*Recall ...*

The processor can only keep a part of a computation; portions must move back and forth as needed.

Same memory holds both instructions and data in a modern stored-program computer.

⇨ **Essential role: hold computation**

✎ Close interaction with processor

✎ Limit computation size

*therefore,*
*3 issues matter*

✎ Limit performance

⇨ **Design goal (ideally)**

Maximum amount of fastest memory

# Memory Solutions

Memory raised substantial challenges due to a growing performance gap between the two parts, essential to performing computations, that must interact closely.

👁

**Primary storage** is an older term for main memory where actively running code is stored in RAM or ROM.

⇨ # CPU-memory technology gap

Not all storage systems are <u>usable as memory</u> (i.e., to hold jobs marked or selected for execution)

⇨ # Revised Design Goal

Maximum amount of fastest memory <u>economically</u>

# Memory: A Deeper Look
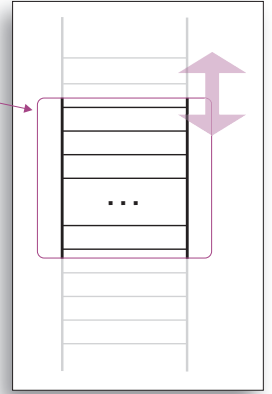# Principle of Locality

⇨ **[Memory] References** 👁

**Memory generally stores program instructions and data, but how are they actually used?**

⇨ ## Memory access patterns

✎ Most recently used

✎ Close locations

*window (high probability) active access*

**Where do these patterns come from? (Next.) How probable? (Later.)**

⇨ ## Mem *reference* locality

✎ Temporal ~ access times

✎ Spatial ~ access locations

**Quiz**
**What's the difference? [memory access context]** *Address, label, pointer,* and *reference*.

Give examples from <u>high-level programs</u> for temporal and spatial locality in data (min 2–3 each). Repeat for instructions.
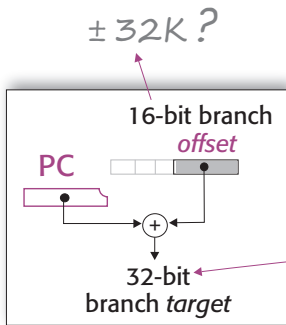
**Exercise**
Lookup how *objects* are typically stored in memory. Suggest locality opportunities.

# **Think about your programs**

✎ Data locality examples  *vars, arrays, data structures...*

✎ Instruction locality examples  *loops, subroutines conditionals...*

**± 32K** *?*

16-bit branch *offset*

PC

⊕

32-bit branch *target*

# **Note MIPS PC-relative addr**
## Store an offset (a difference)

```
beq $4,$5,Exit
```

**Quiz**
Which type of locality is exploited?

# Memory Systems

In 2021–22, as much as 4–5 orders of magnitude (approx. 10K to 100K); see next.

# To sum up

✎ Storage vary wildly in cost

✎ Small active reference window

Only a relatively small amount at a time needs to be referenced very quickly.

Watch keywords: amount (**capacity**), economical (**cost**), and access speed (**performance**).
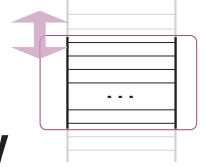
# Hierarchical storage: a solution

Largest amount of economical storage accessed, <u>most of the time</u>, at the speed of the fastest storage
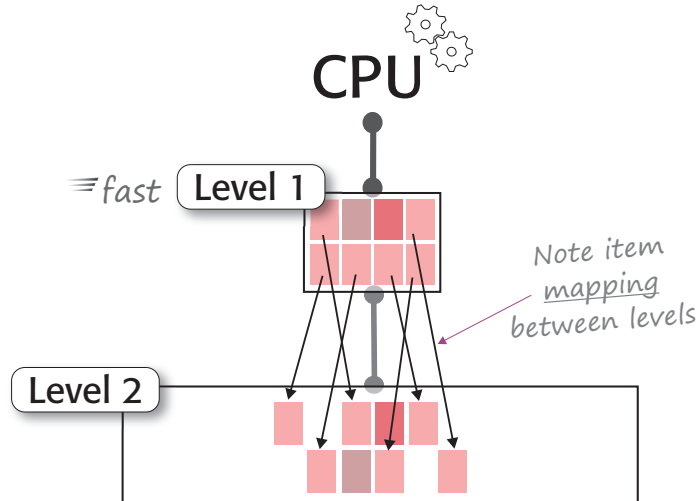
# Memory Systems
# Hierarchical Storage

## ⇨ Inclusion

The fastest memory, typically expensive and therefore small (as much as design economics afford), is directly connected to the CPU.

Typically insufficient, so additional, more economical, but slower memory (too slow to hook to the CPU) may be connected as a secondary back store.
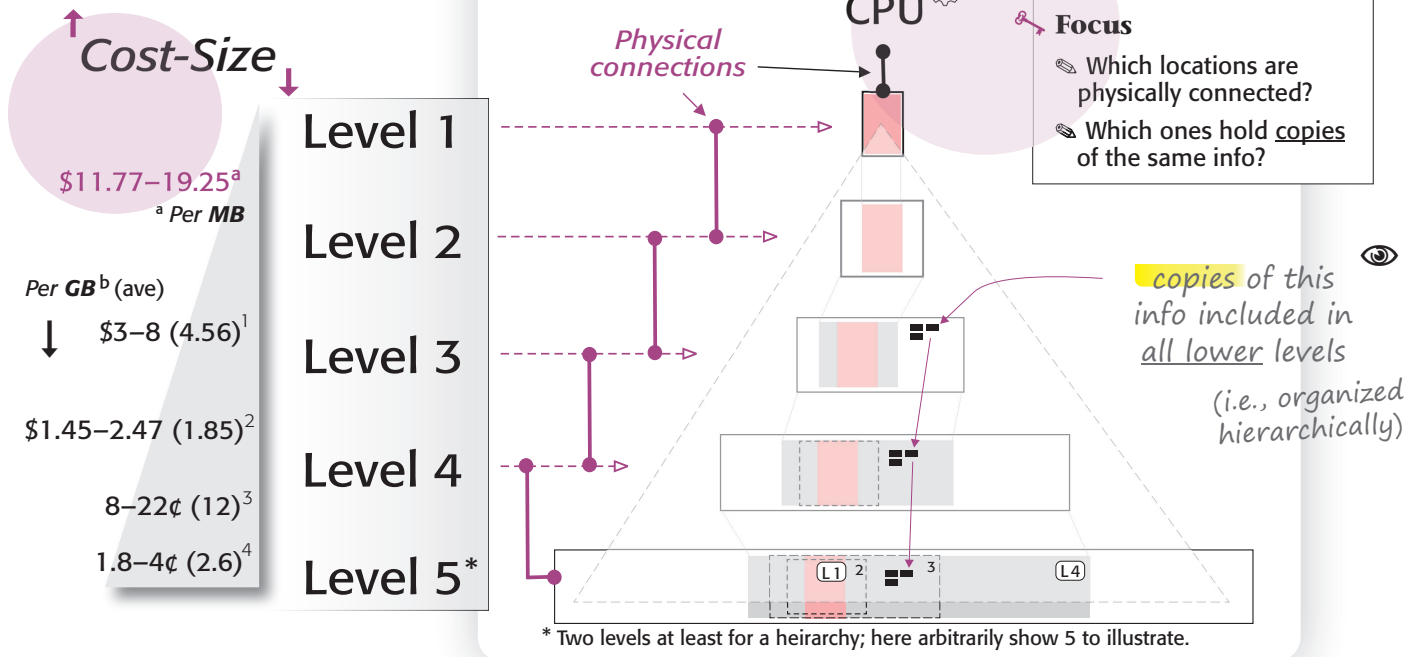
Items in the upper level numbered 1, are <u>included</u> in (a <u>proper</u> subset of) the lower back store level, numbered 2.



CPU

≡ fast  Level 1

Note item <u>mapping</u> between levels

Level 2

# Memory Systems
# A Memory Hierarchy

➩ **Inclusion**

*Cost-Size*

$11.77–19.25[a]

[a] *Per MB*

*Per GB[b]* (ave)

↓ $3–8 (4.56)[1]

$1.45–2.47 (1.85)[2]

8–22¢ (12)[3]

1.8–4¢ (2.6)[4]

Level 1

Level 2

Level 3

Level 4

Level 5*

CPU

*Physical connections*

**Focus**
- Which locations are physically connected?
- Which ones hold <u>copies</u> of the same info?

copies of this info included in <u>all lower</u> levels

*(i.e., organized hierarchically)*

L1  2  3  L4

* Two levels at least for a heirarchy; here arbitrarily show 5 to illustrate.

[a] SRAM/Chips (Cypress/Infineon) 0.45/3/45 ns arrow.com@2021/12/13 • Integrated on-die <$10/MB likely (?) in 2021
[b] Consumer/retail newegg.com@2021/12/14 (best selling) • [1]DRAM DDR3-1333 4/8G DIMM • [2]SSD/3D XPoint (*Intel Optane*) • [3]SSD/NAND 256-1TB • [4]HDD/internal 10-20 TB

# Memory Systems
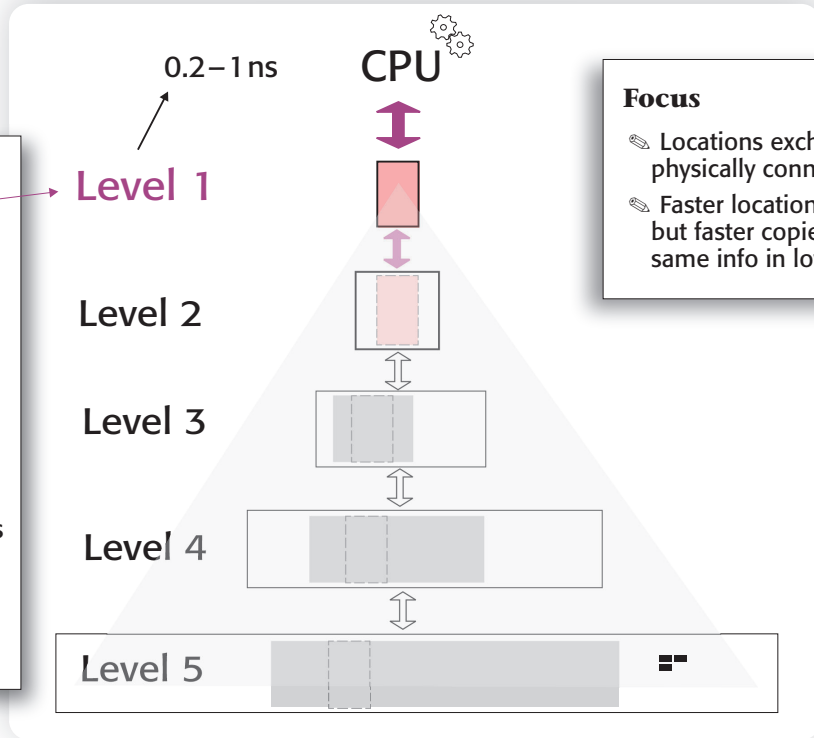# Hierarchy Operation

### ➪ Miss Time *penalty*

Best case scenario:
CPU finds requested
info in Level 1.

*Speed*

0.2−1 ns

CPU

👁
Level access times add up
as info is fetched from
lower levels.

< 1 ns → Level 1

**Focus**

✎ Locations exchange info if
physically connected

✎ Faster locations hold fewer
but faster copies of the
same info in lower levels

Level 2

10s ns

Level 3

10s μ-seconds

Level 4

Least used parts of <u>active</u>
programs; note indicated
items will be very costly
to fetch upon request.

milli-seconds

Level 5

# Memory Systems
# 🕐 Summary

**Quiz**
Even if memory were super fast and cheap, is it a good idea to connect the CPU to all the memory?
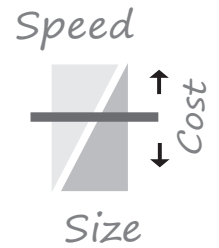
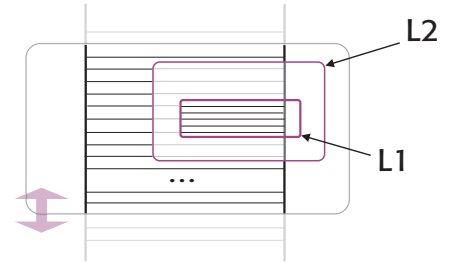**Exercise**
<u>Why</u> store hierarchically (i.e., include addresses in lower levels)? **Hint**: look for answers in reading material; note lower levels are higher-numbered.

⇨ # What is a memory hierarchy?

✎ A speed-size tradeoff

✎ Cost-driven

✎ <u>Access locality</u>-enabled

*Speed*

↑ *Cost* ↓

*Size*

<u>As much</u> of the most actively used memory is kept in the fastest locations (upper part of the hierarchy) close to the CPU <u>as long as possible</u> (not as long as needed, unfortunately, occasionally stuff is evicted prematurely).

⇨ # Why it works?

L2

L1

...

# Memory Systems
# Classic Hierarchy

## ⇨ Coherence

Some logical memory addresses will have <u>many copies</u> of their contents in the physical memory system.
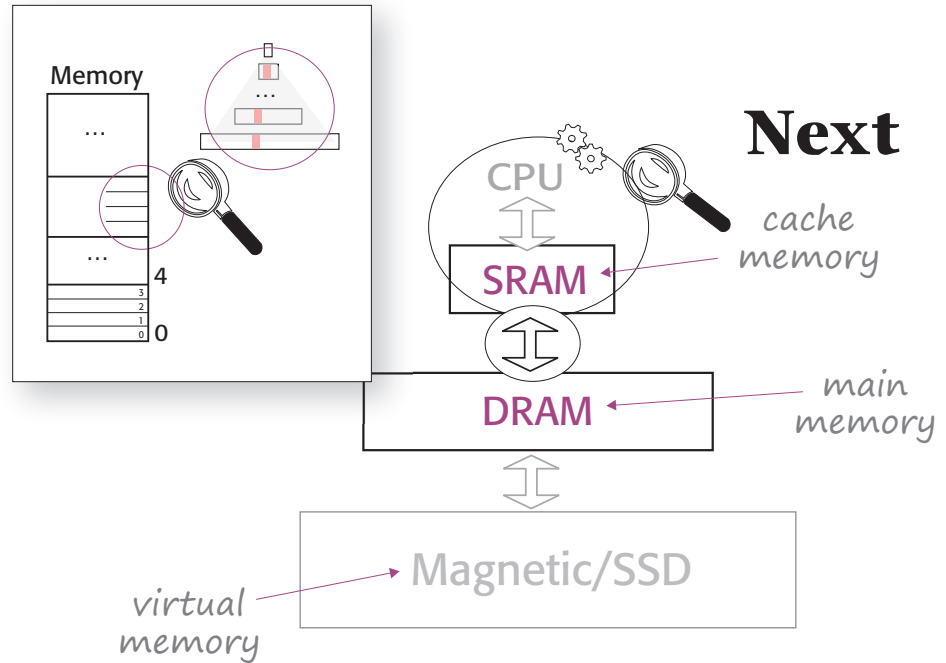
**Quiz** What if the CPU changes its copy? What about a multiple CPU scenario? Should a cache be read-only?

**Static RAM (SRAM)** Faster, expensive volatile bits made from transistors.

**Dynamic RAM (DRAM)** Cheaper, denser volatile bits made from capacitors.

**Magnetic/SSD** Very cheap, persistent bits, electro mechanical magnetic storage or transistor-based solid-state flash memory.

Recently, phase-change solid-state (based on changing resistivity due to crystalline phase change) such as Intel/Micron *3D XPoint* (cross-point).



**Next**

cache memory

main memory

virtual memory