

Energy Berger

energyberger.github.io / CompSCI-630

Systems Lunch @ Fridays - Fall

Random number

TA Abhinav Jangda

Github!

Hot CRP

energy@cs.umass.edu

- Scribes!
 - course load :
 - 1 core / semester for MS
 - ..
 - 1st semester MS / PhD

> 1 core = too much!
 - schedule - yes, we have a midterm & final!
 - reviews
-
- ## FORTRAN (in COBOL)
- 1954
- | | |
|--|--|
| military <ul style="list-style-type: none"> - calculating - artillery trajectories - decryption | first computers
— people! |
| general-purpose computing machine | automatic computers |
| Turing machine
von Neumann machine |] data & code together
Harvard architecture |
| | fetch-decode-execute cycle |

assembly language

- lowlevel
- direct access to memory
- operations not abstract

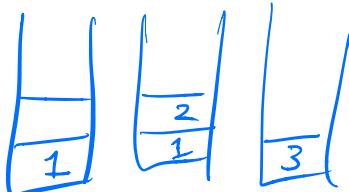
- stack architecture
- register architecture

A, B/C/D
AX BX CX DX

- portability

ASM ≠
~~portable~~

MOV 1, A
MOU 2, B
ADD A, B, C



PUSH 1
PUSH 2
ADD
POP

"Intel" "ARM" (subset of instr.)
JVM, .NET, Python bytecodes,

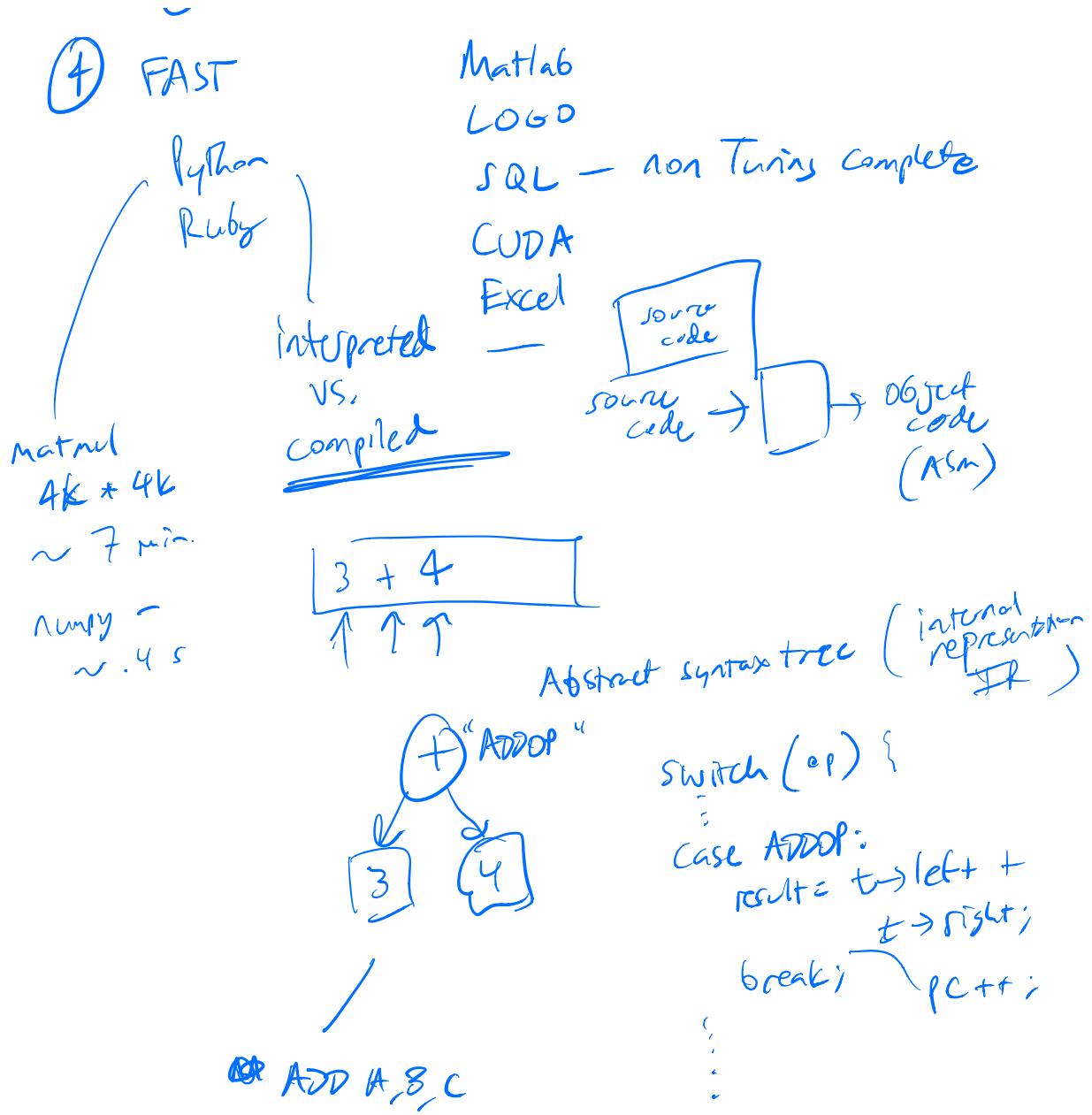
WASM

backward compatibility

- hard to read, hard to debug

- ① PORTABLE
- ② READABLE
- ③ WRITABLE

- COBOL - managers
FORTRAN - scientists } domain-specific languages



programmer time MORE USEFUL than cycles

— modern POU

- productivity
- debugging
- ~ 100x faster!

interpreter → just-in-time compiler (JIT)

Python PyPy
JavaScript V8
 Chrome

WASM

Predicates -
- Parsing
- Compiling
- program analysis
- type coercion
- scope
- object orientation
- modularity
- structured prog.

FORTRAN
1st optimizing compiler

I, J, K integer
X, Y, Z float

DO 10 I = 10 10 *

(was comma) line number

10 CONTINUE

BASIC



$\text{DO } \cup^1 \phi \cup^2 \psi = \cup^1 \phi$

control flow

$\text{DO } \phi \psi = \phi \psi$

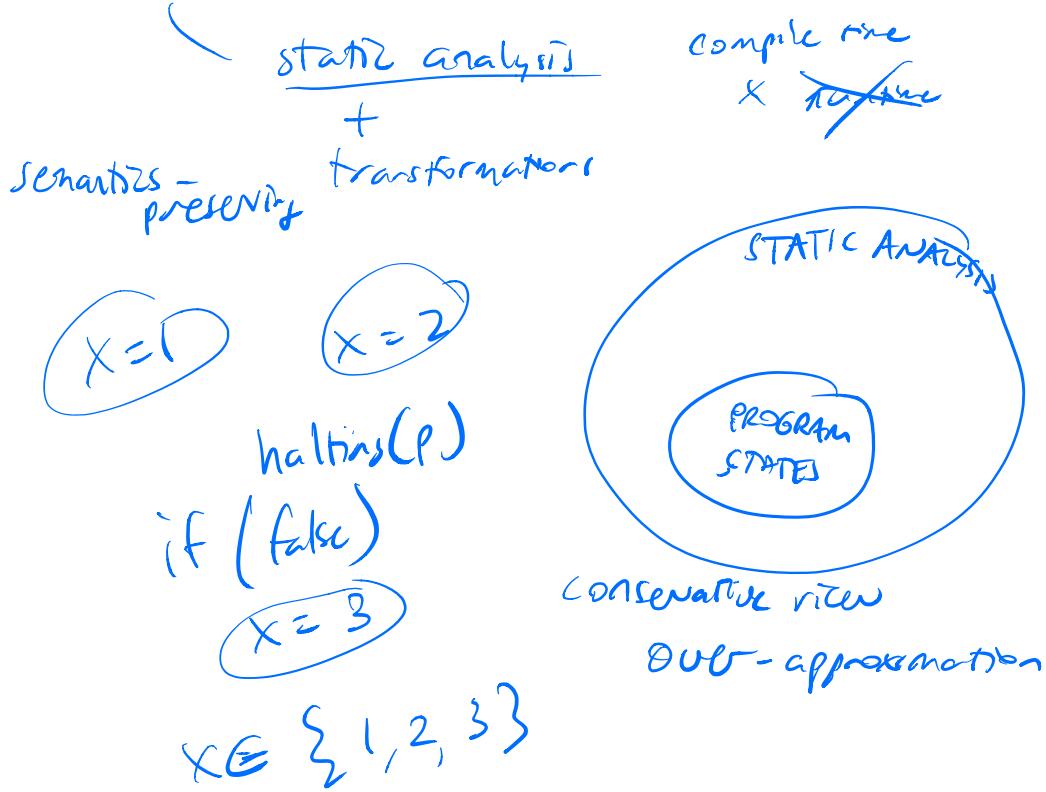
FORTRAN 77, FORTRAN 90 ~

writing

LEGACY
FAST

The Fastest PC today! 2019 !!

optimization



Halting problem

halts (program-text)
→ true if program terminates
 ^{always}
→ false otherwise

p' if $\text{halts}(p)$
 run forever
else terminate

① Constant folding

$$x = 3 + 4$$

↓

$$x = 7$$

* PI * PI

② Constant propagation

$$x = 7 \quad z = 12$$

$$y := x + 3 + z$$

↓

$$x = 7 \quad z = 12$$

$$y = 22$$

③ ``strength' reduction

$$x = x \times x$$

↓

$$x = x * x$$

/ %

expensive
cheap

$$x = x \% 8$$

↓

$$x = x \& 7$$

_____ 111

④ vectorization

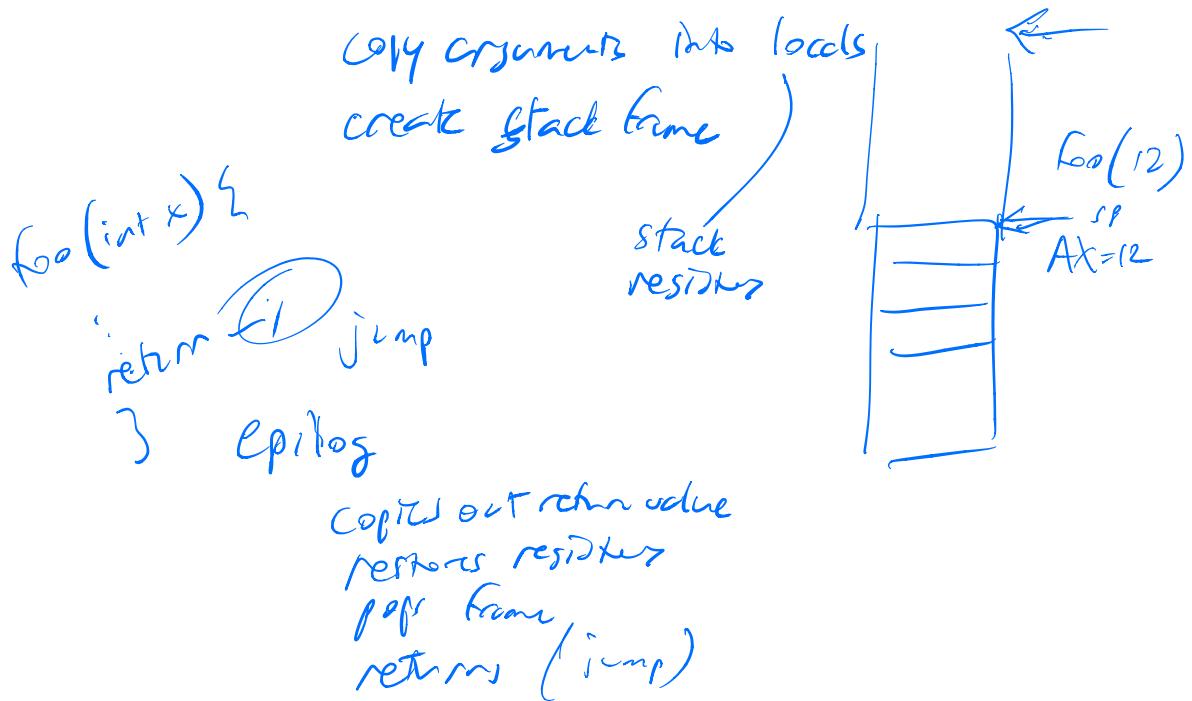
8	12	13	45
↓			
16	24	26	90

```
for (i=0; i<16; i++) {
    a[i] *= 2;
}
|||  

vector ops (4x speedup)
```

⑤ inlining

function prolog

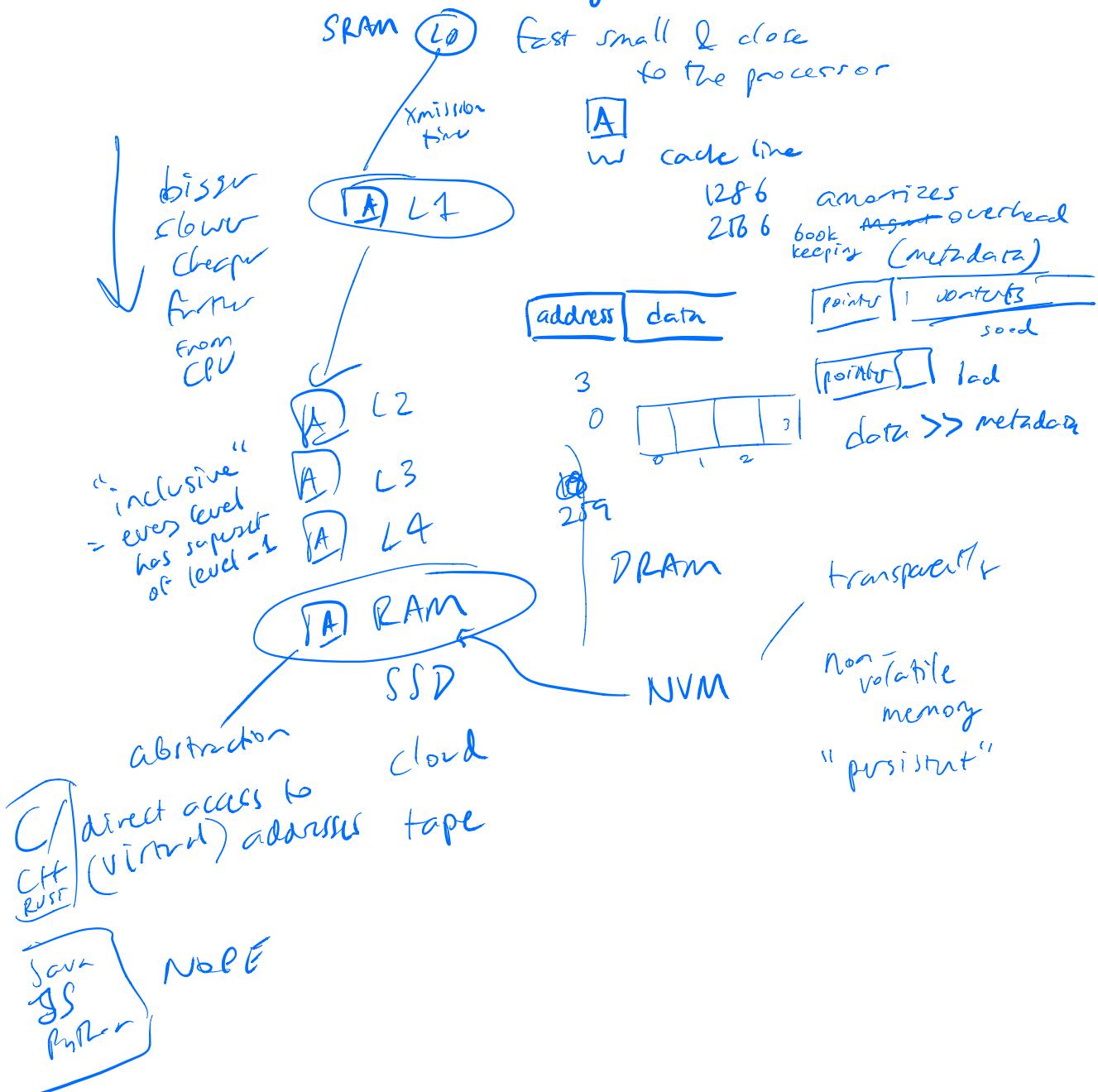


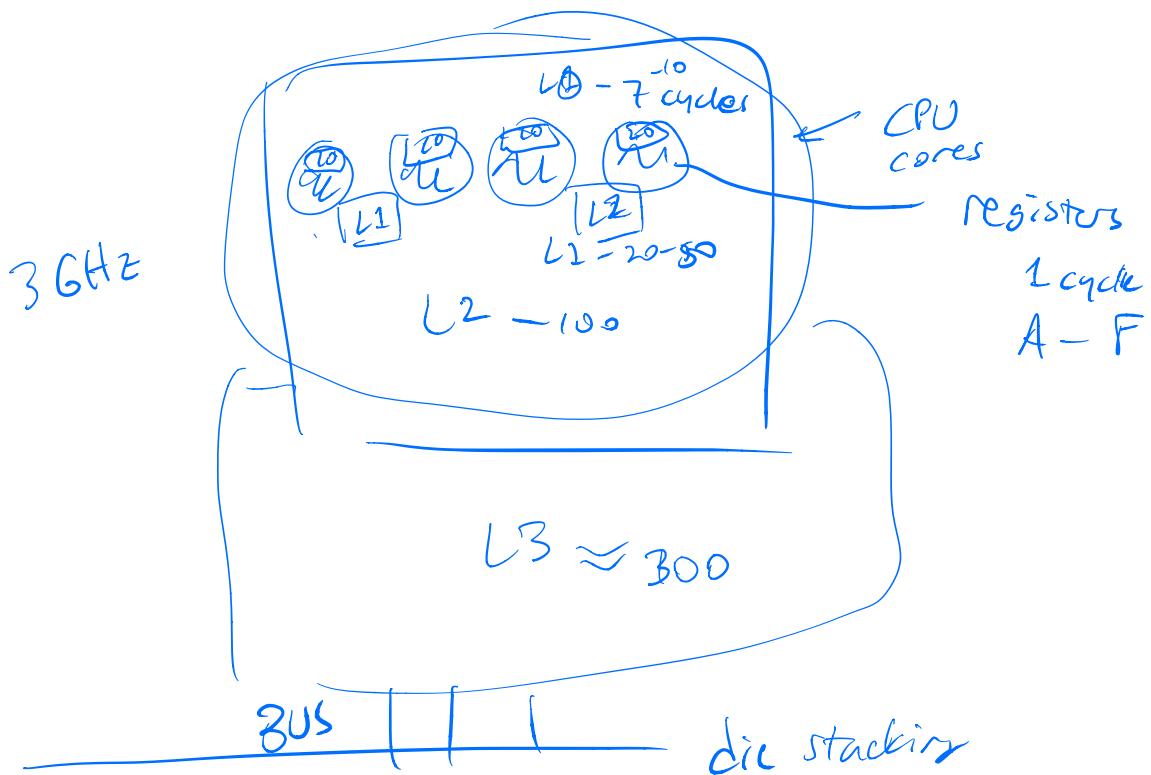
```
int add(int a, int b) {  
    return a + b;  
}
```

foo() {
 z = add(a, b); ↴
}
↓
foo() { z = a + b; }

- Code size
 - I\$ instruction cache

memory hierarchy





RAM "1Gbs"

~~100~~

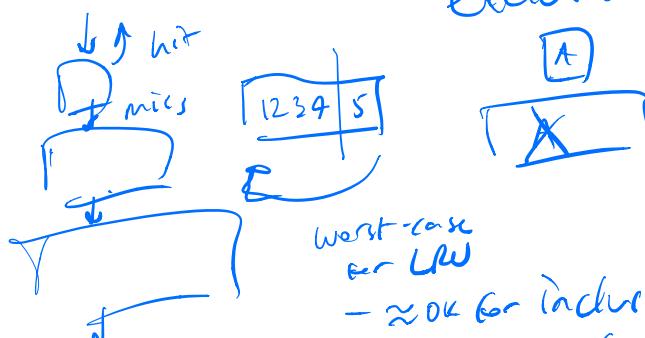
RAM bandwidth

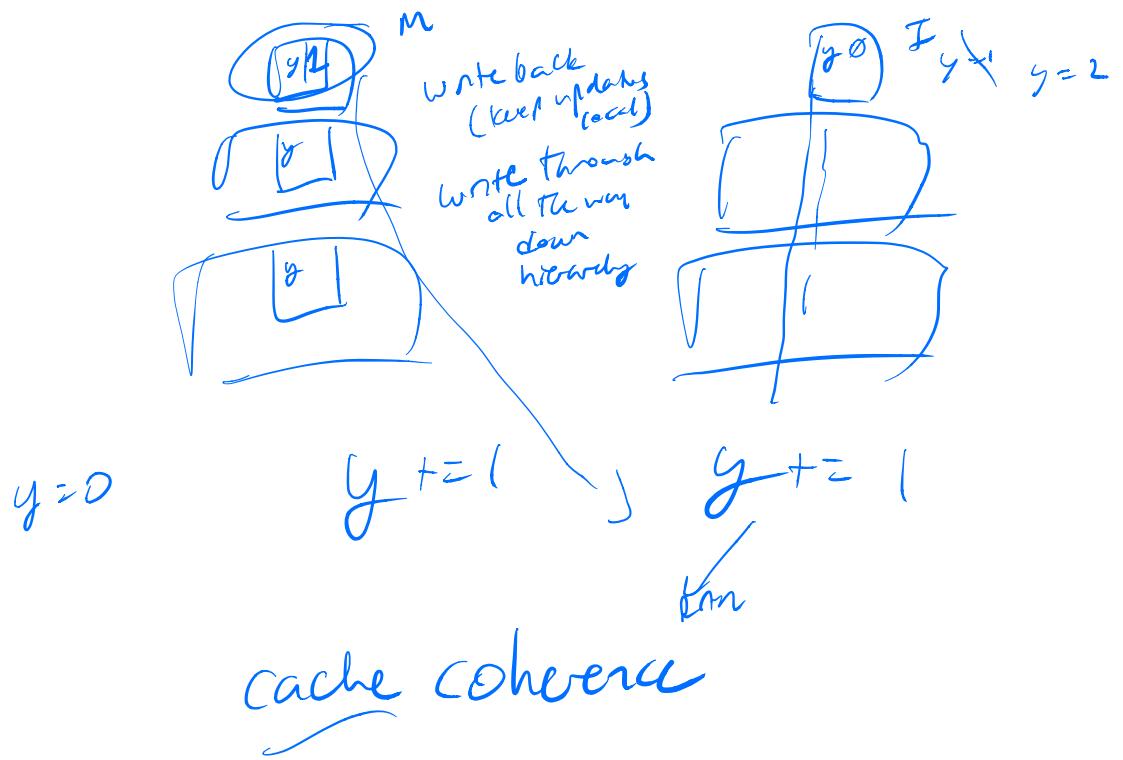
1000
10000

Caches — "replacement policy"

eviction "LFU"

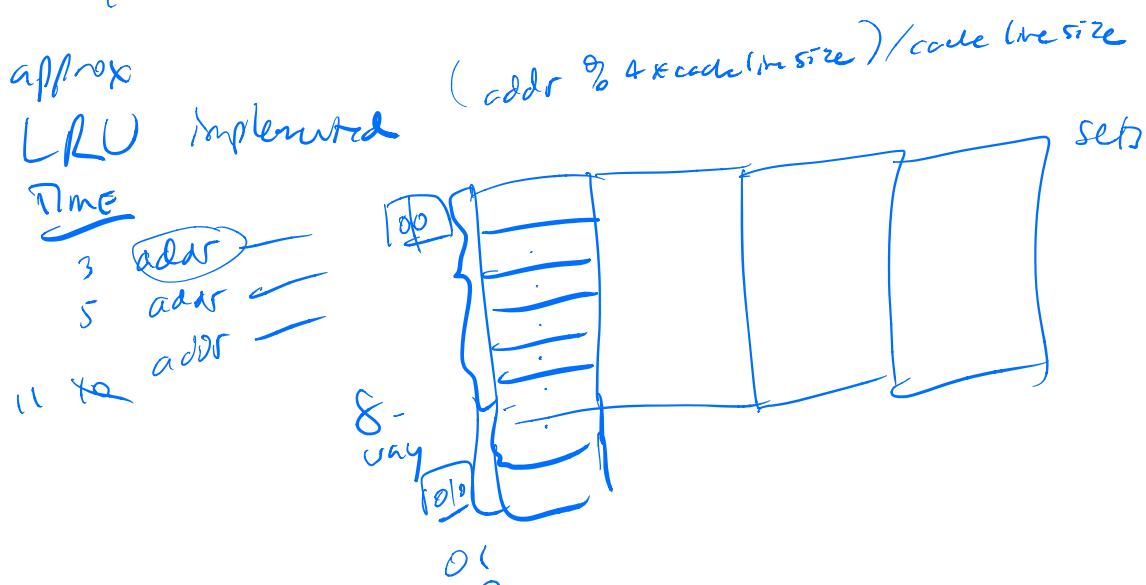
latency
wait time
caches wide latency





MESI protocol

- { Modified
 - { Exclusive
 - { Shared
 - { Invalid
- I changed it
just me
several CPUs have it (modification)
trash



Cache misses

Can be different
types -

Capacity - "hash collision"
Conflict -
Coherence - 3C's model

direct mapped - 1-way

2-way

4-way

8-way

16-way

fully associative - ∞ way

TLB

translation
lookaside
buffer

Cache of
virtual \rightarrow physical addresses

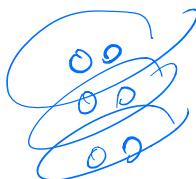
#include < stdio.h >

```
int main() {
    int a;
    printf("a = %p\n", a); 0x10000000
    return 0;
}
```

address

process isolation

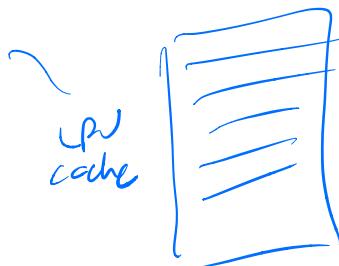
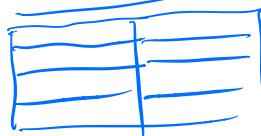
10
11

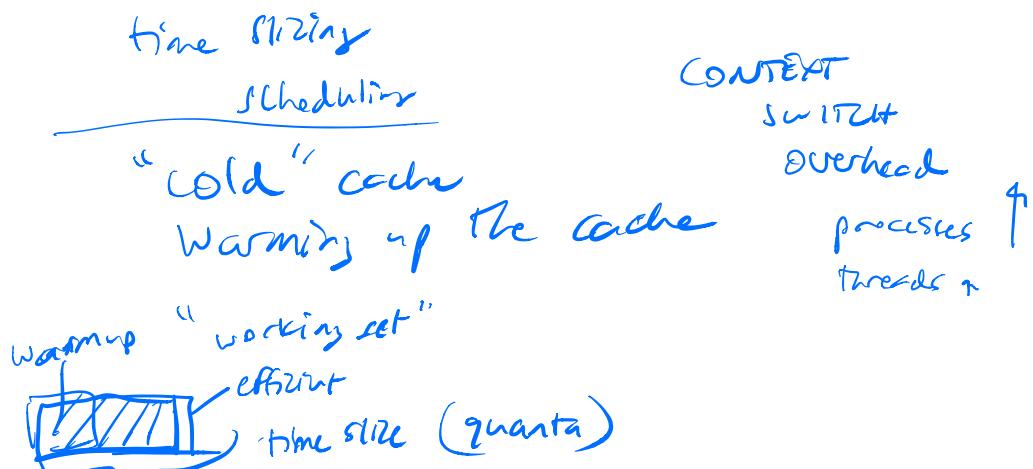
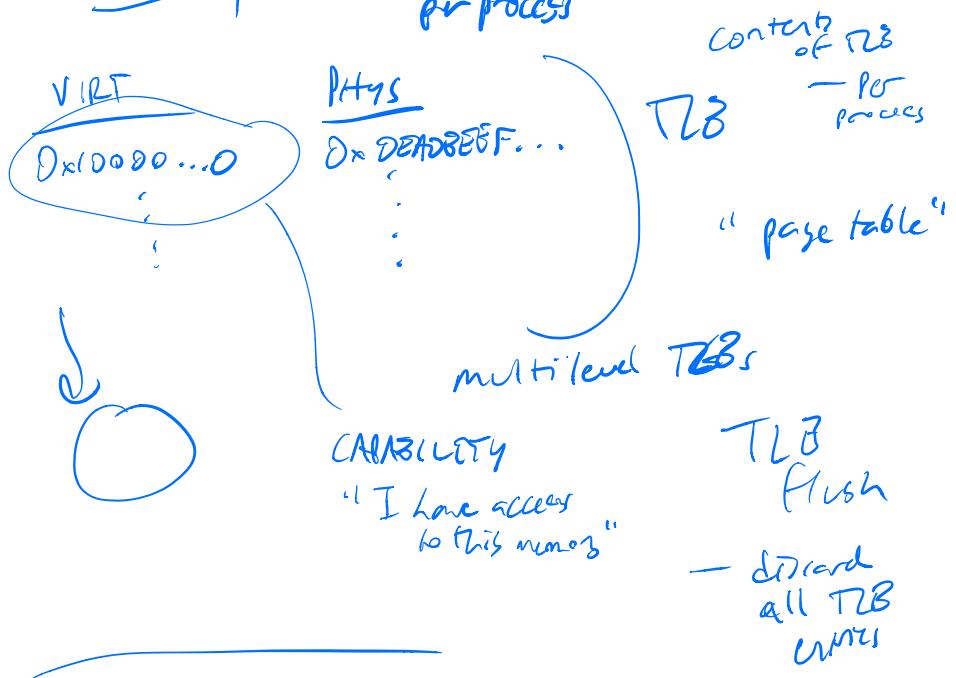
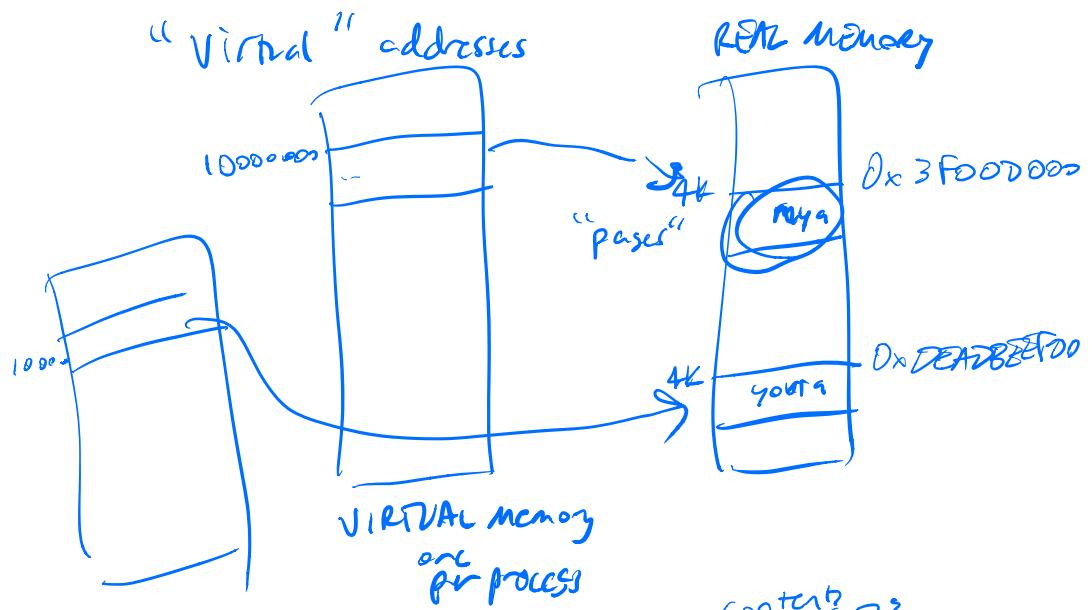


01
10



saturated
counter







anomalous cleanup
across the quantum



REGS < QUANTUM

register allocation

"spill" - copy
contents
of register
to mem "eviction"!

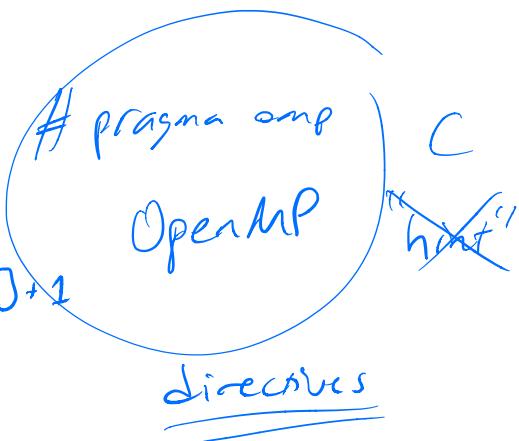
straight-line code
- predictable!

"graph coloring" NP complete!

Why is FORTRAN fast?

~~OMP PARALLEL~~

```
DO I = 1, 100
  DO J = 1, 100
    A[I][J] = A[I][J-1] + 1
  ENDDO
ENDDO
```



automatic
parallelization

"holy grail"

main thread
 S
 ↓
 for ($i=0; i < N; i++$) {
 foo();
 }
 ↗
 are there
 loop carried
 dependencies?

P
 }
 for ($i=0; i < N; i++$) {
 t[i] new thread (foo);
 }
 ↓
 foo₀, foo₁, ...
 [t[i].join()]

$\text{result}(S) \equiv \text{result}(P)$

parallel prog
semantically equivalent
to seq prog.

does `foo()` have side effects?

foo()
 {
 x++
 (_{global})
 }

read x into a register
 increment x
 write the result back

$x = 0$

A
 R₀ ← x
 R₀ ← R₀ + 1
 x ← R₀

$A ; B \Rightarrow z$

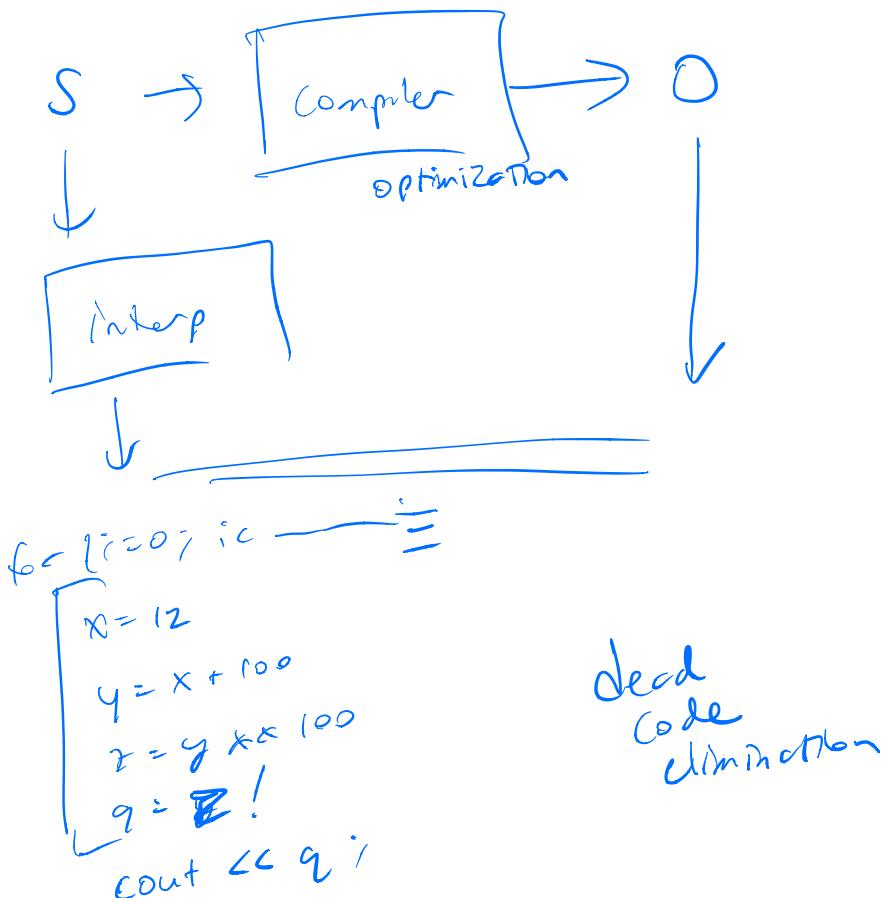
B
 R₀ ← x
 R₀ ← R₀ + 1
 x ← R₀

$$X \in \{1, 2\}$$

non-deterministic!

Race condition

Concurrency errors



```

foo() {
    ~x ~
}
  
```

```

foo (int i) {
    ~i ~
    return ack(i, i+1);
}
  
```

Ackermann

```
foo (int *int i a,) {  
    a[i] = ack(i, i+1);  
}
```

```
foo (int * a, int *int i b) {  
    a[i] = -;  
    b[i] = a[i] + 2;  
}
```

```
int x[200]; int *y = &x[0];  
v = foo(x, y, i)
```

aliasing

$p \{ \quad \}$ alias analysis
 $q \{ \quad \}$ pointer analysis

$p \cap q = \emptyset$ $\text{int } *p; \quad p \in \{ \}$

if ($__$) {
 $p = q;$ $p \in \{q\}$

} else {
 $p = r;$ $p \in \{r\}$

symbolic execution

abstract interpretation \rightarrow
 $\wp \in \{q, r\}$

$$\begin{aligned} x &= 1^2 \\ x &= x - 100 \\ x &= \dots \end{aligned}$$

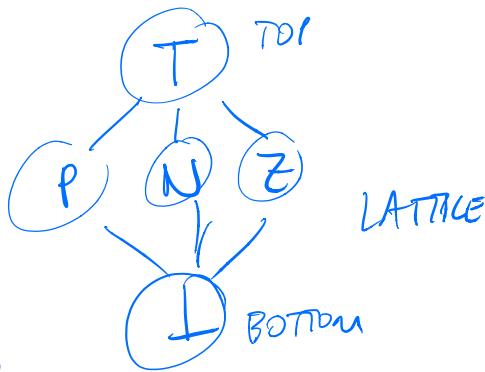
if ($x > 0$) {
 |
 }

$x \in O, A, P$

$x = \dots \quad \{O, N, P\}$
 $x = x - 1 \quad \{ \}$
 does this run over?
 $O \Rightarrow O$
 $N \Rightarrow P$
 $P \Rightarrow N$

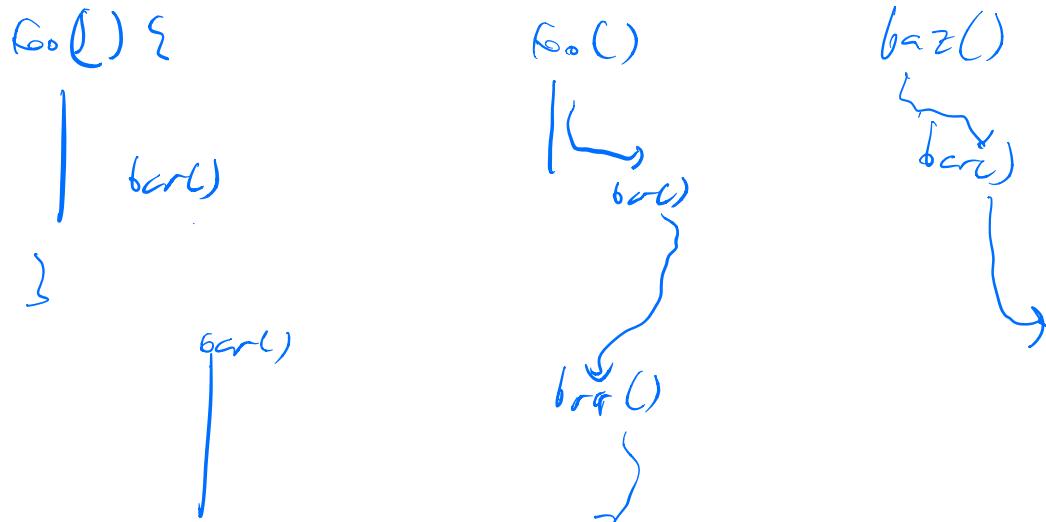
CONSERVATIVE APPROX
 PROG BEHAVIOR

$\{x : T\}$
 if (\leftarrow) {
 $x = -1$
 then
 $x = 1$
 }
 $\{x : \perp\}$



$$\begin{array}{c} \text{if } x : P \\ \text{then } x : P \\ \hline x : P \end{array}$$

intraprocedural vs. inter procedural



flow sensitivity



Context sensitivity

Steensgaard's algorithm

PCR sensitivity

false positives

precision - recall

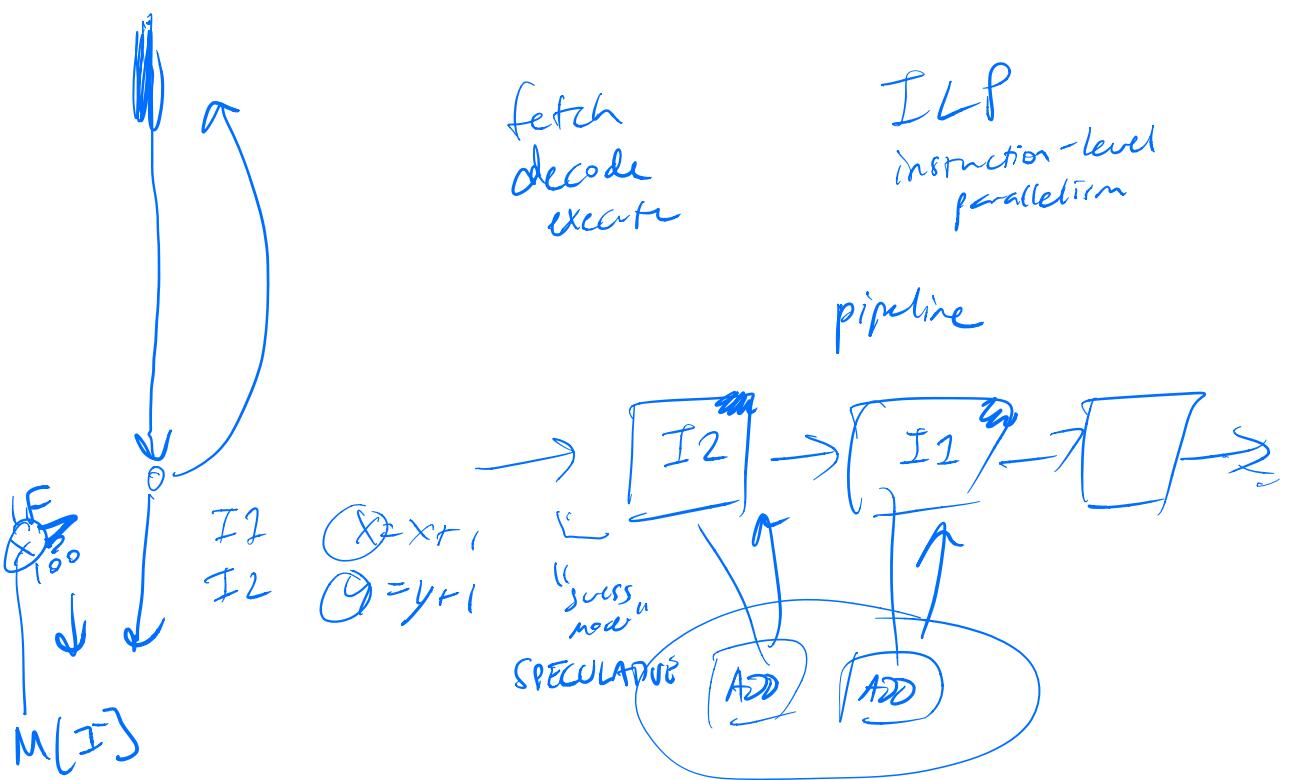
false negatives

tradeoff

$$V \sim \{\emptyset\}$$

-W

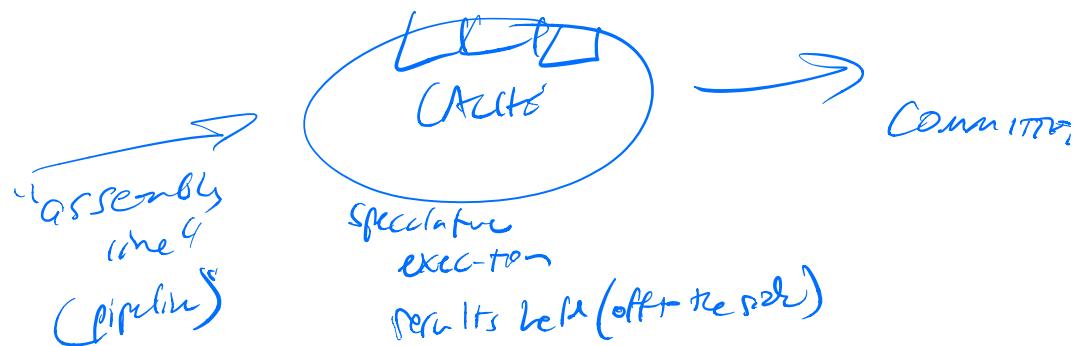
if ($x == \underline{\hspace{2cm}}$)
 $y = x / @ V$



Pipeline stall

"bubble"

"BACK EDGE DATA"



Speculative execution
hide memory latency

SPECTRE / Meltdown

timing channel

Covert channel

side channels

last branch taken

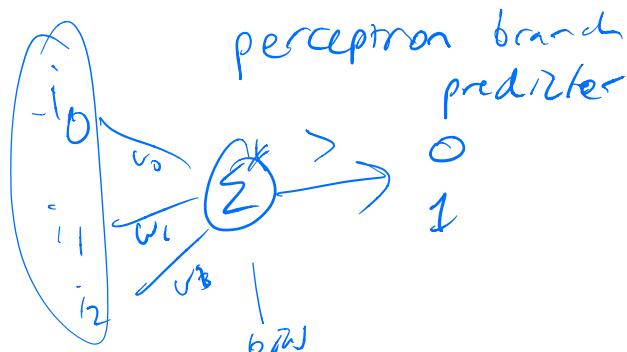
85-90% effective

PHAT

Modern

branch predictors

99.9% effective



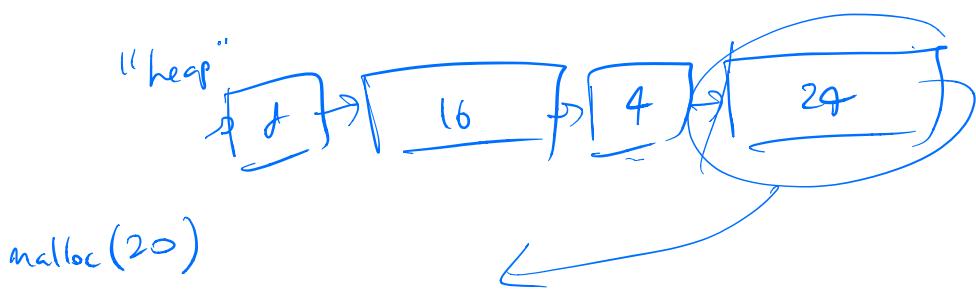
dynam. memory allocation

globally

stack variables

dynamic memory

"the
heap"



Object 24

First-fit → SPACE INEFF

best-fit

SPACE
EFFICIENT

slow ...

$$n \cdot O\left(\log \frac{M}{m}\right) * n$$

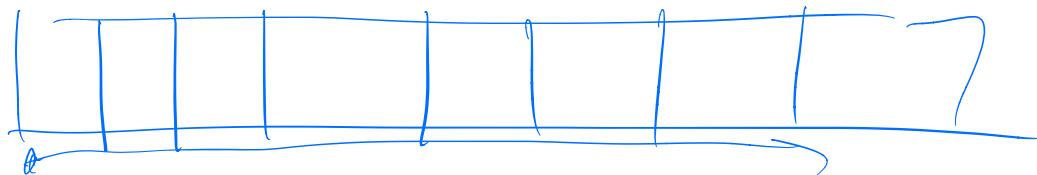
M - biggest
m - smallest

$$O(M \times m)$$

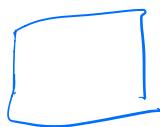
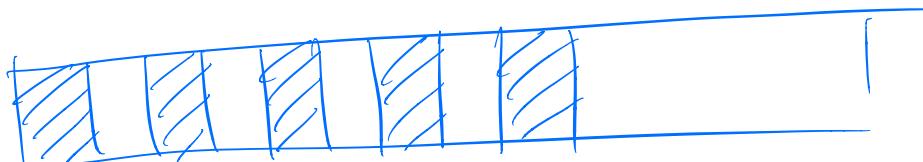
best fit

worst-case
"fragmentation"

$$= \frac{\text{mem consumed}}{\text{mem requested}}$$



D → D → D



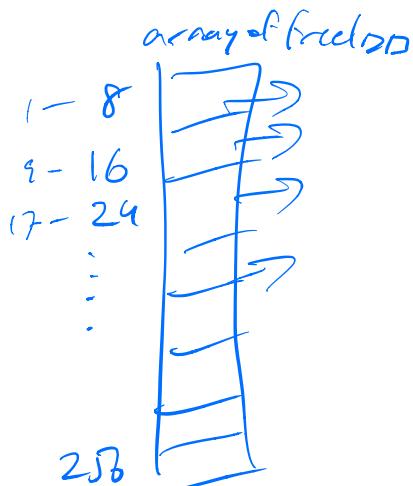
Compaction
NOT IN C/C++

(L1E')

Python, Ruby

"custom allocators"
"ad hoc allocators"

malloc - "too slow"
Reconsidering Custom Mem. Allocation



Py-object
- malloc

if ($sz \leq 28$) {

use small obj.

else

use malloc

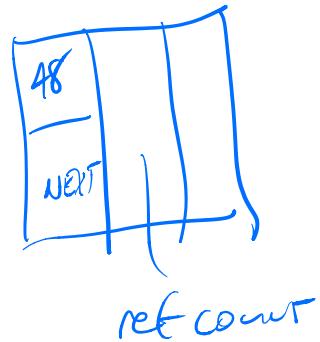
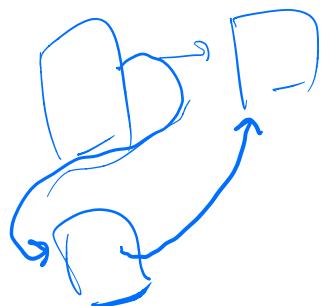
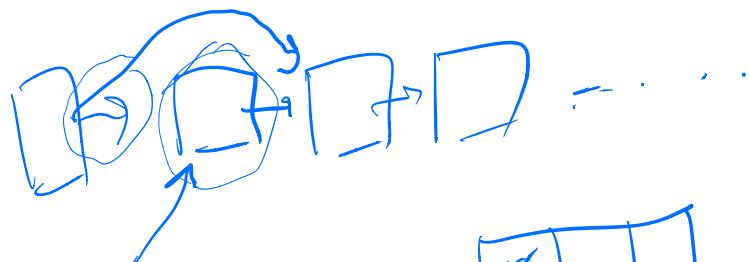
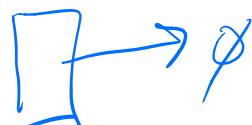
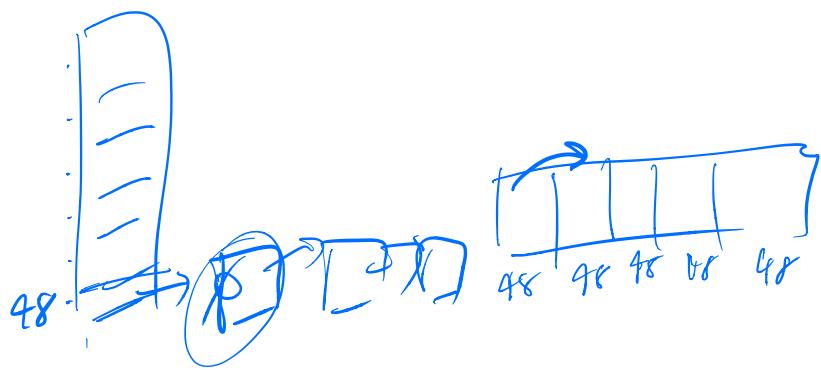
int ≈ 32 bits C Python
long ≈ 64 bits 48 24/28

"a"
1/16 50

{ "a": }
~~2~~

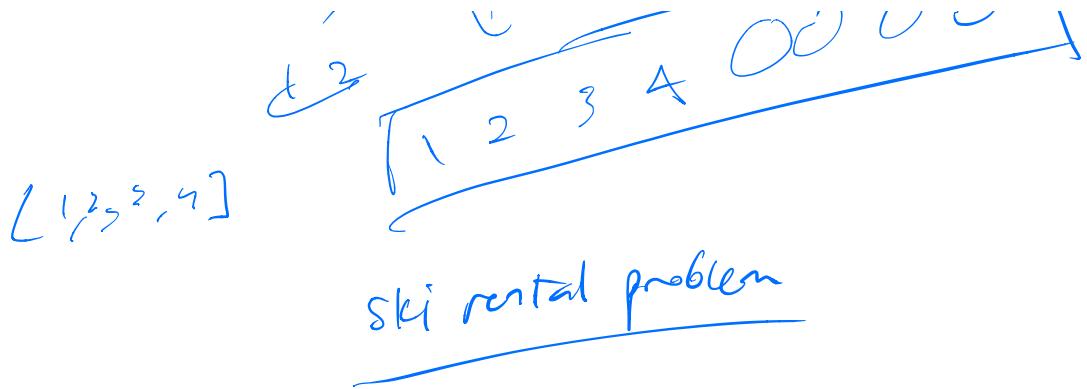
240

60 20 75 100 140



$$x = [1, 2, 3]$$

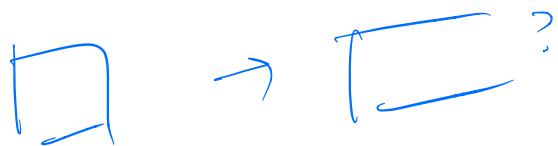
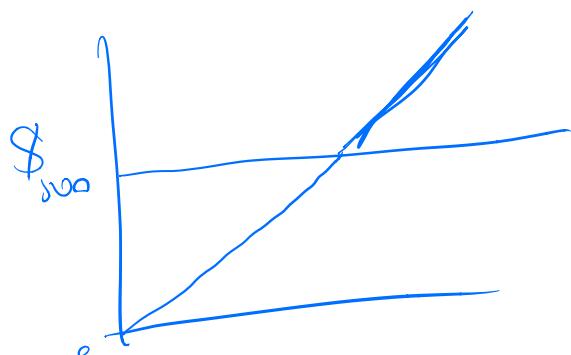




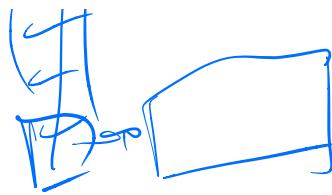
\$50 rent

\$500 buy

rent until
you reach
cost of buying



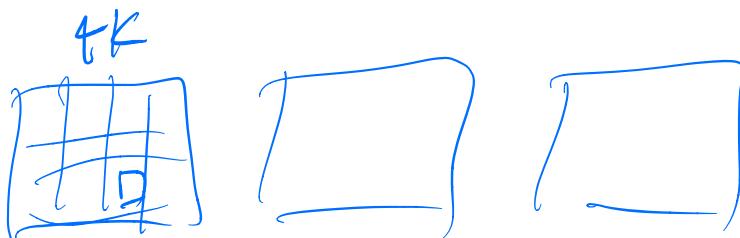
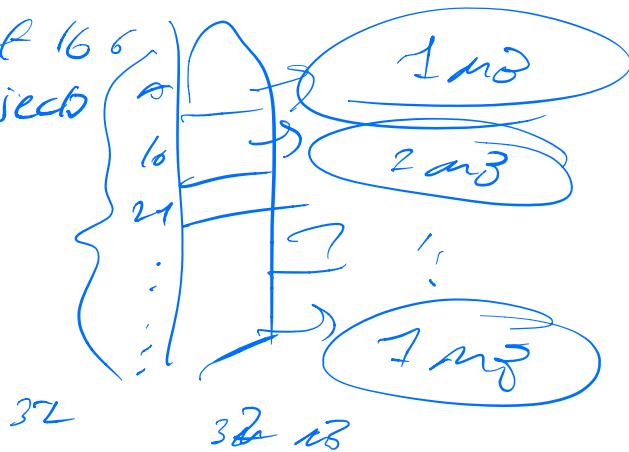
Segregated
size
classes



allocate 1MB of 8 64B objects
free them all

allocate 1MB of 16 64B objects
free them all

actual
MMT 1MB
32x



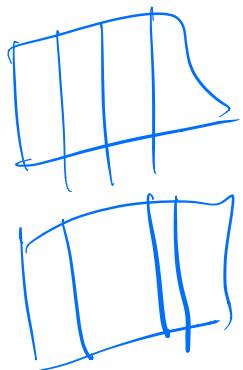
OS passes

Space-time tradeoff

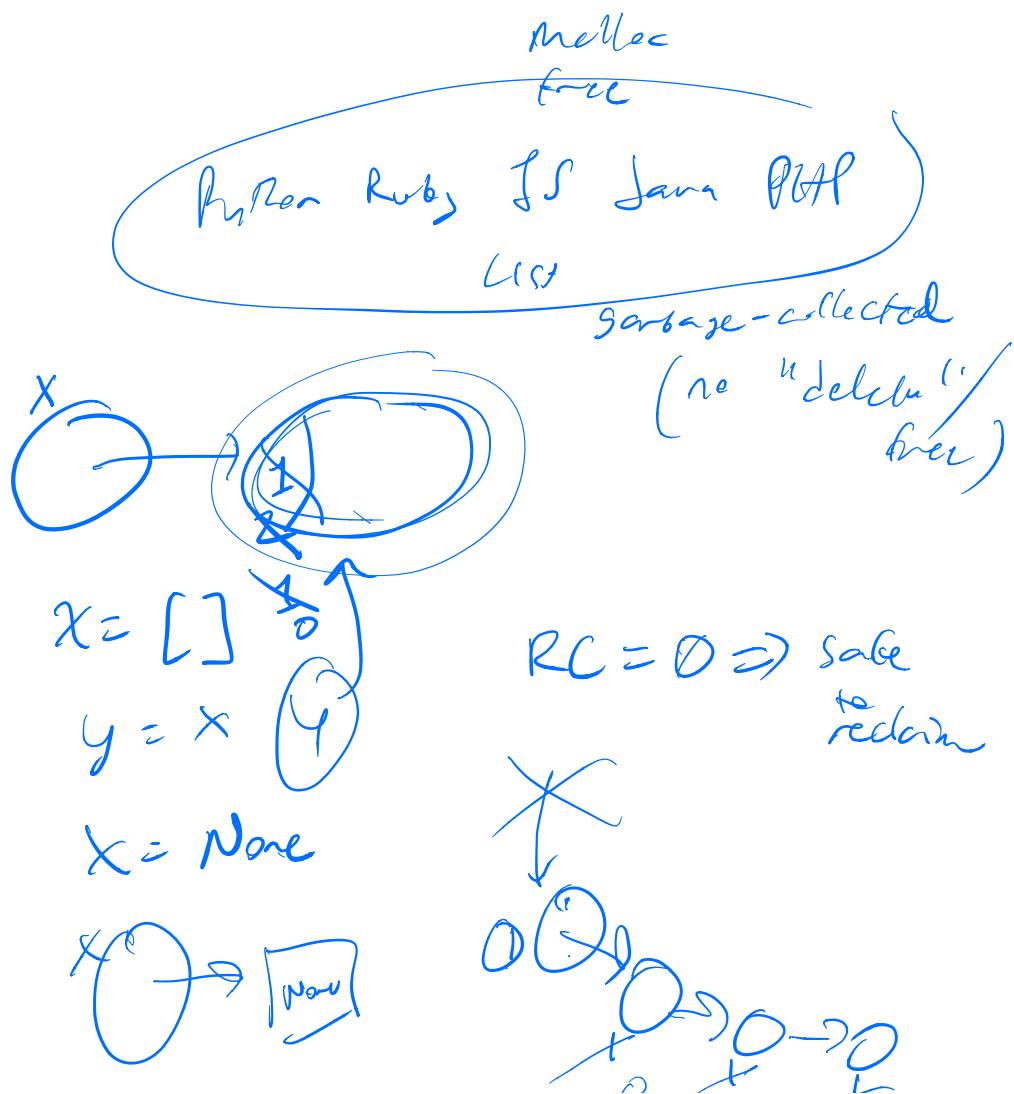
$$f(a) =$$

$$\begin{aligned} f(0) &= a \\ f(1) &= b \\ f(2) &= c \end{aligned}$$

0	a
1	b
2	c



Reference counting



$y = \text{None}$

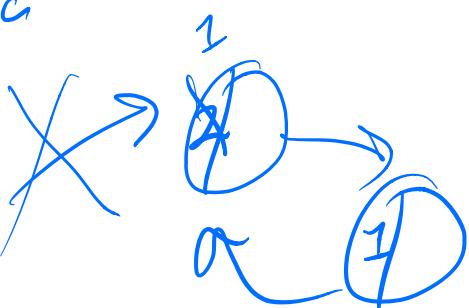
$\leftarrow \circ \rightarrow$



doubly-linked list

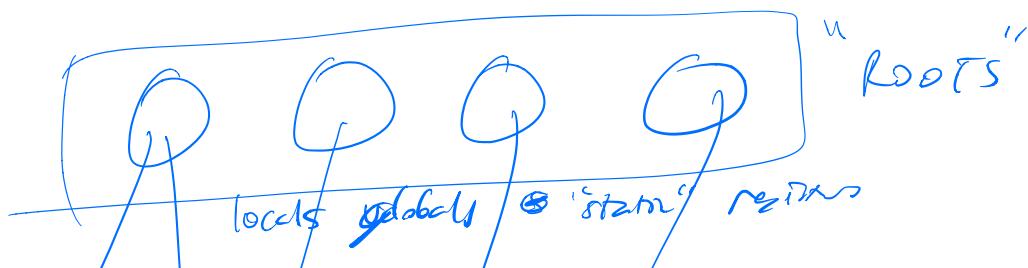
$a = \{\text{'key'}: 0\}$ CYCLE

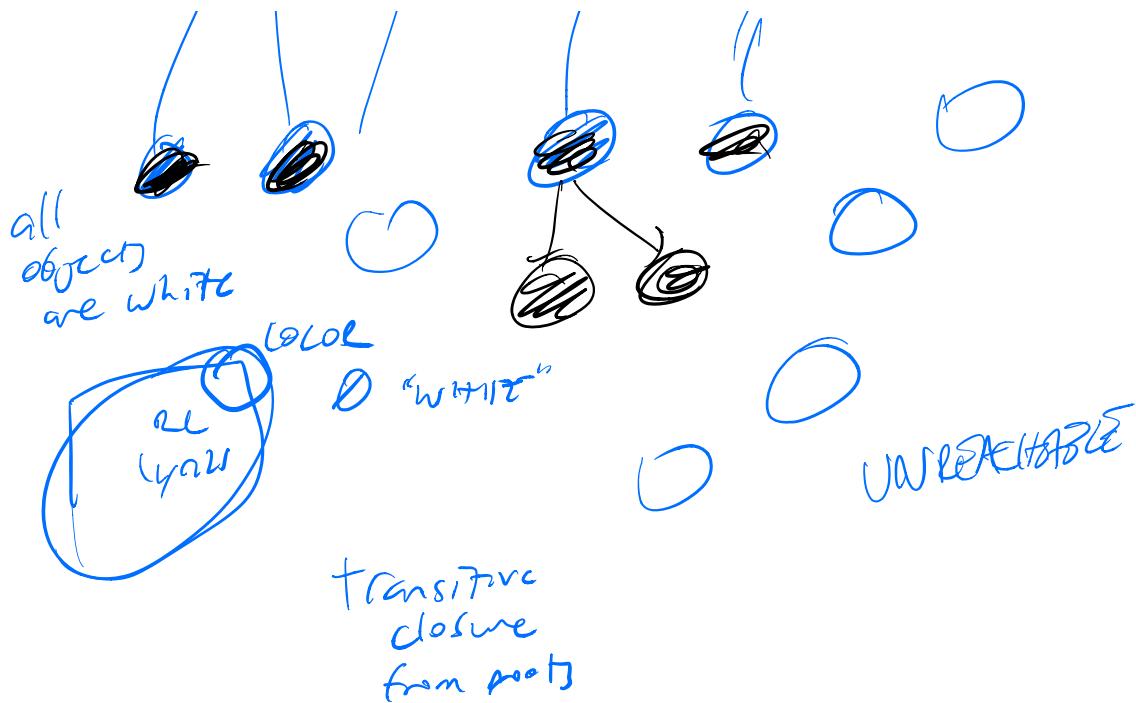
$a[\text{'key'}] = a$



HC
(1956) Incomplete
GC algorithm

Mark-sweep (1978)
"generational" GC (1983?)





David Unser

smalltalk self

Jö Holzle
 David Hamburger

Ole
 Azum polymorphic
 oriented

dynamic

Luc
 Bak



Google

→ V8

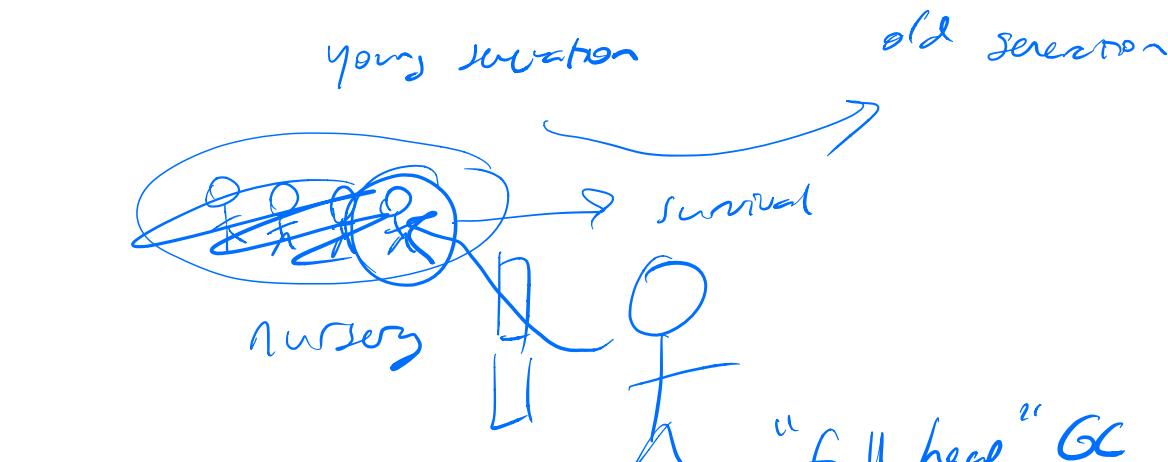
GC algorithms

JIT
 compilation

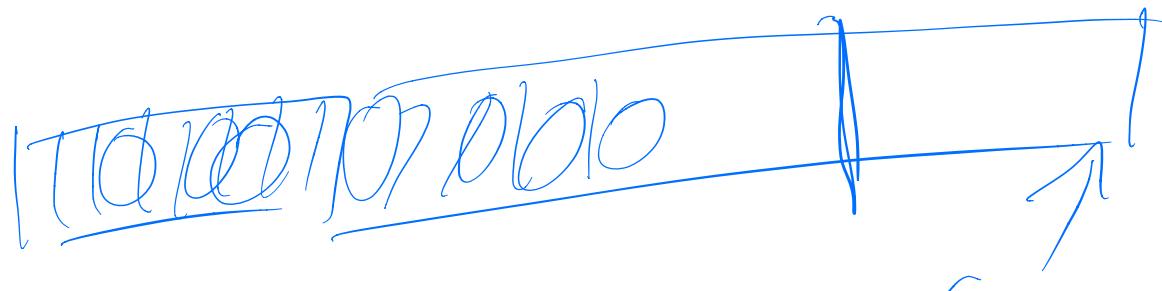
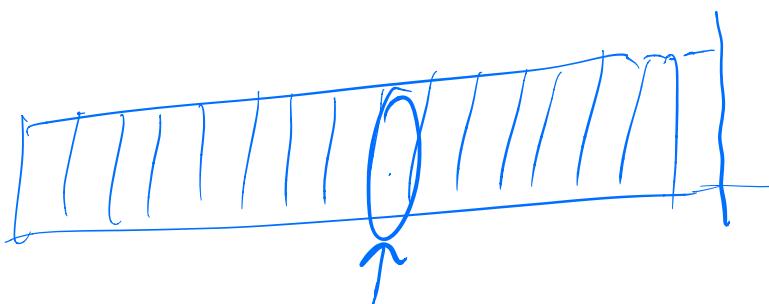
dynamic
 languages

polymorphic
 inline
 cache
 "Pic"

most objects die young



gen GC →
takes pressure
off the full
heap GC

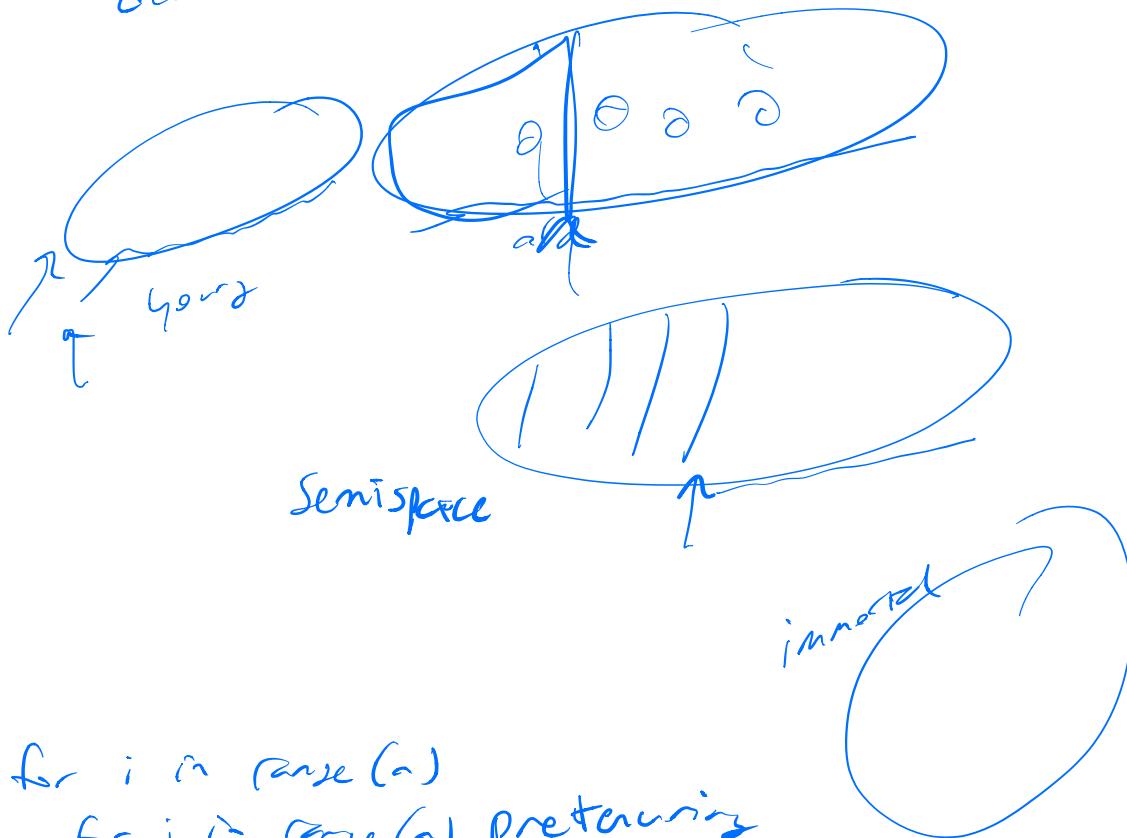


GC

$\sim 2 \rightarrow 5x$

"live
objects"
"dead"

Quantizing the cost of
GC vs. Explicit deallocation



for i in range(a)

 for j in range(n) pretenuring

 for k in range(n)

$$C[i][j] += \delta A[i][k] + \delta[k][j]$$

• pyc

PyPy