Application Architecture

Lecture 5 Srinivas Narayana http://www.cs.rutgers.edu/~sn624/553-S25



Components of an Internet Service

Routers

Storage

App compute and communication patterns

Modularized applications

Endpoints

Interconnect: Routers

Data Center

Servers

Offline and Online components



Review: Web server design

• Process other requests while waiting for one to finish



process socket





listen()

accept()

recv()/send()/..

bind (IPaddr_B, port_B)

Review: Parallelism

- Process requests in parallel with other requests
- One design: multiprocessing/multithreading (MP/MT)





recv()/send()/..

Can block if any of the requests block (asynchronous IO support can be incomplete & complex)

Avoid overheads of multiple processes and threads

Using parallelism + concurrency

Asymmetric Multi-Process Event Driven (AMPED)



Flash: An efficient and portable Web server



Microservice Architectural Pattern





In short, the microservice architectural style is an approach to developing a single application as a **suite of small services**, each **running in its own process** and communicating with lightweight mechanisms, often an HTTP resource API. These services are **built** around business capabilities and independently deployable by fully automated deployment machinery. There is a **bare minimum** of centralized management of these services, which may be written in different programming languages and use different data storage technologies.

-- James Lewis and Martin Fowler (2014)







Cost of communication: Performance

Profiling a warehouse-scale computer (Google). ISCA'15.

Cost of comm: Hotspot spreading

A. NGINX Saturation

Deathstarbench. ASPLOS'19.

Cost of comm: high level failure handling

Are microservices always ideal?

- Just an architectural style. Look at solving problems first
- How to evolve the splitting of components?
 - Refactoring microservice interfaces later isn't easy
 - Interface changes need buy-in from multiple dev teams
 - Components should compose cleanly in the first place
- How to design apps?
 - Monolith first, or microservices from the beginning?
- Testing, Observability, Deploy automation
- How significant are dev coordination overheads?
- Complexity

Partition-Aggregate

Processing interactive search queries

Web search: some numbers (circa 2003)

- 10s of terabytes of web corpus data
 - Read 100s of megabytes per query
- 10s of billions of CPU instructions per query
- Data accessed depends on the query; hard to predict
- All results to be returned to users within (say) 300 milliseconds
- Cannot process on a single machine within acceptable time

Example: Google search architecture

Web search for a planet, MICRO'03

Quick Review: Compute & Memory Org

Measurements from one (index) server

- Not too fast single-threaded
 - Data dependencies
 - Branches often mispredicted
- Small instruction memory footprint
- Data locality within a block, but not across blocks
- Numbers not much better on a newer architecture
- Can't drive high single-threaded performance Use parallelism

Characteristic	Value
Cycles per instruction	1.1
Ratios (percentage)	
Branch mispredict	5.0
Level 1 instruction miss*	0.4
Level 1 data miss*	0.7
Level 2 miss*	0.3
Instruction TLB miss*	0.04
Data TLB miss*	0.7
* Cache and TLB ratios are pe	er
instructions ratirad	

monucions icincu.

Web search for a planet, MICRO'03.