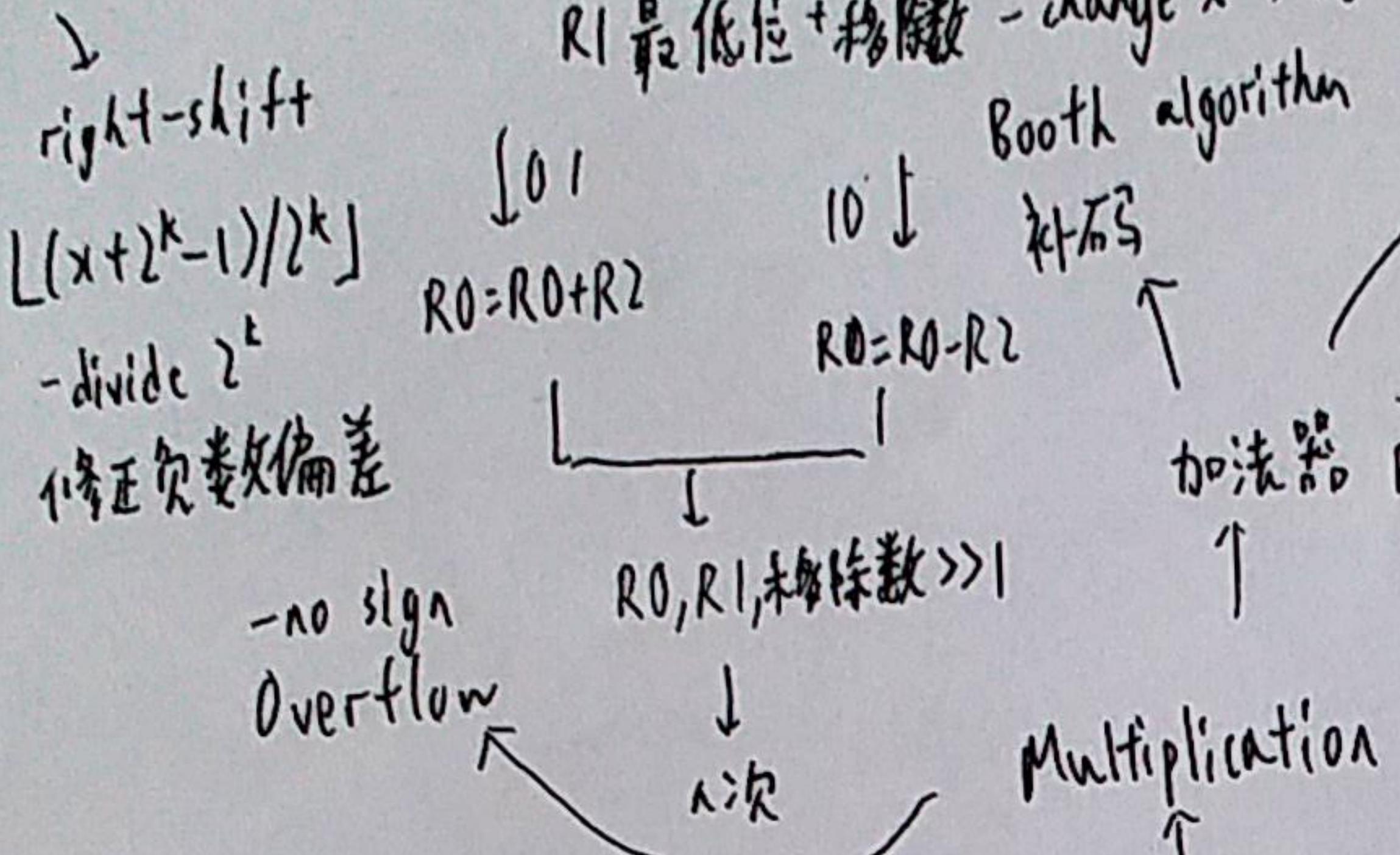


signed bit-shifting

left-shift  
works for + -



throw away higher bits  
\* may cause data loss

Type conversion

signed extension  
前面加符号位 ← 位扩展 Bit Extension

0扩展 Unsigned  
前面加0

$$[X]_{\text{补}} = \begin{cases} X, 0 \leq X \leq 2^n \\ 2^{n+1} + X, -2^n \leq X < 0 \end{cases}$$

$n$ : 模 Modulo

$$[X]_{\text{补}} - [Y]_{\text{补}}$$

11

$$[X]_{\text{补}} + [-Y]_{\text{补}}$$

\* can be directly calculated

\* 补码 ← Two's complement

1) positive:  
unchanged

2) negative:  
取反加1

last 1 unchanged,  
reversed rest of front

Coding → 原码

反码

One's complement  
- first bit as sign  
(0:+ ; 1:-)

binary

Signed magnitude  
- first bit as sign  
(0:+ ; 1:-)

/  
has 2 0s

算术

arithmetic shift  
if sign bit=1

add 1 as highest

bit

same sign addition

different sign subtraction

↑

检测

溢出 Overflow

↑

加法器 / 减法器

→ use xor gate to

determine sign

将 k 设为 00/11, if 01,10 overflow

或

进位位相异或 Cn xor Cn-1 (↑R)

→ 符号位判断 (个量)

↑

加法器 / 减法器

→ use xor gate to

determine sign

Addition /

Subtraction

↑ not

Logical

逻辑

→ || or

→ && and

~ one's complement 取反

↑ ↑ ^ XOR & > ^ > | > & & > ||

& ← Bitwise → | or

按位

(x 2<sup>k</sup>)

左移 << k

right-shift

(divide 2<sup>k</sup>)

left-shift

- logical = arithmetic

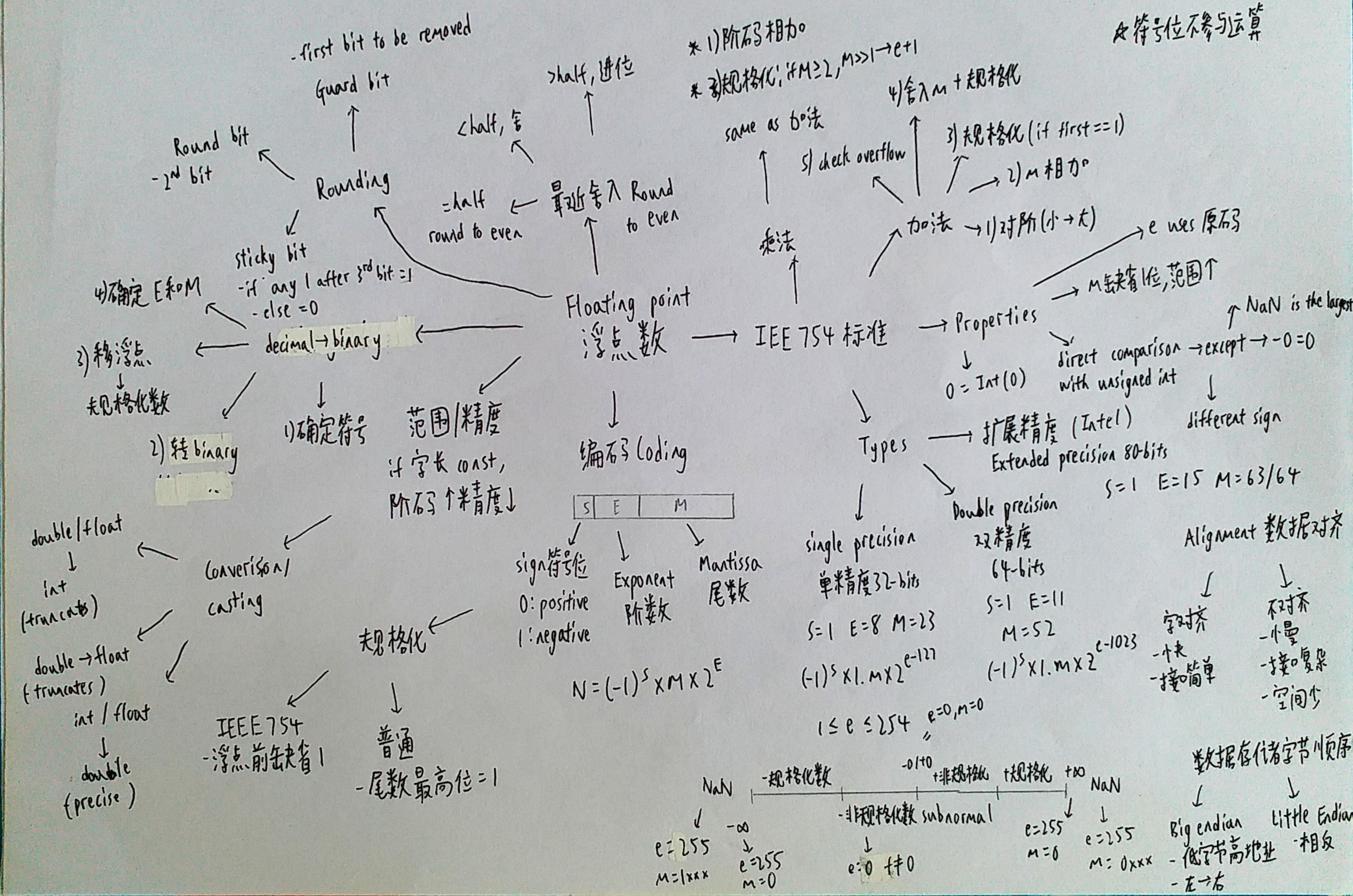
- may overflow

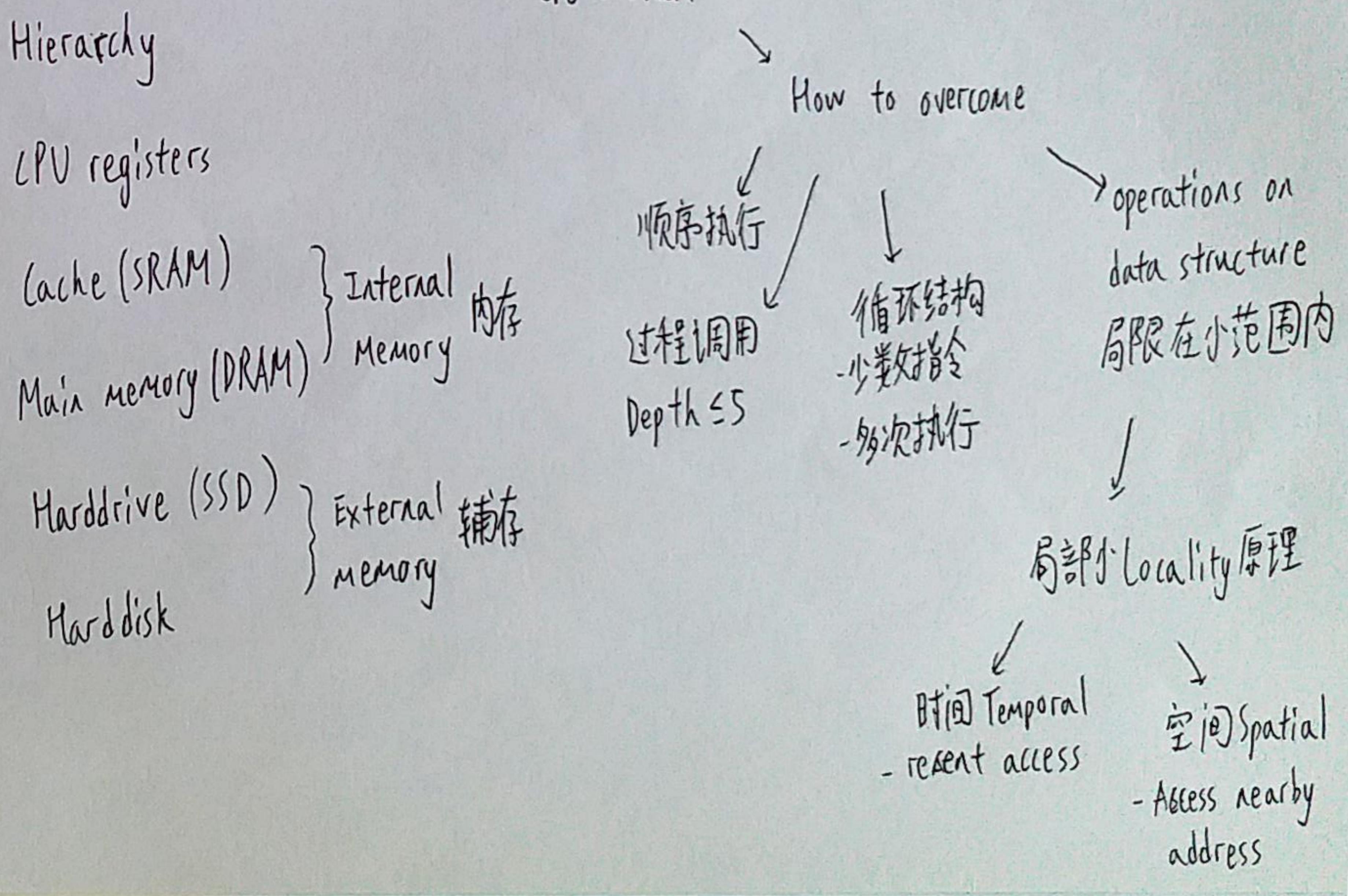
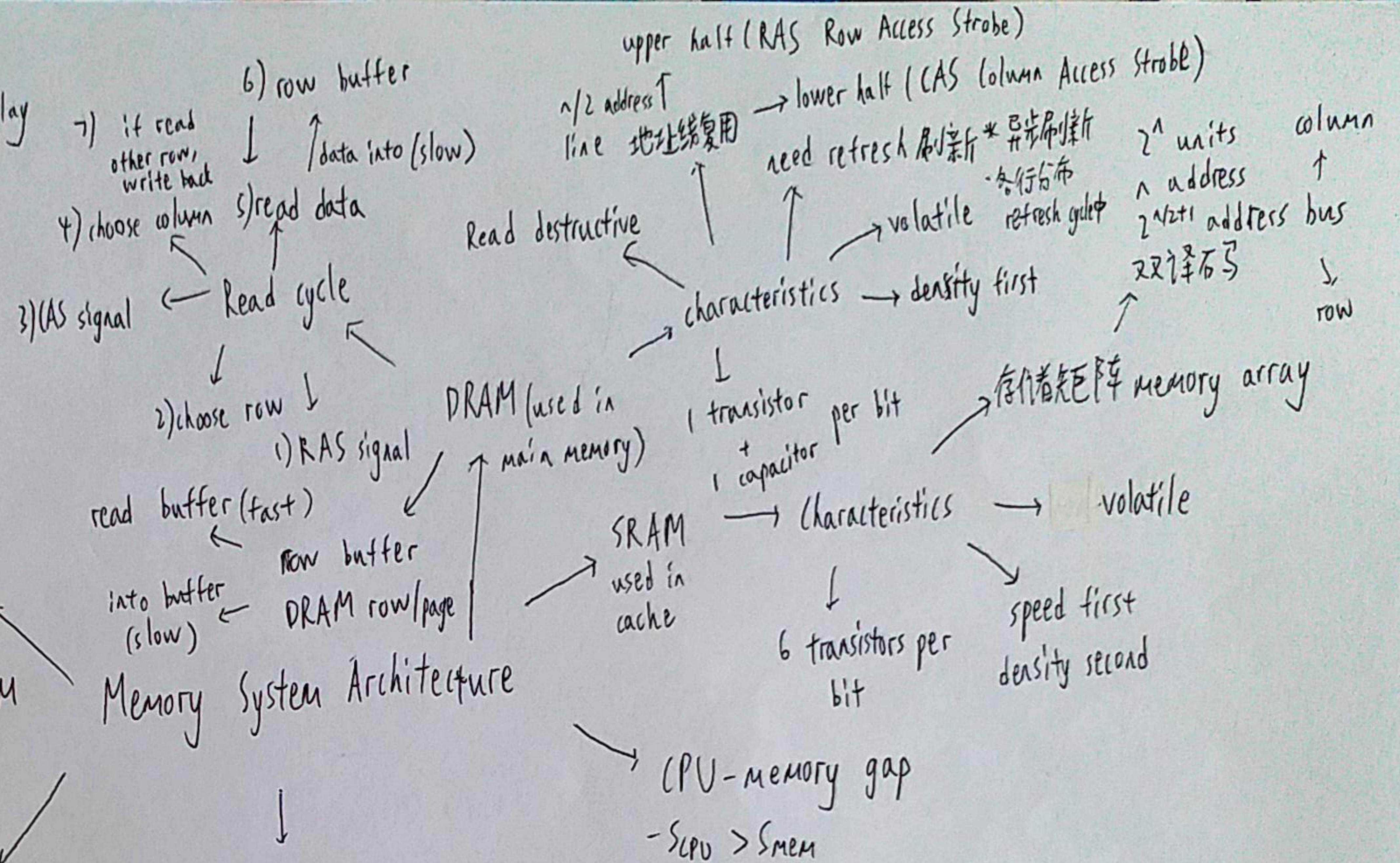
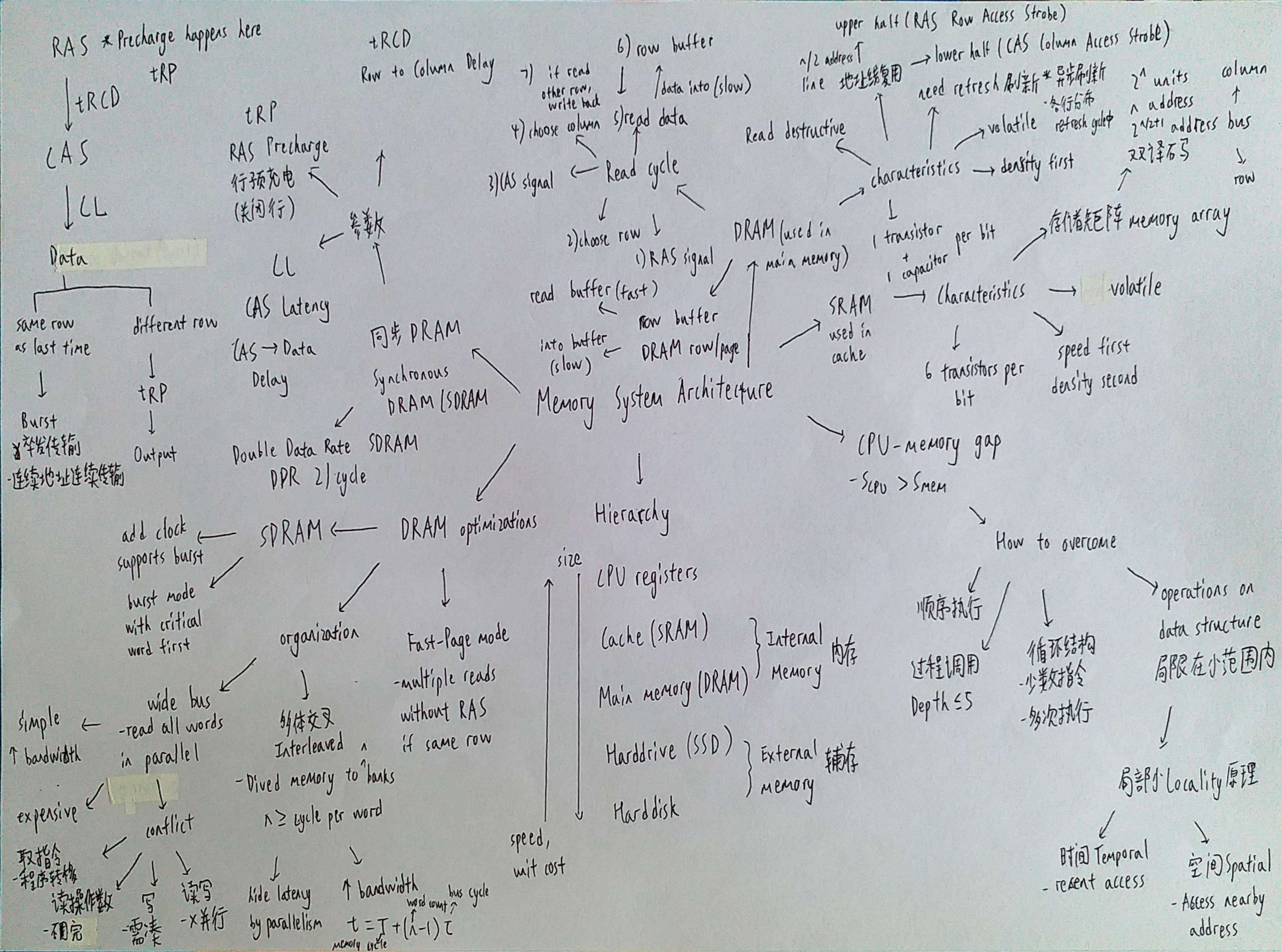
逻辑

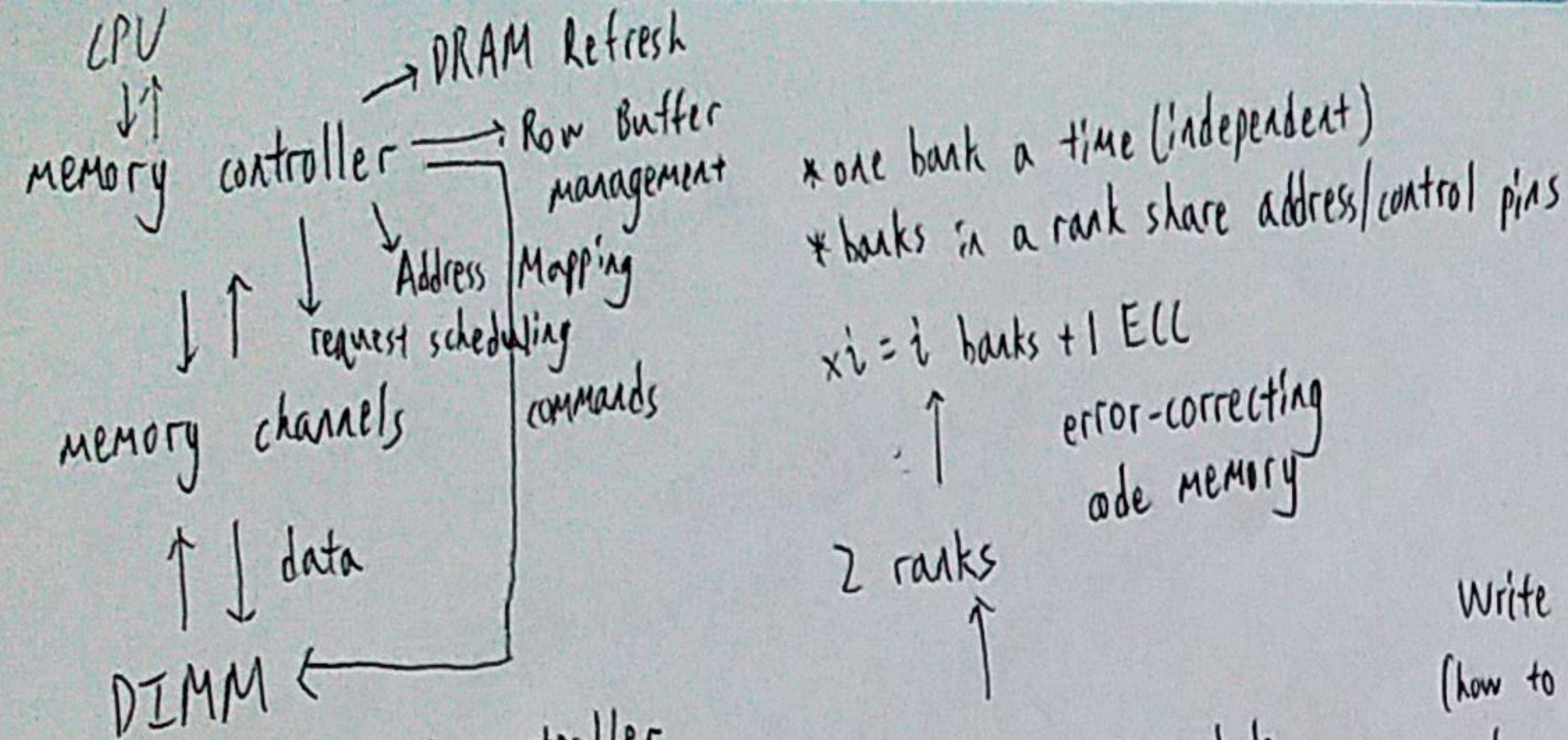
logical shift

shift in zeros

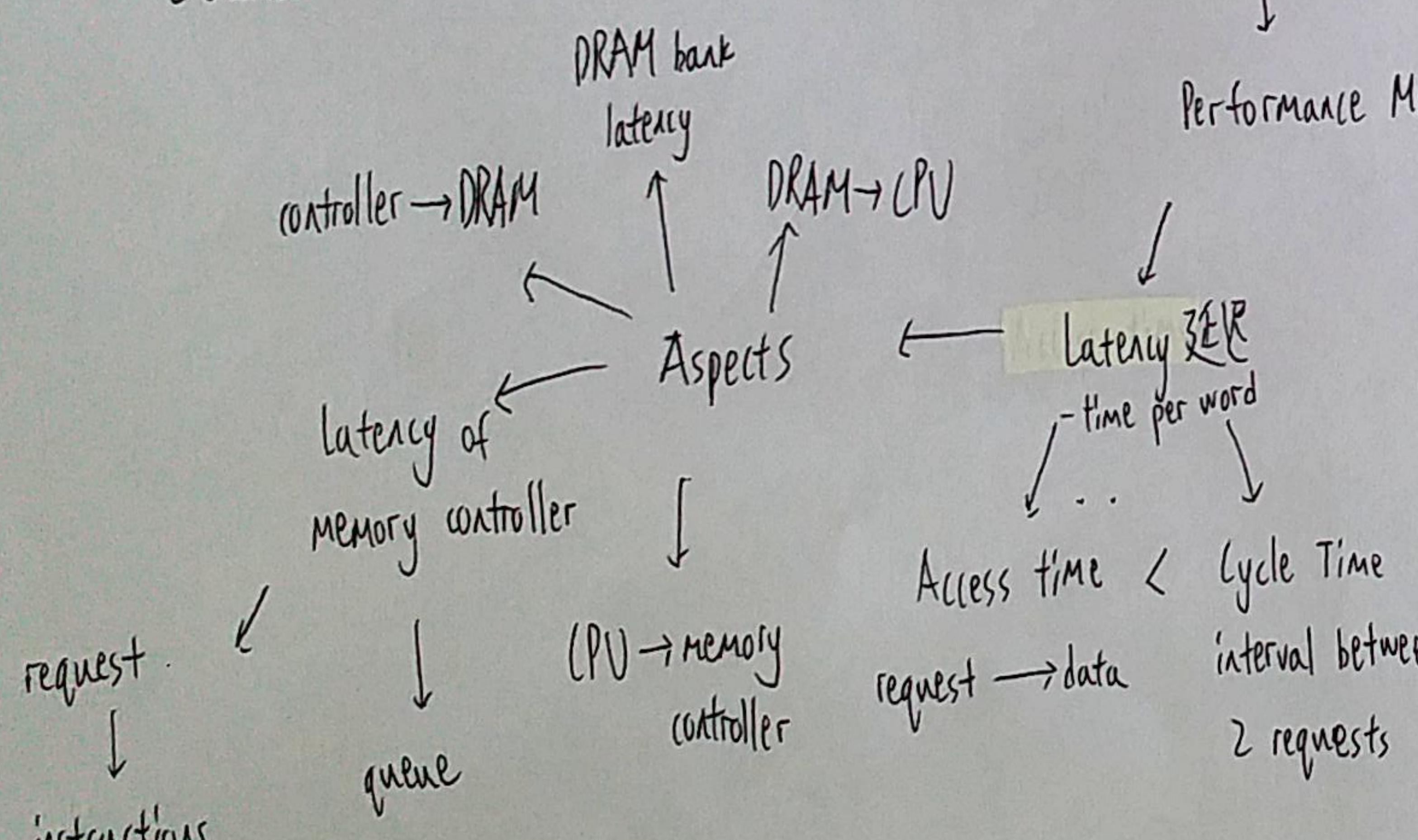
no matter what







DIMM  
 | controller  
 | channel  
 1 controller  
 2 channels ← memory channels ← DRAM organization  
 ↘      ↓  
 2 controllers      ↑ channels → ↑ bandwidth  
 2 channels



$t_a = H \times t_A + (1-H) \times t_M$  :  $H = \frac{\text{cache access no.}}{\text{cache access no.} + \text{memory access no.}}$

Bit extension  
(extend data width)

- parallel memory array

address

↓

no. of chips =  $\frac{N_{CPU}}{N_{MEM}}$

total      chips      decoder  
Tl            II             $k \rightarrow i$

↓      ↓      ↓

address line       $\log_2 N_{CPU}$        $\log_2 N_{MEM}$       total  
directly connect      to every chip      - chips

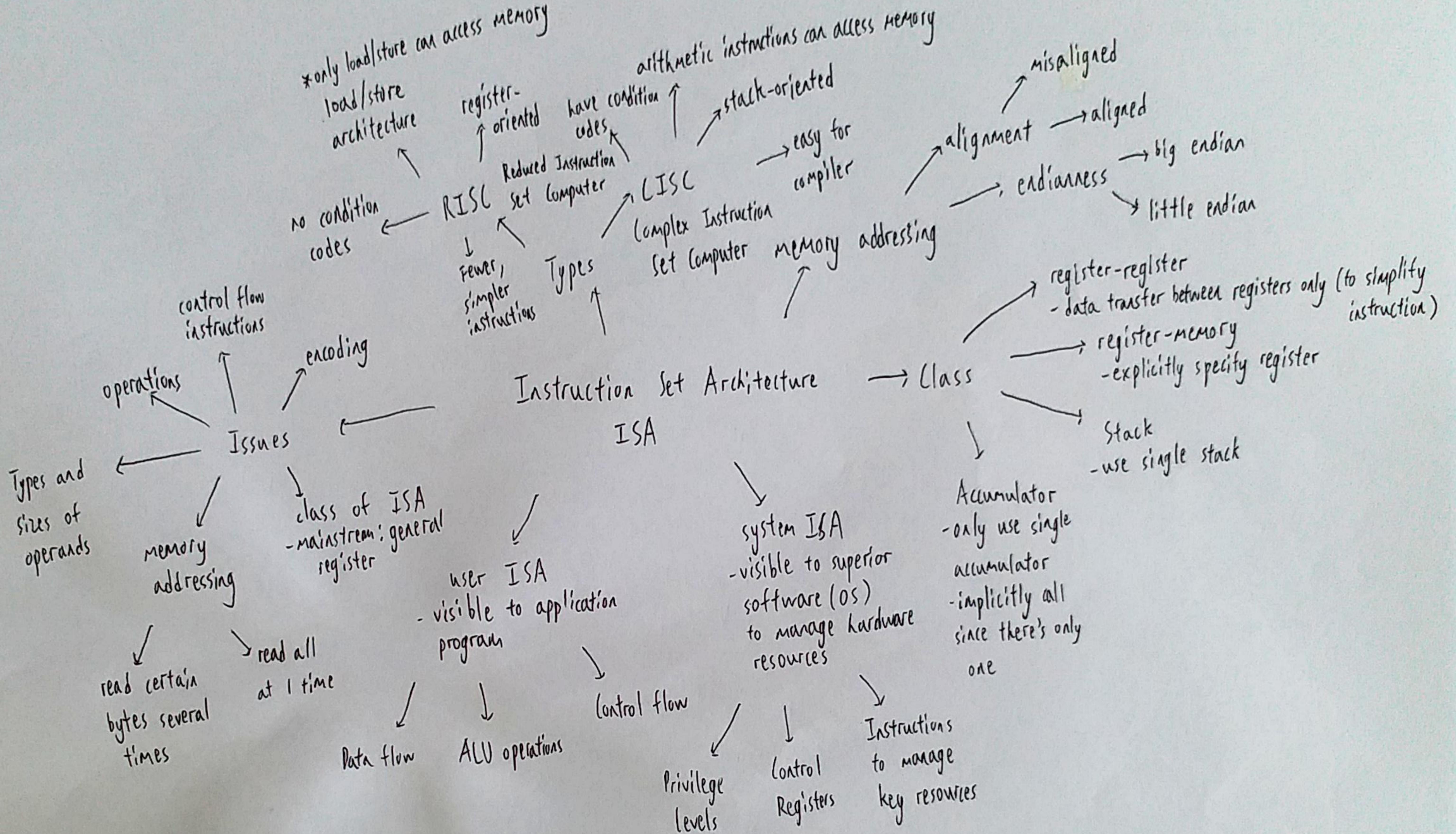
↓

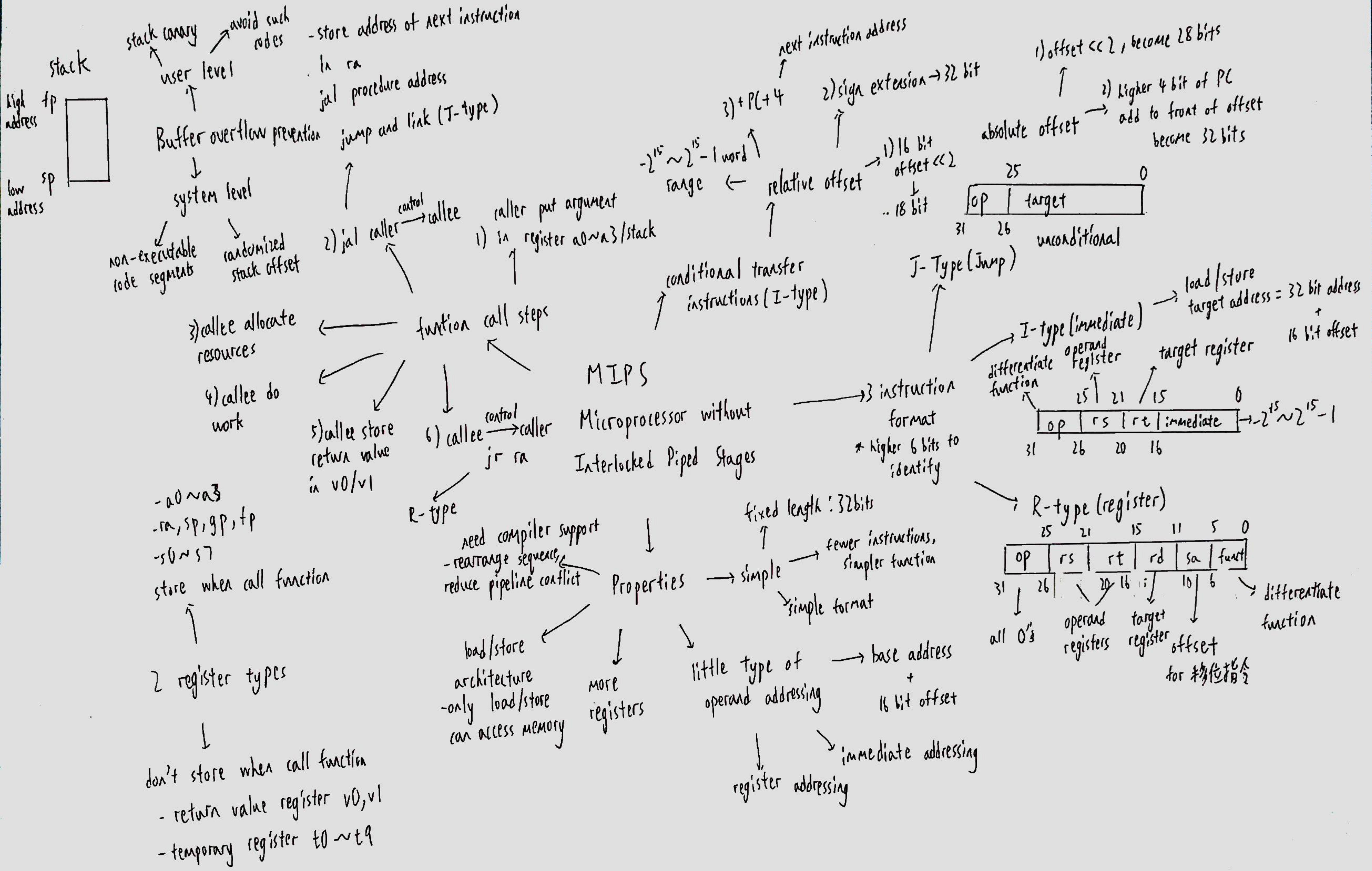
a lines      no. of chips

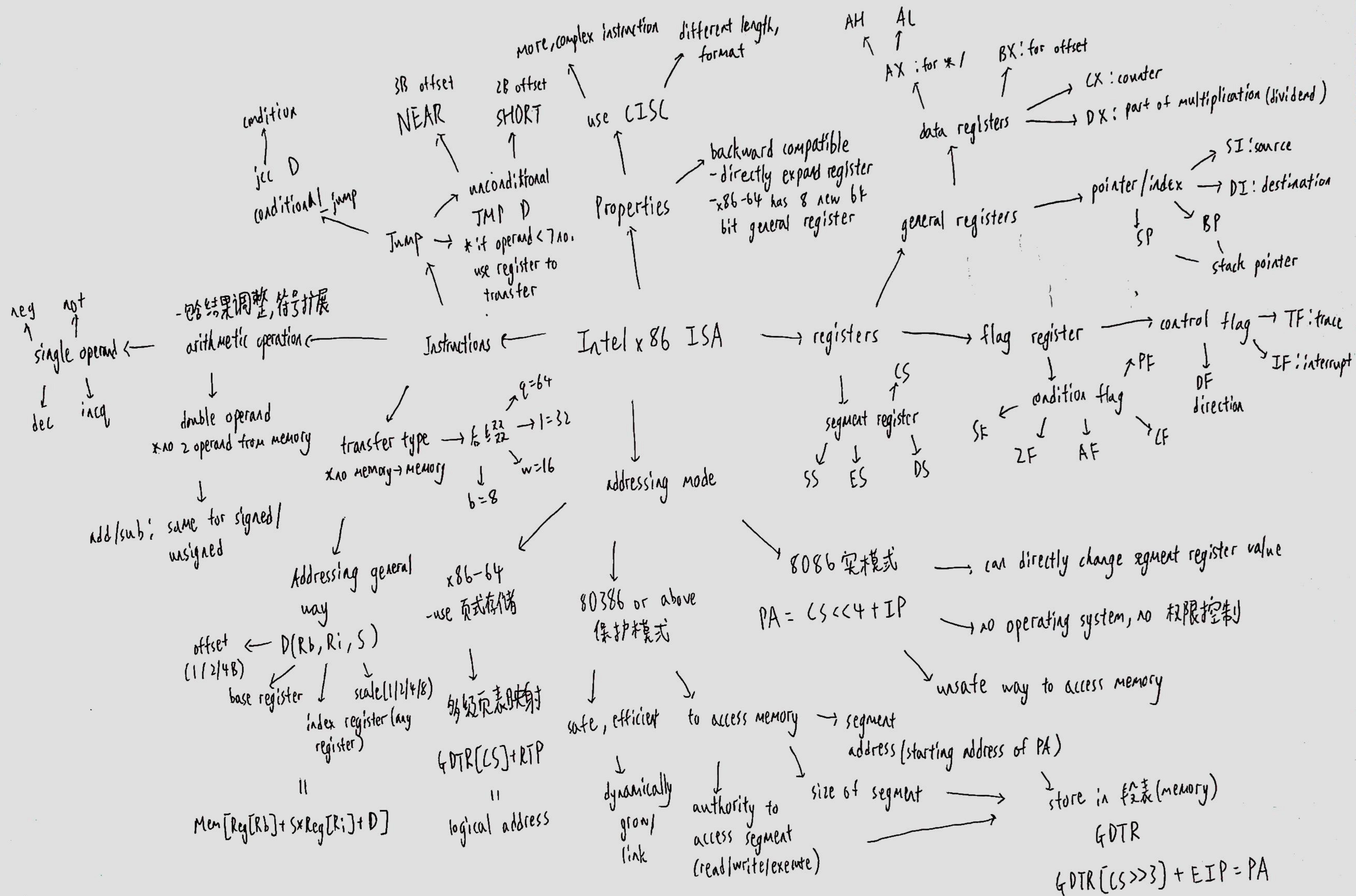
↓

outted      - chips

chip      =  $\log_2 N_{CPU}$







路↑, 组↓, 关联度↓

\* if 组 = 1 → fully-associated

关联度

路↓, 组↑, 关联度↓

\* if 组 = 1 → direct map

- directly change

memory

No-write-allocate security

Data → cache write-allocate \* performance  
change in cache 写分配法

good for 连续访问

hit  
write-back  
- write row in cache

dirty bit  
- check it  
changed

if cache refresh  
- directly write  
memory + cache

main memory

分 M 区

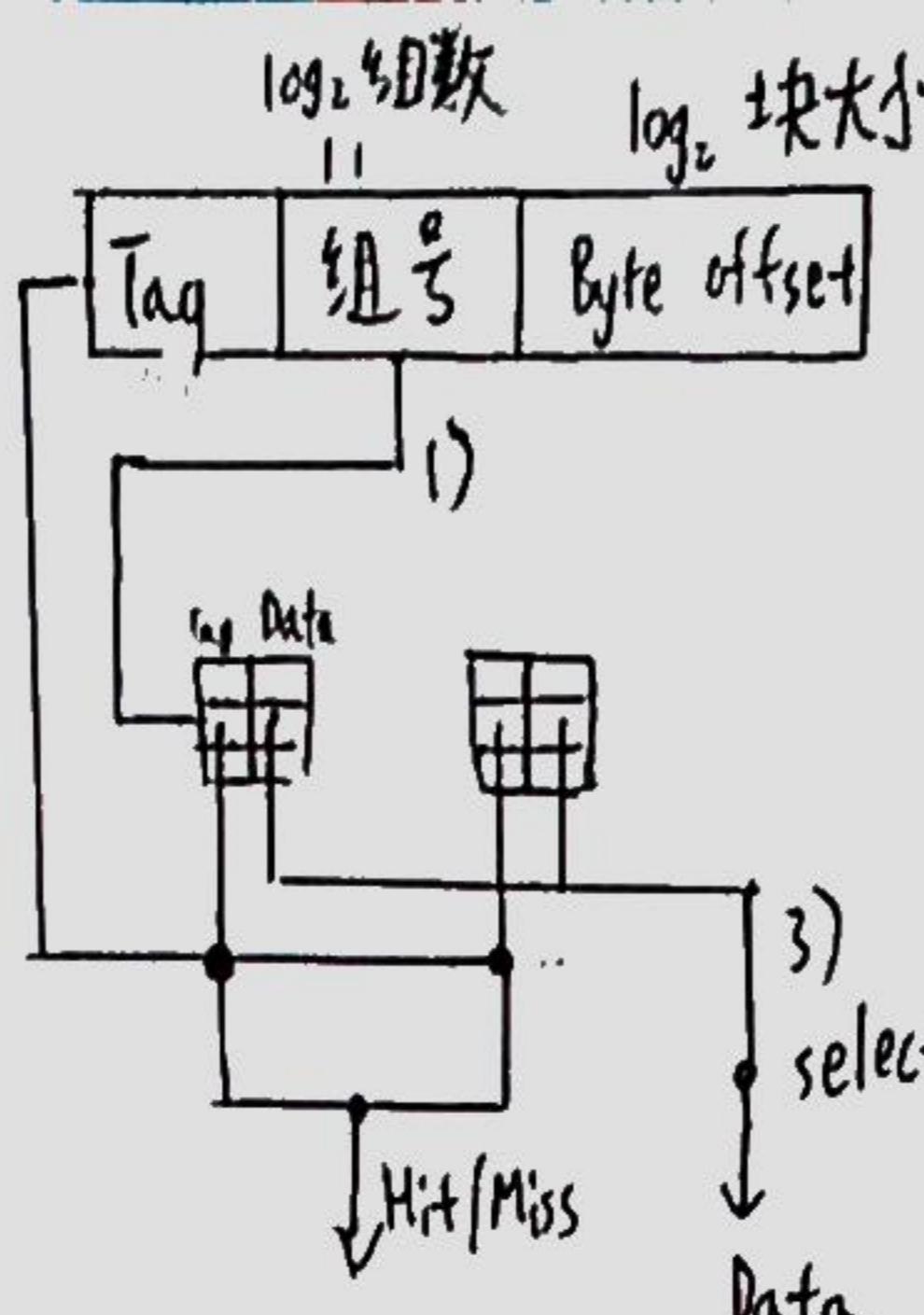
区分 N 组 (direct map)

组分入块

fully-associated

$$\text{行数} = \frac{\text{Total capacity}}{\text{Data bus}}$$

$$\text{组数} = \frac{\text{Data bus}}{\text{way}}$$

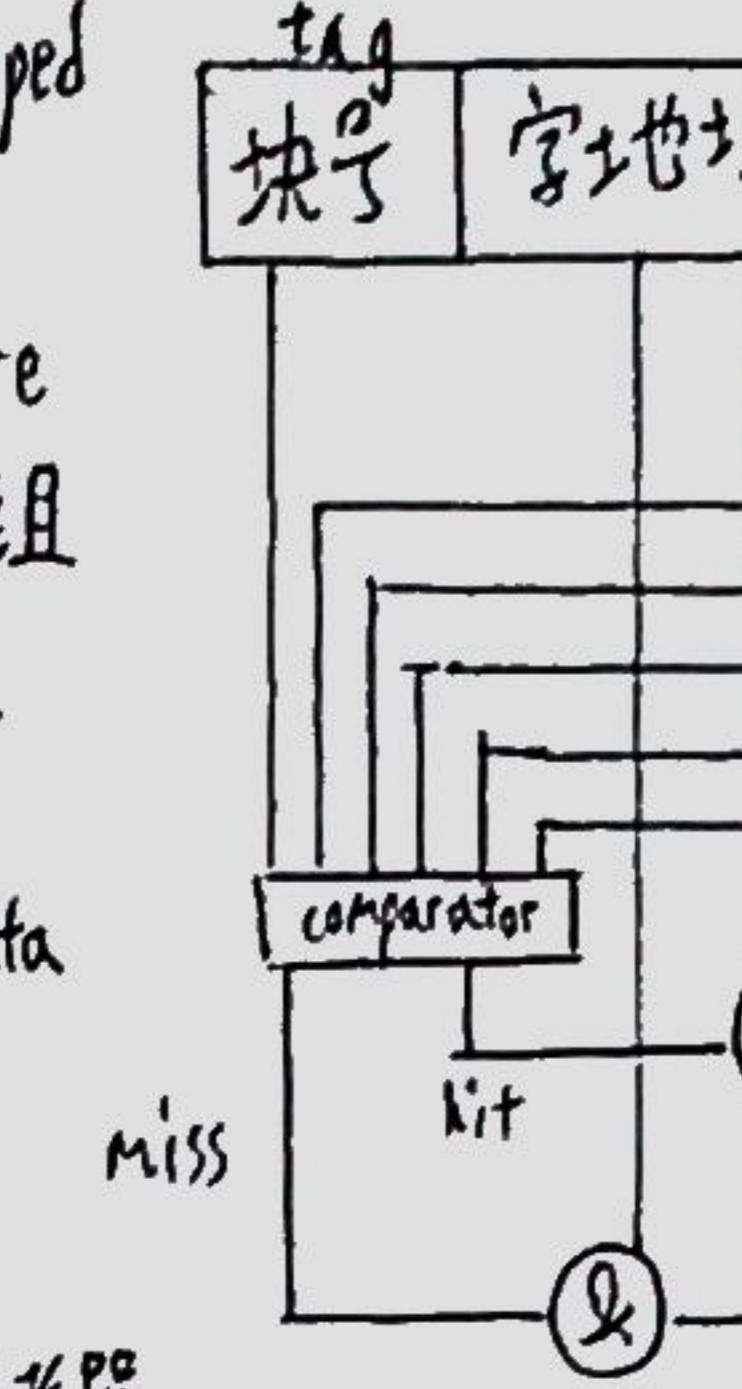


1) 组号 Direct mapped to cache

2) Tag → compare with rows in 组

3) tag same, hit

N-way → 1 data



cache

memory

conflict only if full

↑ cache use rate

one-to-many

→ advantages

hit rate ↑

→ disadvantages → hardware cost ↑

complex ↘

Main memory

divided to N 区

(same size with cache)

Cache Write Policy

- uses binary tree O(n)

LRU

Cache 替换策略

LRU  
least recently used

FIFO

随机替换 - faster

直接 - fast

完全 / set associated

hit rate ↑

hit rate ↓

硬件实现

$\log n! \approx O(n \log n)$

状态位

- 记录访问顺序

cache access time

Memory → Address Mapping (Cache)

upper address into storage

Tag blockid

区号 块号 字地址

valid bit

Direct Mapped 直接

→ Main memory

divided to N 区

(same size with cache)

Cache

Tag

in cache

not in cache

Memory

advantages ↘

fast

simple

disadvantages ↘

conflict ↑

low hit rate

↓ cache use rate

1) use block id to find which block

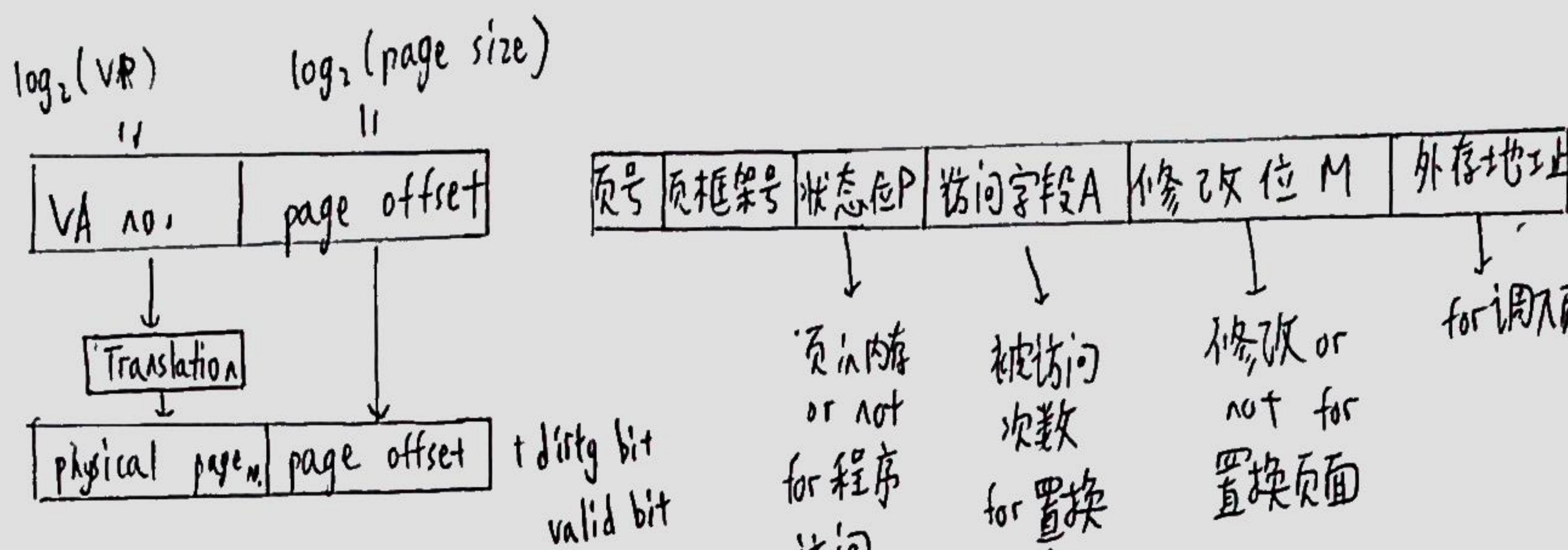
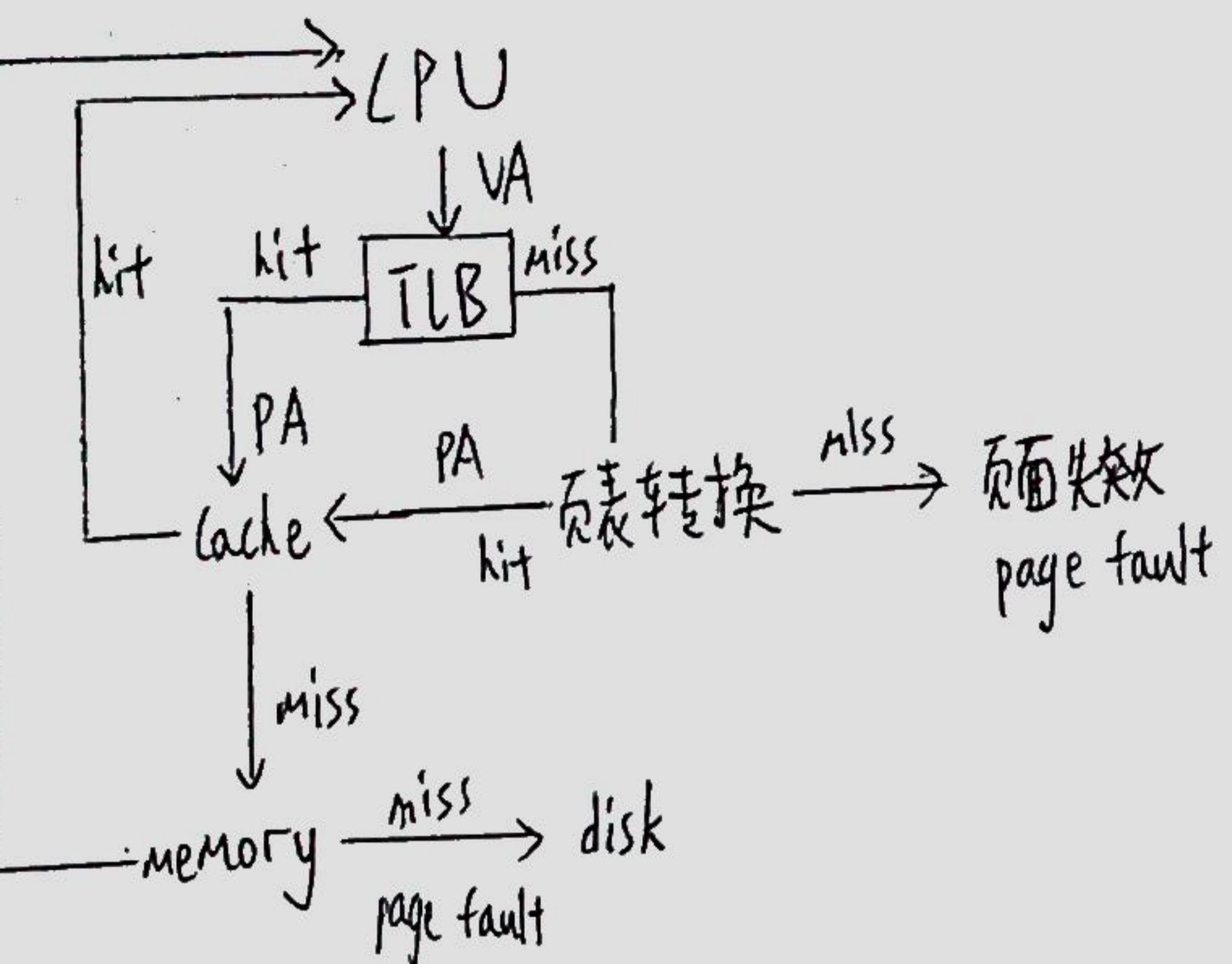
block id =  $\log_2(\text{行数})$

2) use tag to find which line

tag =  $B \text{ addr} - (\log_2(\text{行数}) + \log_2(\text{word in block}) + 2)$

3) use offset to find which word

offset within B



缺页中断请求  
used page not in memory

run program

### 地址转换策略

一级  
指向地址 ← 多级页表 ← 页表机制  
二级  
→ 页地址

→ 请求分页技术 → page into memory

Inclusive caches  
simplifying finding cache  
block by other entity  
Multi-level  
caches  
wastes capacity  
low-level has copy of higher-level  
only stored in 1 level  
capacity ↑  
exclusive caches  
direct mapped  
set associative  
need uniform block size  
addressing  
VI-PT  
Virtually Indexed  
physically Tagged  
no aliases  
VI-VT  
virtually indexed,  
virtually tagged  
hit time ↓  
contain synonyms  
no aliases

Translation  
Lookaside  
Registers TLBs & faster than  
cache

Types → fully associative  
typically  
Cache 取代石子  
direct mapped  
set associative

思路  
COMMON case  
use locality  
faster  
e.g.: 内层循环  
- translate before access  
repeatedly use same variable  
temporal

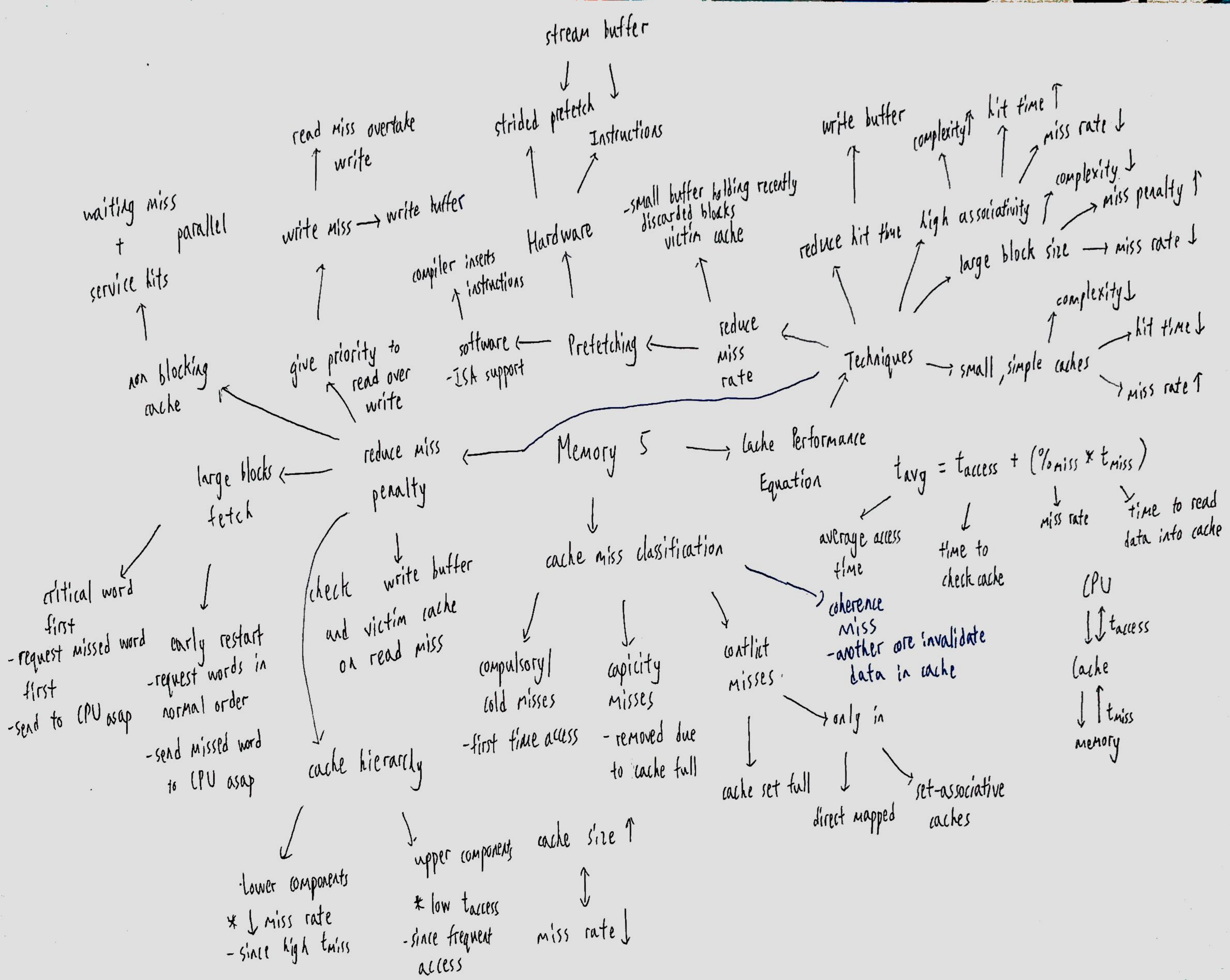
program address space  
Why?  
How? → put others in disk  
program > memory  
only put code and data used now in memory  
dynamically exchange during run time  
memory can't fit all program in  
spatial

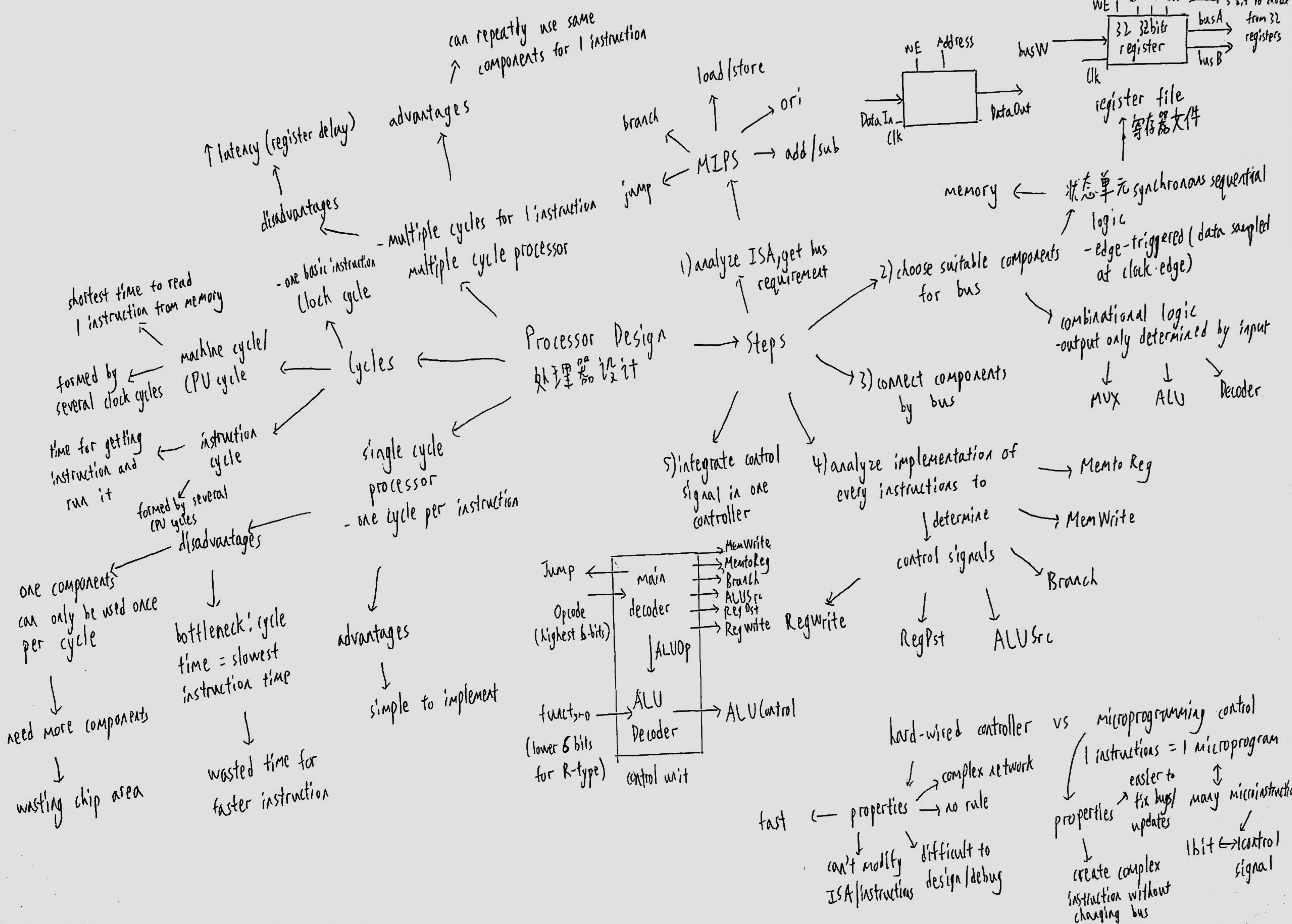
why

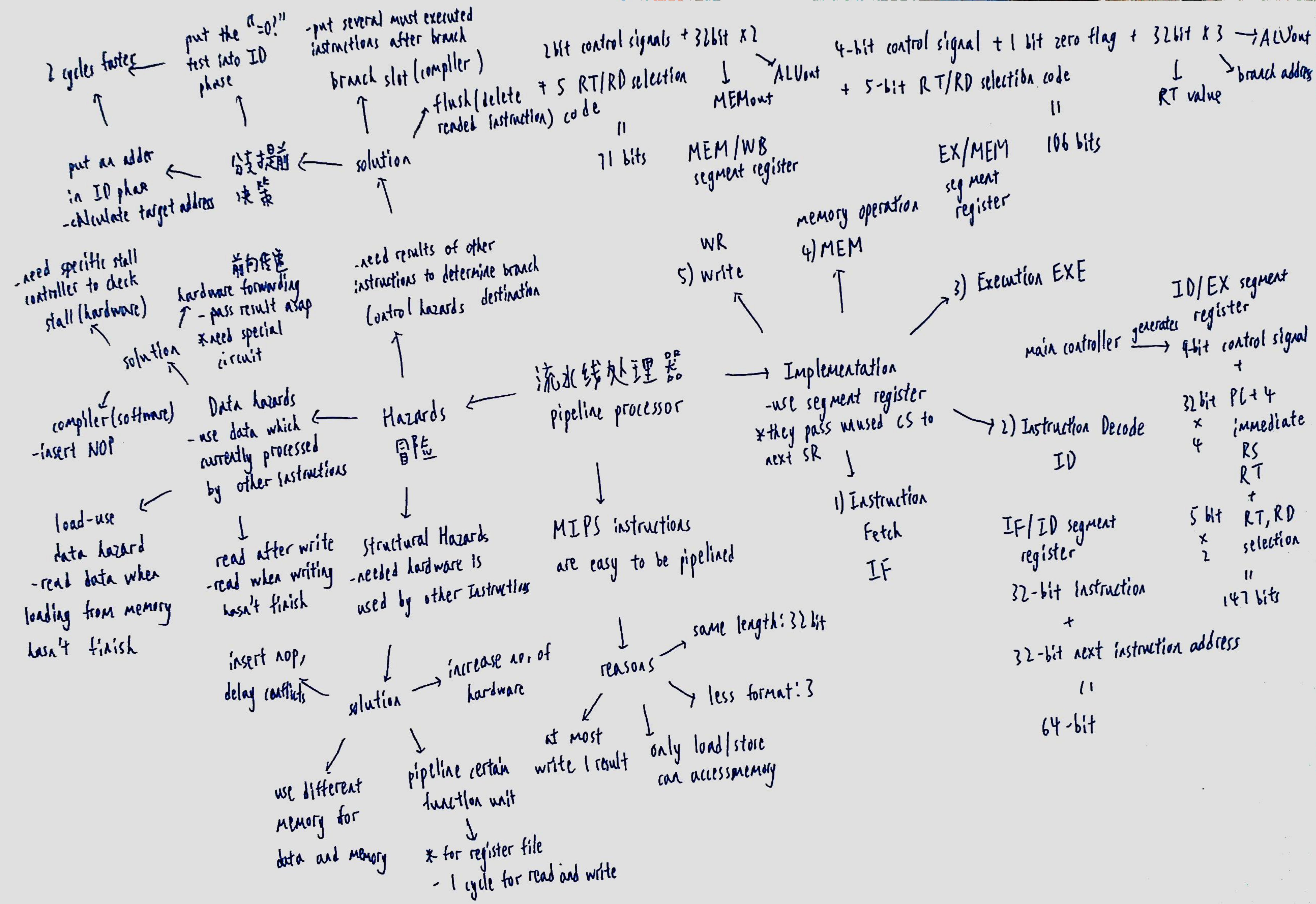
reduce

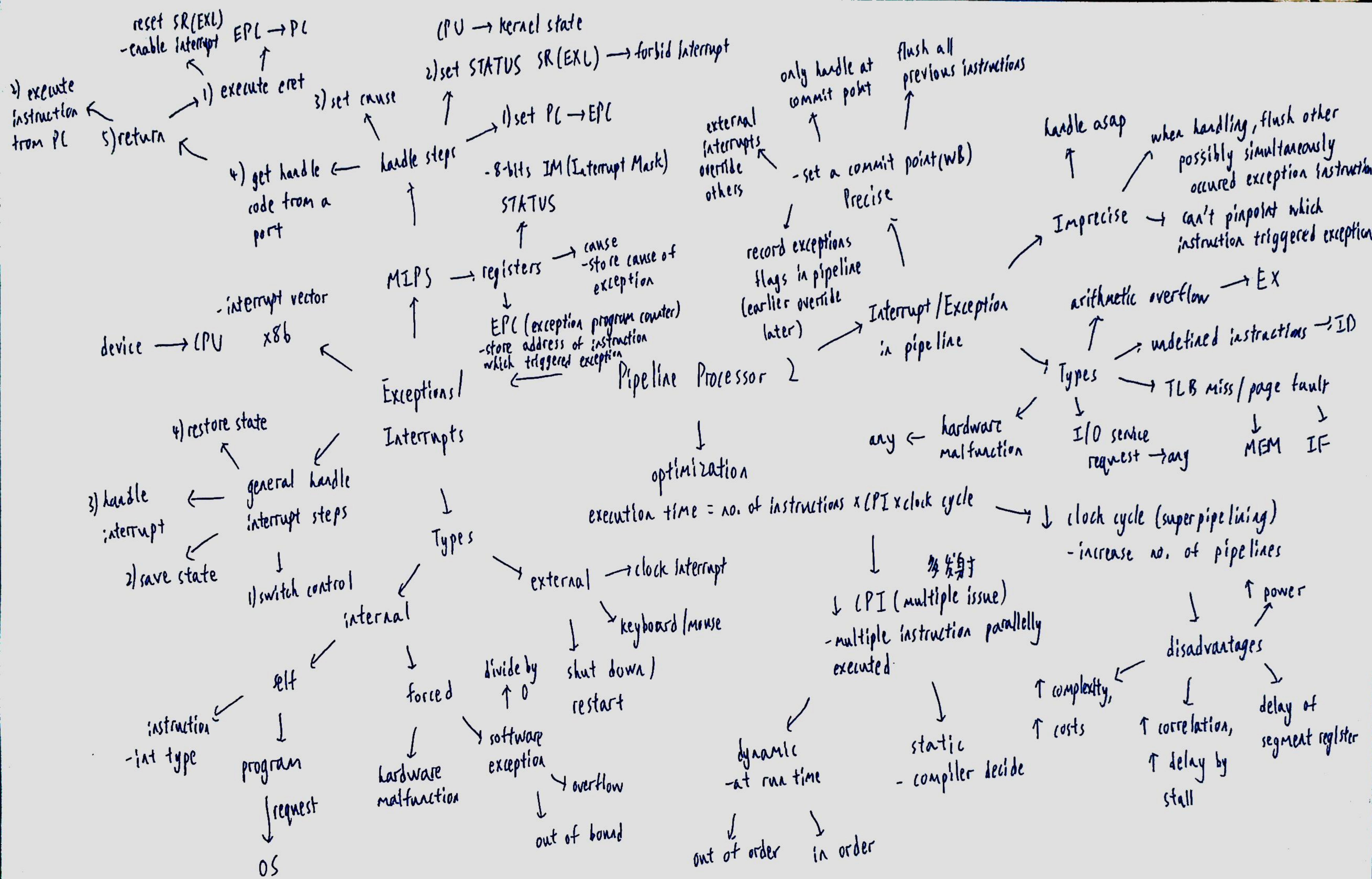
unused  
fragments

PI-PT  
Physically indexed,  
physically Tagged  
hit time ↑  
temporal  
spatial









Fetch → Decode → Rename → Dispatch

In-order front end

↓  
Window of instructions

Issue → Reg-Read → Execute → Writeback

Out-of-order execution

↓  
In-order commit

BTB + BHT

Control Hazards

Use branch prediction

dynamic  
-out-of  
-order

Branch Target Prediction

precise interrupt

NOP

true

false

register renaming  
RS > registers  
better than compiler  
renaming

WAW 矛盾  
write after write

rename registers  
as value/pointer  
to RS  
avoids WAW  
WAW

every Function Units have buffer  
reservation  
Tomasulo Algorithm station RS  
high performance without  
special compilers

use common  
Data Bus (DB) to  
broadcast results

↓  
Stages

- 1) Issue: get instruction from FP Op Queue  
-if RS free, rename registers
- 2) Execution: operate on operands  
-if not ready, wait (DB)
- 3) Write result: write to (DB), free RS

Return Address Stack RAS

Branch Target Buffer Multiple Issue Processor

BTB

-store history as a  
table

\* works for  
direct transfer

Very Long Instruction Word

VLIW

-pack 2 or more instructions into 1 VLIW

expand loop  
↑ Instructions  
↑ ILP

one instruction stuck, whole process  
stuck

lock step

target code  
not compatible

limitations

expand code

complex compiler

↑ storage

↓

Precise Exceptions

-Reorder Buffer ROB

speculative Tomasulo

- 1) Issue: send instruction and reorder buffer no.

- 2) execution: same

- 3) write result: write to (DR and ROB)

- 4) commit: update register with reorder result according to ROB

initially choose  
BHT  
chooser  
after threshold,  
choose fshare  
fshare  
PC xor BHR

↓  
BHT

Branch History Table BHT

↑ (calculate branch target)  
6 bits (MO)

Branch Outcome Prediction

cost

simple  
hardware

change instructions  
order

register renaming

↓

correlation

↓

correlation

Explicit Register Renaming

-use a physical register

file > no. of registers

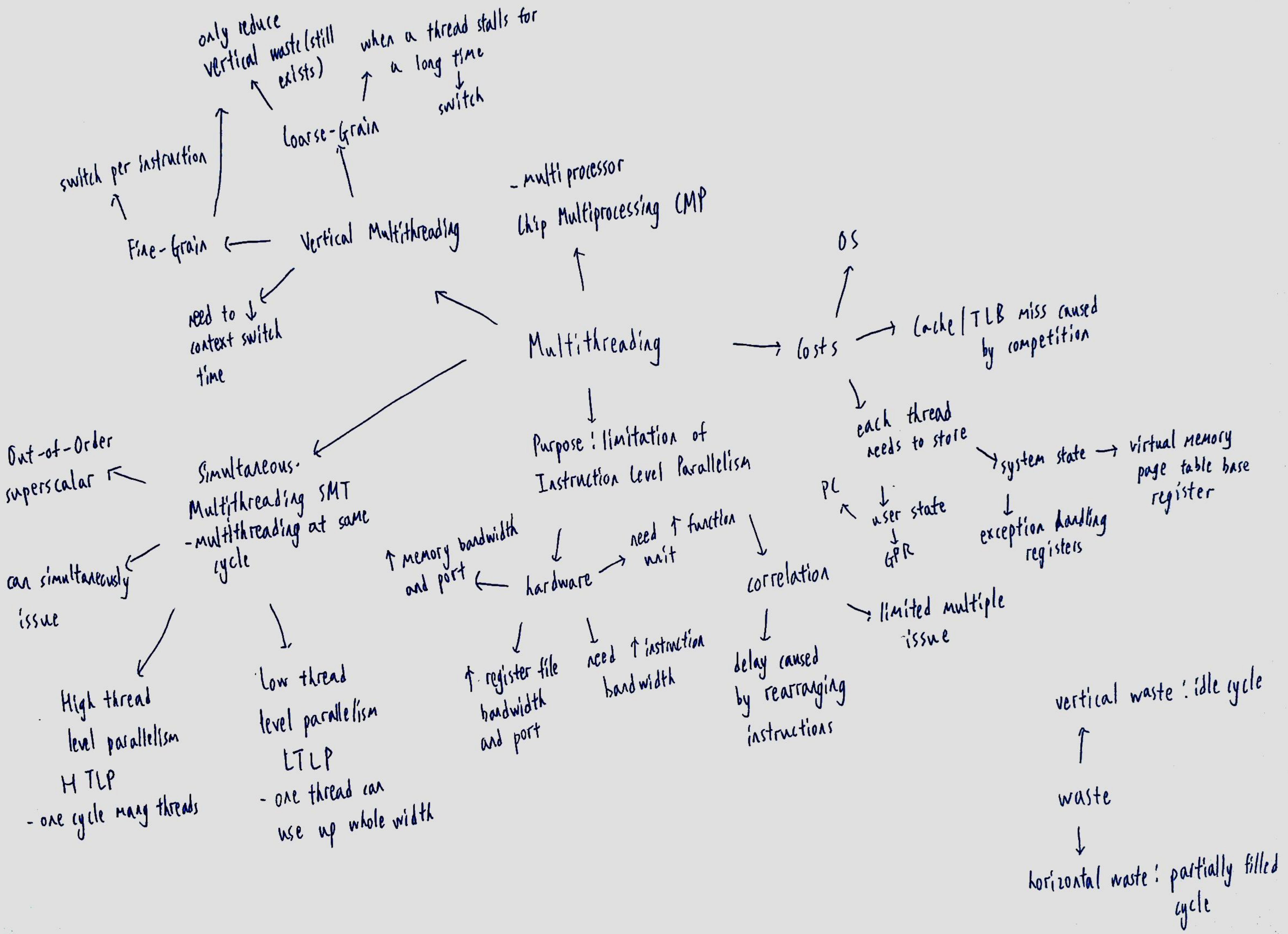
-map ISA register → physical register

precise interrupt

allow data fetch  
from single register  
file

decouple renaming

from scheduling



- one or more CPU share same block,  
writing different part of it

False Sharing

Cache consistency

- shared cache + private cache

Modified/Exclusive/Shared/Invalid

- 1 copy  
- up to date data

consistency protocol

MSI

Modified/Shared/Invalid

- can read/write  
- 1 dirtied copy  
- can read  
- multiple clean copy

Valid/Invalid VI

- can't read/write  
only allows one copy  
exists

context-switching

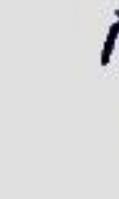
多路复用同处理器  
multiplexed-pipeline

1 processor 1 pipeline

Consistency

- unordered load/store

Release consistency



Memory consistency

→ Total Store Order TSO  
- store in order

allow reading  
store buffer  
to bypass value

↓  
Sequential consistency SC  
- load/store order same as  
code for every CPU

↓  
store buffer  
- put all store into  
buffer to hide  
write latency

Explicit parallelism  
显式并行性

multi-threading

↓  
hardware multithread  
硬件多线程

1 pipeline

multiple PC and  
register files to  
save thread status

switch thread

throughput ↑  
↑  
advantages

↓  
pipeline use rate ↑

multi-core → multi pipeline

↓  
1 memory