

*North pole* ☺

**Core**

Address
Data: 1 cycle

$ 10s of KBs

Address
10 cycles

$ 100s of KBs

Address
30 cycles

$ 1000s of KBs

Multi-level cache hierarchy

How many levels ?

Total latency < DRAM latency

*South pole* ☺

# Takeaway

**L1 $**

Latency and bandwidth (multiple ports)

**L2 $**

**L3 $**

# Takeaway

**L1 $** | Latency and bandwidth (multiple ports)

**L2 $** | Latency

**L3 $**

# Takeaway (Do not forget the word microarch.)

**L1 $**    Latency and bandwidth (multiple ports)

**L2 $**    Latency

**L3 $**    Capacity

# PAUSE

# Accessing a cache

One byte

**Core**

# Bytes to blocks (lines)



One byte

One line

**Core**

Typical line size: 64 to 128 Bytes

CASPER

# A bit deeper: 1024 lines each of 32B

4 GB DRAM

One byte

**Core**

Address (32-bit)

Line 0

One line

Line 1023

# A bit deeper: 1024 lines each of 32B

4 GB DRAM

One byte

Core

Address (32-bit)

Line 0

One line

Line 1023

Line number (index): 10 bits

Byte offset   (offset): 5 bits

CASPER

19

# Direct Mapped Cache

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

Line 0

Line 512

**byte**

.
.
.

Line 1023

# Direct Mapped in Action



31                               14                  4                 0

| Cache Tag | Cache Index | Byte Select |
|---|---|---|

Ex: 0x50              Ex: 0x01         Ex: 0x00

Valid Bit              Cache Tag                         Cache Data

| | Byte 31 | .. | Byte 1 | Byte 0 | 0 |
|---|---|---|---|---|---|
| 0x50 | Byte 63 | .. | Byte 33 | Byte 32 | 1 |
| | | | | | 2 |
| | | | | | 3 |
| : | : | | | : | |
| | Byte 1023 | .. | | Byte 992 | 1023 |

CASPER

21

# What if we have multiple ways?



Tag | Index | Offset

Set 0
Tag
Tag

Set 511
Tag
Tag

One byte

Line 0

byte

Line 1023

CASPER

# 2-way associative in action



CASPER

# Knobs of interest

Line size, associativity, cache size
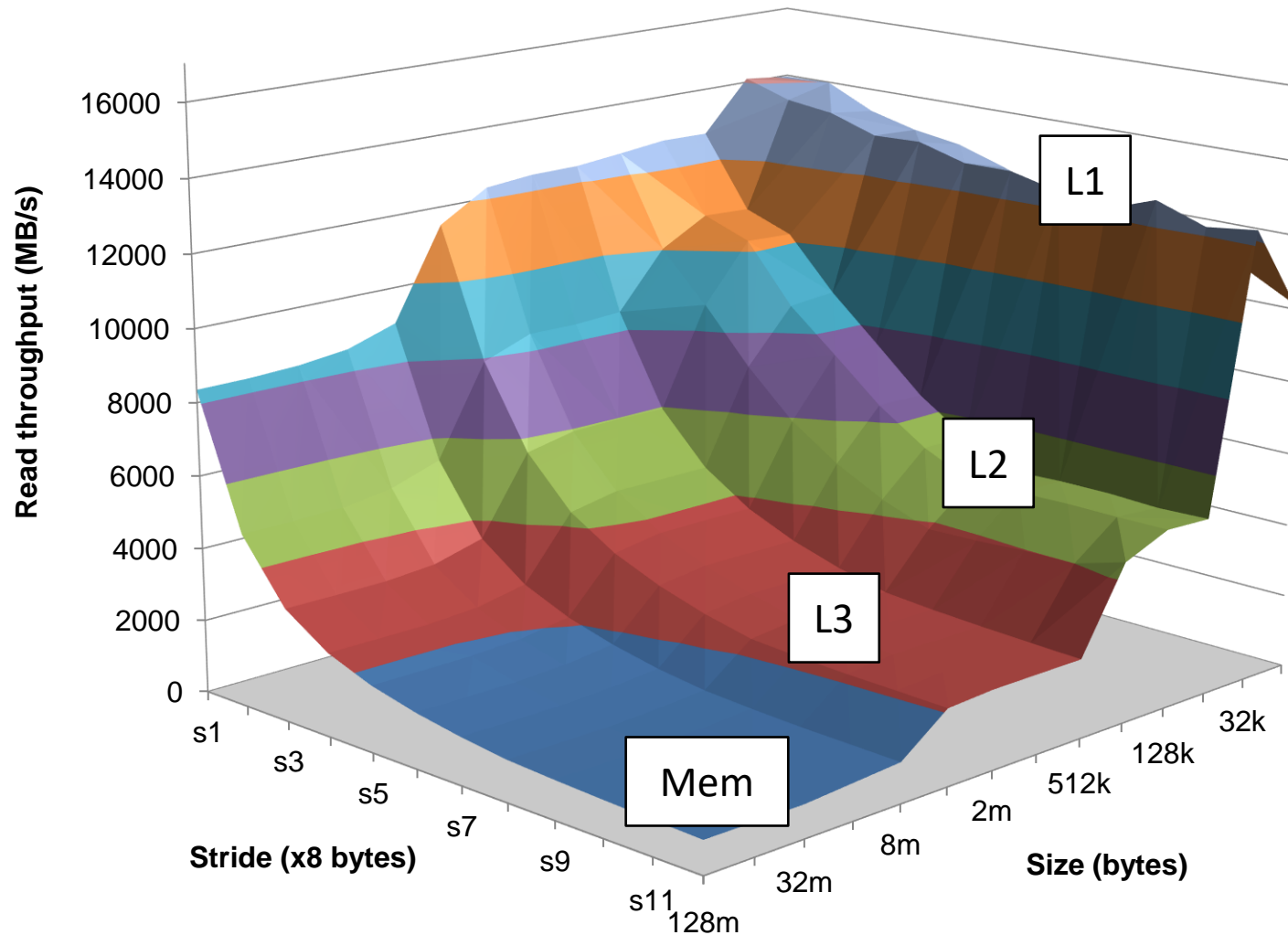
Tradeoff: latency, complexity, energy/power

Tips: Think about the extremes:
Line size = one byte or cache size
Associativity = one or #lines
Cache size = Goal oriented: latency/bandwidth or capacity

# Memory Mountain



Think about it, before you write your program:
Assignment-I is live

CASPER

# Into the Real World

sudo dmidecode -t cache

cat /proc/cpuinfo

getconf -a | grep CACHE

lscpu

Wiki chip: https://en.wikichip.org/wiki/WikiChip

Perf tool: https://perf.wiki.kernel.org/index.php/Main_Page

sudo perf stat -e cache-misses ….

PAUSE