

grammars / parsing / compiling
 optimization recursively enumerable

Chomsky hierarchy context-sensitive

context-injective "free" CFG

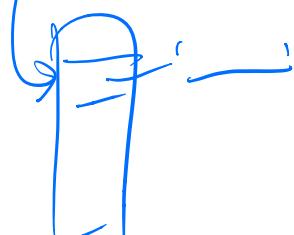
|
regexp

ALGOL-60

((((3))))

Backus-Naur form BNF

$\text{expr} ::= (' \text{expr}') \quad |$
 number |
 $\text{expr op expr} \quad |$
 unaryop expr

"Var" "parameter"
 42 h₁₅


backtracking ~ exponential time

N tokens 2^N

text

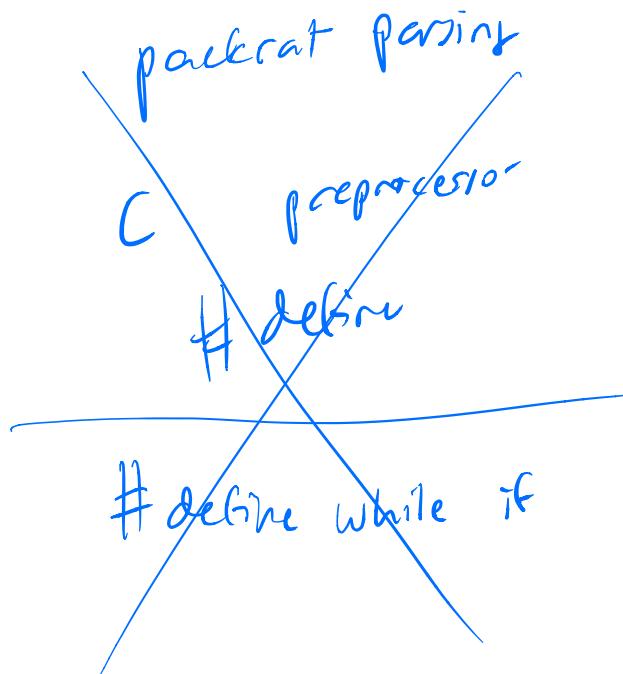
lexer → tokens → lex
 parser → parse flex

Symbol table

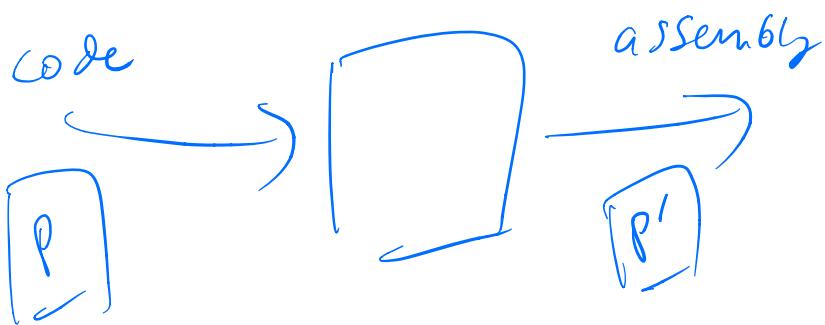
grammar → Yacc → parse
 bison

Java...ANTLR

LL(1)
LALR(1)



optimization



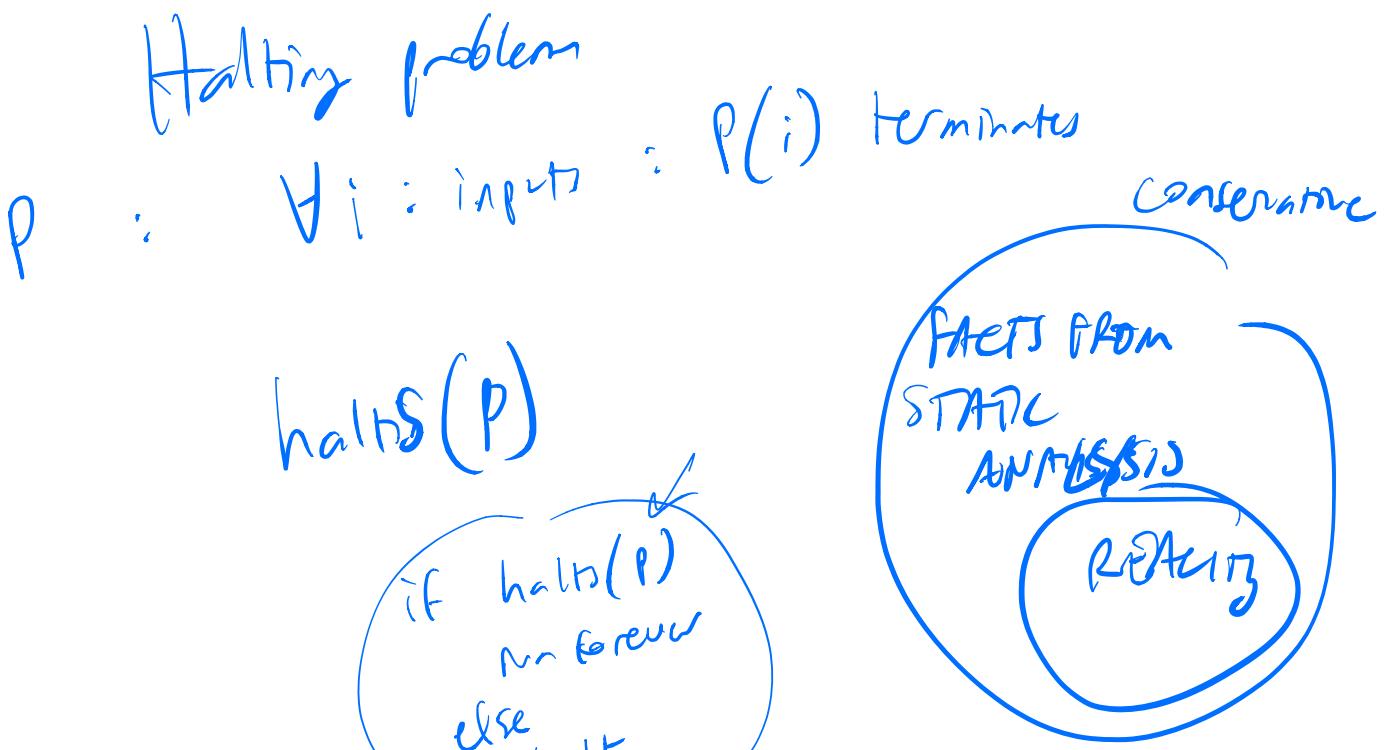
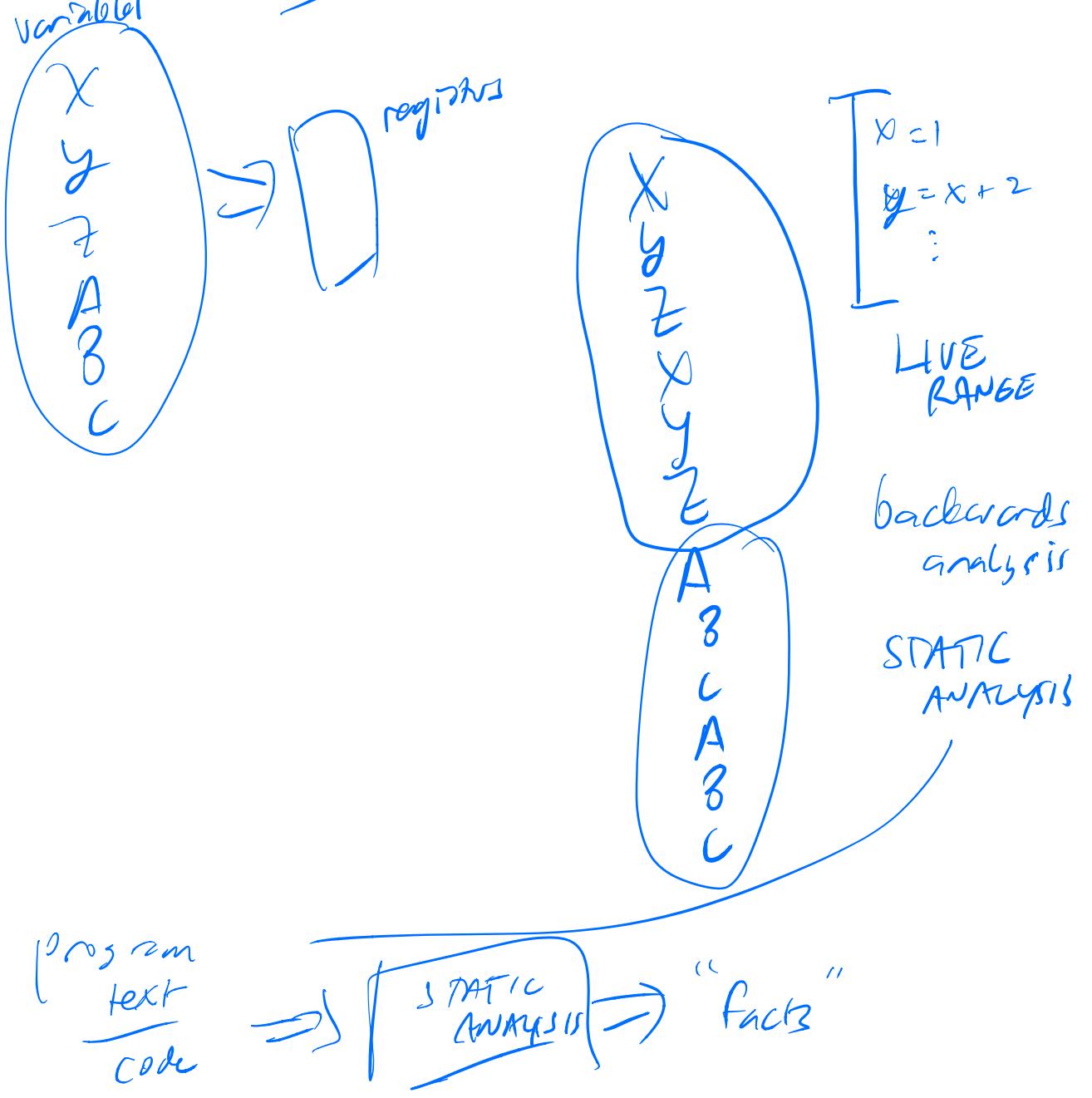
$\forall \text{ inputs} : P(i) = P'(i)$

semantic-preserving

registers CPU
memories - cache

memories hierarchy





STRAIGHT LINE CODE

Loops

Recursion

FUNCTION CALLS

FIRST-CLASS FUNCTIONS

EVAL EVIL

POINTERS

$x = [1, 2, 3]$

$y = [2, 2, 3]$

$\text{INT} * x$

$\text{INT} * x[i]$

ALIASES
ANALYSIS

POINTER
ANALYSIS

for
 $x[i]$ $y[i]$
 $x+i$ $y+i$

$f(x, y)$

$x: \{A\}$ $(x_1: \{A, Z\})$
 $y: \{B\}$ $(x_2: \{A, Y\})$

$x = A$
 $y = B$
 $i = C$
 $x = D$
 $x = E$
 $x = F$
 $x = G$
 $x = H$
 $x = I$
 $x = J$
 $x = K$
 $x = L$
 $x = M$
 $x = N$
 $x = O$
 $x = P$
 $x = Q$
 $x = R$
 $x = S$
 $x = T$
 $x = U$
 $x = V$
 $x = W$
 $x = X$
 $x = Y$
 $x = Z$

register allocation
= graph coloring
NP complete

strength reduction

$$\begin{aligned} X &= \text{pow}(3, 2) \\ X &= 3 \times 3 \\ X &= 9 \end{aligned}$$

3^2

$$\begin{aligned} Y &= 12 \\ X &= 12 \\ Z &= 13 \\ Z &= 2 \\ &\vdots \end{aligned}$$

Constant propagation

Inlining — exposes optimization opportunities

```
int add(int x, int y) {  
    return x + y;  
}
```

reduces
for call
overhead

foo()

```
int a = 12;  
int b = 13;
```

```
add(a, b);
```

closed-world

Whole-program

$$int c = a + b \\ = a + b \cdot 2^5$$

P, analysis
modular
analyses

FORTRAN - OPTIMIZING COMPIRATION

LISP - INTERPRETERS

switch (-) {

case ADDONE:

~~var(p) = var(p) + 1~~

break;

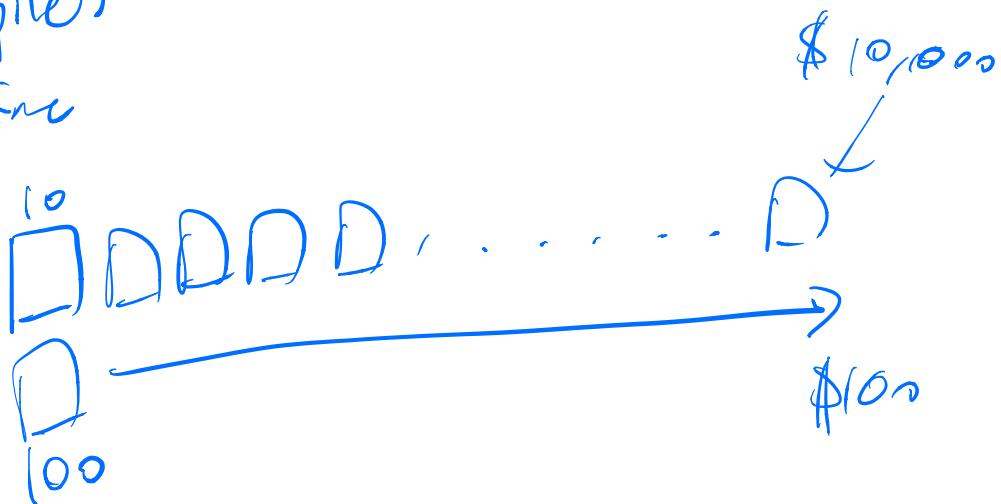
;

JIT Compiler

just in time

rental: \$10

ski purchase: \$100



ski rental problem

rent until you spend \$100
then you buy

rent-case: Spend \$200
could have spent \$100

Intros

Fair JIT - -00
-01
-02

Jim Larus
LARUS

Dumb

Slow code

:

Combiner

smart
Fast code

parse

internal representation

AST abstract
syntax tree

homogeneity

LISP

lambda calculus

functional programming

function - centric

functions first-class objects

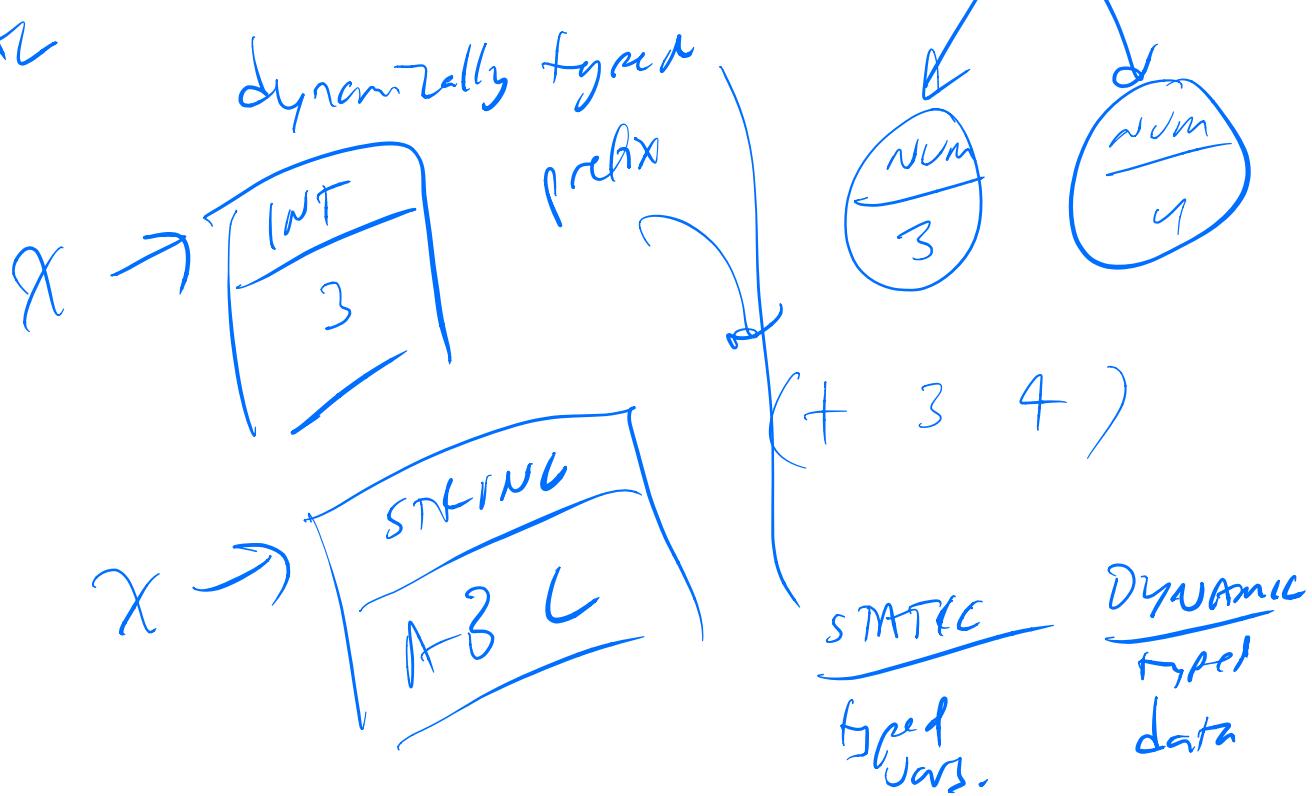
program as data

3 + 4

garbage collection

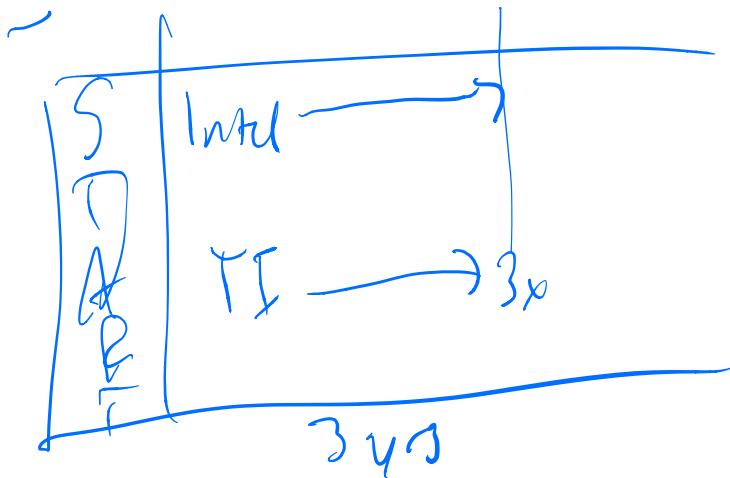


EVAL

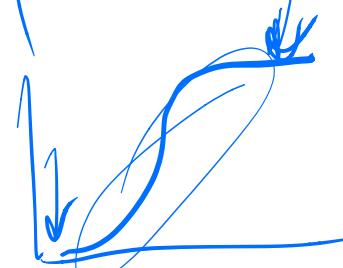


LISP machines 70s - 80s

TI Explorer



DB machines
Computer
Network
Neural net



Moor's Law
#gates in area
doubles
~ 18 mos

Demand scaling
↑ density \Rightarrow ↑ clock speed

Moore's Law

compiler optim.
double perf
every 18 YEARS

Self
Chambers Unkar
Hölzl
Jeff Dean

GC

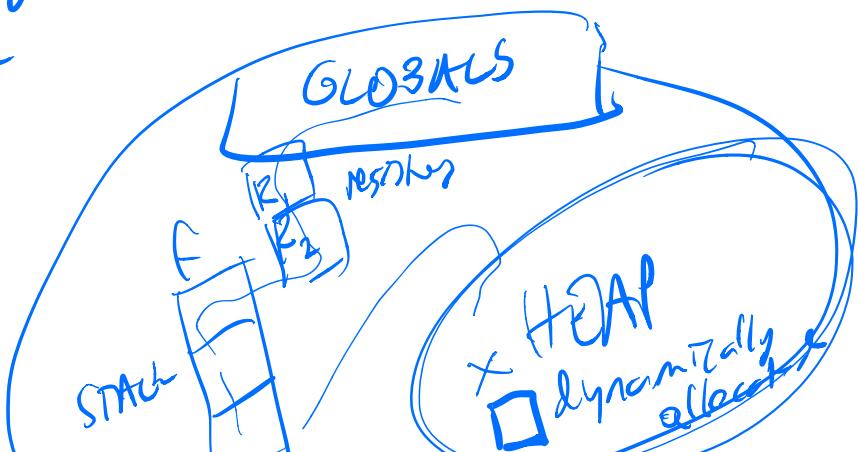
Mark-sweep

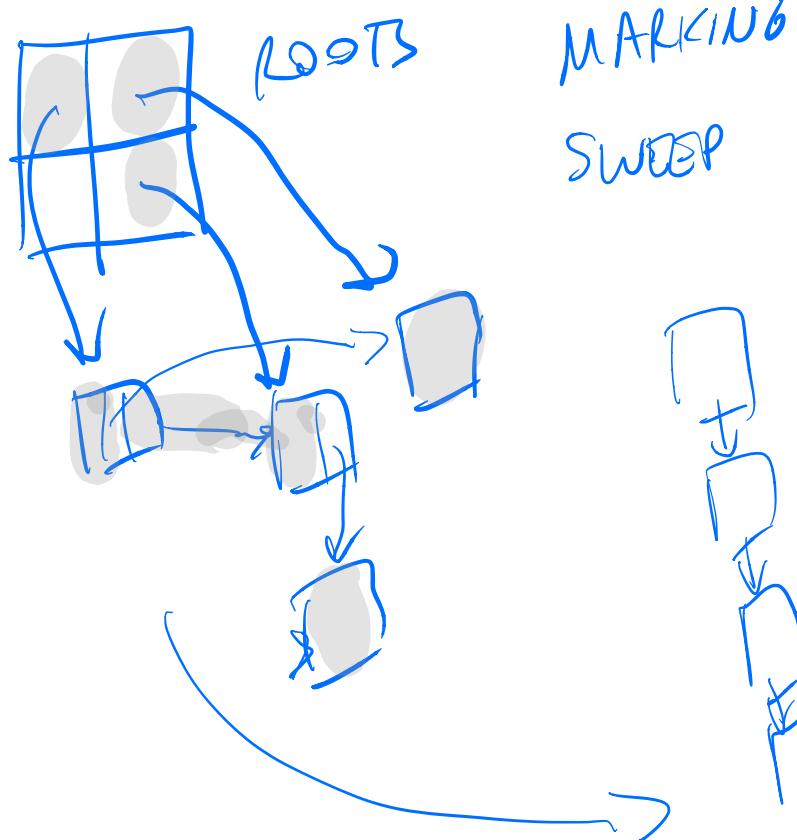
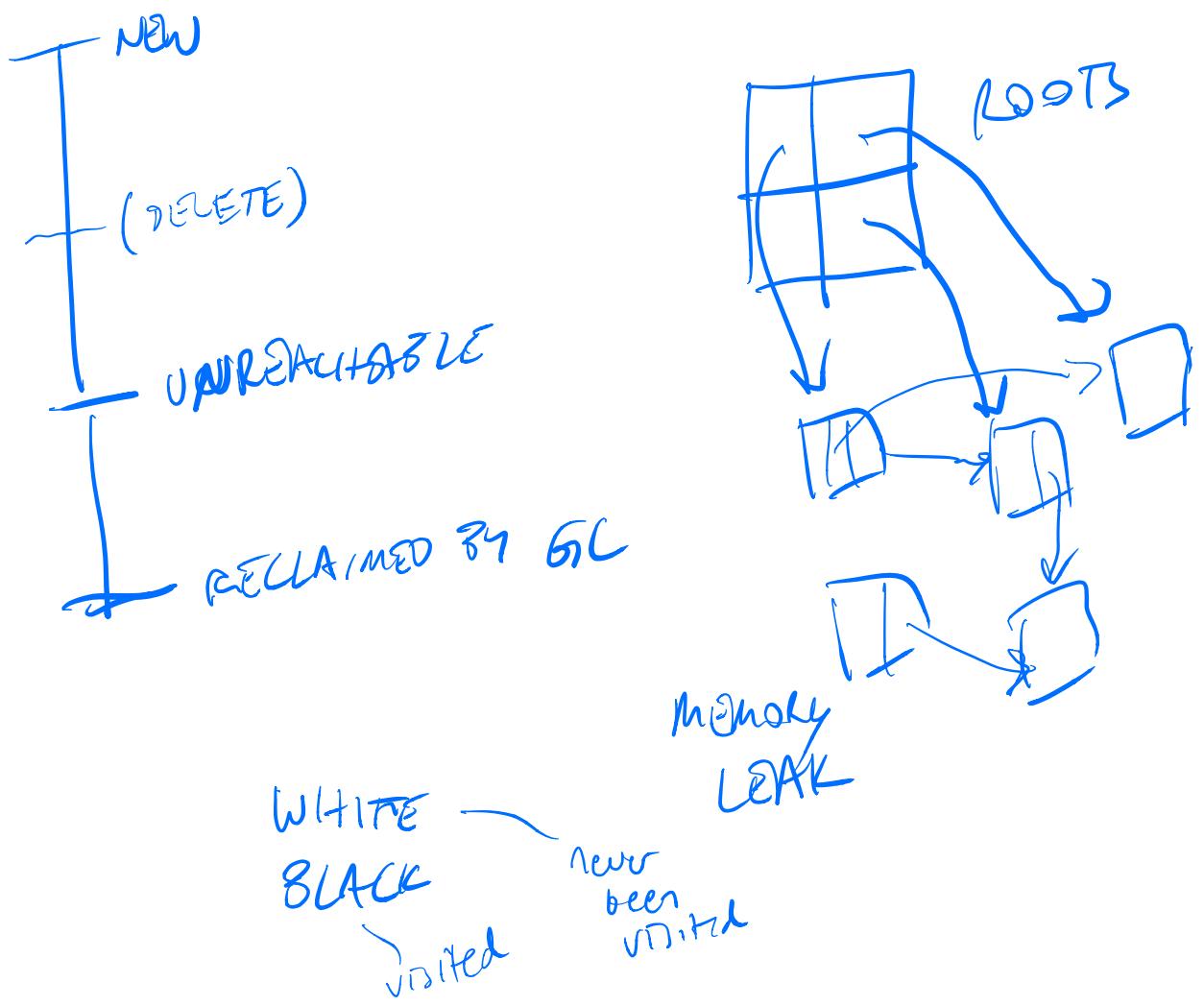
reachability

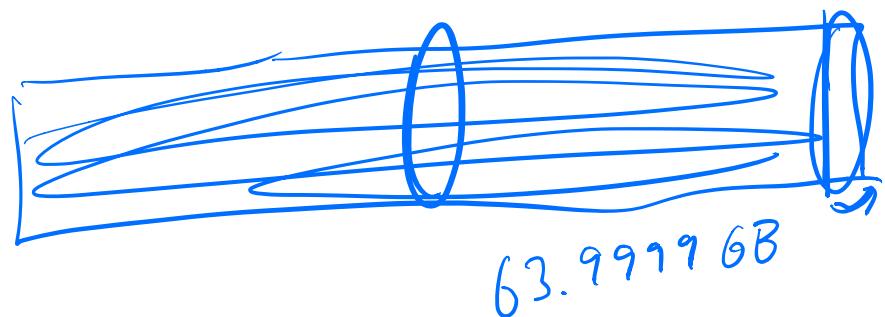
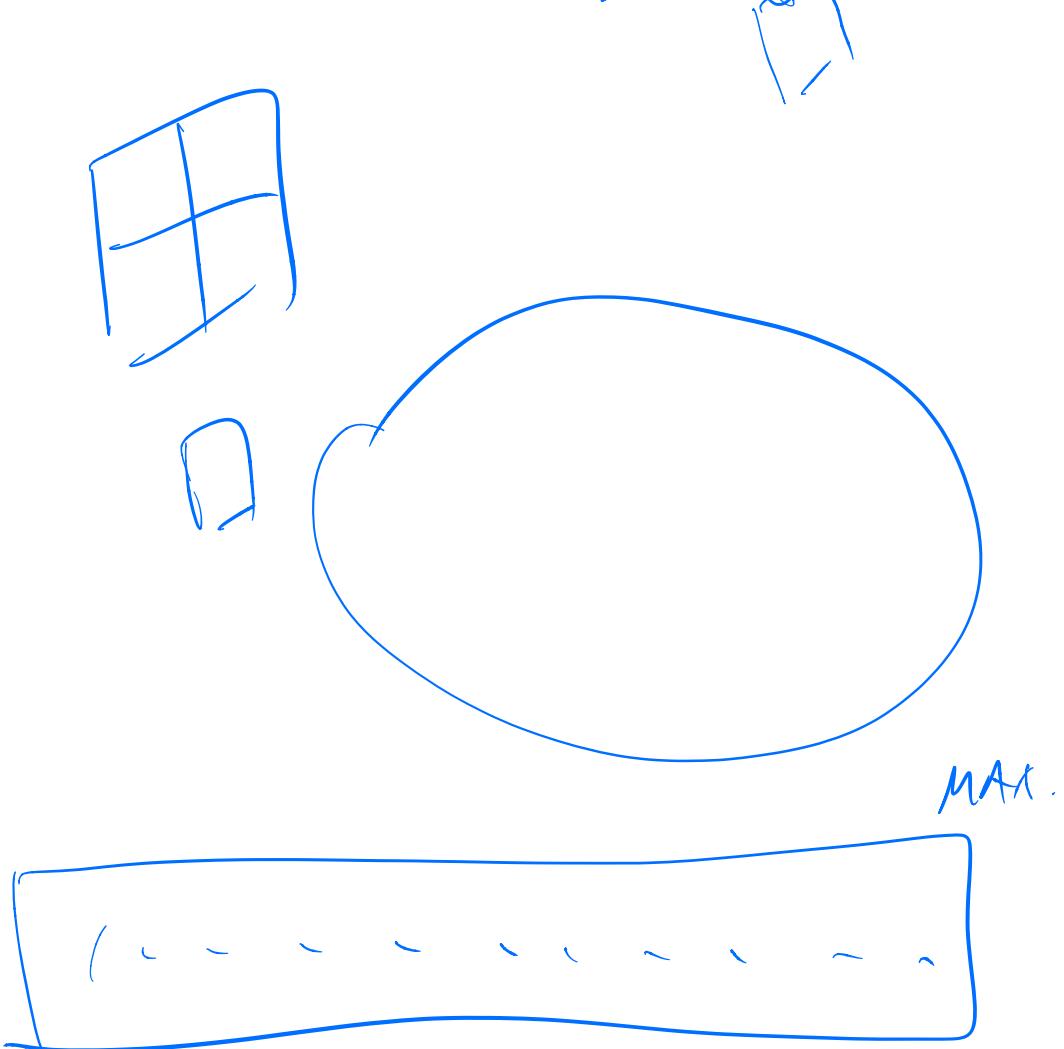
explicit mem mgmt
malloc/new
free/delete \Leftarrow object "dead"

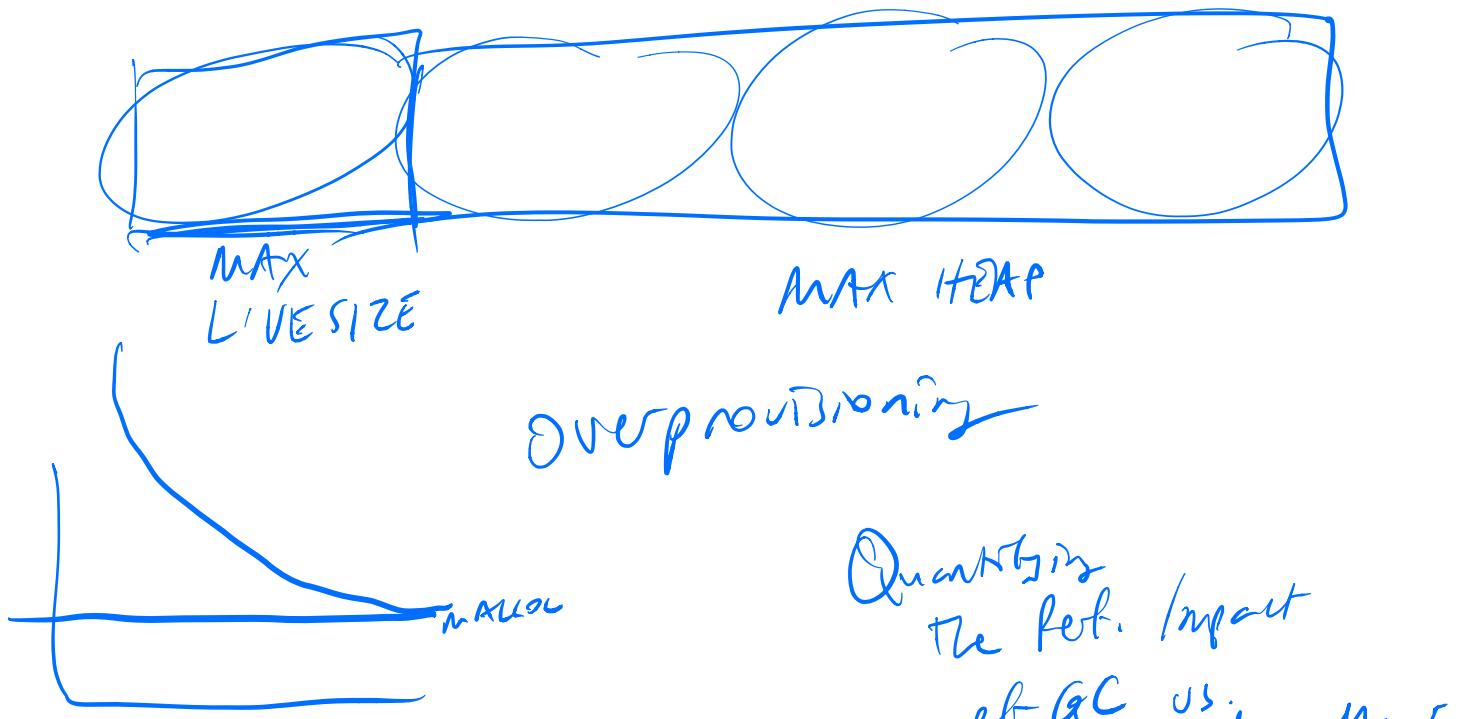
reachability \approx
Oracle

PROGRAM









MULT
FACTOR
OF
HEAP

$\sim 3x - 5x$

more RAM

mallay/
free

Quantifying
the perf. Impact
of GC vs.
Explicit mem Mgmt

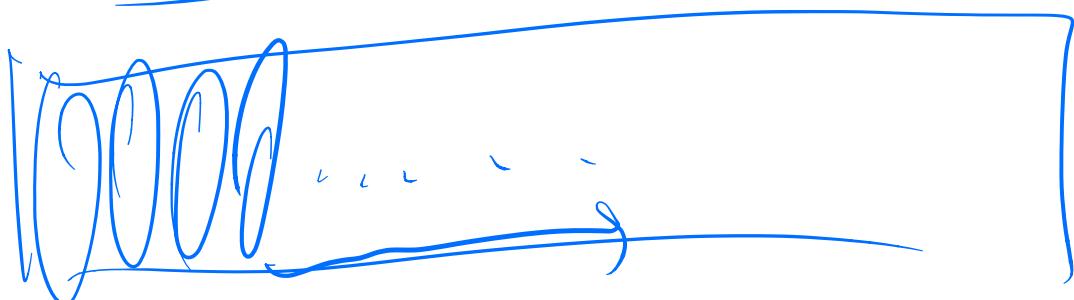
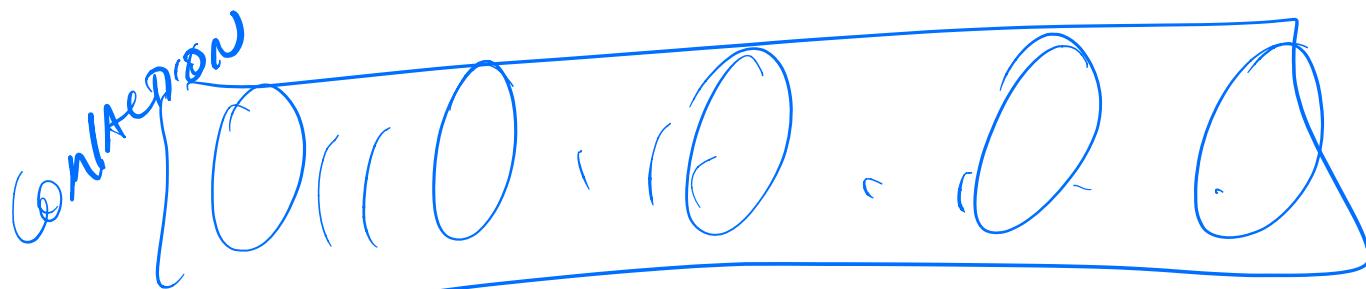
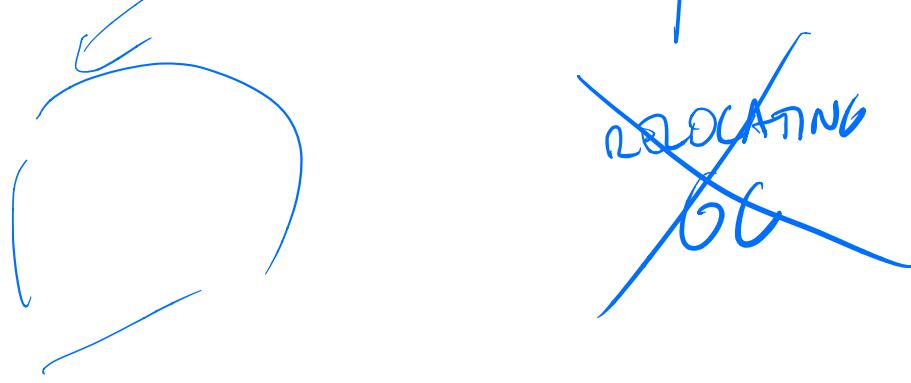
UNCOOPERATIVE GC
Hans Boehm

PRECISE GC

Duck Test

Conservative GC

3.99×10^{15}



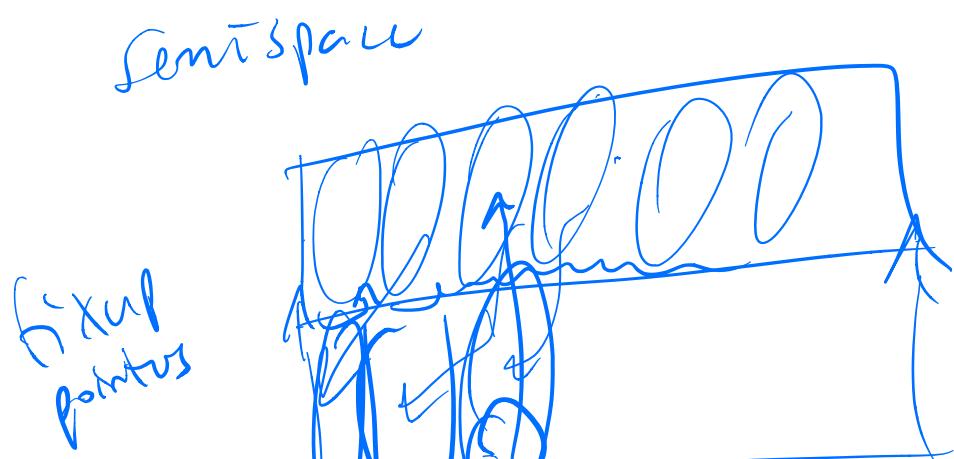
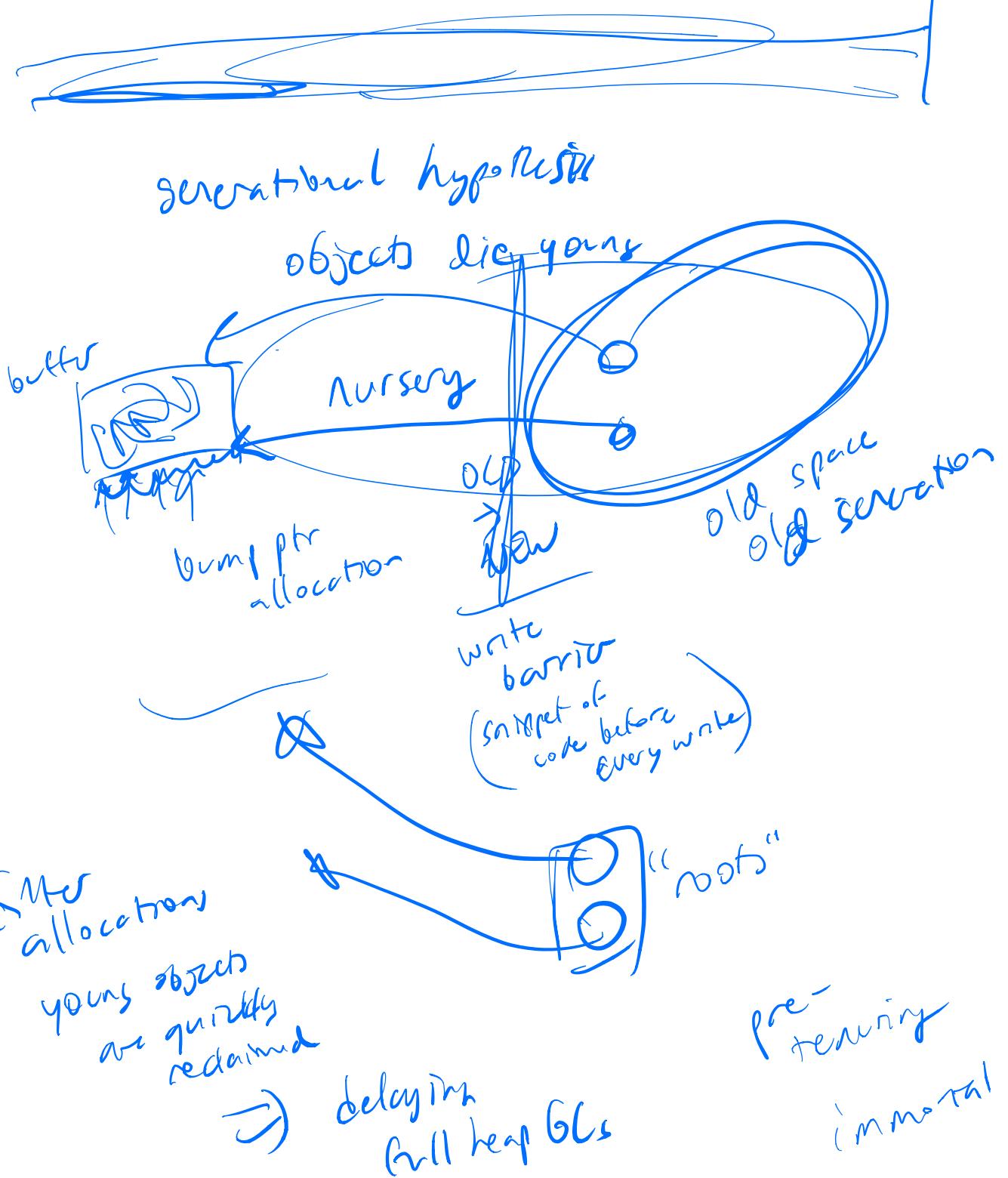
MARK-SWEEP - COMPACT

reclaiming

"scanning"

MS
MSC

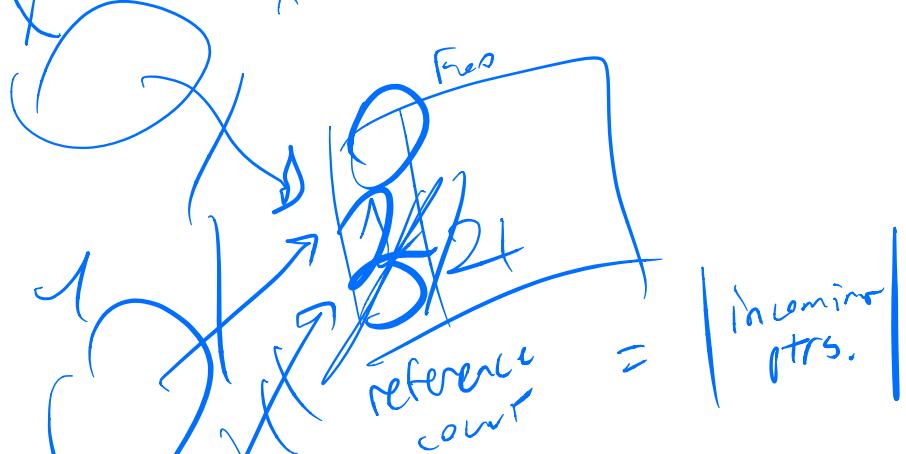
generational GC



~~new~~ ~~ptr~~

reference counting

$x = \text{new } \text{Fox}$



x

$y = x$

$z = x$

$z = \text{null}$

$y = \text{null}$

$x = \text{null}$

$\text{ptrupdate}(x, \text{null})$

If $*x = \text{null}$ do nothing

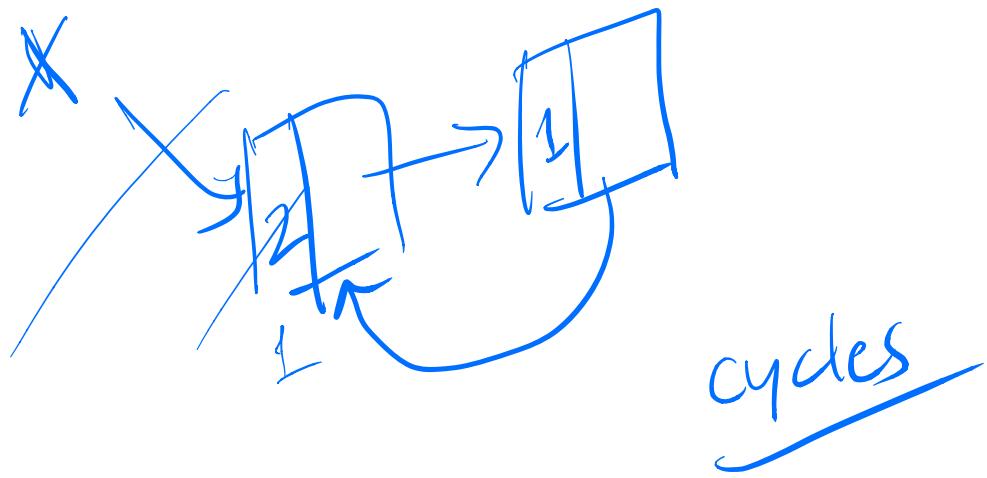
else $\text{ref}(*x) --$

$\text{ref}(q) +=$

If ($\text{ref}(*x) == 0$)

delete $*x$

cx
smart pointer



$X = \text{null}$



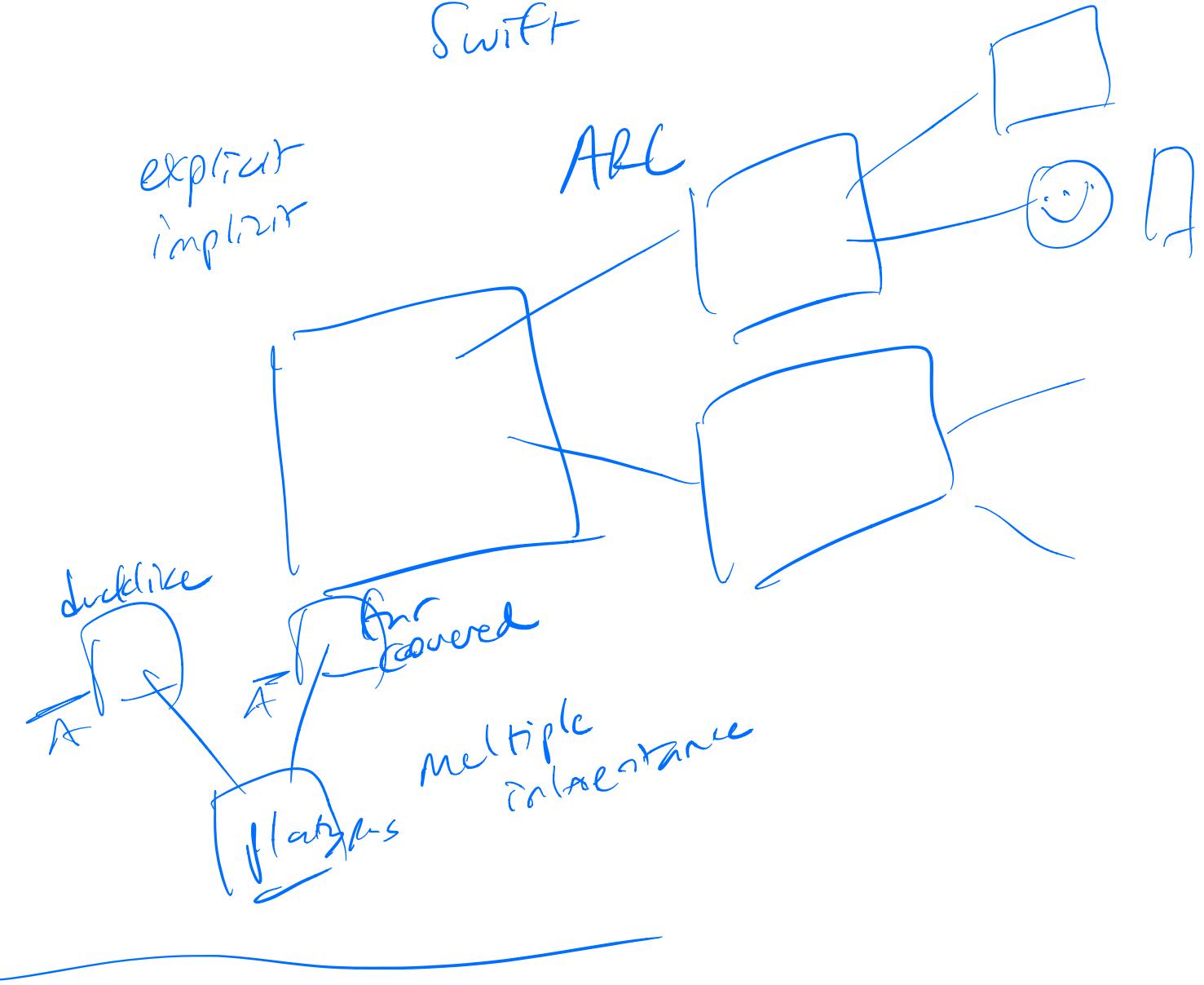
cycle collection
 \approx GC

region pool allocation areas

Finalizer RAII

delete obj
 class ~Obj ~Obj() ~

LC → Objective-C
Swift

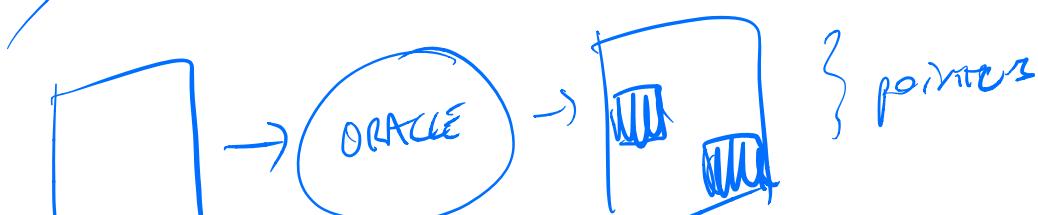


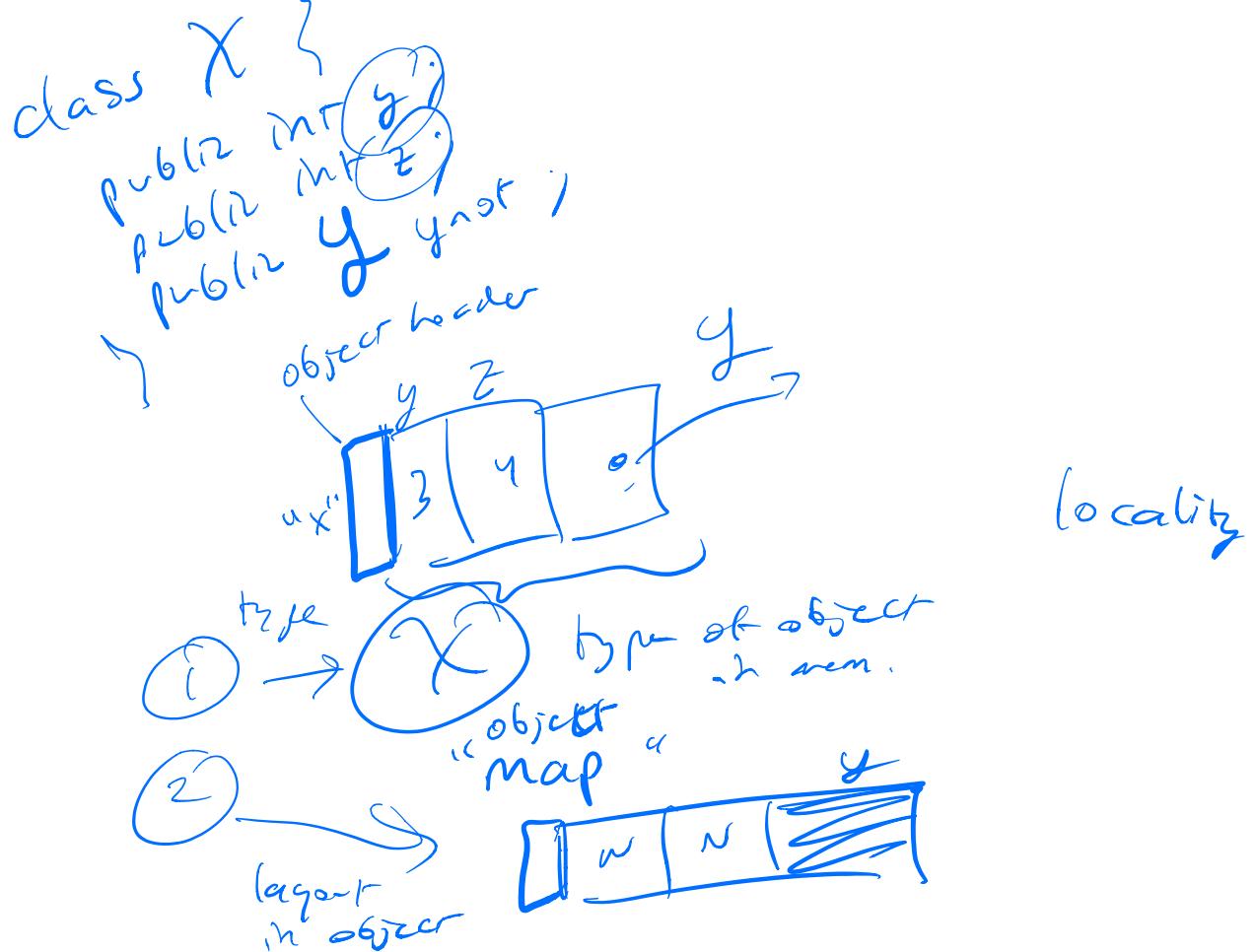
Boehm collector

conservative GC

cannot distinguish ptrs. from values

precise GL
- can identify ptrs.





managed languages

garbage collection

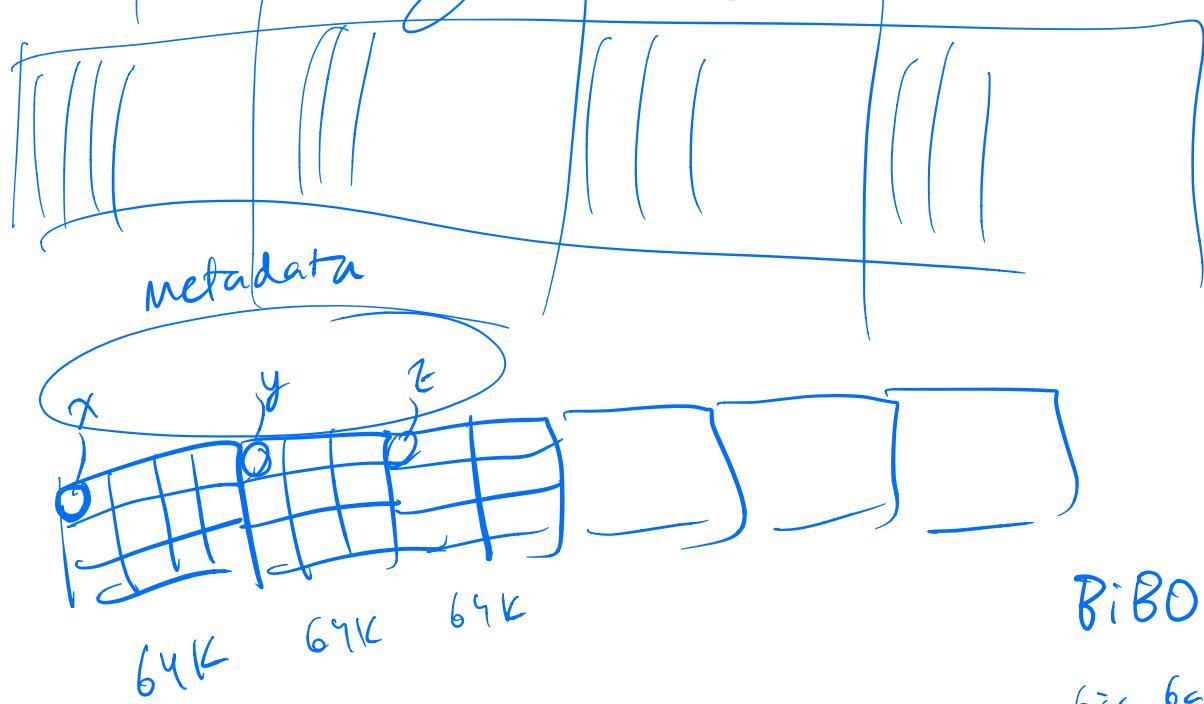
type safety -

bounds checking

errors - uninitialized reads



X, Y, Z, Ynot



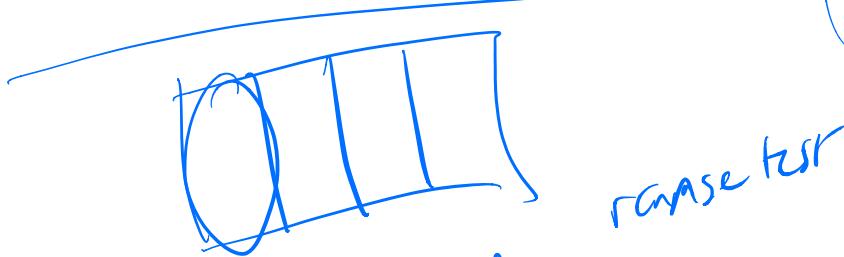
6 bits
of
pages



threaded
free list

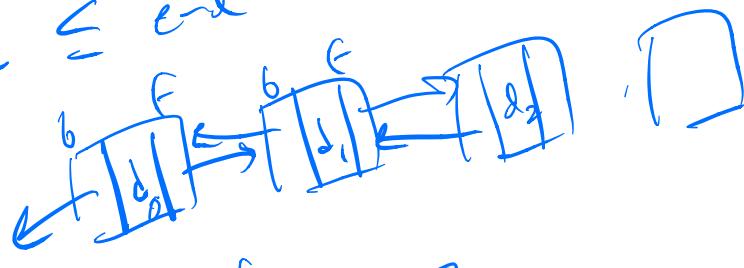
Morris
worm

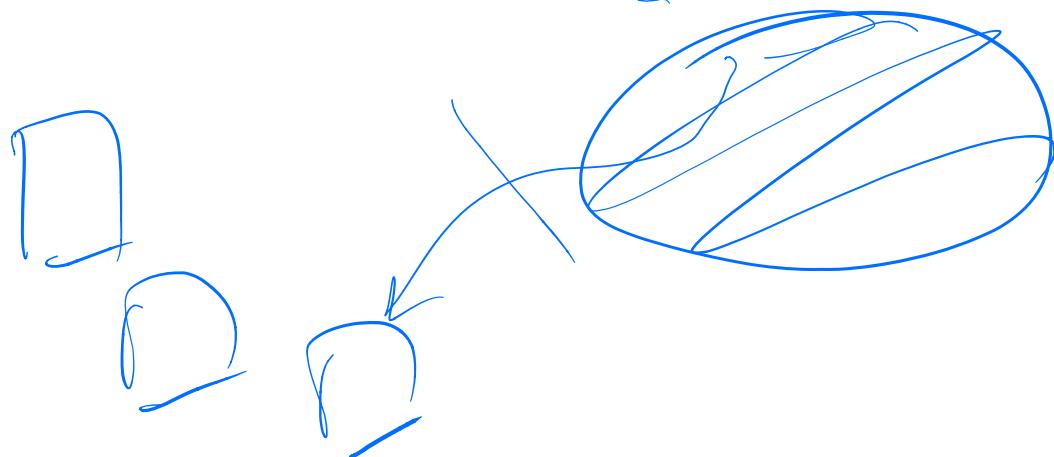
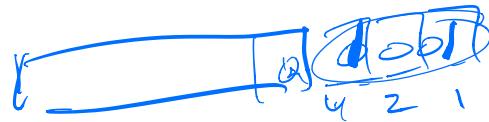
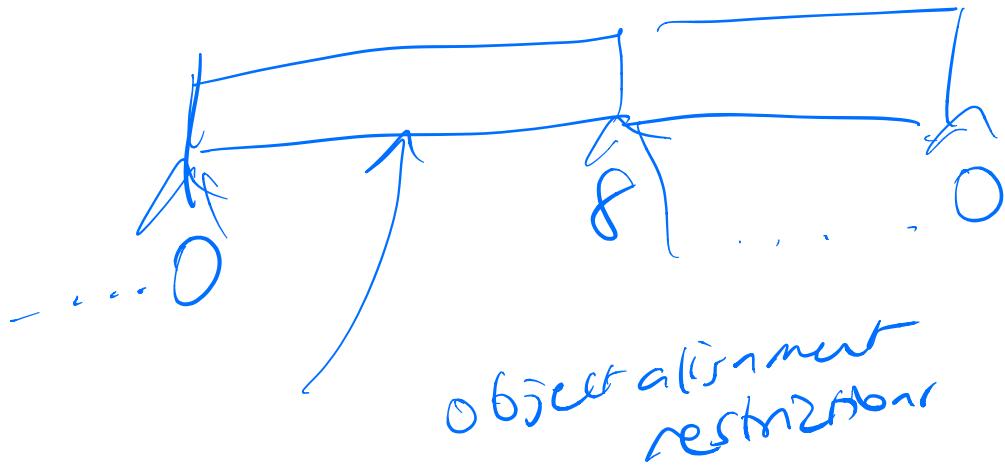
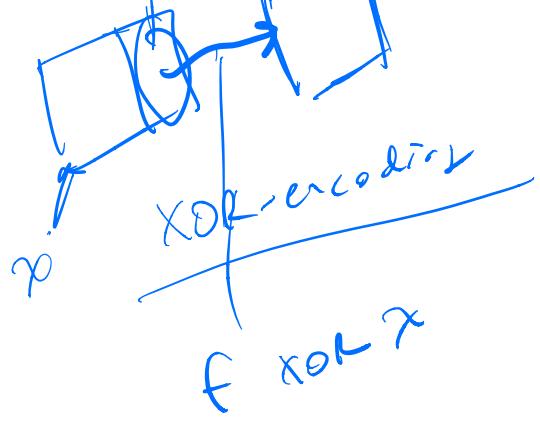
DDoS
malware



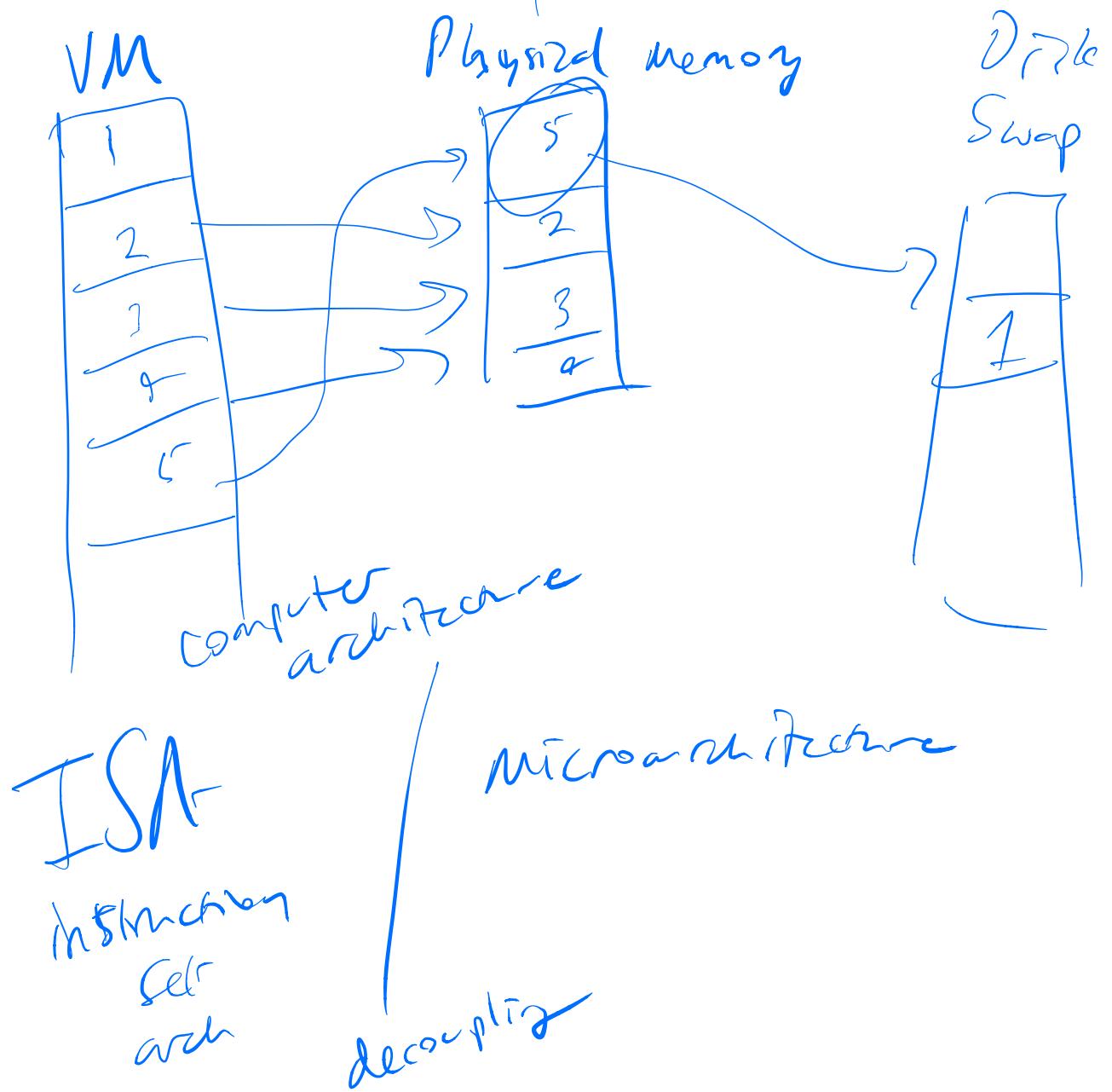
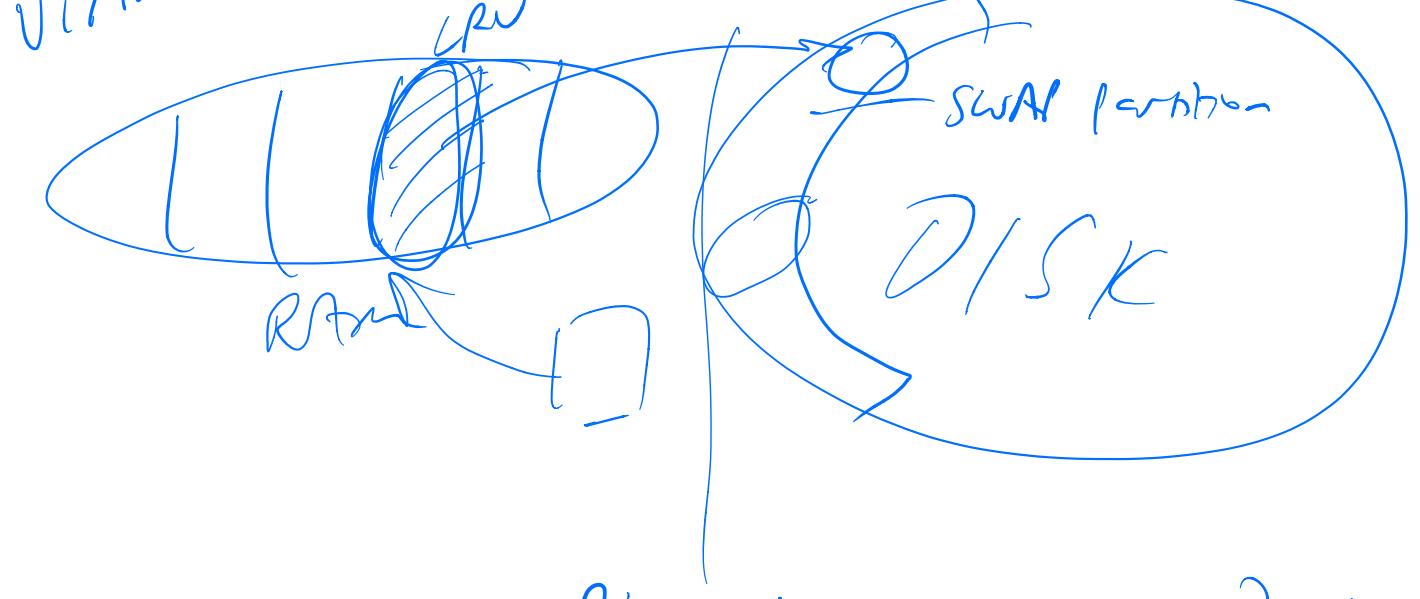
reassembler

base $\leq \alpha \leq$ end



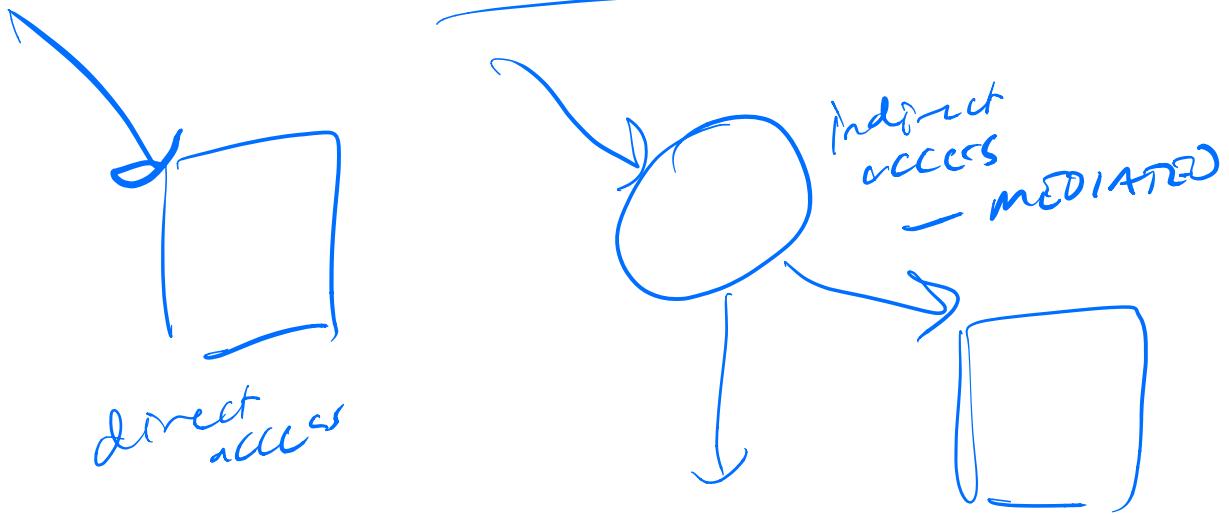


Small memory passing



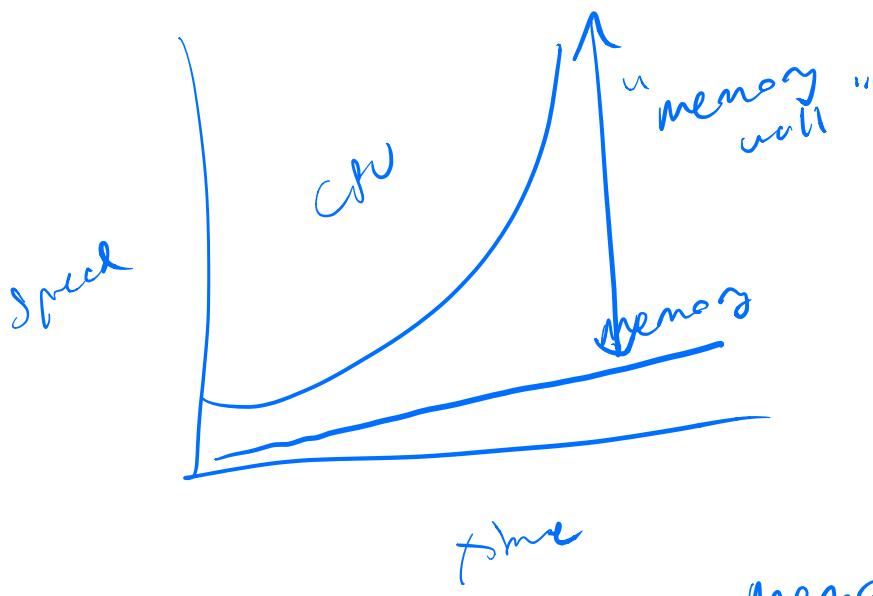
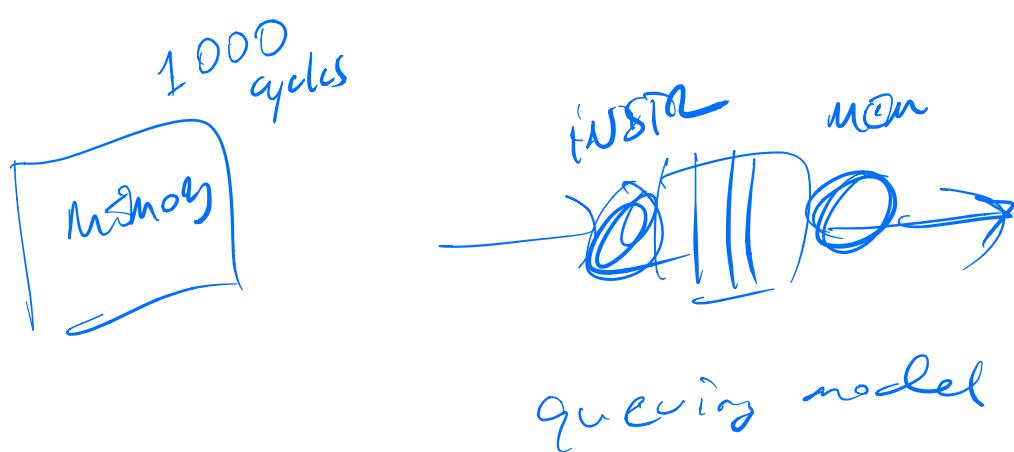
Virtualization
memory hierarchy
caches

"all problems in CS
are solvable
by adding
one level of indirection"



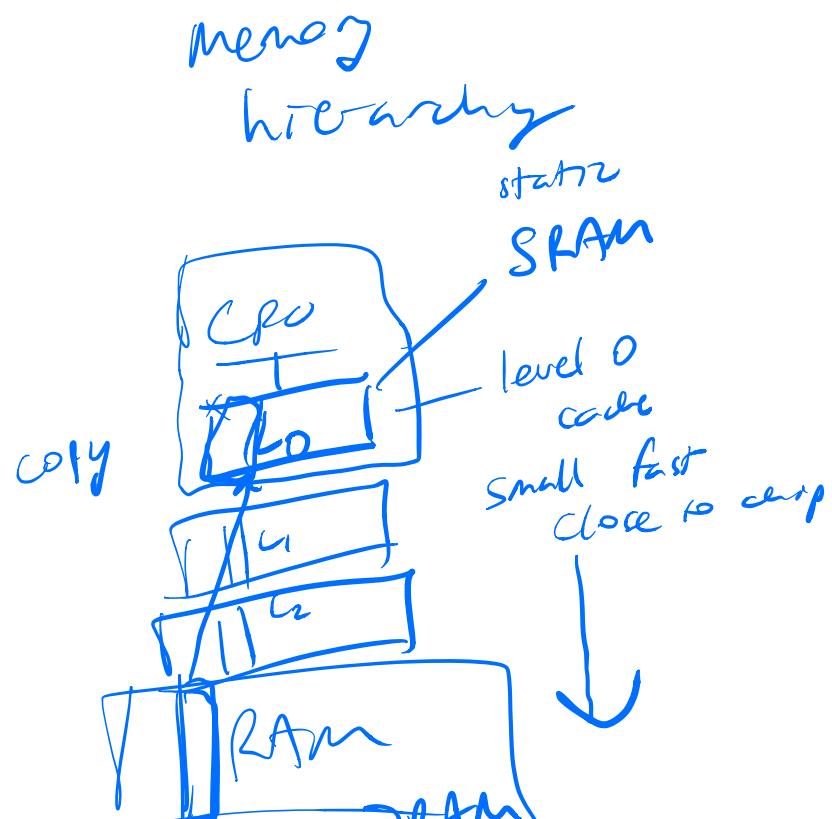
Meredith Rosenthal
Carl Waldspurger





ASIC
grows every
SOLVER iteration

ADD MOVE



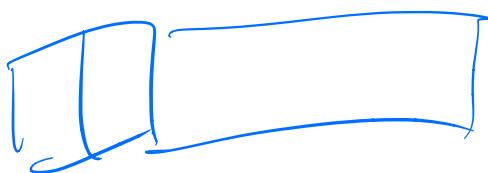


~~Dynamic~~
cache line
 $32b \rightarrow 128b \rightarrow 256b$

inclusive - copy in all levels
exclusive - copy in one level

"LRU" eviction

Cache coherence protocol

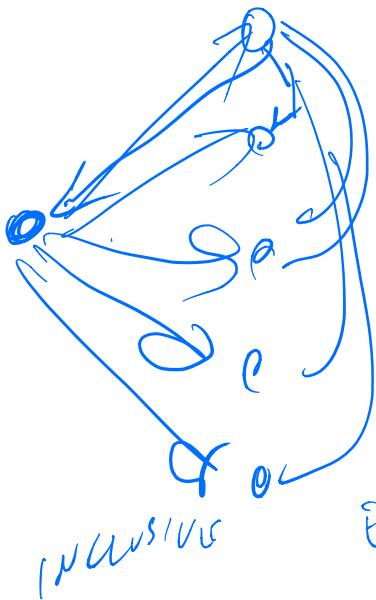


MESI



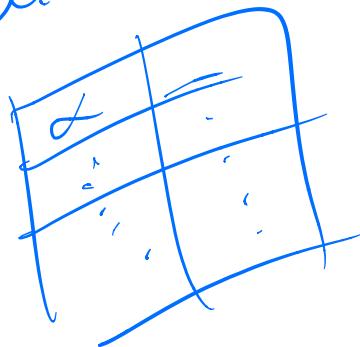
scalability w.r.t N
⇒ execution time

$\sim \text{polylog}(N)$



$$N(N-1) \approx N^2$$

directed-based scalable

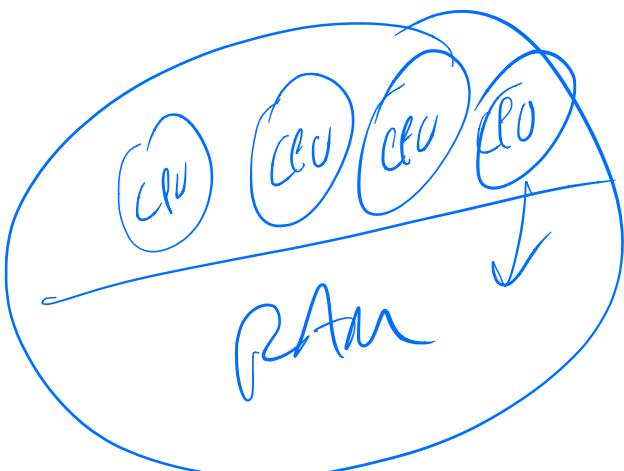


$\sim \text{hash table}$



Cache levels

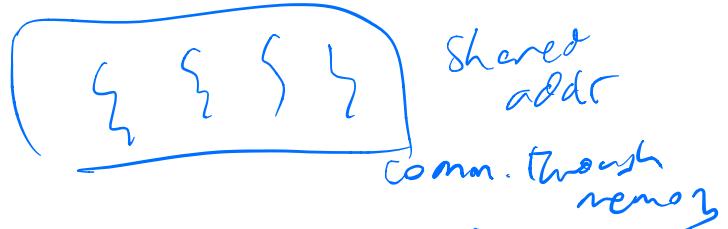
Andahl's Law



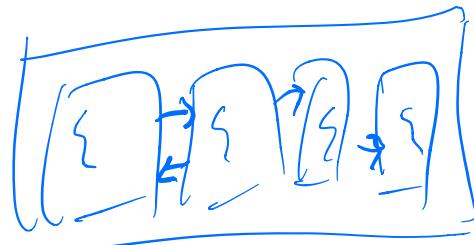
automate parallelization



Multithreading



Multi processing



```
for (100,000) {  
    x[i] ← x[i] × 2;  
}
```

```
for i = ... 25K  
    x[i] ← x[i] × 2  
    x[i-1] × 2
```

lock join
parallelism
loop-carried
dependencies

```
x[0] ← x[0] × 2  
x[1] ← x[1] × 2
```

transform & eliminate
(e.g.
carried
dependencies)

Polyhedral analysis

$$x[1] \leftarrow x[0]^3 \cdot x^5$$
$$x[2] \leftarrow x[0]^4 \cdot x^4$$
$$x[i] \leftarrow x[3]^2 \cdot x^2$$

$x[y[i]]$

DSL

constrained computation
 \rightsquigarrow easier to parallelize

manual parallelization

dynamically allocated objects
reversion
side effects
complex expr.

(libraries)

mult(x,y)

manually written

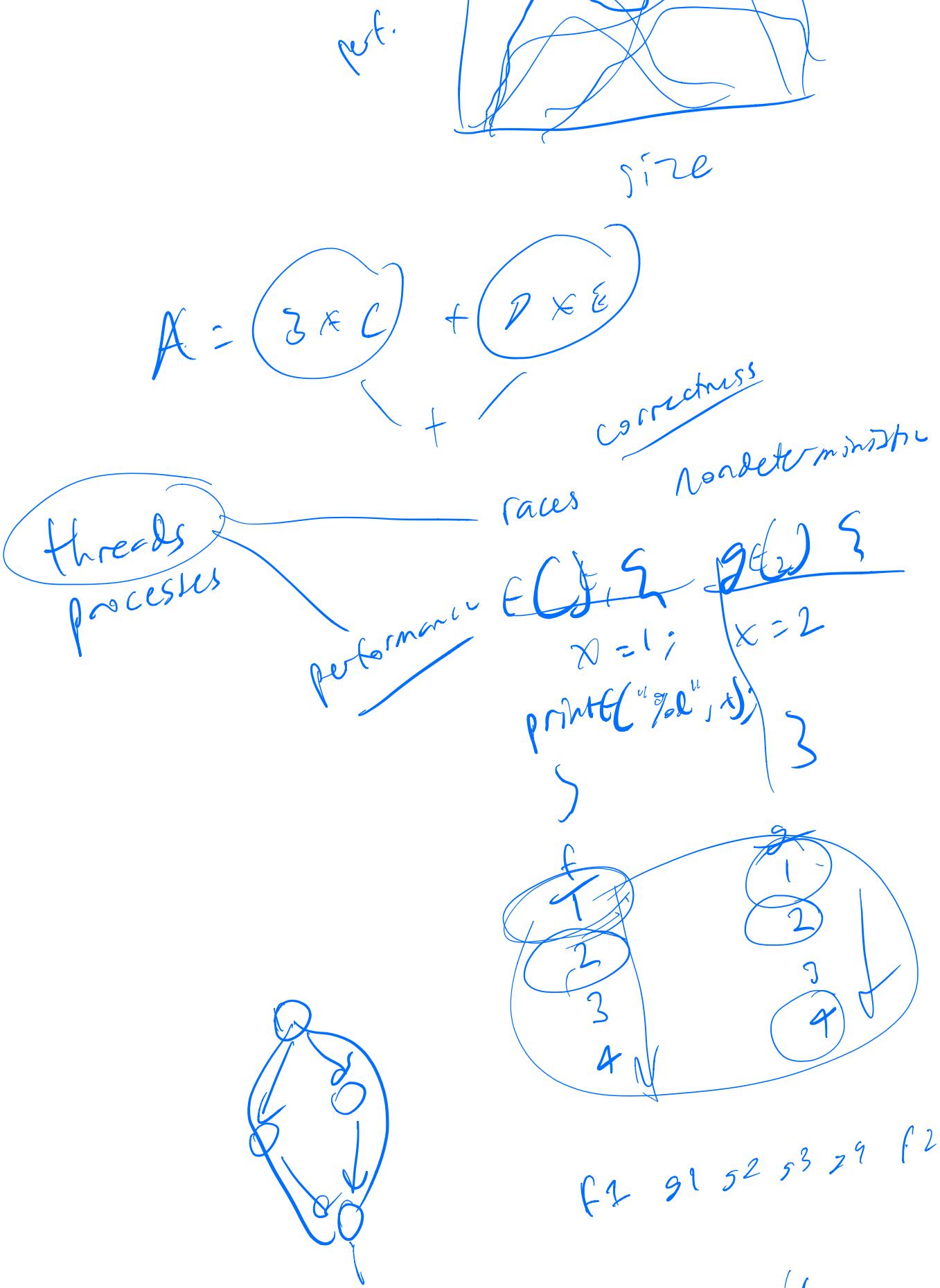
BLAS

Basic Lin Algs Subroutines

ATLAS

FFTW

program synthesis



$\sim T^{\frac{2^u}{I}}$ " 1

Synchronization
Contention (cache line level)
memory level

SEQUENTIAL part

(Thread t1 = new Thread(f, args);
Thread t2 = (f, args);
(t1.wait();
(t2.wait();

$T \ll P$
 $T \gg P$
 $T \approx P$

insufficient # of threads

too many threads



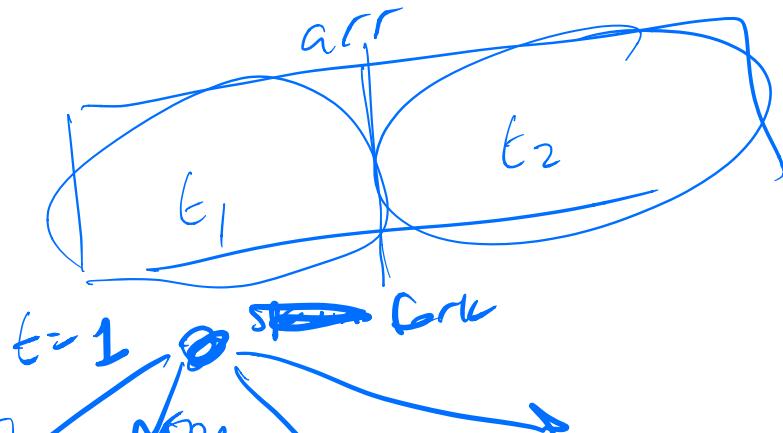
too fine-grained

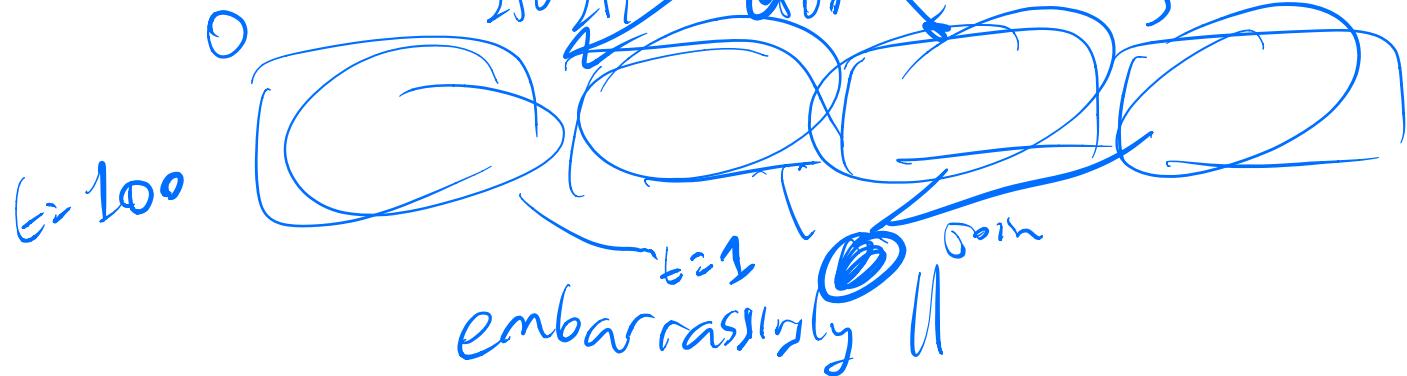
granularity



coarse-grained

$T \approx P$





T_1 time with 1 processor (serial)

T_p time with p "goal": $T_p \approx T_1/p$

T_∞ "critical path"

$$T_p \approx \frac{T_1}{p} + O(T_\infty)$$



$$\Omega_2: \frac{.9}{2} + .1 \\ .55 \sim 2$$

$$f: \frac{q}{4} + \dots$$

.325

~ 3
~ 5

2
3
5

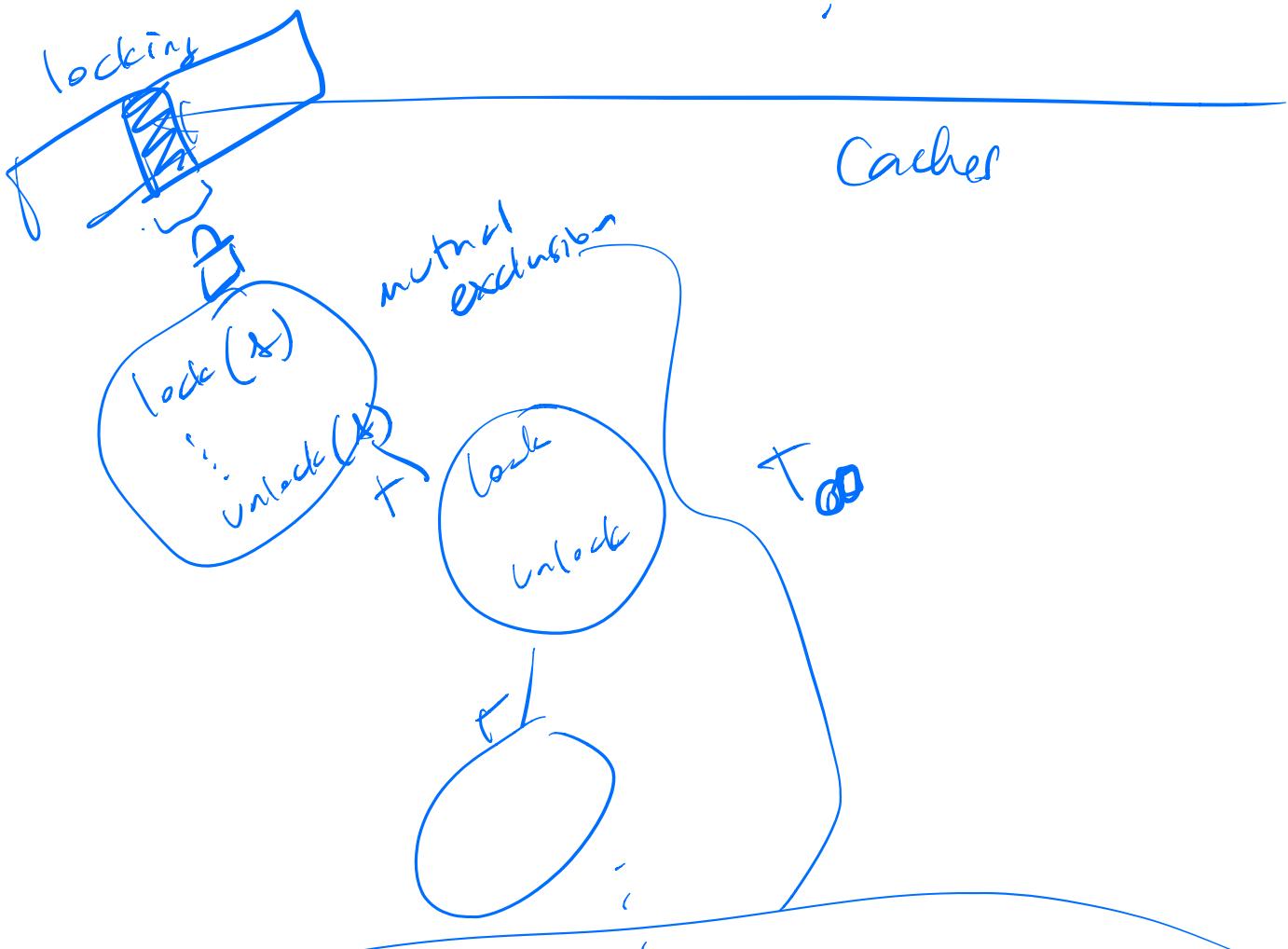
...

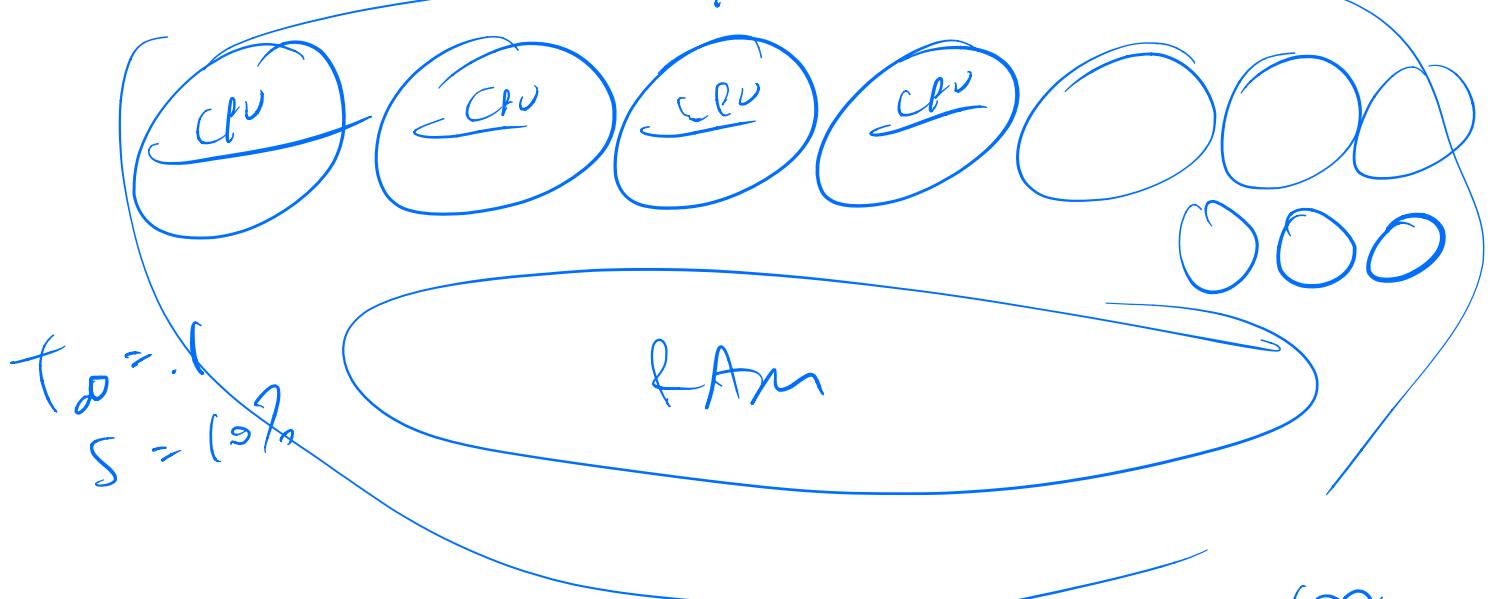
10

$$\bar{T}_{\infty} = .1 \approx 10x \text{ speedup}$$

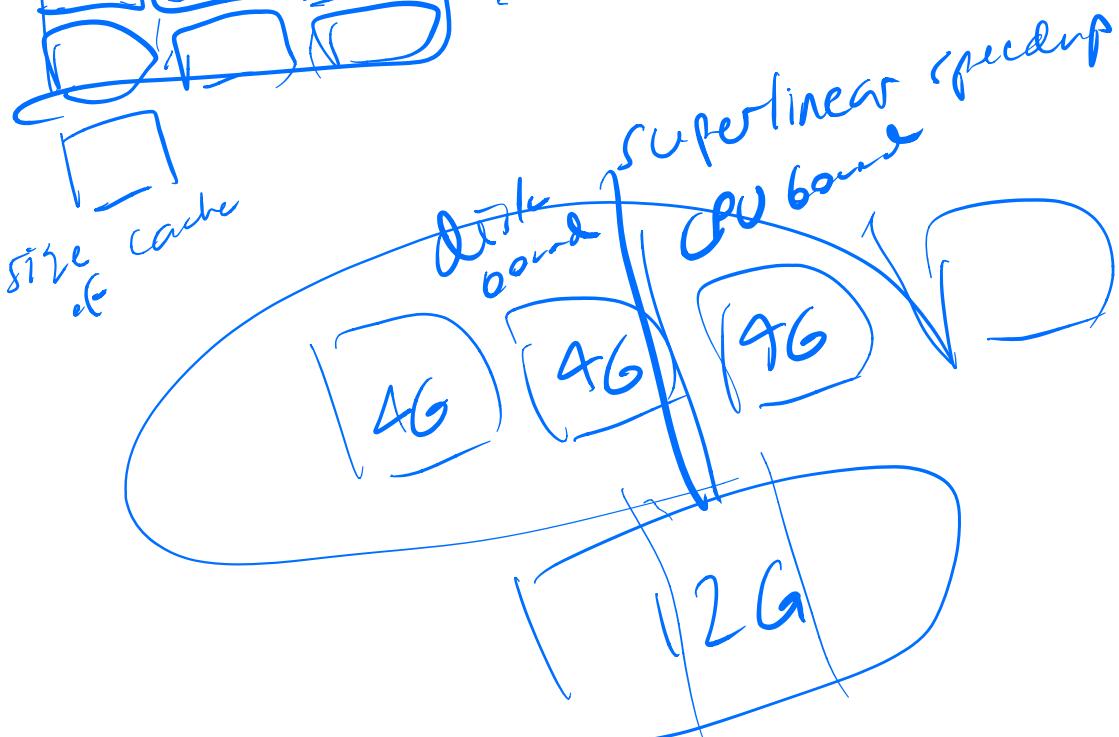
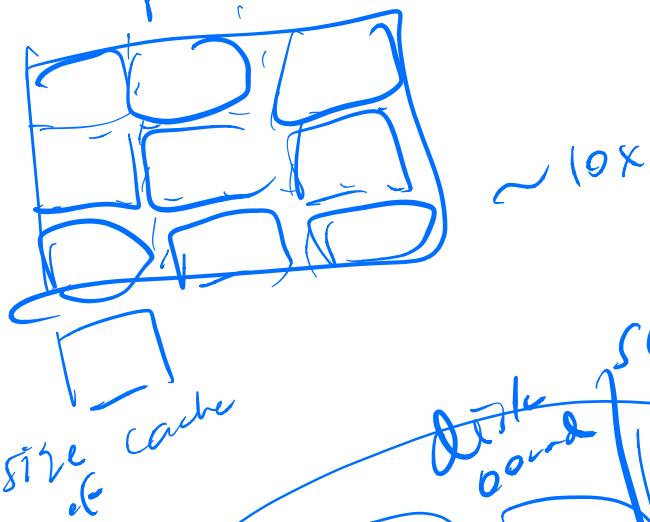
$$T_{\infty} = .01 \times 100x$$

;

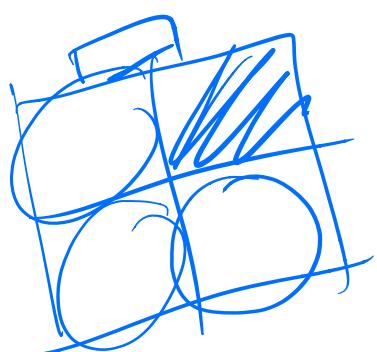
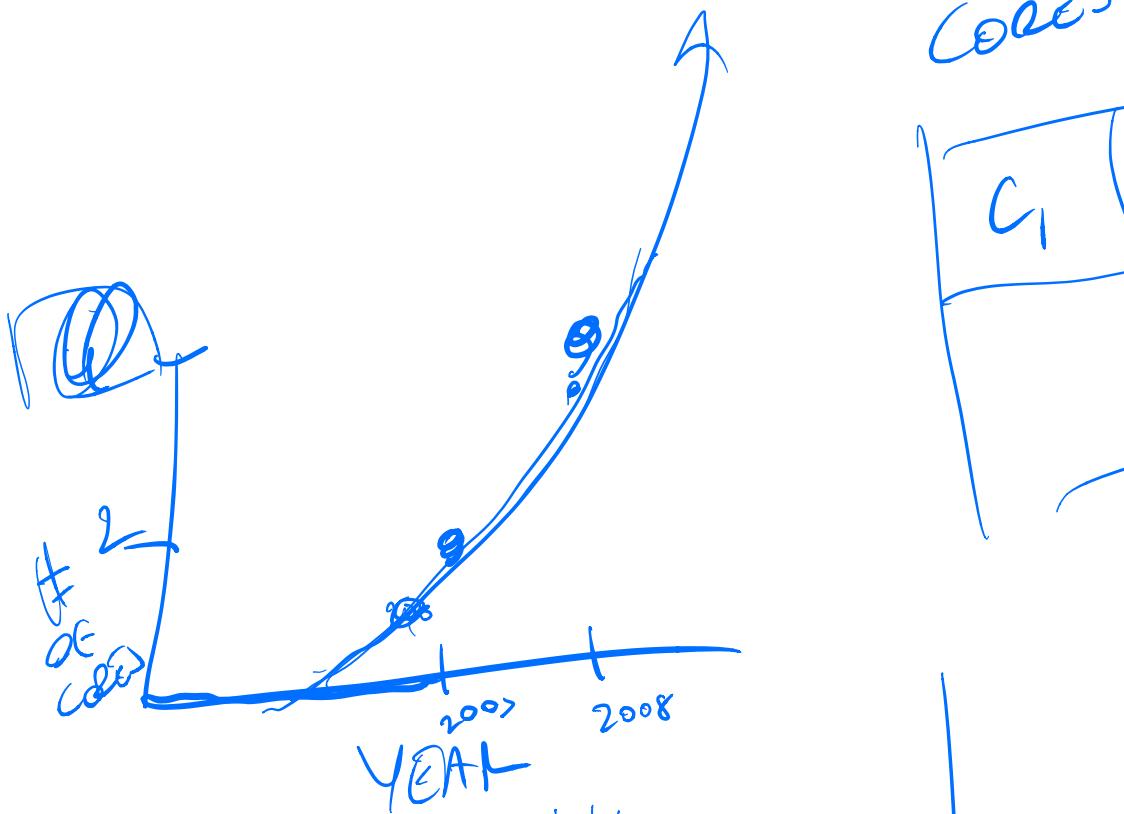
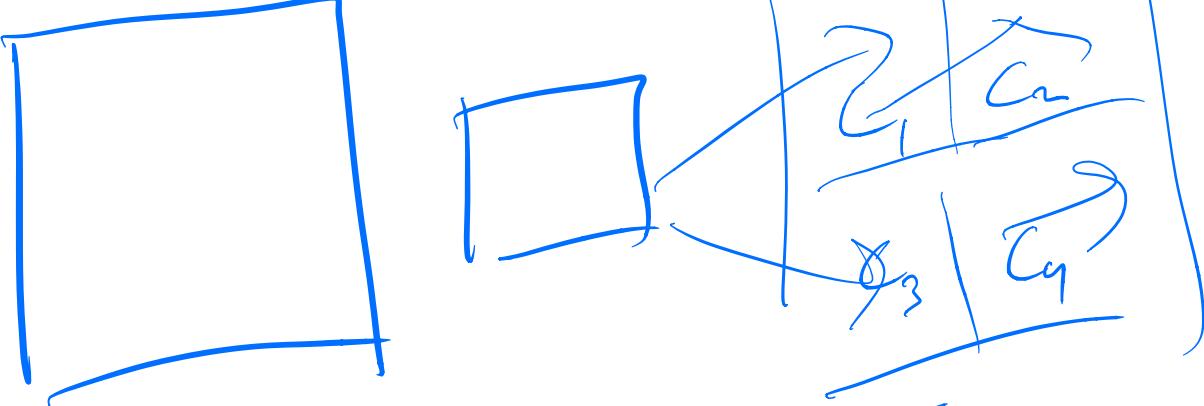




problem size $\Delta \text{dist} \leq 10x$ I see 100x



12G



Cache misses

3 core
5 core
7 core

Caches
dataflow
prefetching

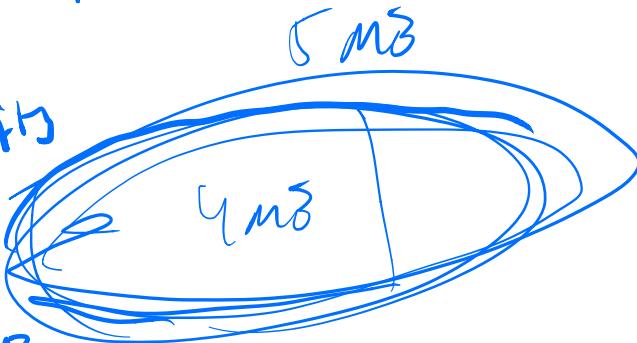
CAPACITY

COMPULSORY

First access

CONFLICT MISS

Lru feature
associativity
(full)

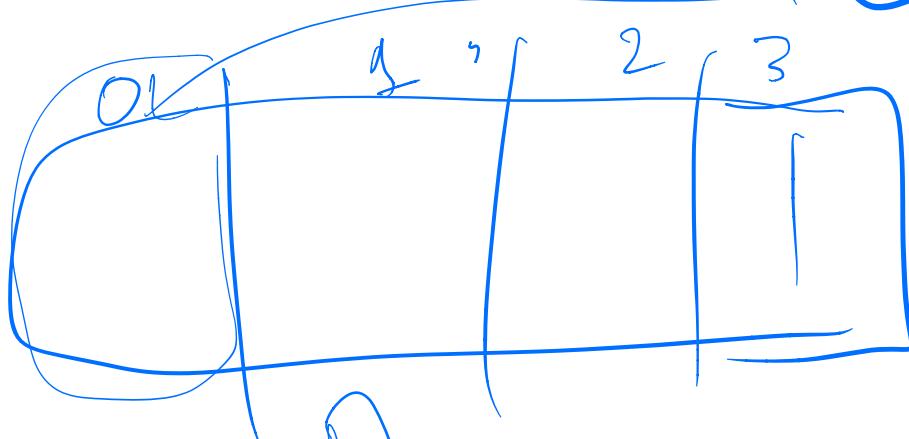


COLLISION MISS



0	0	0	0
0	0	0	1
0	0	1	0
⋮			

saturating counter



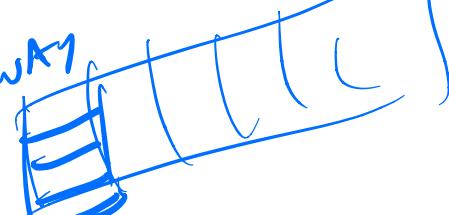
sets

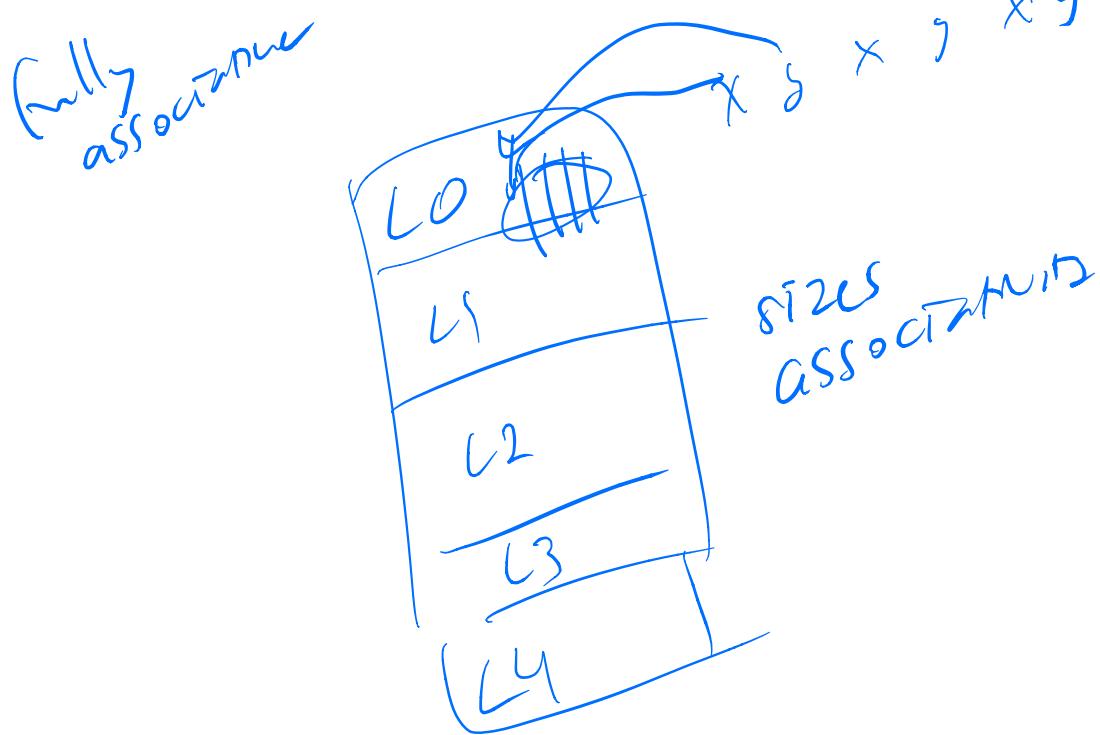
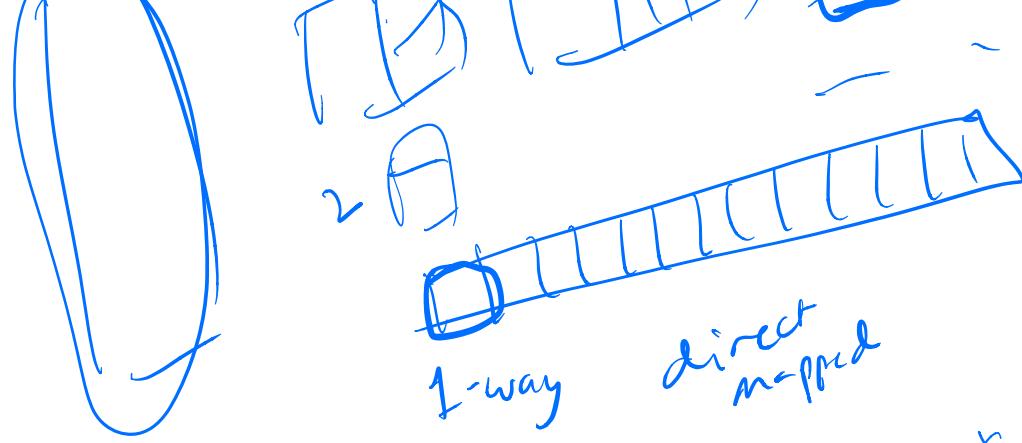
8

4

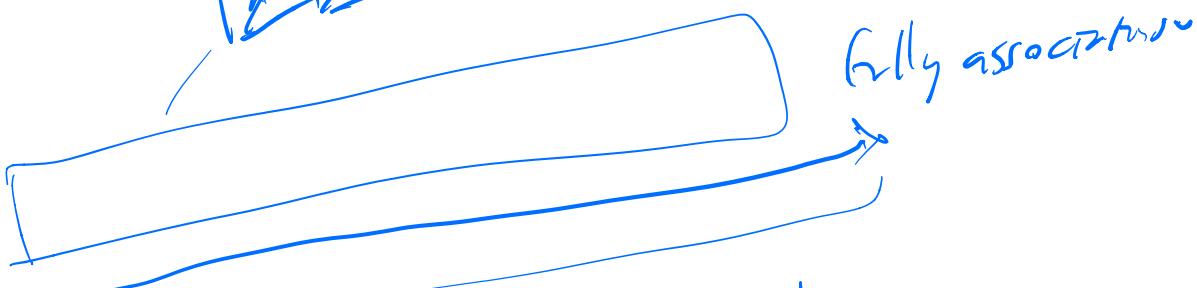
set - association

WAY

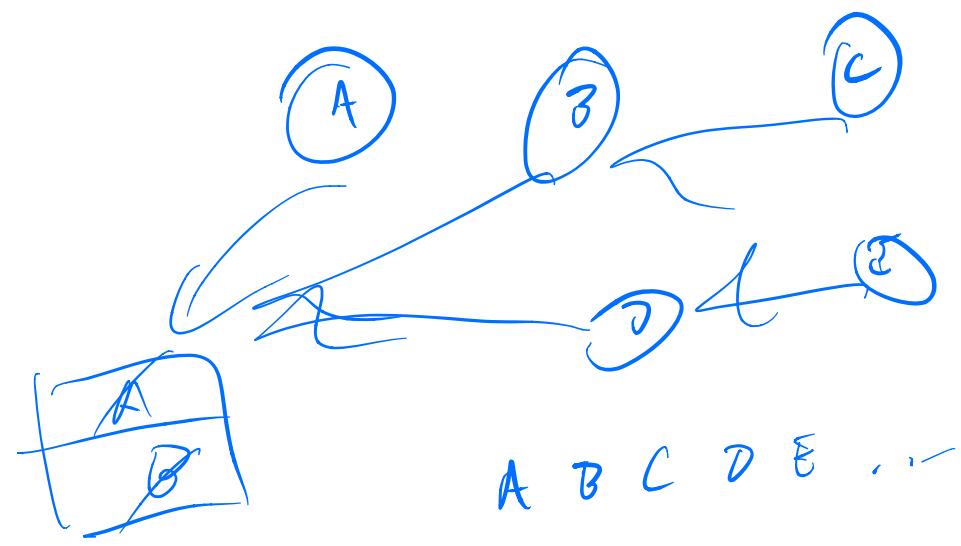




$\text{Gr}(i=0; i < 1025; i++) \{$
 $x[i] = \sim i$
 }

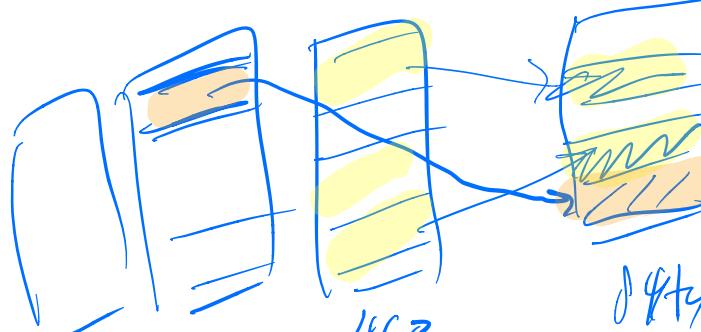
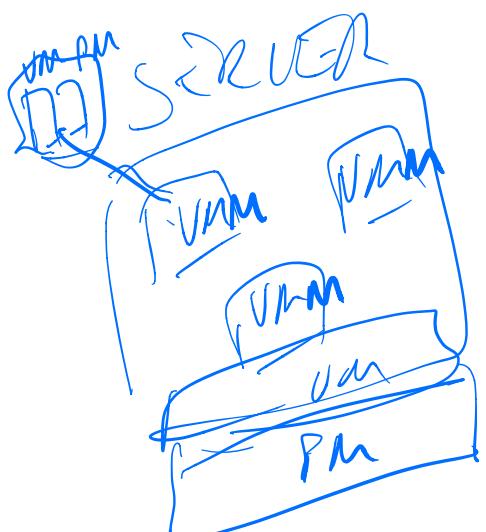
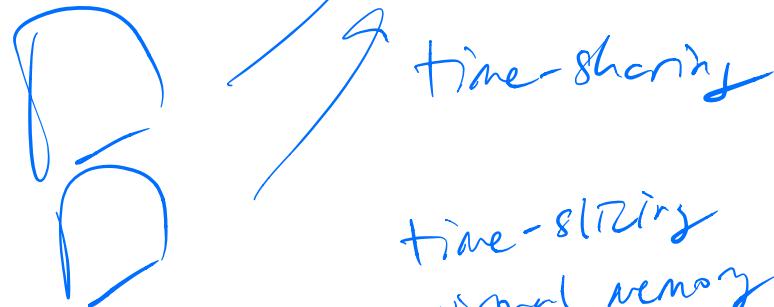
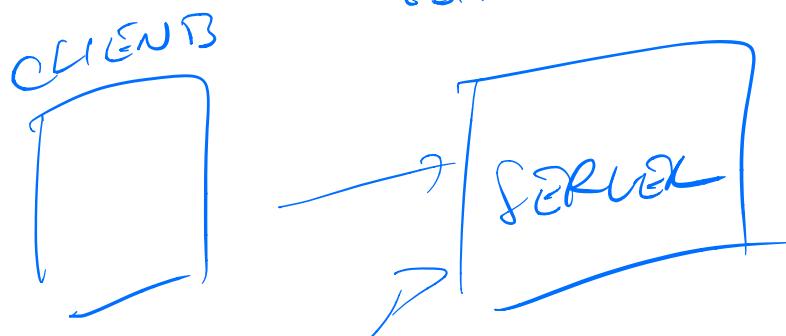


pathological
worst-case



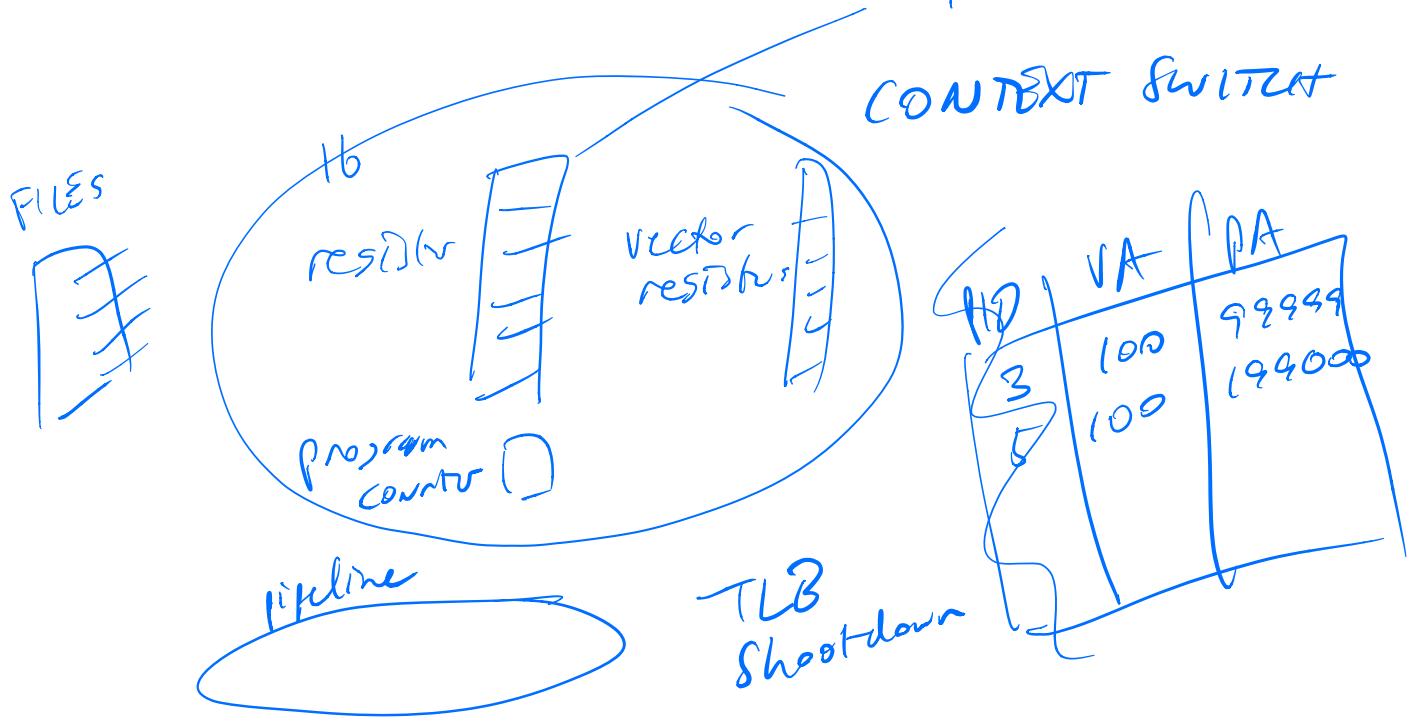
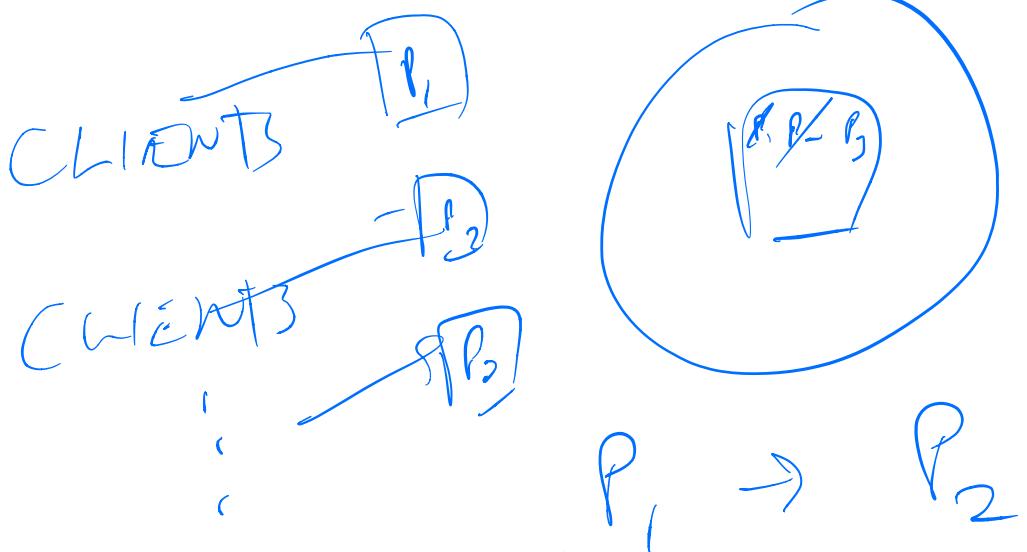
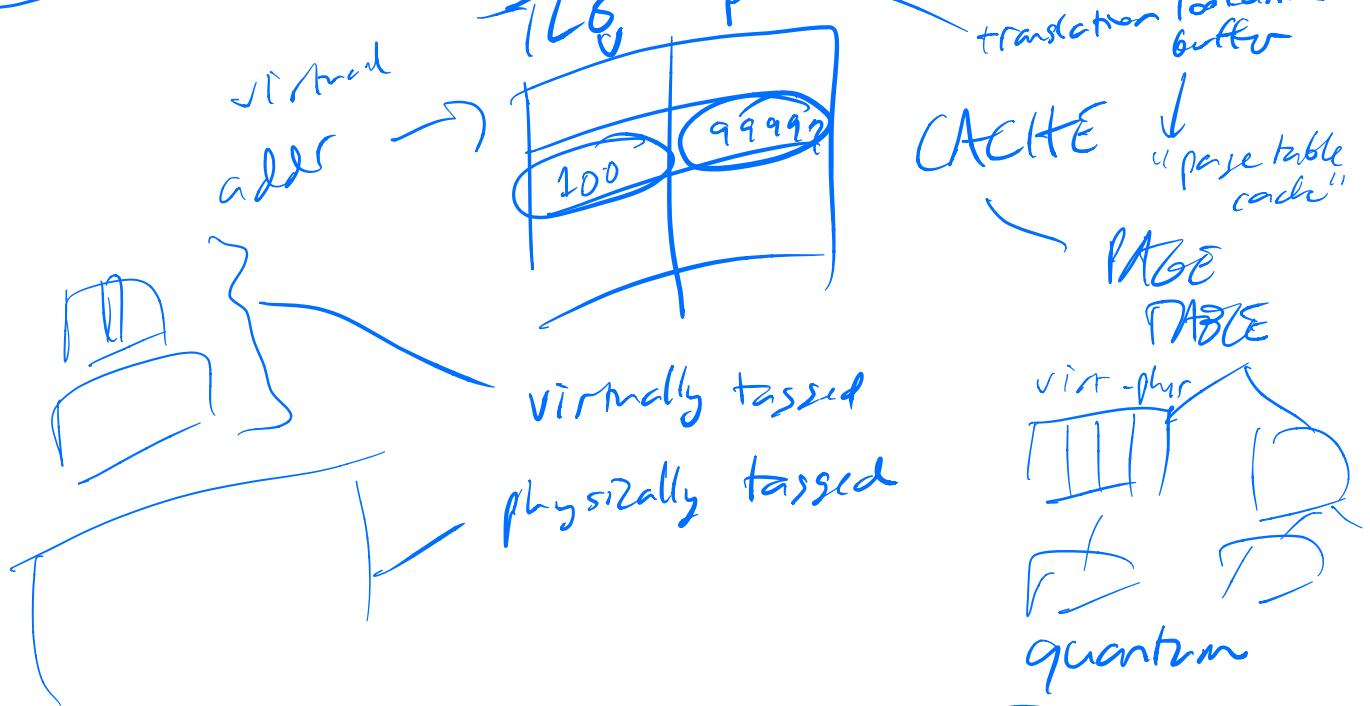
interlude

software architecture

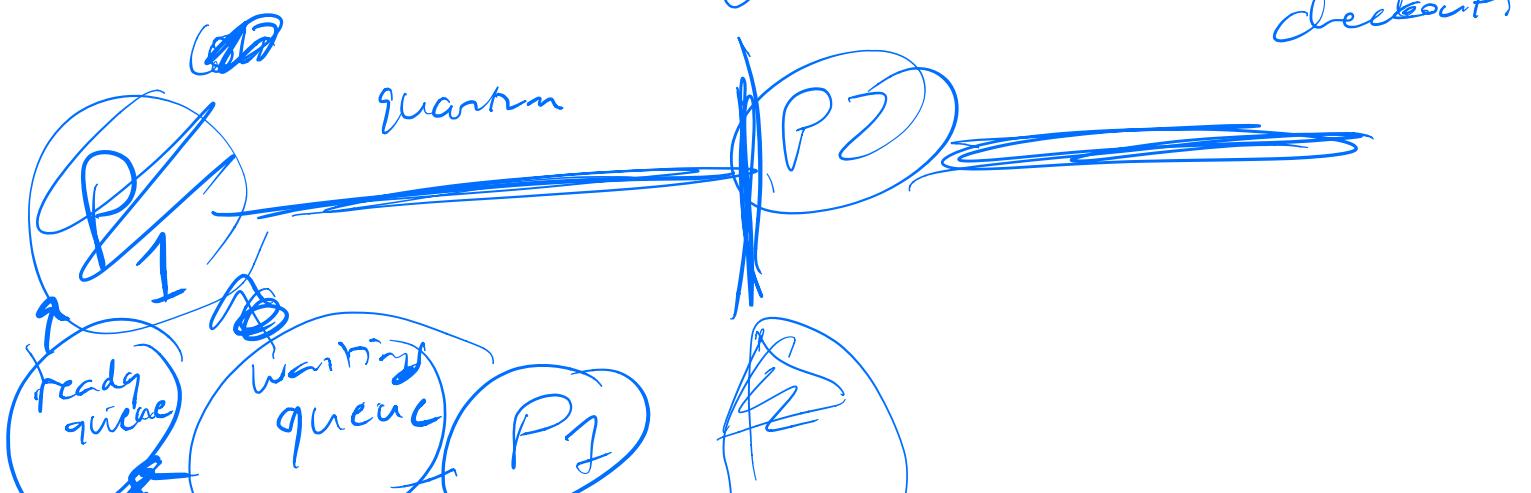
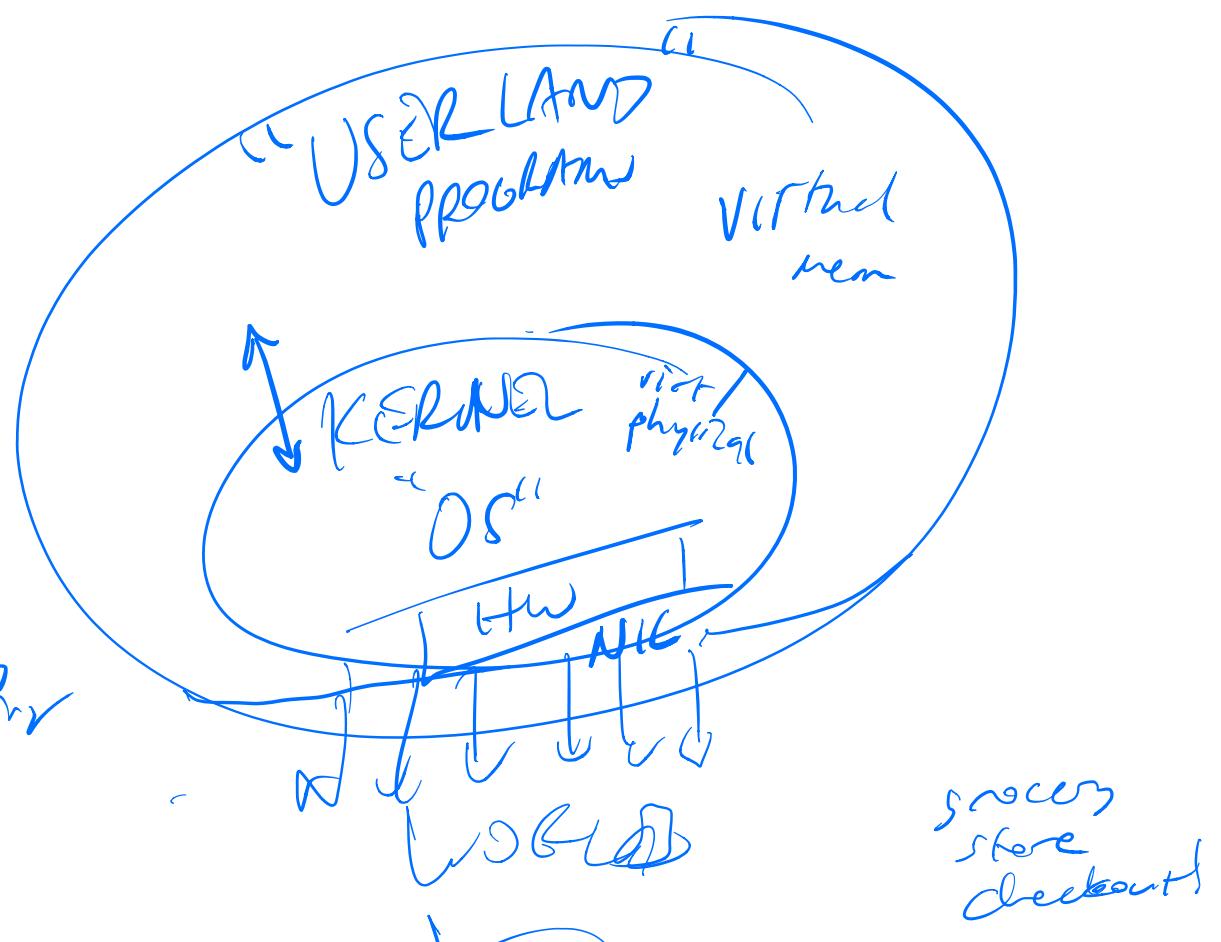
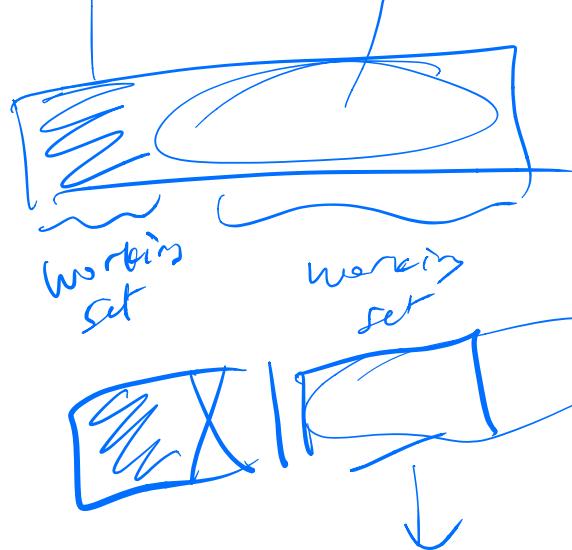
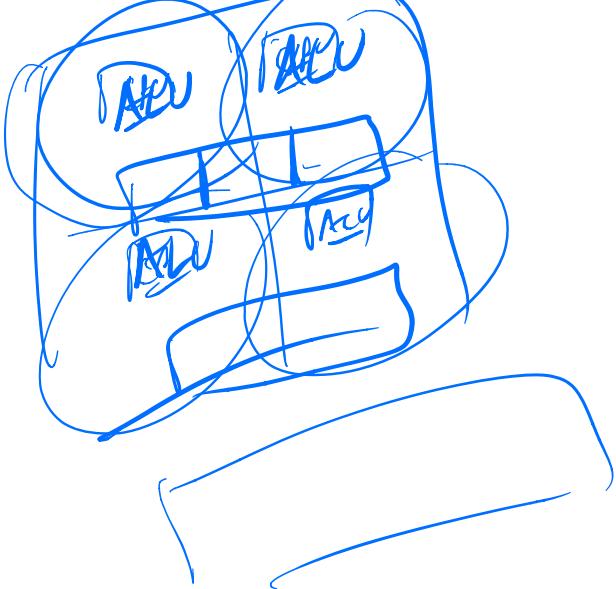


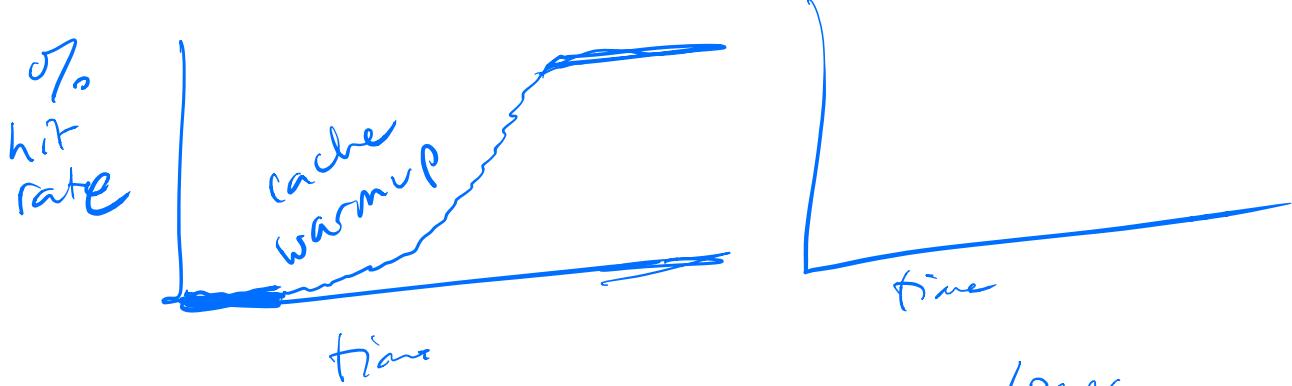
Physical
memory

1 terabyte

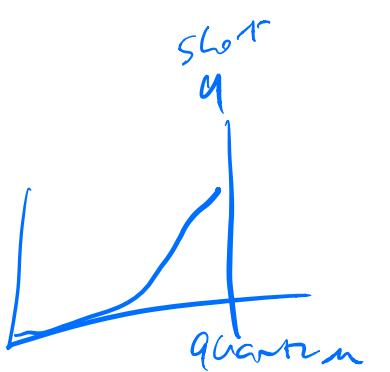


CORE₁ CORE₂

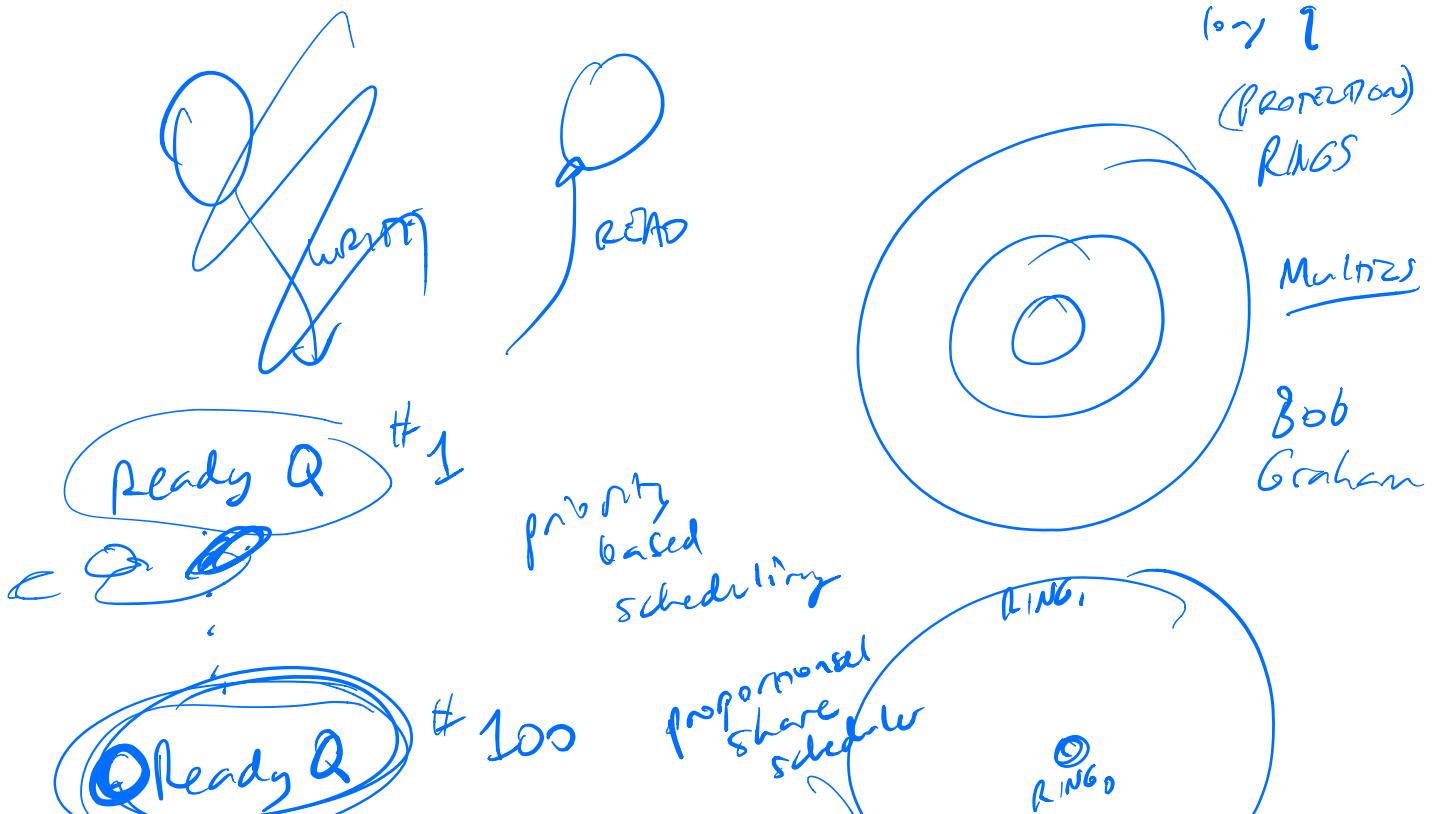
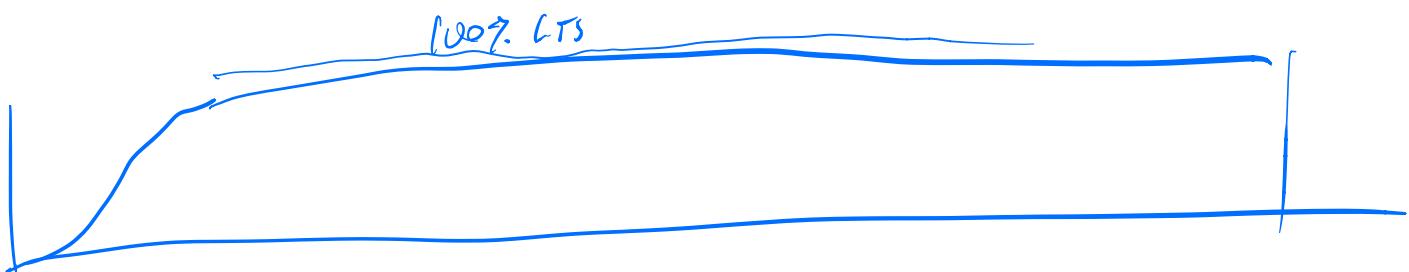


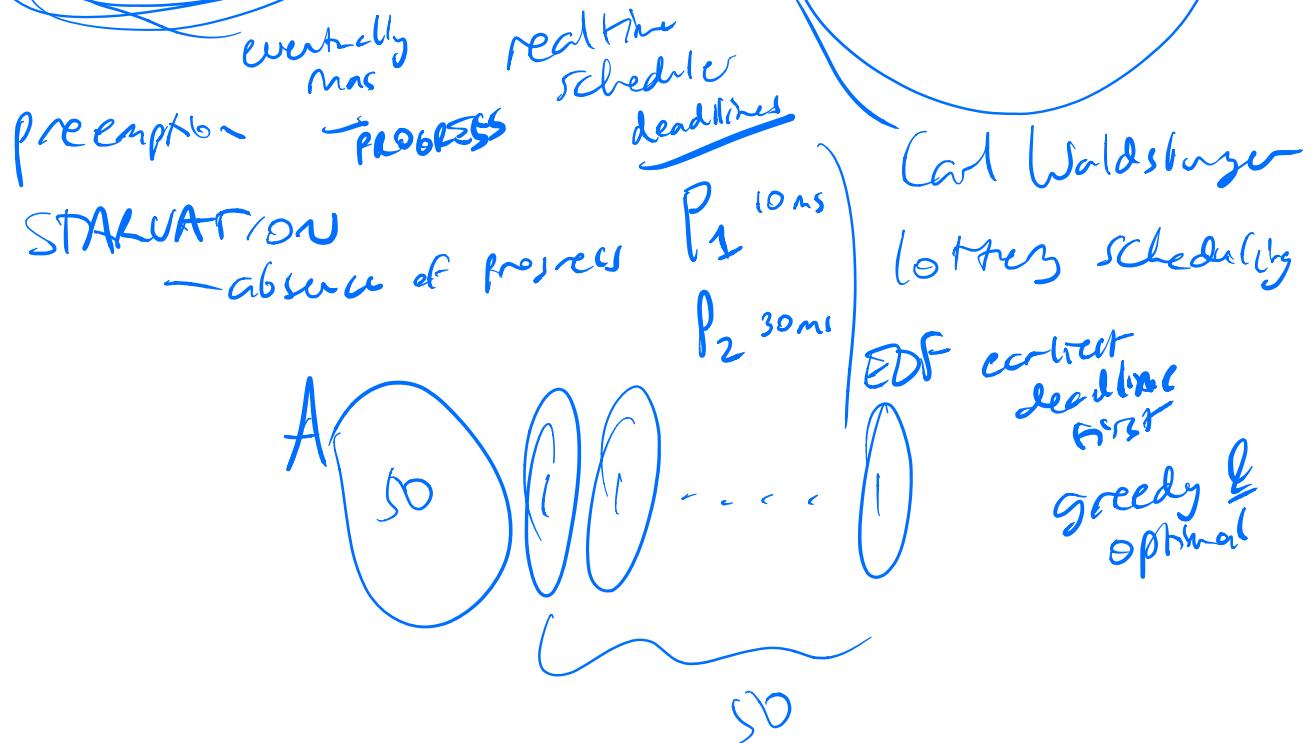


better responsiveness - lower latency
quantum - small

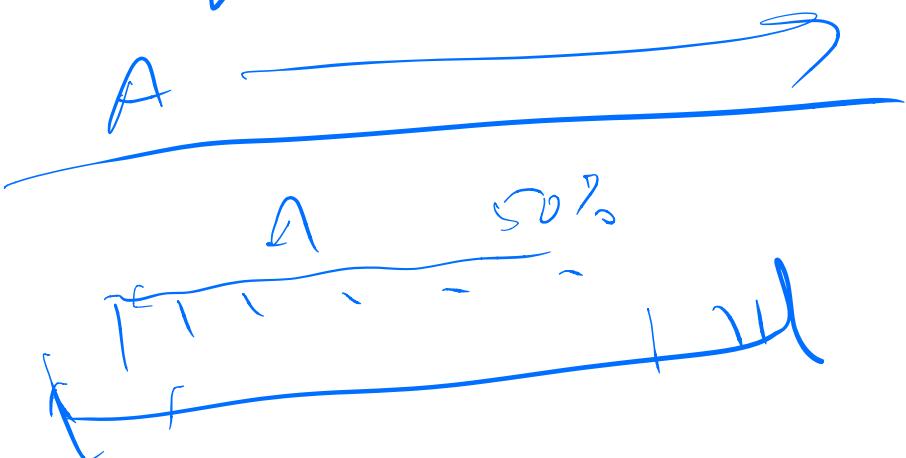


bigger throughput "Faster"
quantum - large



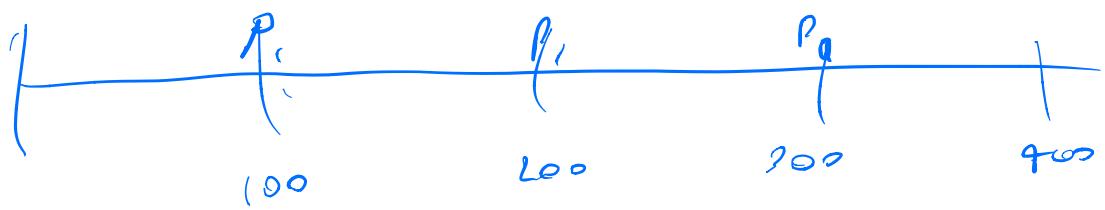


$$\Pr[A \text{ wins lottery}] \approx .5$$



$$\delta_i \in S \quad \text{corr}(\delta_i, \delta_j) = 0$$

$$\delta_0, \dots, \delta_{n-1} \rightarrow \delta_n$$



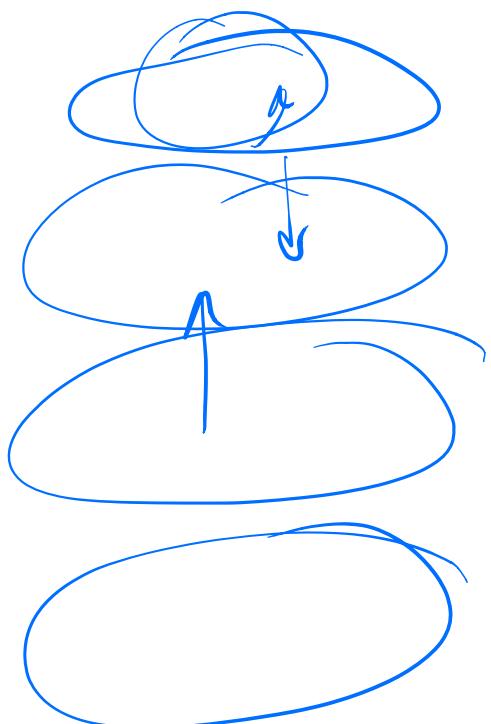
|||||

uniform sampling

|||||

requested: 1

actual: random sampling



1

priority

2

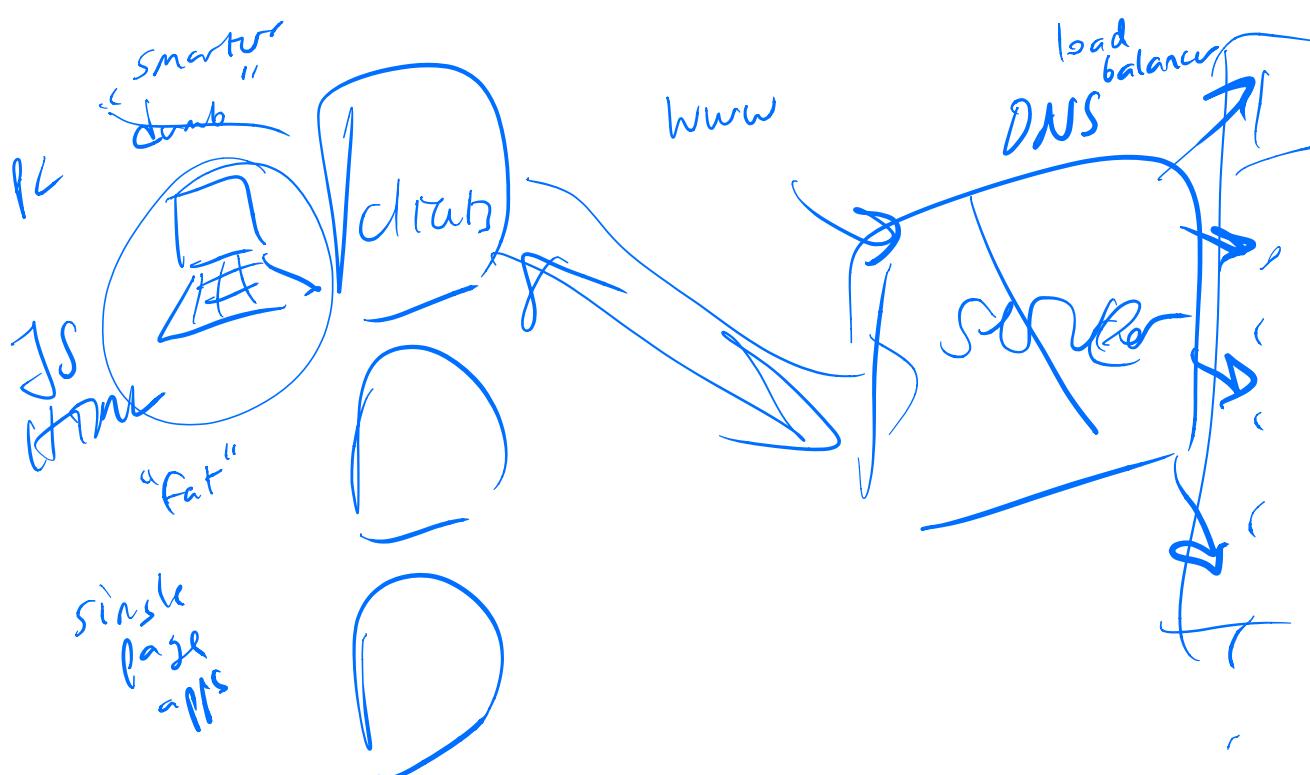
100

CPU utilization : 100%

QoS guarantees
quality of service

access speed Gbit/s availability
latency 1000 ms > 99.9%

$$\Pr \left[\text{10% wire every 100ms} \right] \geq ?$$

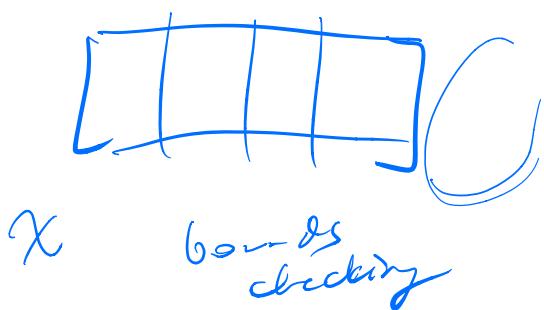


memory leak

Bleak

replacement mem allocator for Rust
(in rust)

"safe" language



mem algmt "safe"
no dangling pointer errors

Rust → no GC

explicit
malloc
fast, clear
but no
control

OWNERSHIP
TYPE

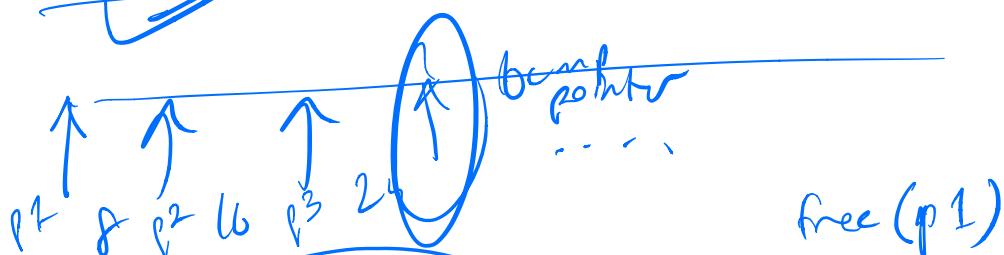
GC
safe
"fast"
mem leak

smart pointers
CTT

"met" UNSAFE
traits

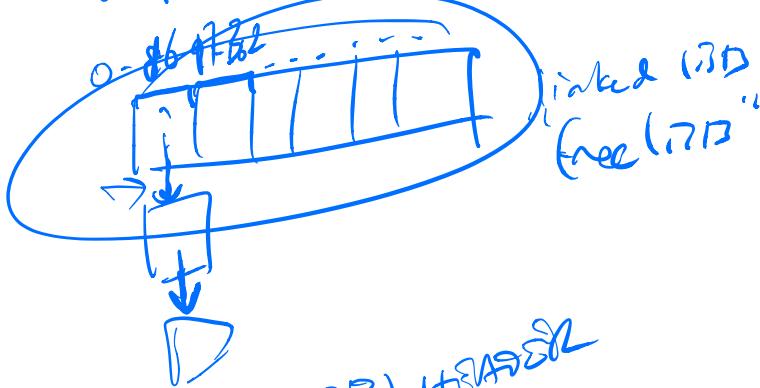
malloc
free
set size

↓
my malloc
my free
my set size



THREAD
SAFE

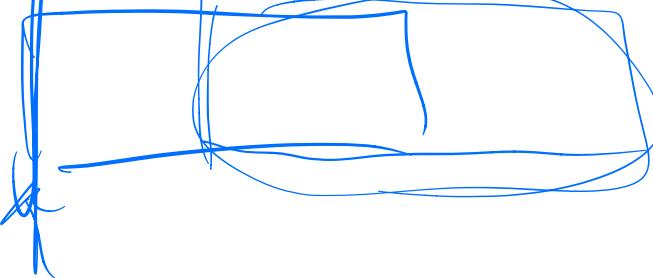
free (p_1)



16
32
64
;

power of
two size
classes

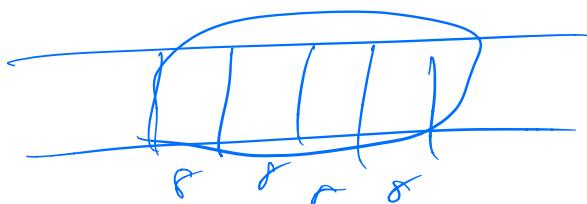
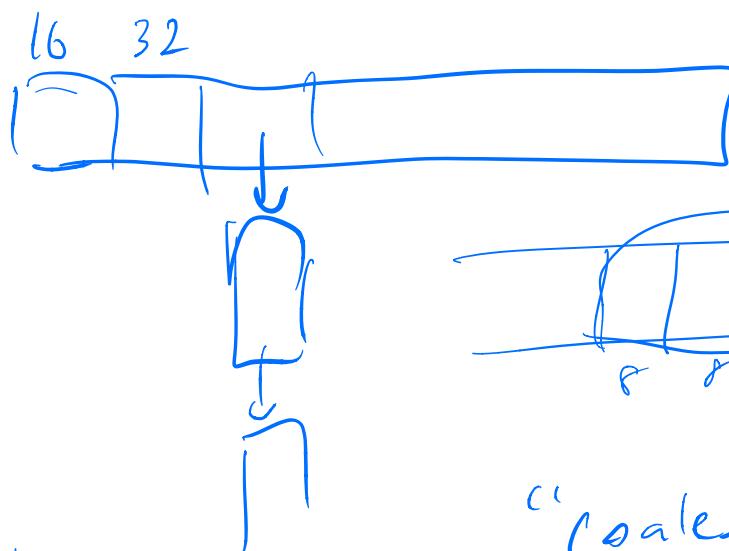
17
 $2^i + 1 \dots 2^{i+1} \dots$
 2^{i+2}



all objects
must
be
"aligned"

SEGREGATED
FITS
(BEST-FIRST)
ALLOCATOR

fit size 0 8
6 byte 0

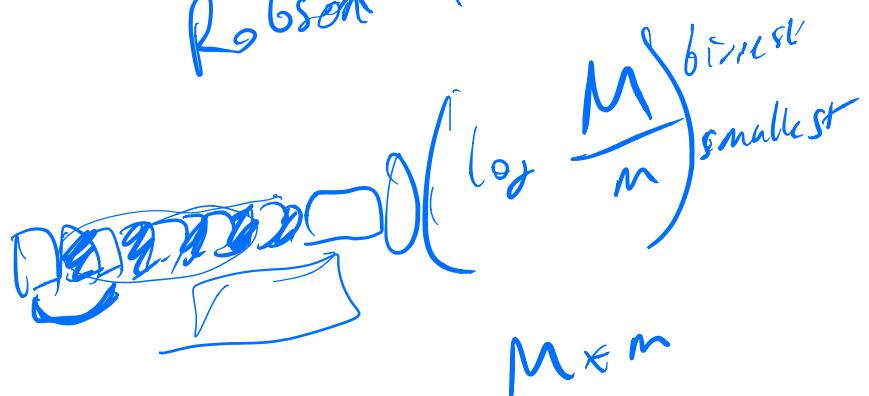


Paul
Wilson
92/94

"Coalescing"

sums
of
uniprocessor GL
mem allocator adjacent small freed objects →
One big freed object

Robson (972)



memory

power-of-two size classes

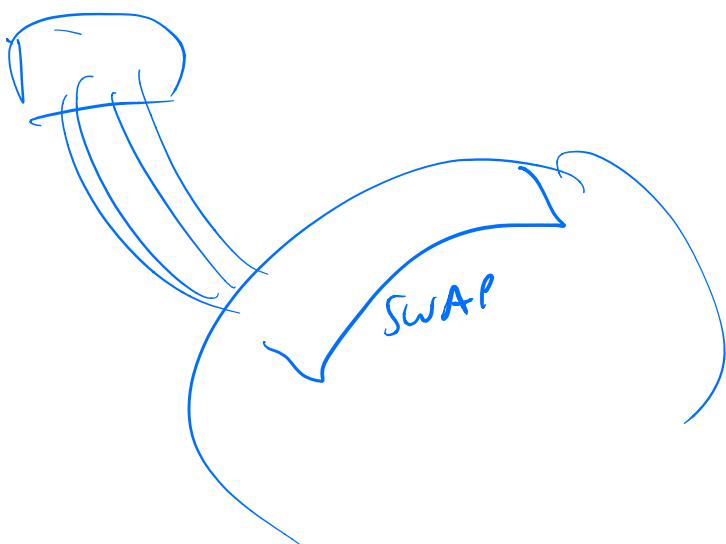
~~integers~~

arrays (stacks)

64

bit

systems



OS —
has 1 exabyte

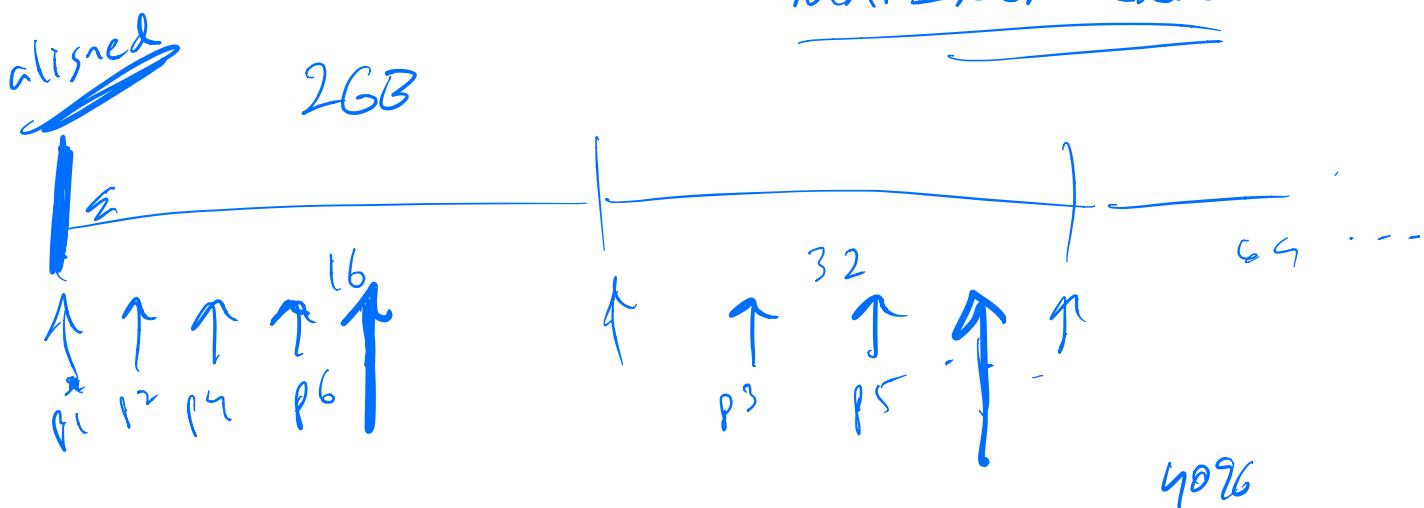
Memory mapped

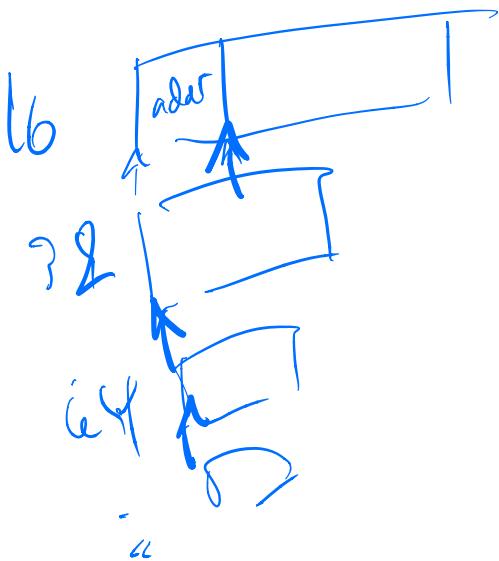
Madville
(DONTNEED)

"mmap"

Readable
writable
Anonymous
private

MAP_NORESERVE



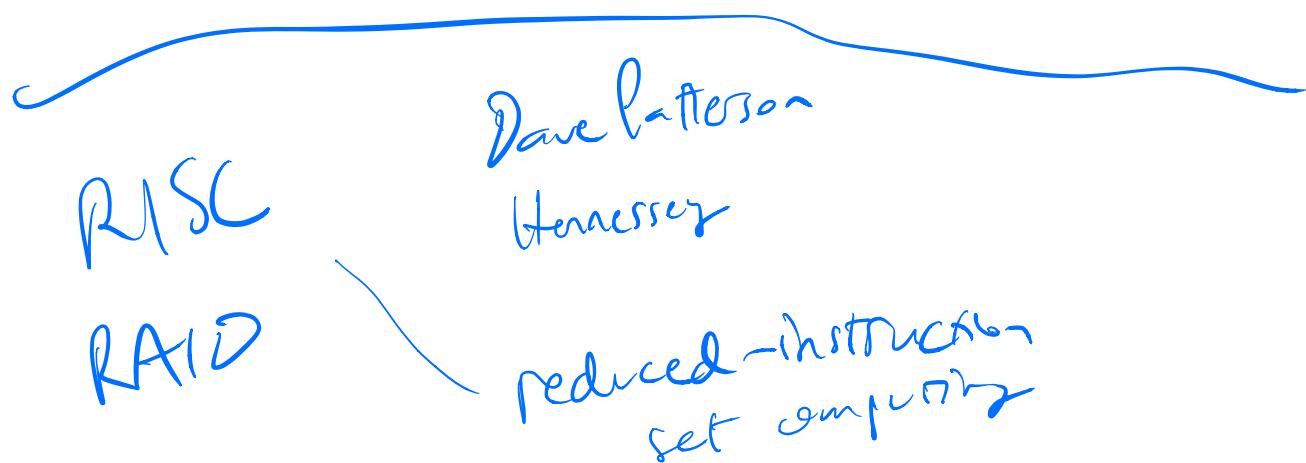


malloc: If stack ! empty
pop stack & ret ptr
else memos 0 to
bump memos & return

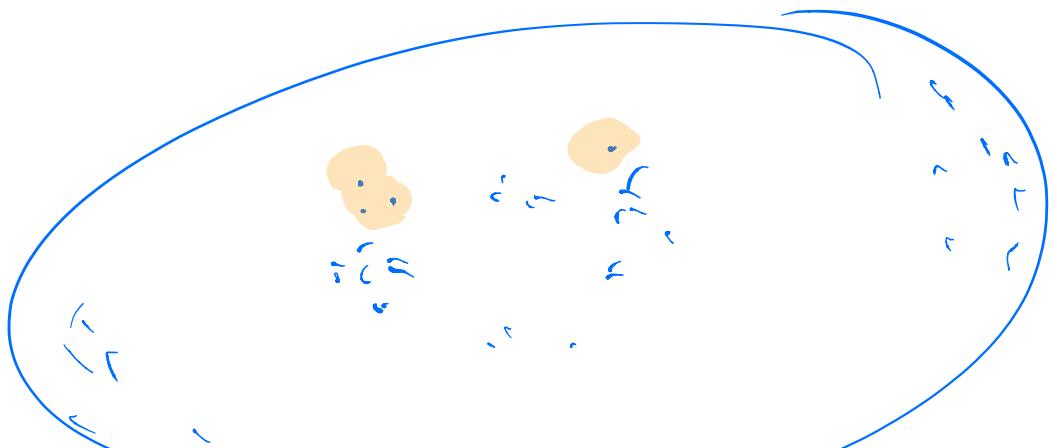
free: push ptr onto stack
inUse
maxAllocated

"Heap"

zero on-demand page



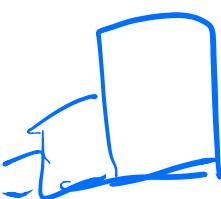
every new instruction for life



CISC

Conventional instruction set computing

~~SIMD for compiler writer~~



~~RISC~~

RISC-II

Kosta Asanovic

ARM

power
energy consumption

very
few
addressing
modes

MOV
ADD

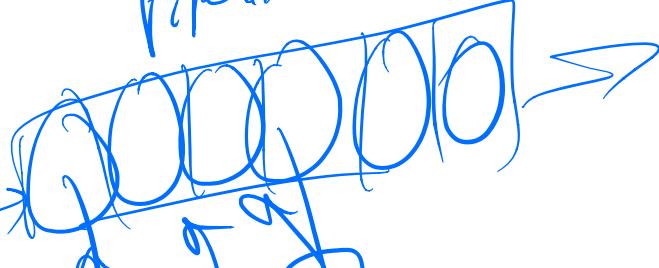
same size

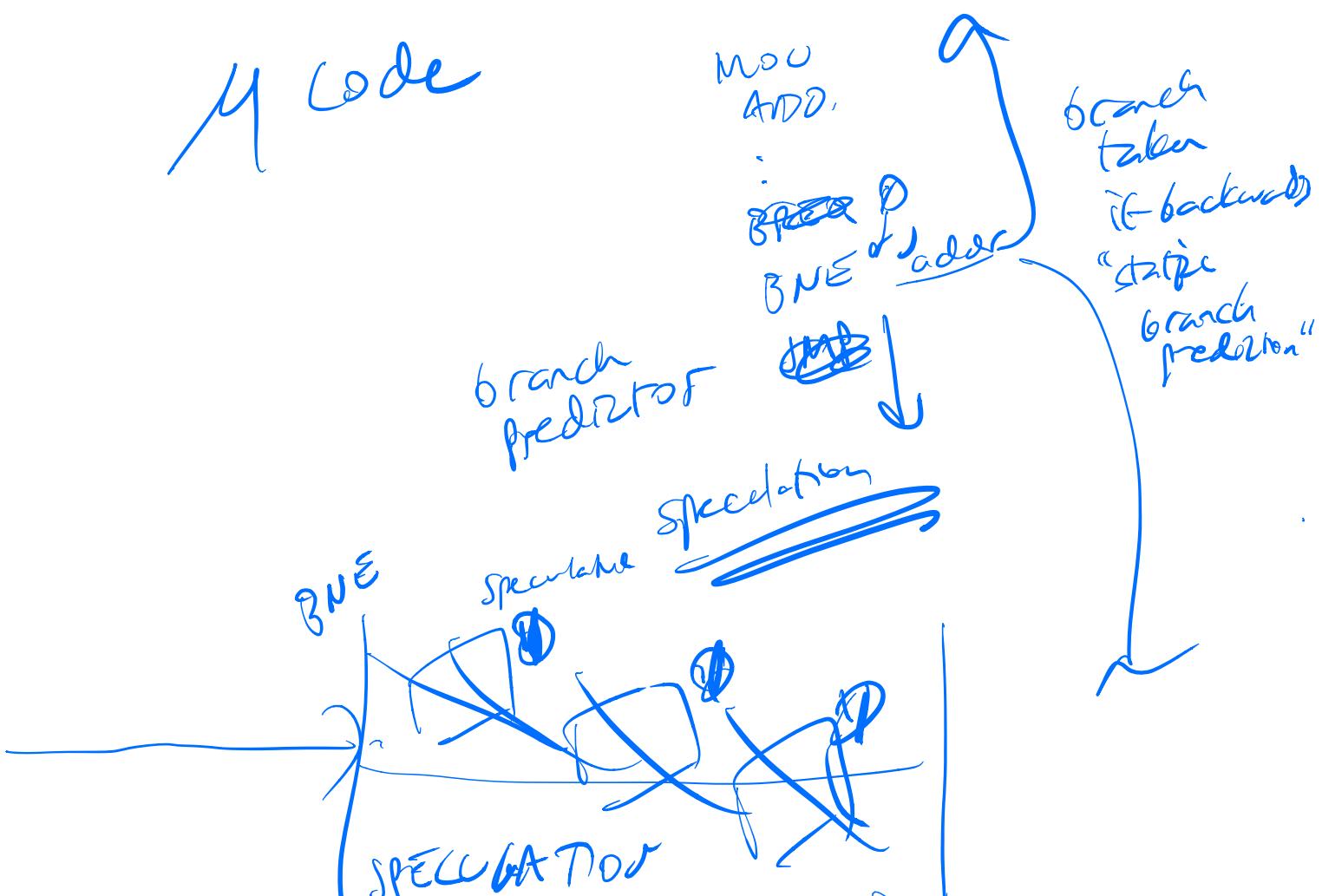
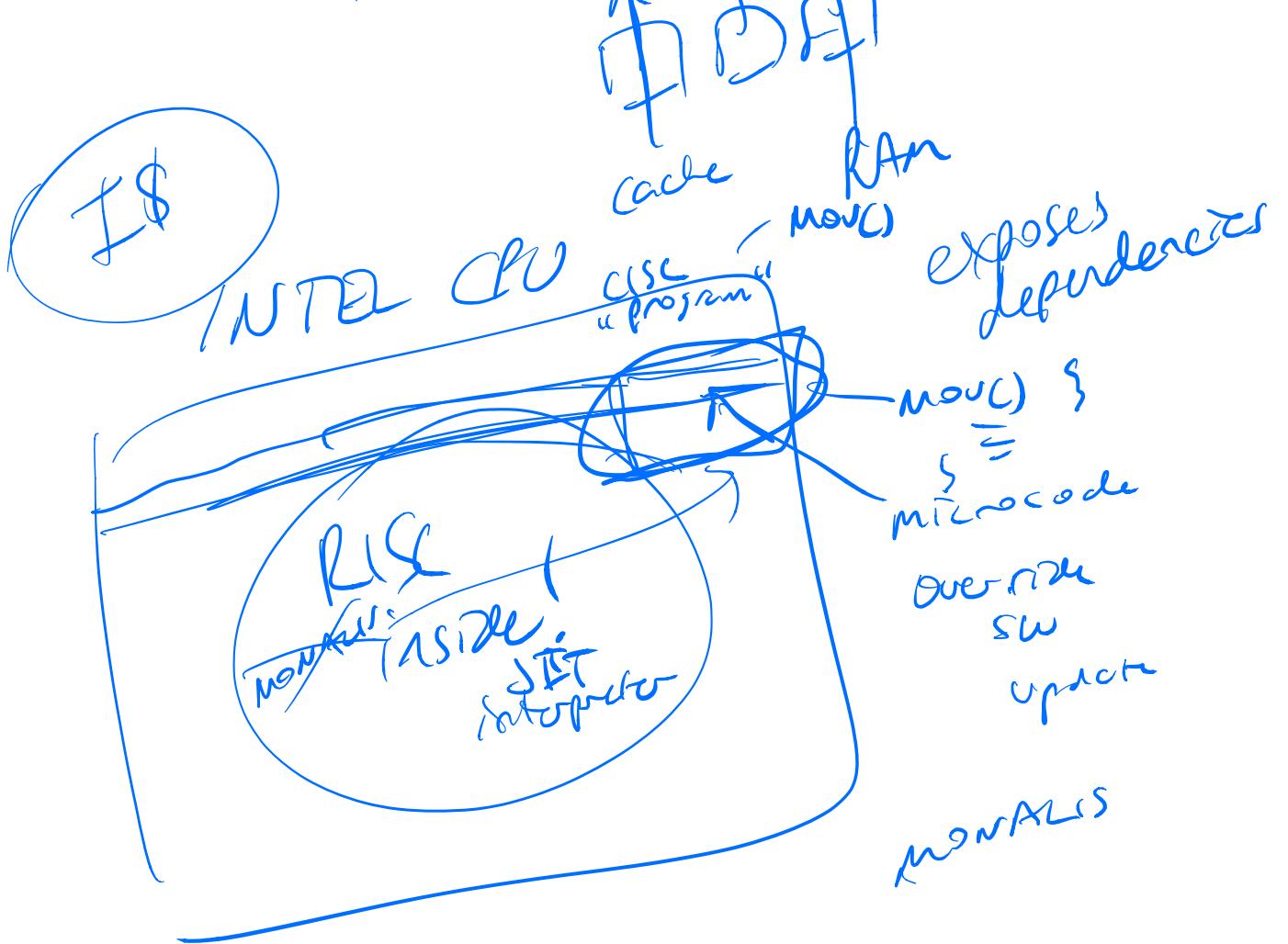
SSE
SSE2
SSE3
crypto
MMX
TSX

fetch
decode
execute

ROP

pipeline





speculative
state

~~yes!~~
~~No~~

