

CS 3214 lecture #26

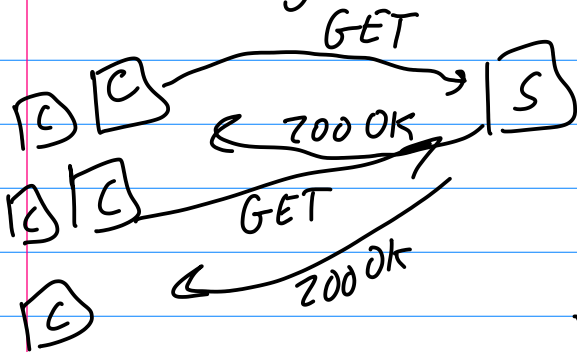
"web servers: threads vs. events"

HTTP

+ cloud computing

SPOT survey - 14%

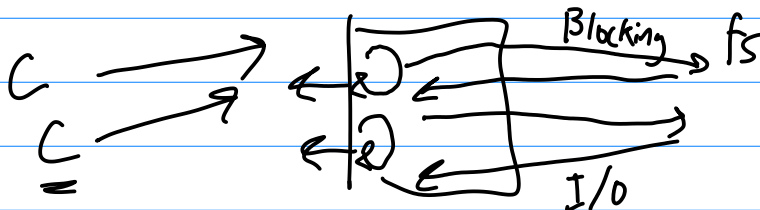
Concurrency -



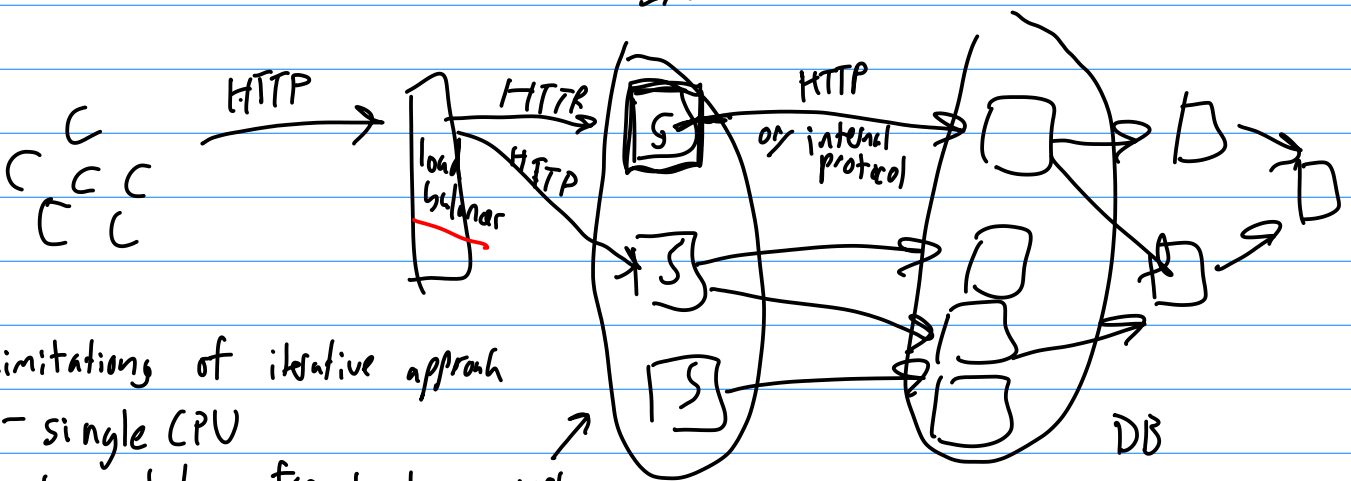
accepts client connection
→ processing
sends response

"iterative server"

- one req at a time (sequentially)
Serially



low utilization



limitations of iterative approach

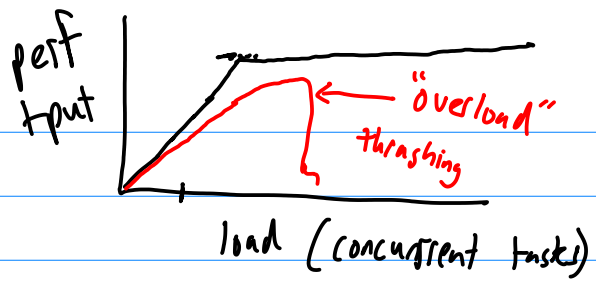
- single CPU
- high latency for client
- low utilization

web servers

back-end services

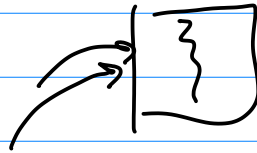
"Microservices"
"service mesh"

handle multiple requests concurrently
1978
- threads / processes (apache)
- event-based (nginx)

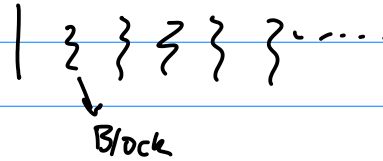


Thread-based approach

- underlying system to maintain context + state



- PROS {
- Familiar linear control flow (blocking is fine)
 - UNIX makes it easy
 - isolation (process-based)



pool



fixed amount of thread
how to pick n?

- CONS {
- lots of clients = LOTS of threads
 - state transitions, ctxt switches, mode switch
 - concurrent programming is hard
 - tuning (# threads)

Event-based model

- reorganize program into event handlers
- application manage client states
- application decides what runs when
- Nothing blocks

→ how to avoid blocking?

→ how to manage state

How to write event-based program

- identify where program would block
- where do events come from
- maintain client state

parsing HTTP header.

read()

[GET /api/foo HTTP/1.1
: . . .



How to identify events?

- new request arrived
- new data rec'd
- connection closed
- ...

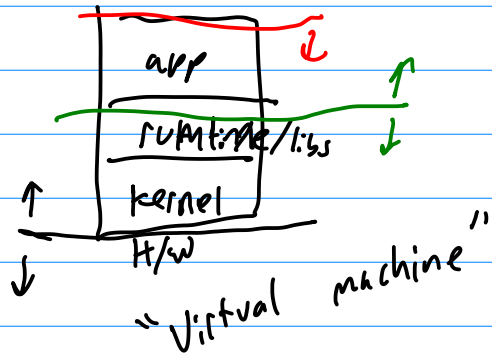
fd's ← which fd's have data pending?
or: will a read block on this fd?
will a write block on this fd?

- Non-blocking mode - EWOULDBLOCK
- select ← tell me when fd's are ready
poll
- epoll ← fastest, most features (on Linux)

libraries: libevent, libuv, put abstraction layer
↑
node.js

High-level language support: async/await in Javascript
↳ make it look a bit more thread-like

- + speed,
- + no need to tune # threads
- not easy to program



Cloud computing ←

X as service

- infra - User s/w kernel, runtime, apps / cloud manages power, cooling, etc. ex: EC2
- platform - user s/w above a runtime ex: Heroku
- software - entire application in cloud

IaaS PaaS SaaS

deployment -

- cloud
- hybrid cloud
- on-prem

defining characteristics

- elasticity: can scale up & down instances on cloud
only pay for what you use (sometimes autoscale)
granularity ↓
- multi-tenancy
- isolation

[cloud, VMs, containers]