

Define Cache

**Understand
how cache
works**

**Write cache-
friendly code**



Cache Memory

~

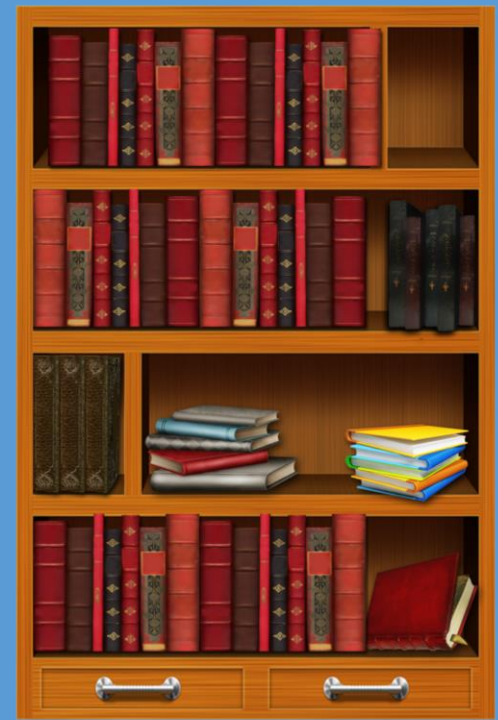
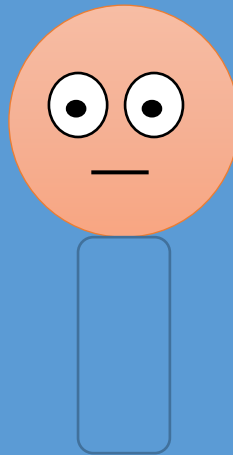
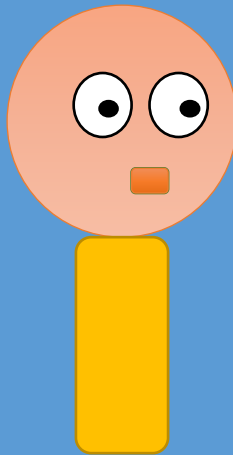
- ① Cache Concepts
- ② Cache Organization
- ③ Writing Cache-friendly Code

~

Cache Concepts

Memory

CPU



Cache





CPU-Z



CPU

Caches

Mainboard

Memory

SPD

Graphics

Bench

About

L1 D-Cache

Size

32 KBytes

x 4

Descriptor

8-way set associative, 64-byte line size

L1 I-Cache

Size

32 KBytes

x 4

Descriptor

8-way set associative, 64-byte line size

L2 Cache

Size

256 KBytes

x 4

Descriptor

4-way set associative, 64-byte line size

L3 Cache

Size

8 MBytes

Descriptor

16-way set associative, 64-byte line size

Size

Descriptor

Speed

CPU-Z

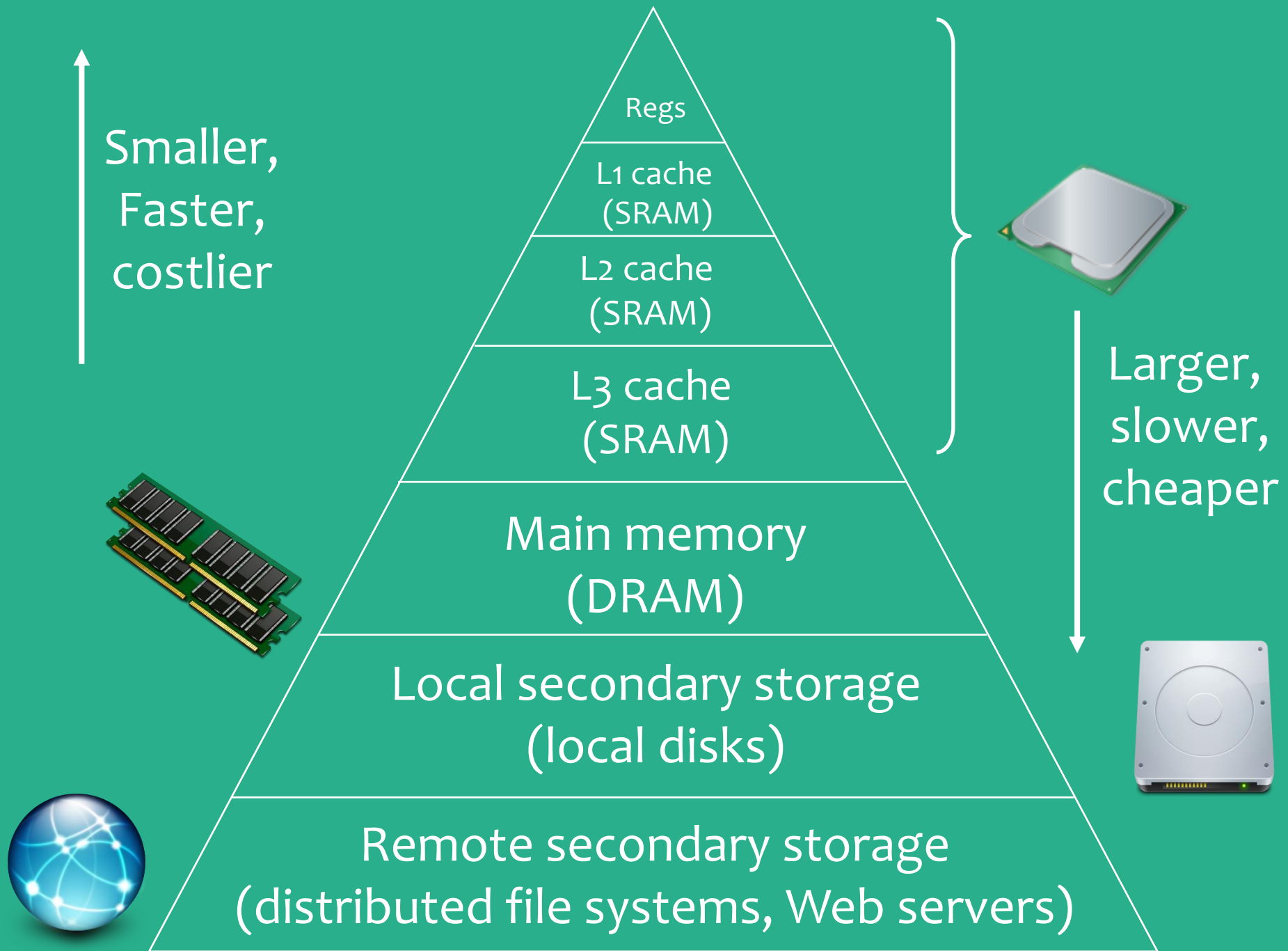
Ver. 1.78.1.x64

Tools

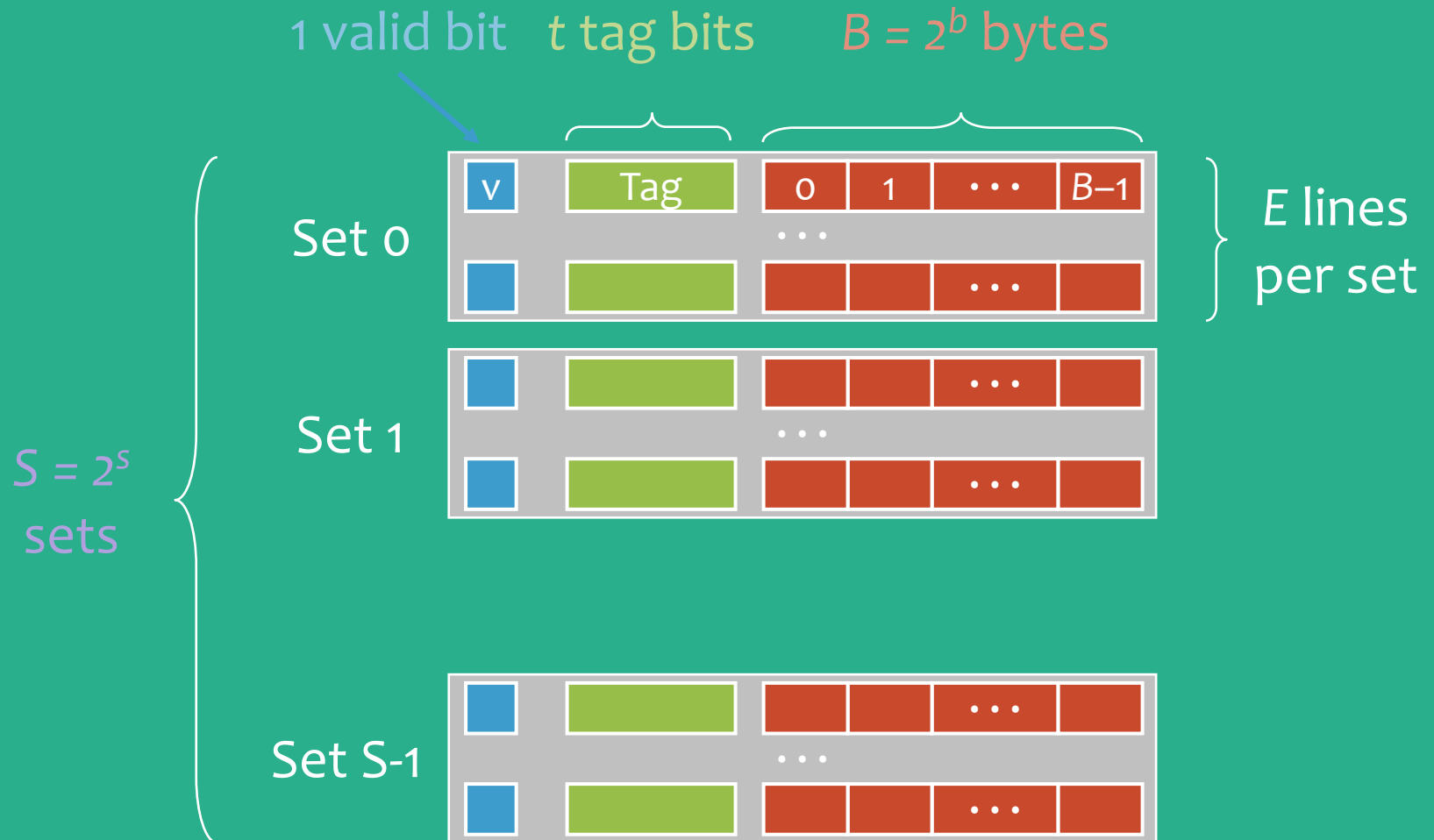


Validate

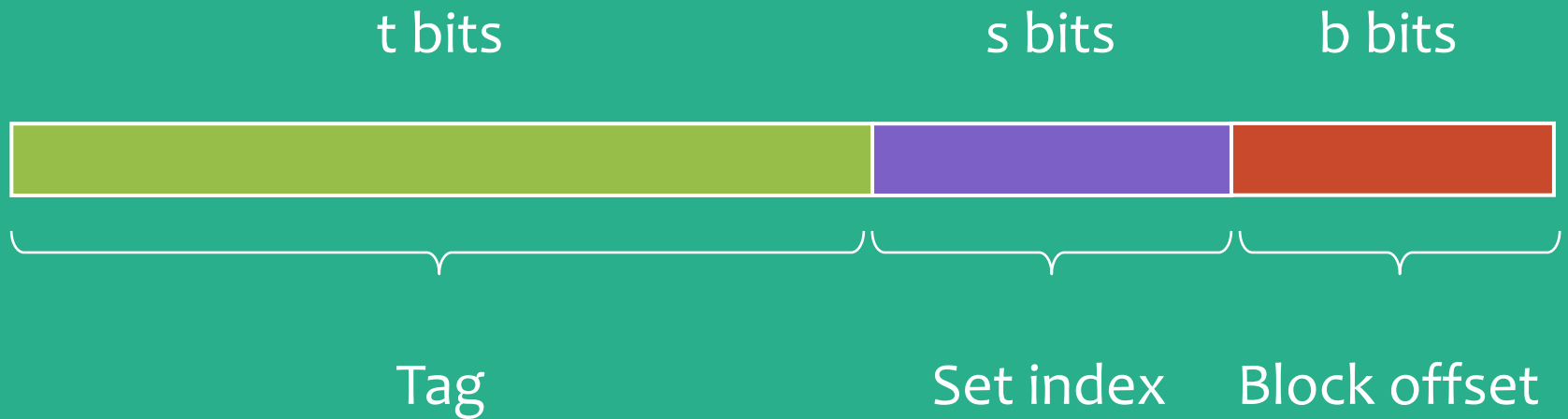
Close



Cache organization



Partition of Address



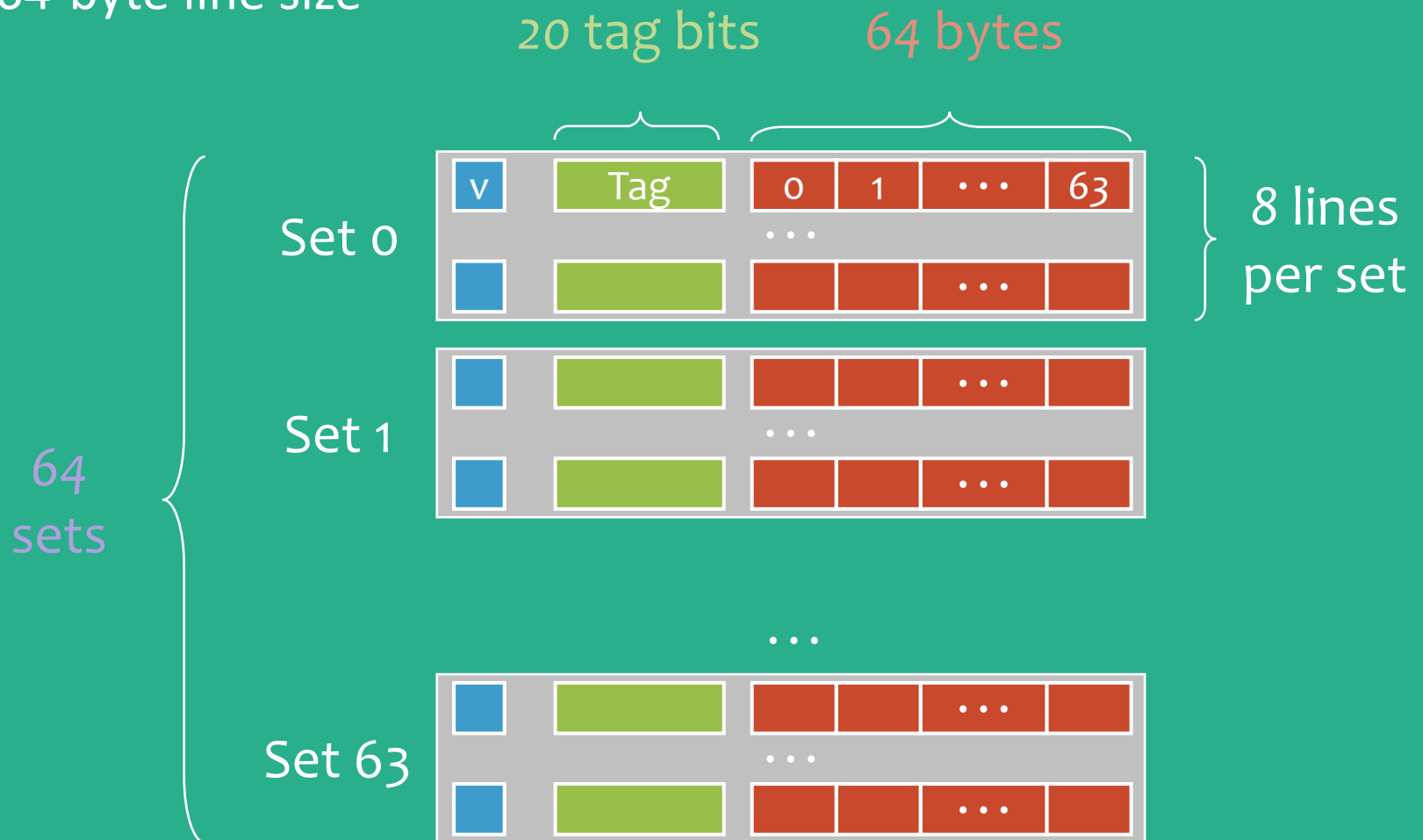
Memory address has 32 bits, determine S, t, s, b?

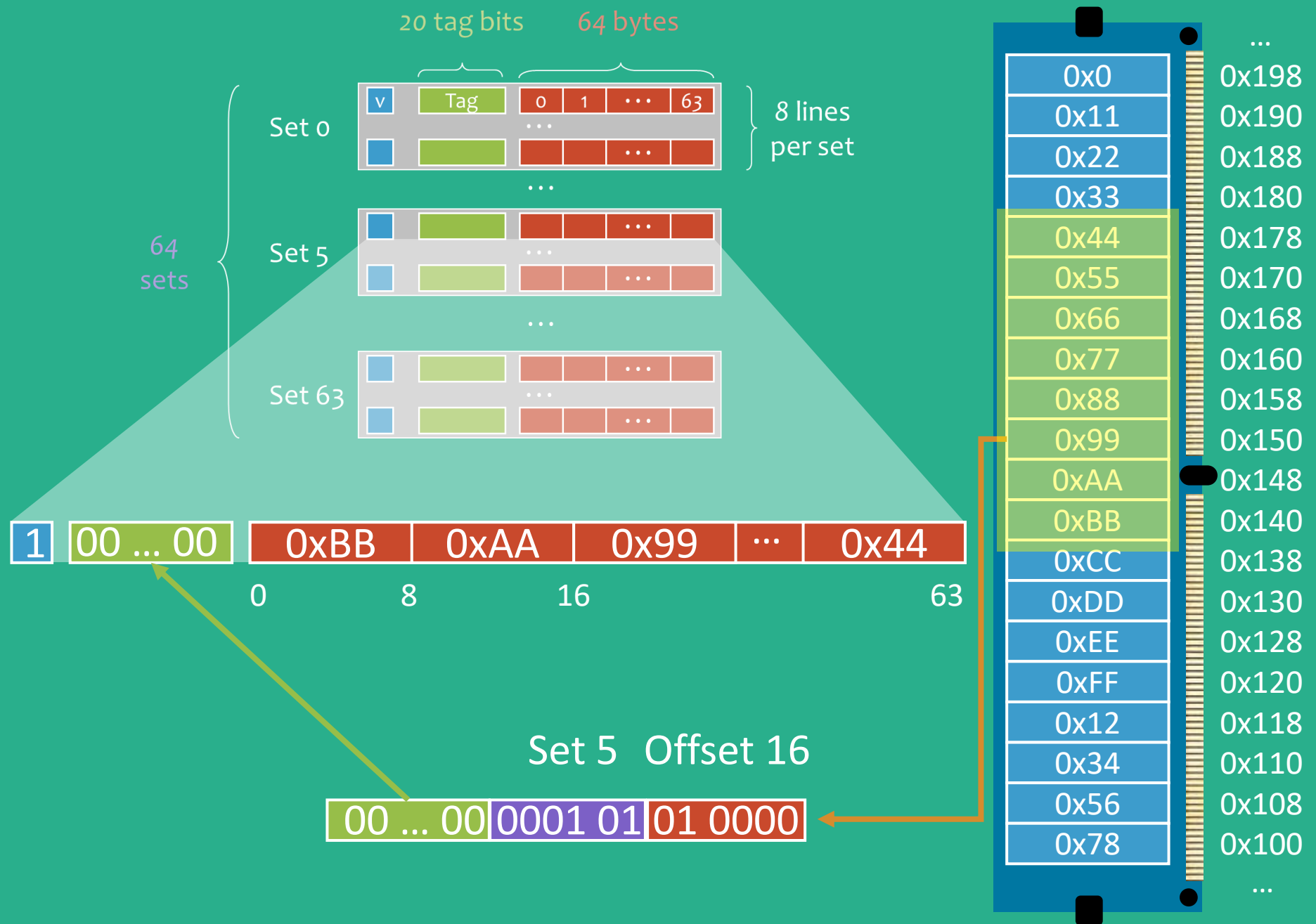
C	B	E	S	t	s	b
1,024	4	1				
1,024	8	4				
1,024	32	32				
32KB	64	8				
246KB	64	4				
8MB	64	16				

32 KBytes

8-way set associative

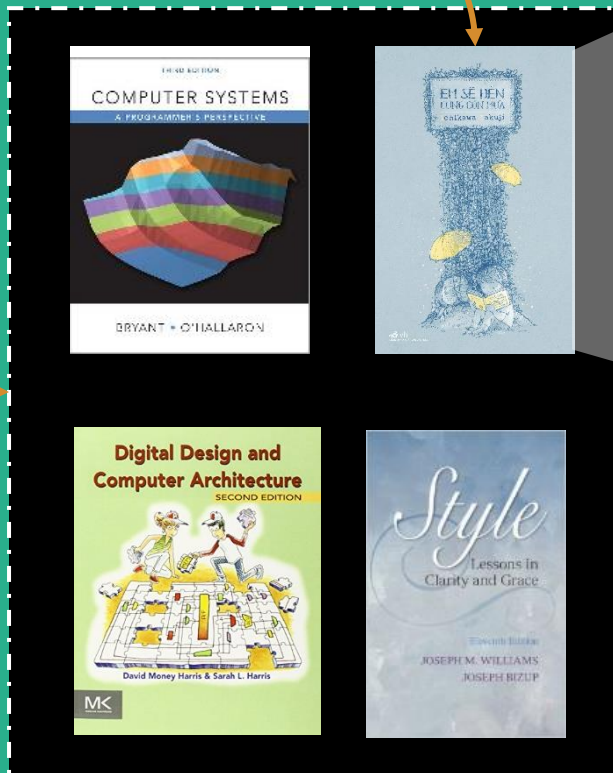
64-byte line size



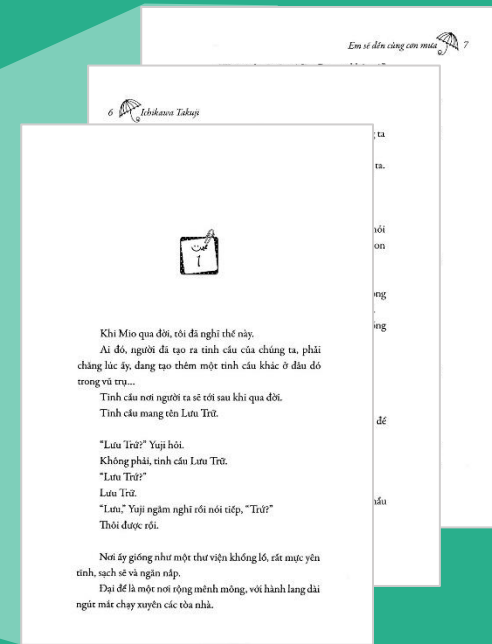


Line (Book)

Set



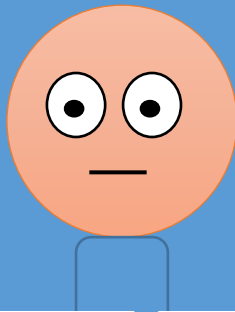
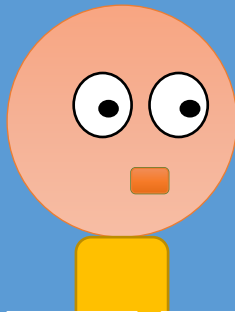
Block (Pages)



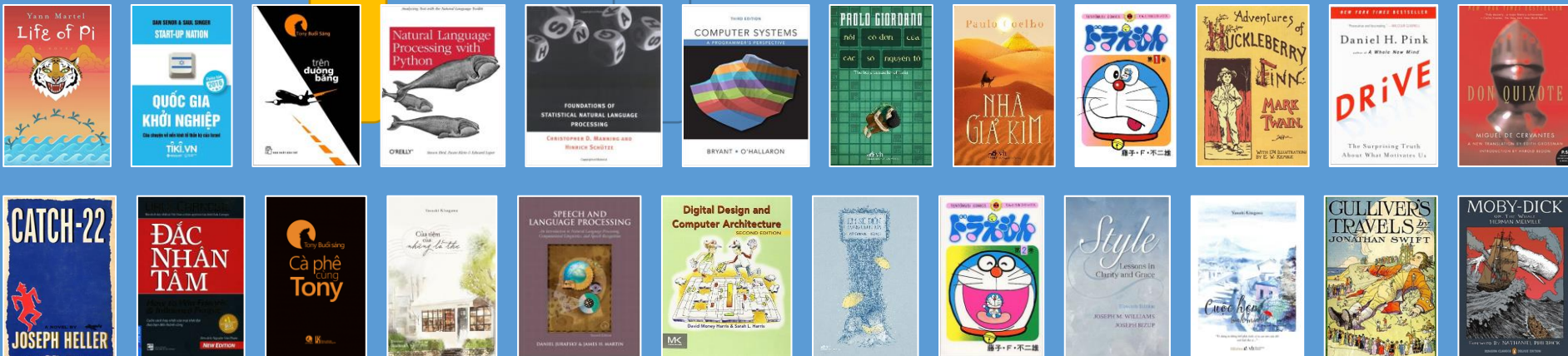
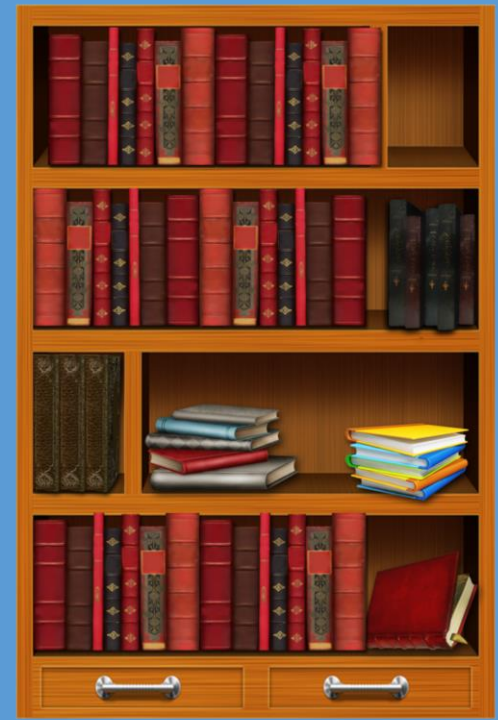
Cache Concepts

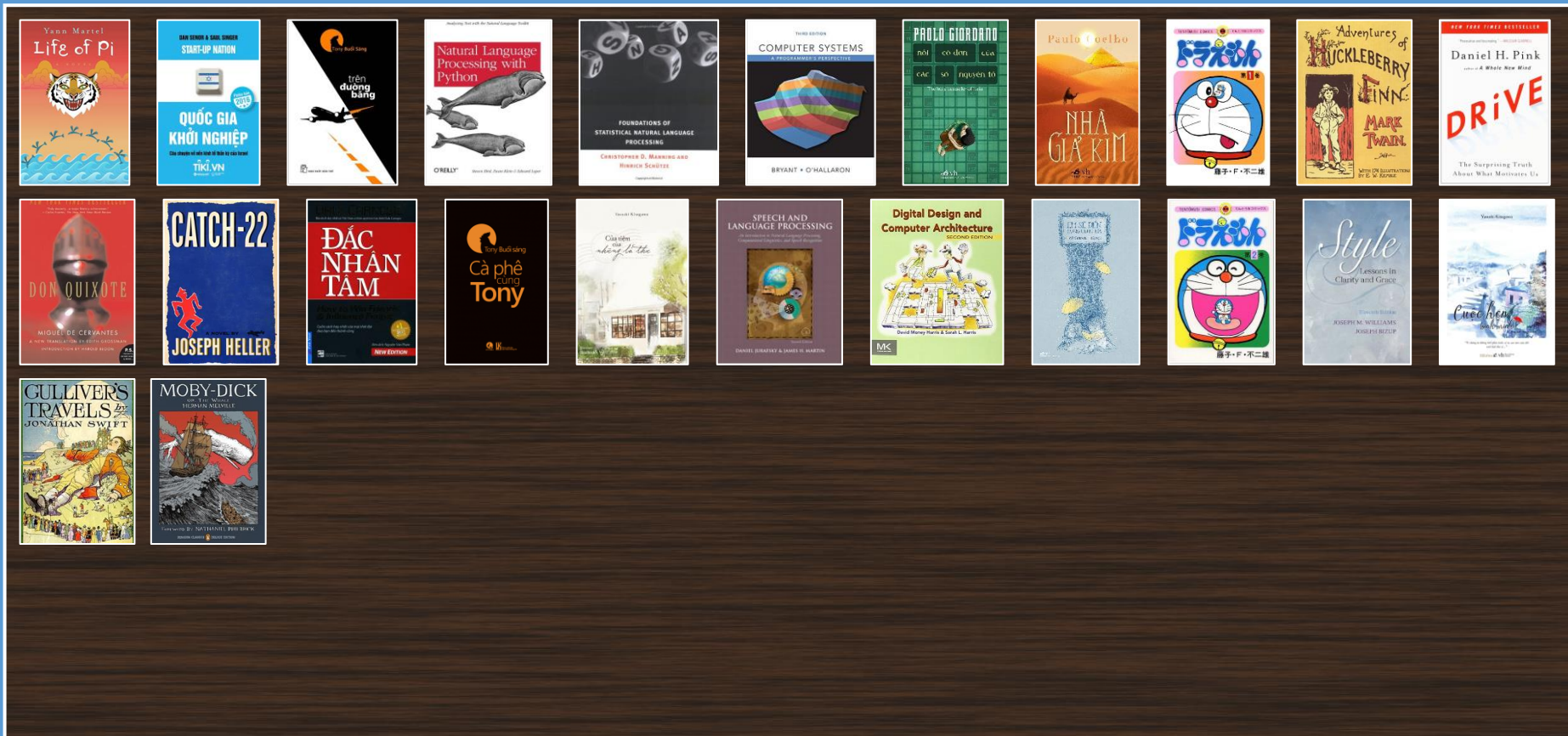
Memory

CPU



Cache



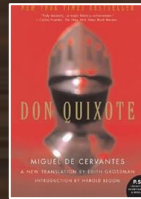
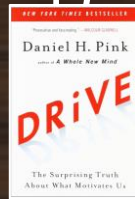
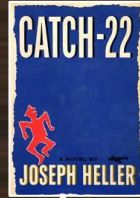
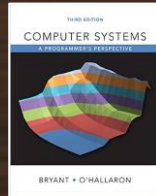


Put any book anywhere

→ hard to find

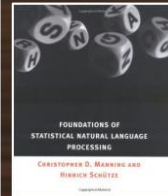
A

B



E

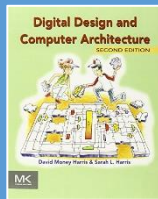
F



...

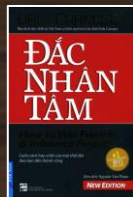
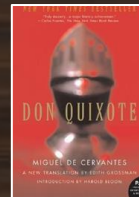
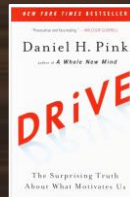
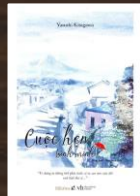
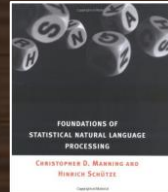
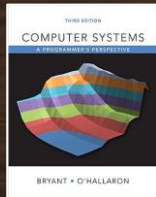
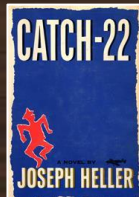
Y

Z



conflict

→ easy to find

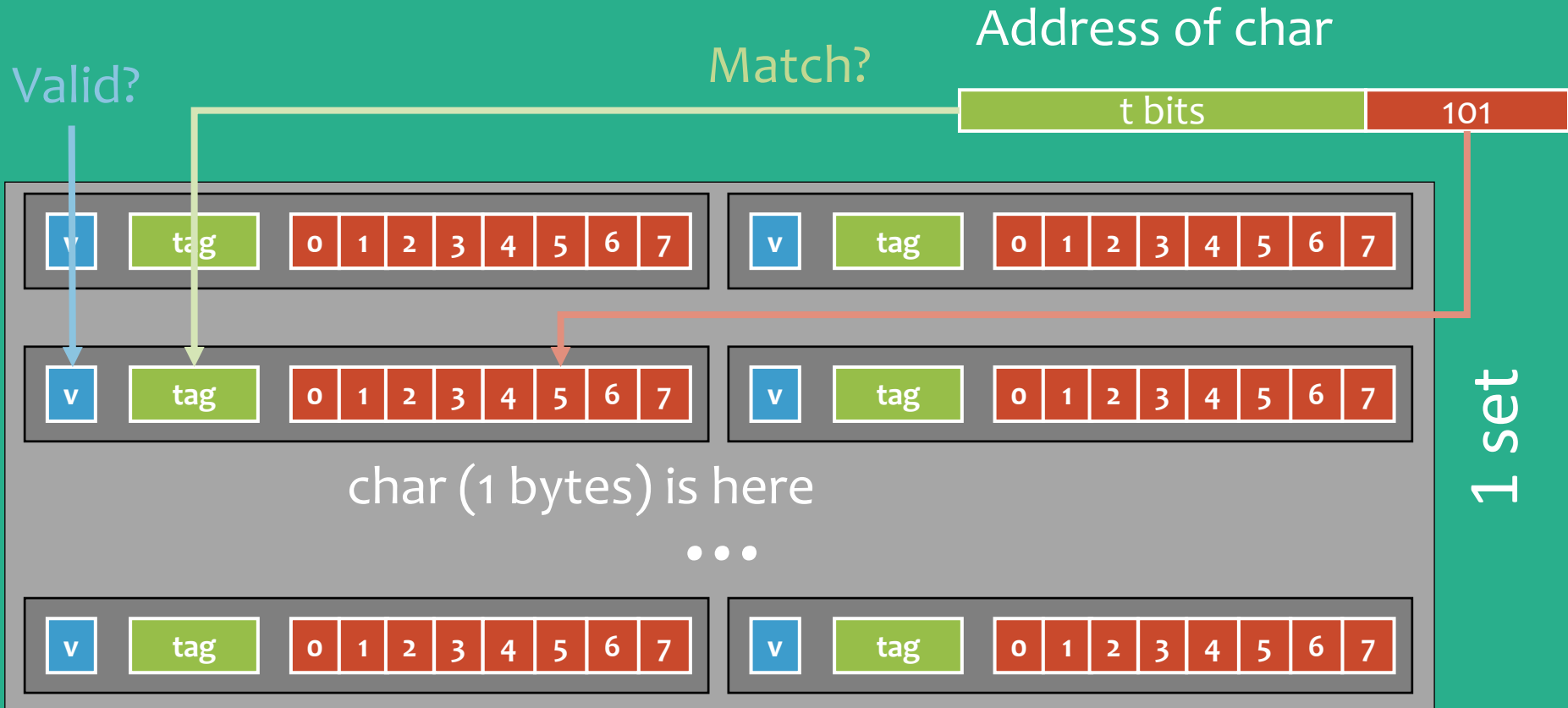


A or B or C

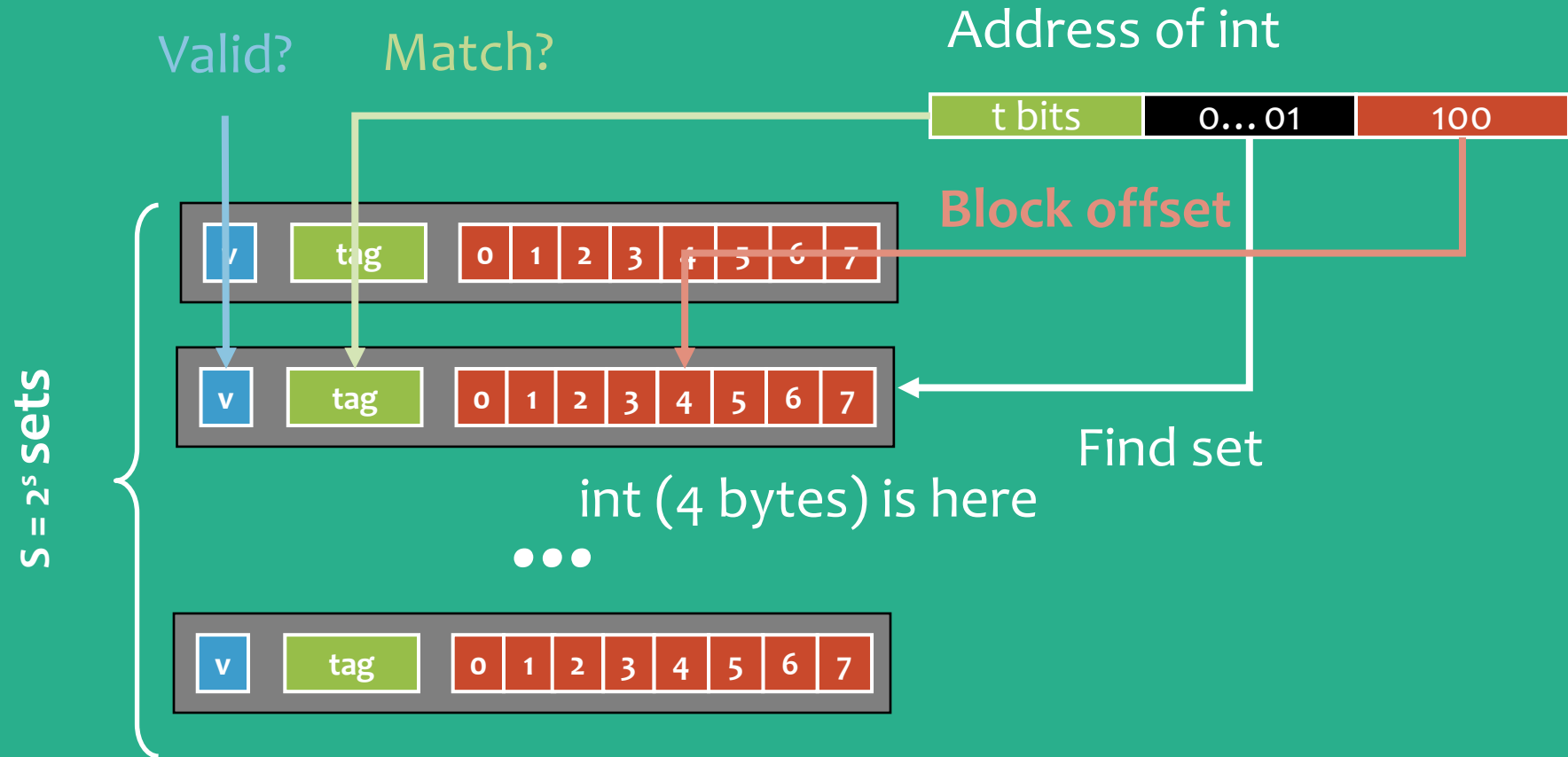
D or E or F

Y or Z

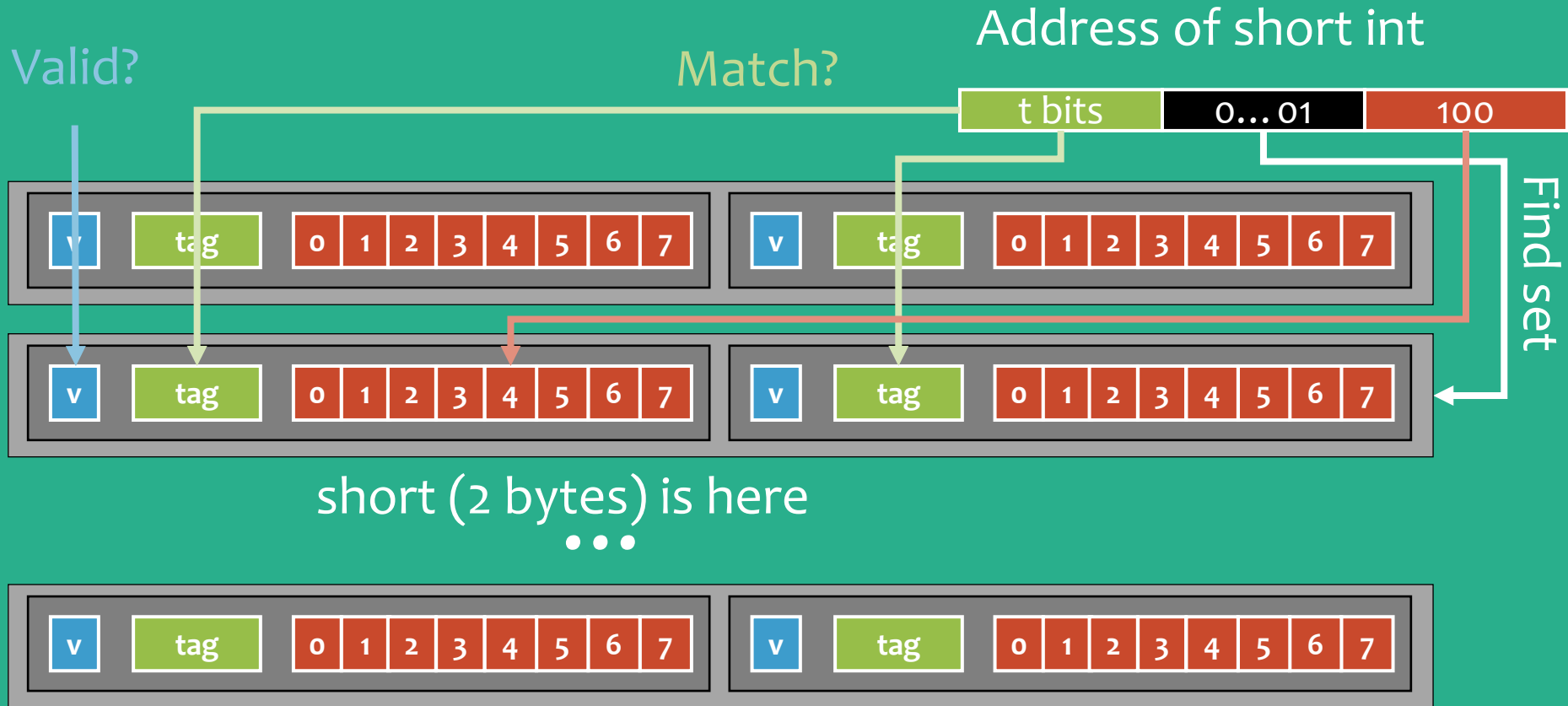
Fully-Associative Cache



Direct-Mapped Cache

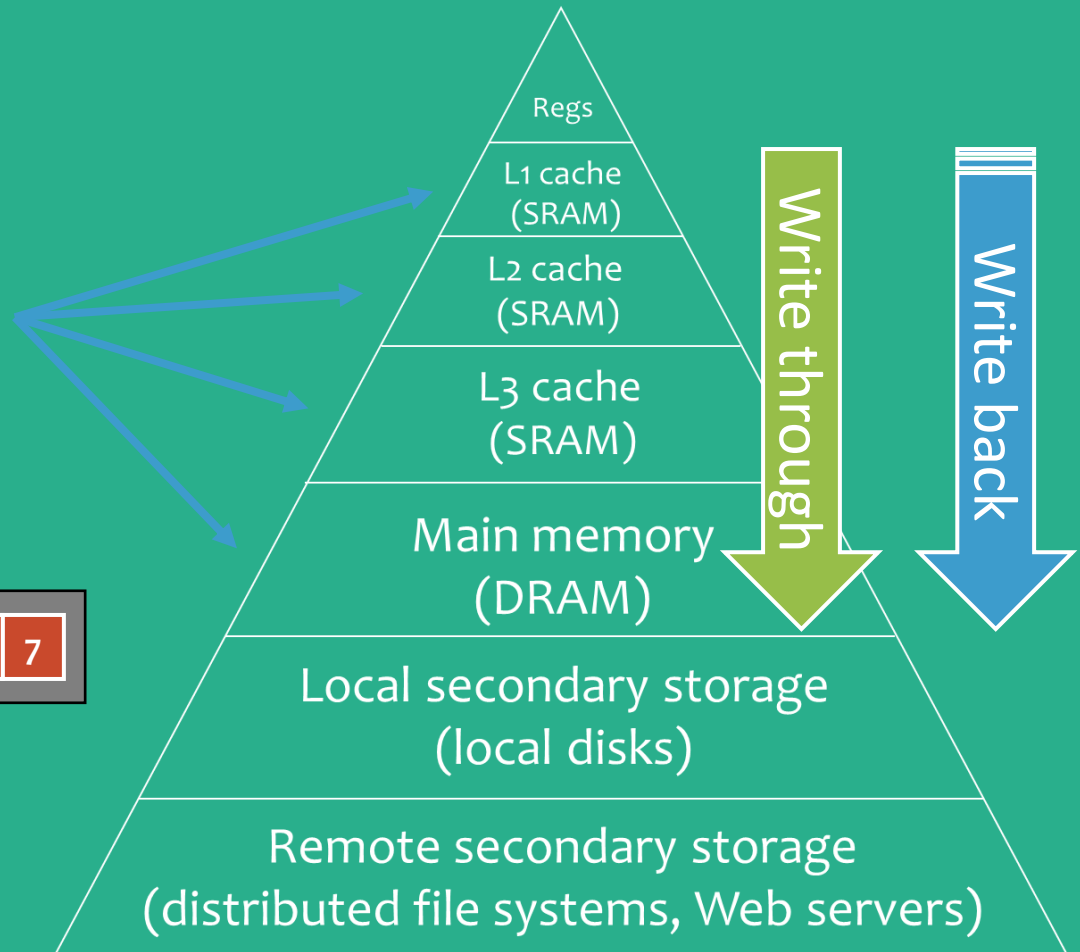


2-Way Set-Associative Cache



Valid bit

multiple copies



Summary

- Cache Concepts
 - Cache Hit
 - Cache Miss
- Cache Organization
 - Direct-mapped
 - E-way Set Associative
 - Fully Associative



Gene Myron Amdahl

formulating Amdahl's law

“

$$S_{\text{latency}}(s) = \frac{1}{(1 - p) + \frac{p}{s}}$$

”