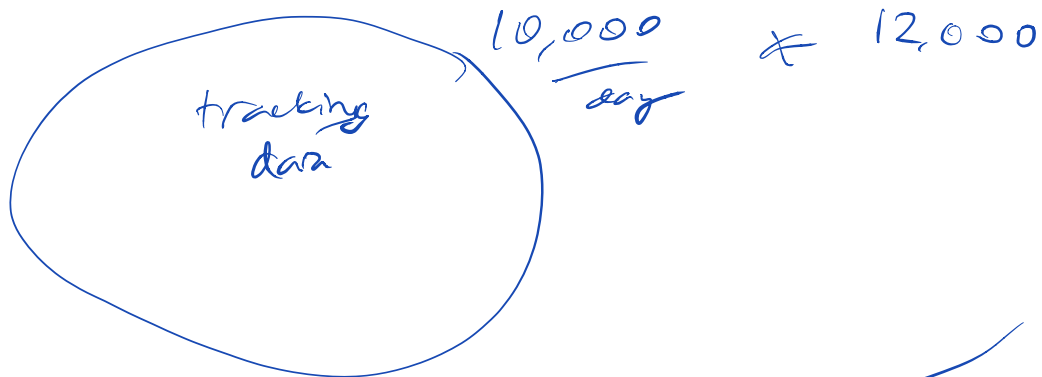# Performance

time is money

lots of data $\Rightarrow$ process $\Rightarrow$ info

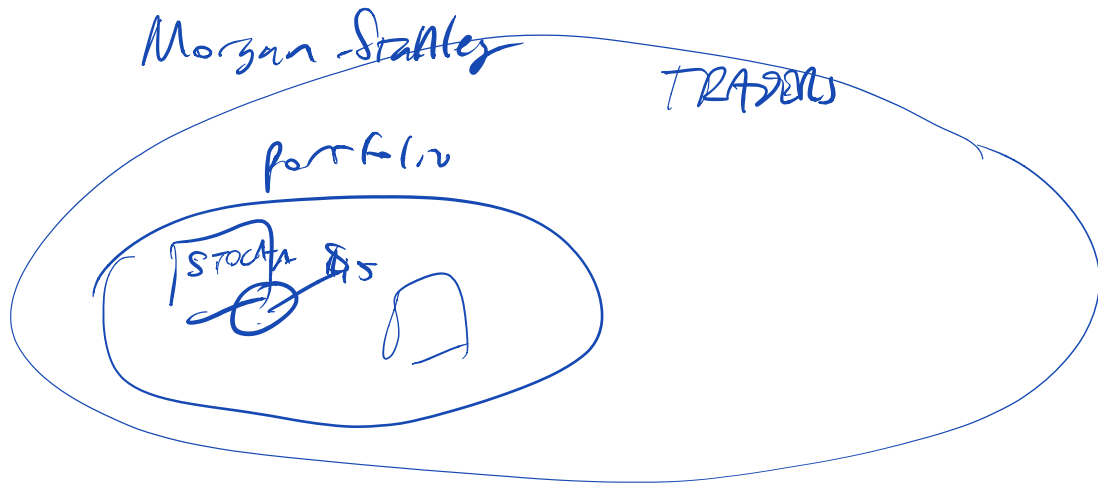insights

limited time budget — 1 day, 1 week

MORE data
$\Rightarrow$ better info

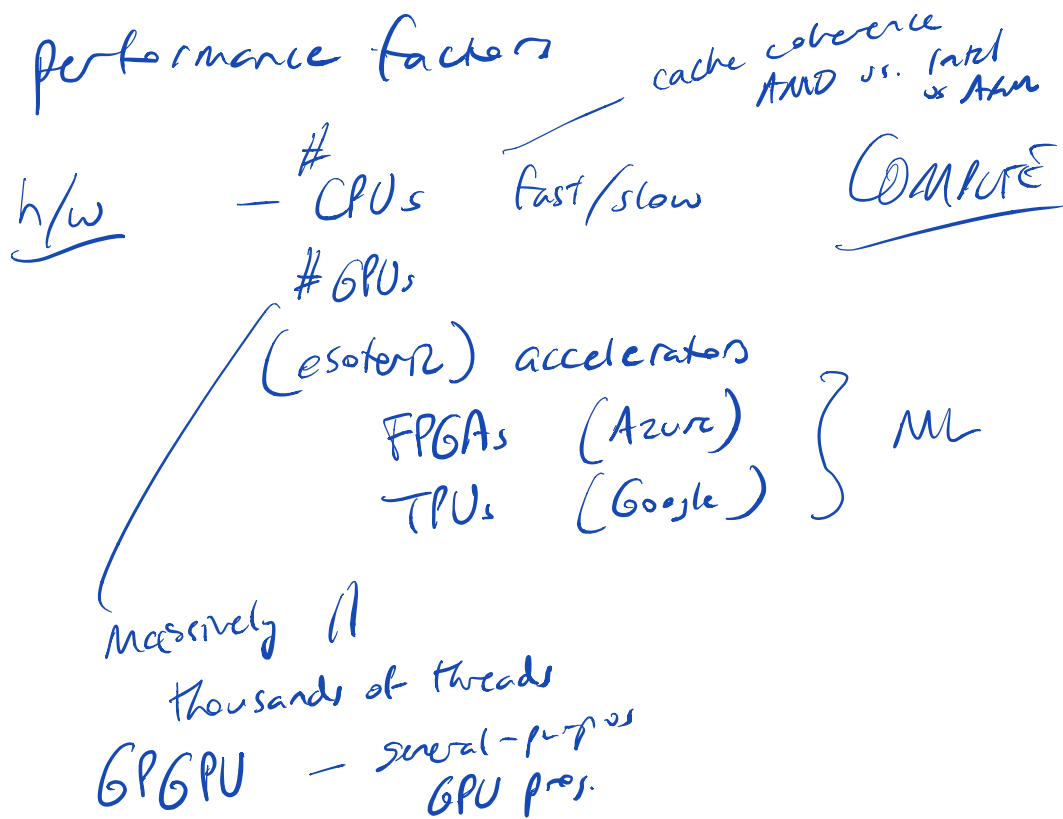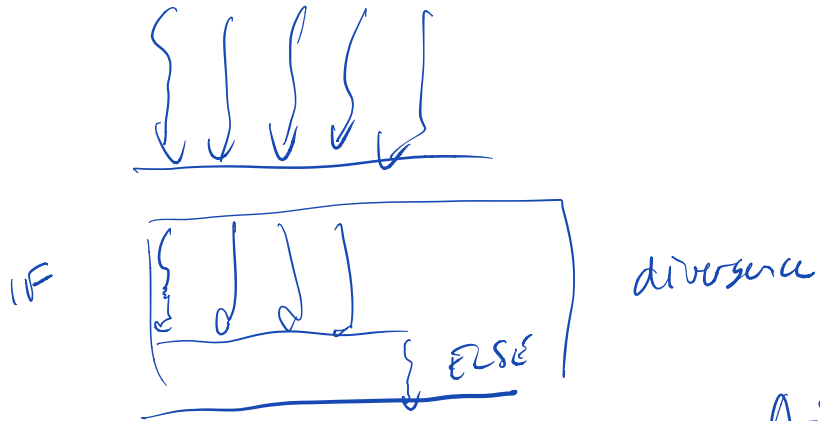Wal-Mart — every single sale $\times$ # of WalMarts

$$\frac{10,000}{day} \quad \times \quad 12,000$$

tracking data

ANALYTICS

$\leq 8$ hours

OPTIMIZATION
$\Rightarrow$ $$

Order
planning (logistics)

Morgan-Stanley

TRADERS

portfolio

STOCK $s

risk models (overfit)

time — many

---

Performance factors

cache coherence
AMD vs. Intel
vs. ARM

h/w — #CPUs    fast/slow    COMPUTE

#GPUs

(esoteric) accelerators

FPGAs (Azure)  } ML
TPUs (Google)

Massively ||
thousands of threads

GPGPU — general-purpose
GPU prog.

IF { ELSE

divergence

ASICs
application
specific
IC

CUDA   Nvidia
OpenCL

CPU ⟶ GPU    COMPUTE

MEMORY

—— NETWORK

topology
(hops)

router

router

— latency
— throughput    high bandwidth
InfiniBand

—— DISK

$\Big($ — HDD      latency   seeks
— SSD  —   capacity  $\Big)$ bandwidth

RAID   increase b/w

— MEMORY
     speeds      bus ~~data~~ freq.
                        $\rightleftharpoons$ b/w
     Capacity  ~ MORE = BETTER
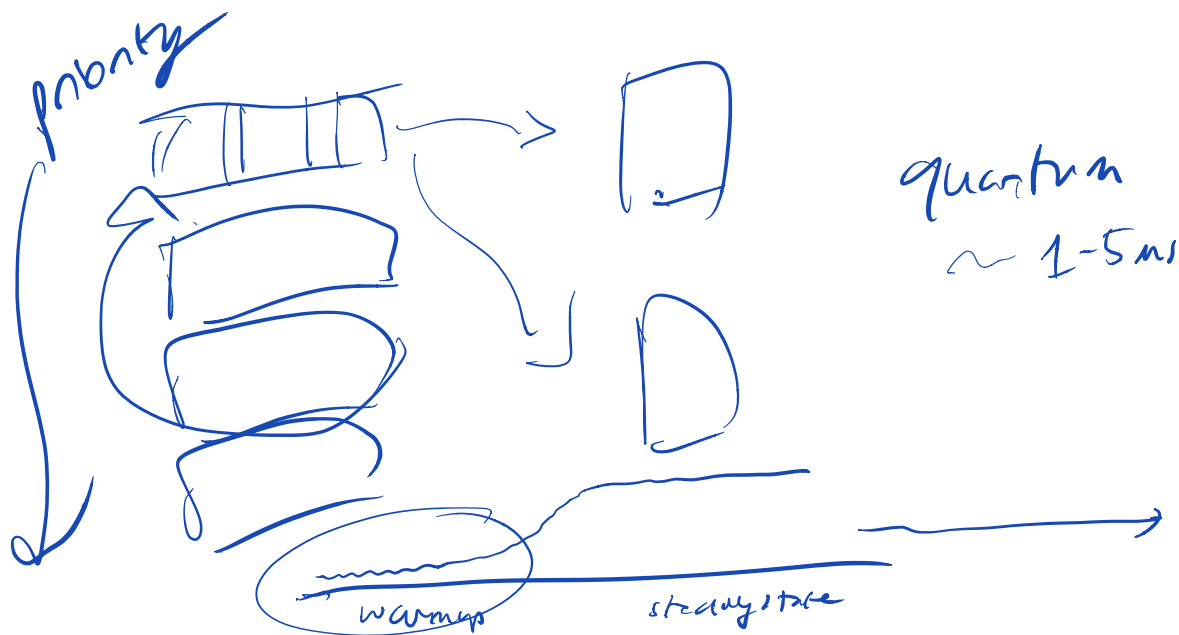     Cache (SRAM)   bigger

lowend  small codes      big codes
        2.4 GHz           ~ 4 GHz

SW $\rightarrow$ perf.  — $O(n) \rightarrow O(1)$
     algorithms
PL — GC  ~ "managed"
     JIT compilation — Java / Javascript
                        warmup period
OS —
     File systems    throughput
     Networking stack
     Scalability of OS
     optimized OS
     Scheduler

priority



quantum
~ 1-5 ms

warmup        steady state

design choices

optimize for Thruput
OR latency

(disk I/O) ⬜ ⬜ ⬜ ⬜ ⬜ interactive
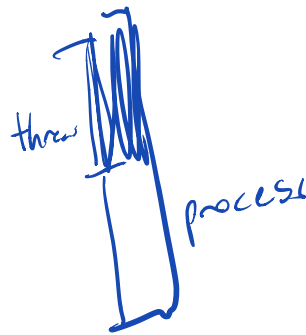
⬜⬜⬜⬜ batching

MT      node.js      1 thread concurrent I/O

MP      [ node.js
        [ node.js
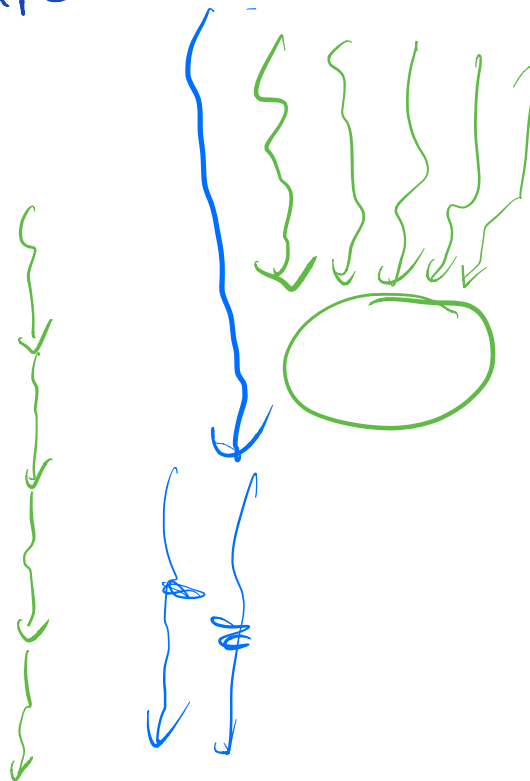          ⋮

— fault tolerance

MP good └ it no "shared state"
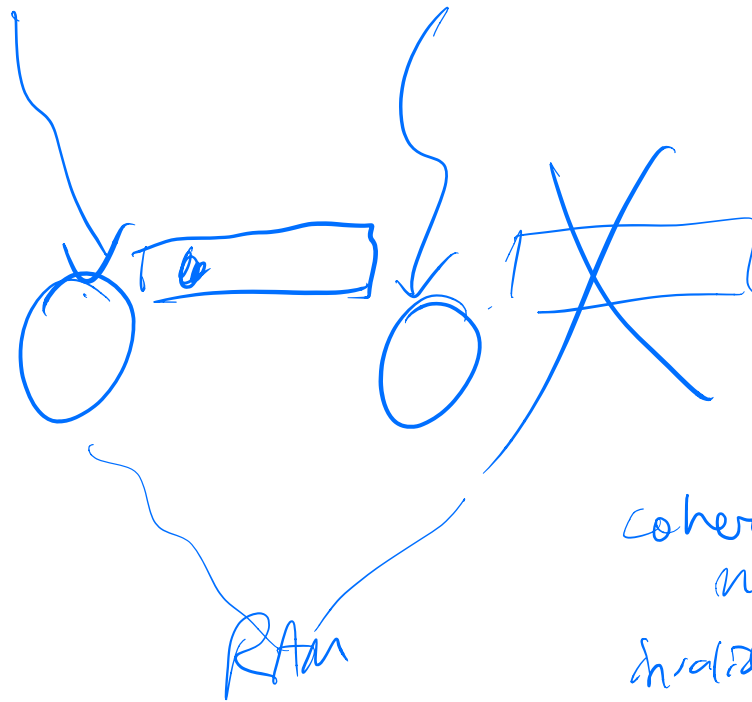MT good — lots of shared state
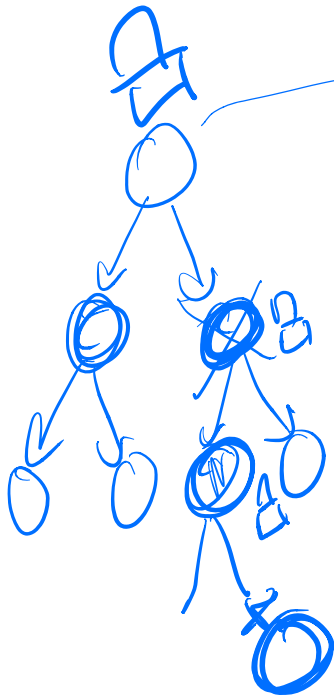
MT [ Apache — a patchy sever
   [ nginx

throw | | | |
         process

MT
painful
(correctness)

CACHE

coherence
miss

invalidation

RAm

hand-over-hand
locking

CONTENTION

abstractions

APIs

fut. ultimate leaky abstraction
_____

⟶ workload - dependent