

Engineering Robust Server Software

Server Software

Servers Software

- Servers accept requests from clients
 - Exchange information (take requests, give responses)
 - Generally do much of the "computing"
- We'll start with two example categories
 - Unix Daemons (sshd, httpd, ...)
 - Server side code in websites (Django)
- So what is so special about server software?
 - Why is it different enough to be in the course title?

Most Code You Have Written

- Run on input, get output
 - Then done
- Error?
 - Print message and exit
- Run by you
 - Trusts user
 - On one computer...
- Deals with one input at a time
 - Serial code
 - Don't care about performance

Servers: Different

```
while (true) {  
    .....  
}
```

- Run "forever"
 - Implications of this?

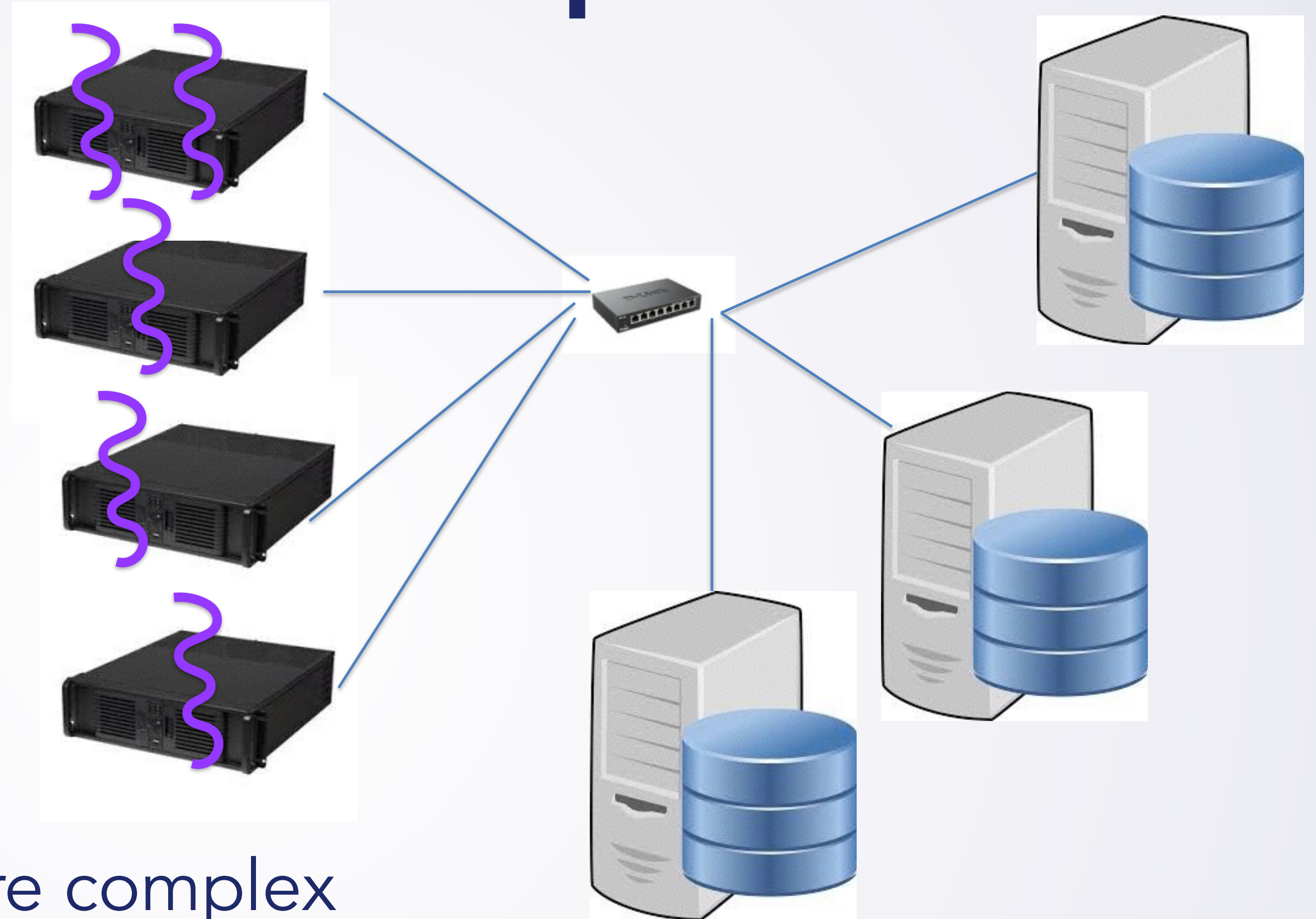
Run Forever

- Resource (memory, file descriptors,...) Leaks: Unacceptable
 - Restart Chrome every week b/c memory leak? Annoying
 - Restart Google every 5 minutes b/c memory leak? No way..
- Then again...
 - DukeHub has a memory leak
 - Solution: restart every so many requests.
- But you all are pros at writing leak-free code

Run Forever

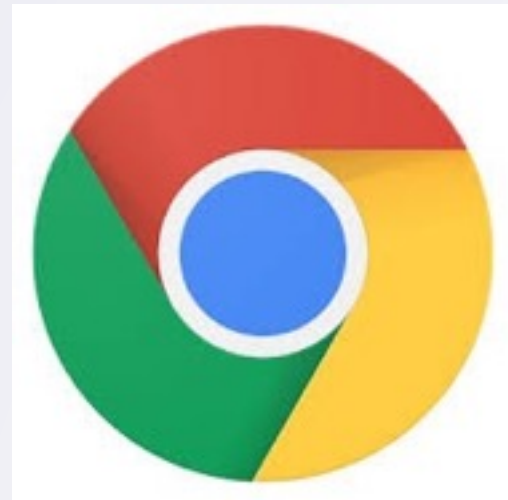
- How to handle errors?
 - abort? No way.
- Report and keep going
 - Need to keep handling other requests
- Log:
 - Nobody is watching terminal.
 - Want admins to know? Need log files (/var/log/...)
- Inform user
 - Send (informative?) error response

Maybe more complex?



- Many server systems: more complex
 - Introduce more complexities in terms of running "forever"

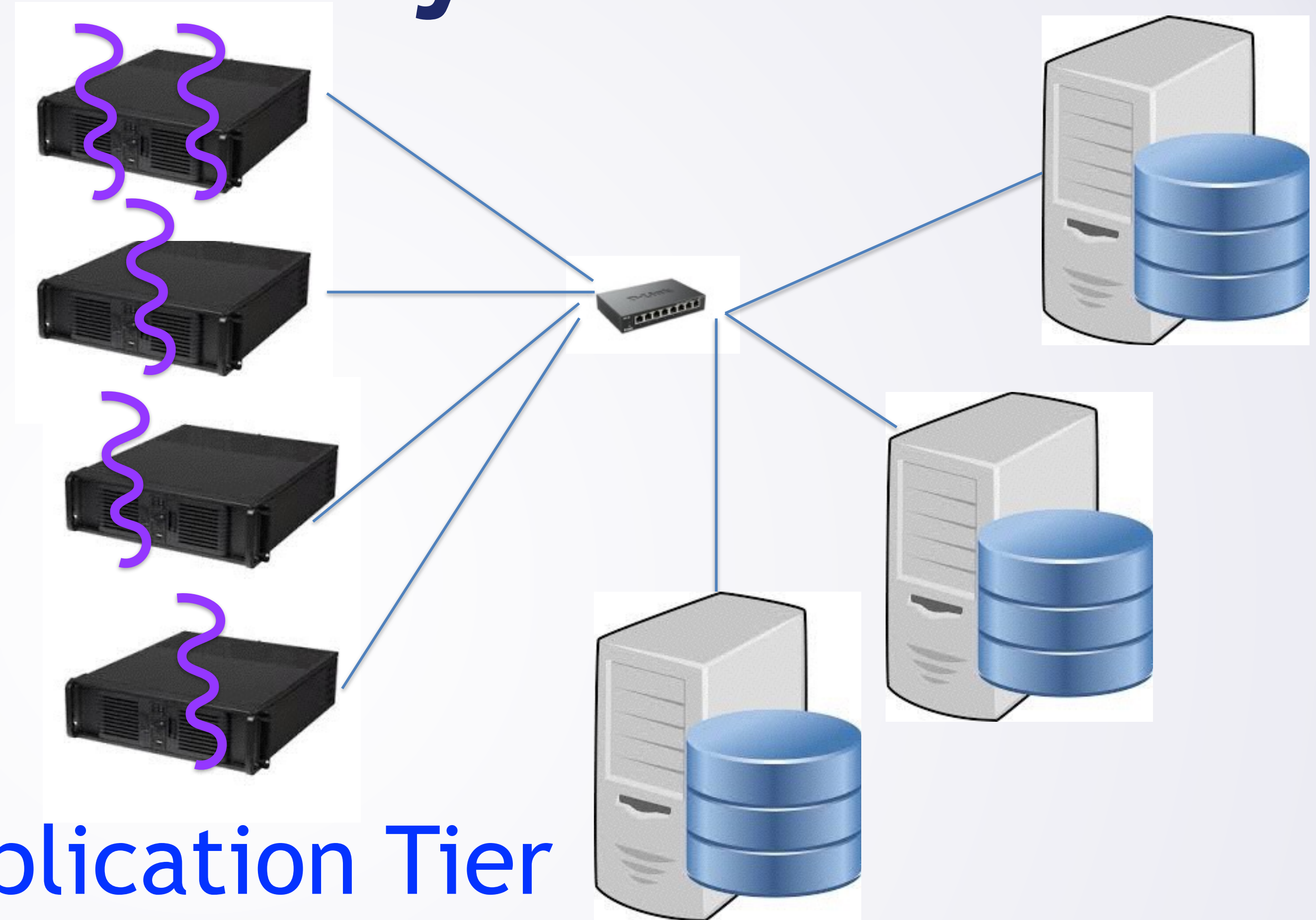
Three Tier System



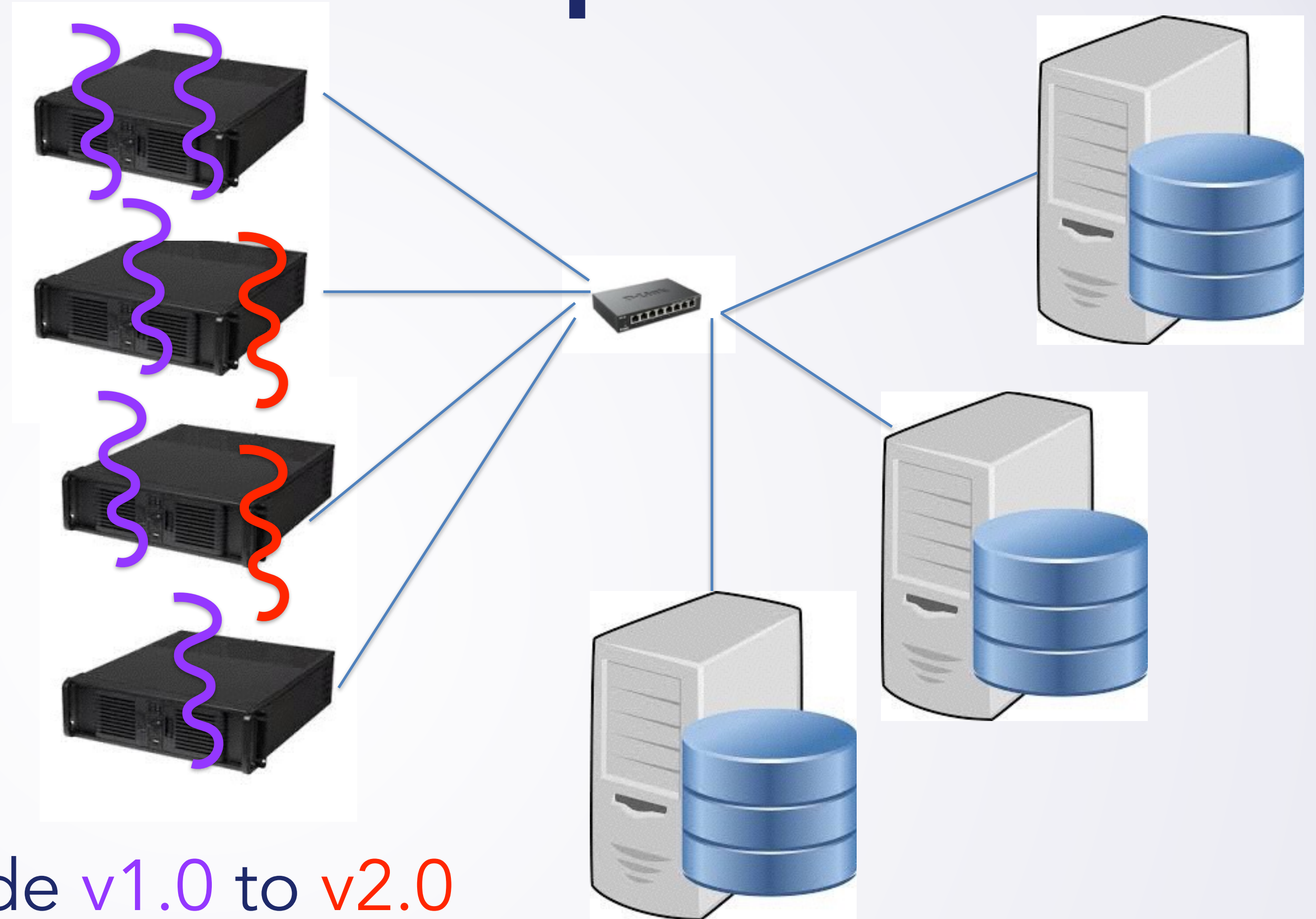
1. Presentation Tier

2. Application Tier
(Business Logic)

3. Storage Tier

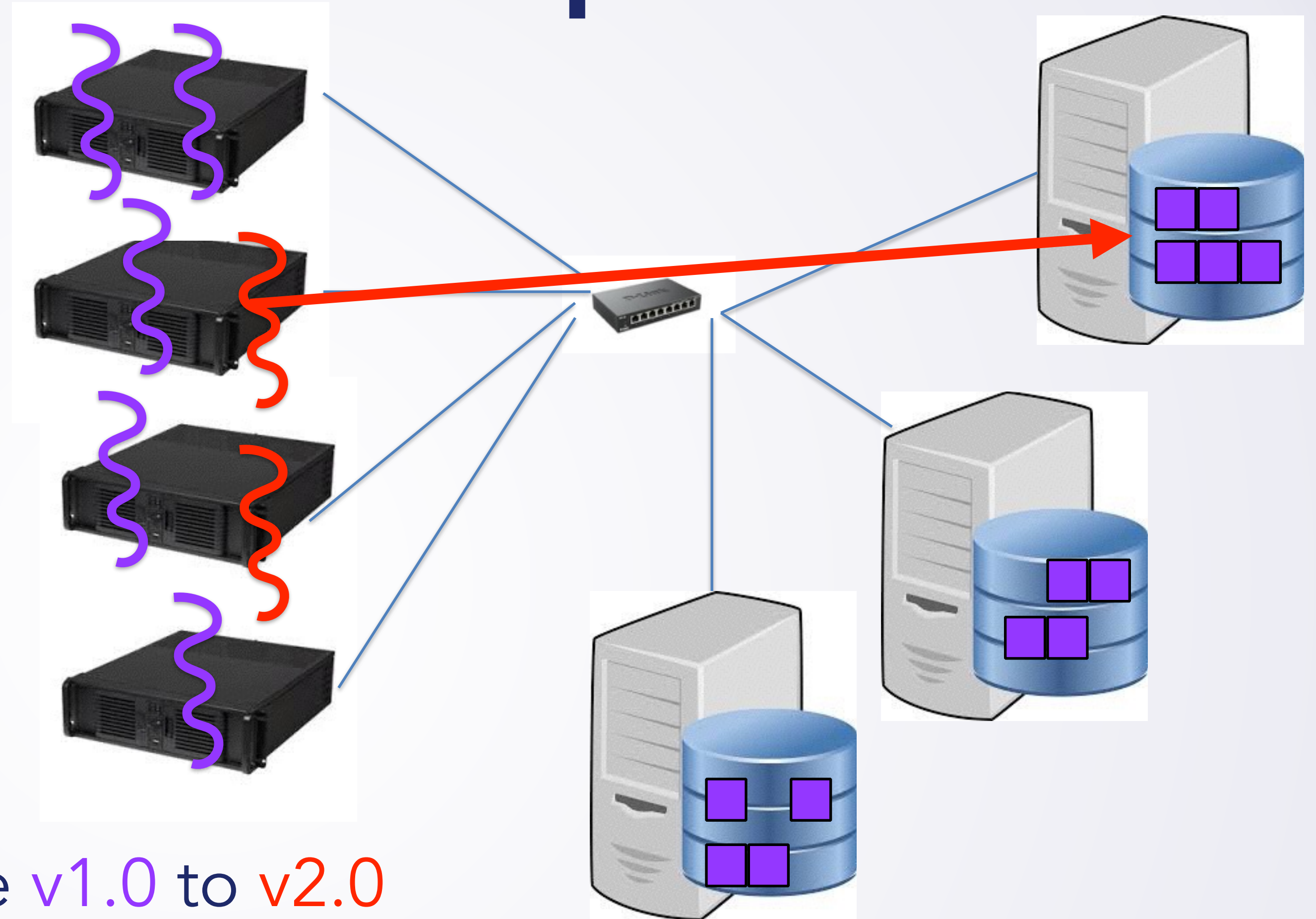


Maybe more complex?



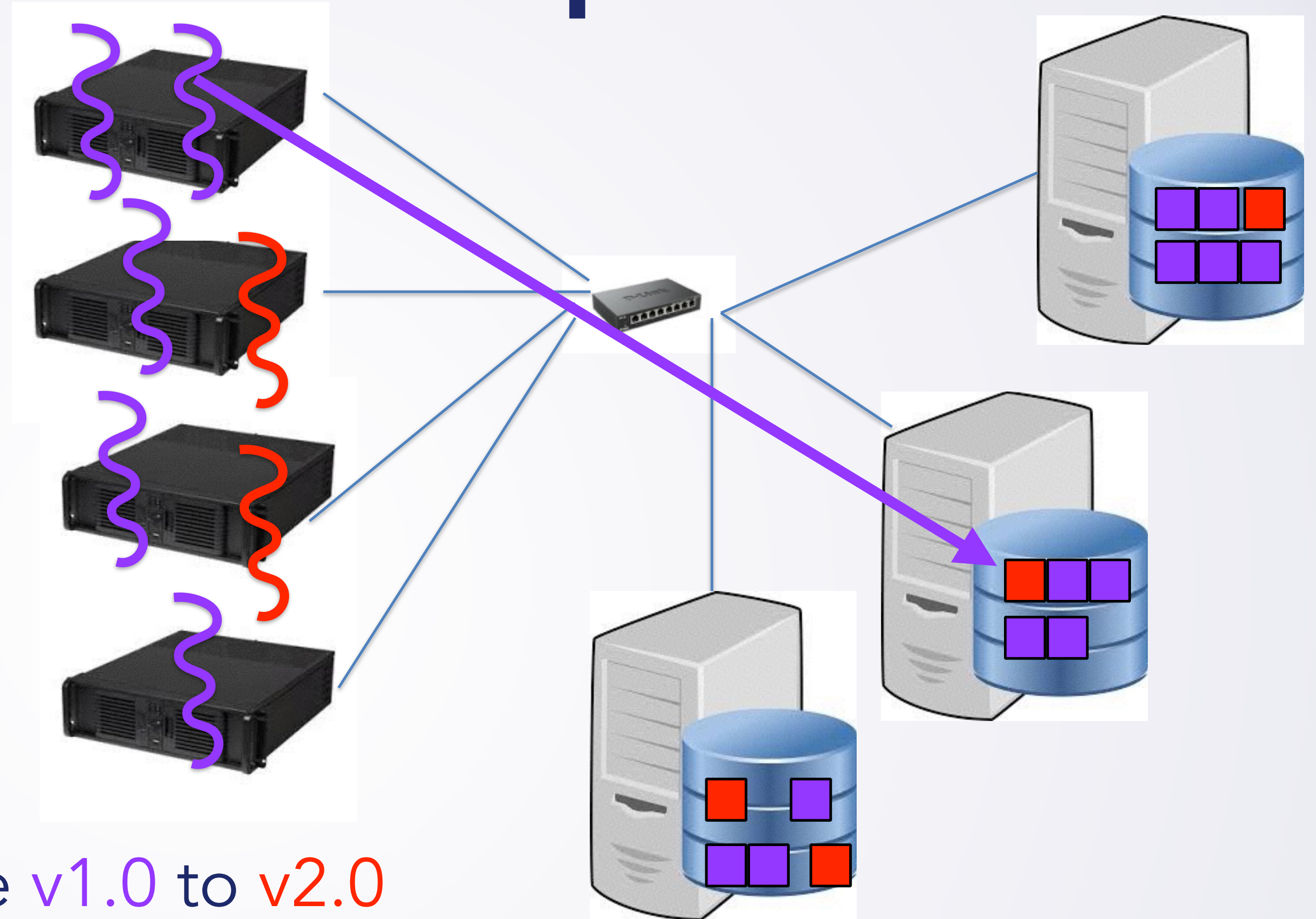
- Maybe we want to upgrade v1.0 to v2.0
 - Difficulties?

Maybe more complex?



- Maybe we want to upgrade v1.0 to v2.0
 - Difficulties?
 - Version 1.0 data: accessed by v2.0 software..

Maybe more complex?



- Maybe we want to upgrade v1.0 to v2.0
 - Difficulties?
 - Version 2.0 data, accessed by version 1.0 software...

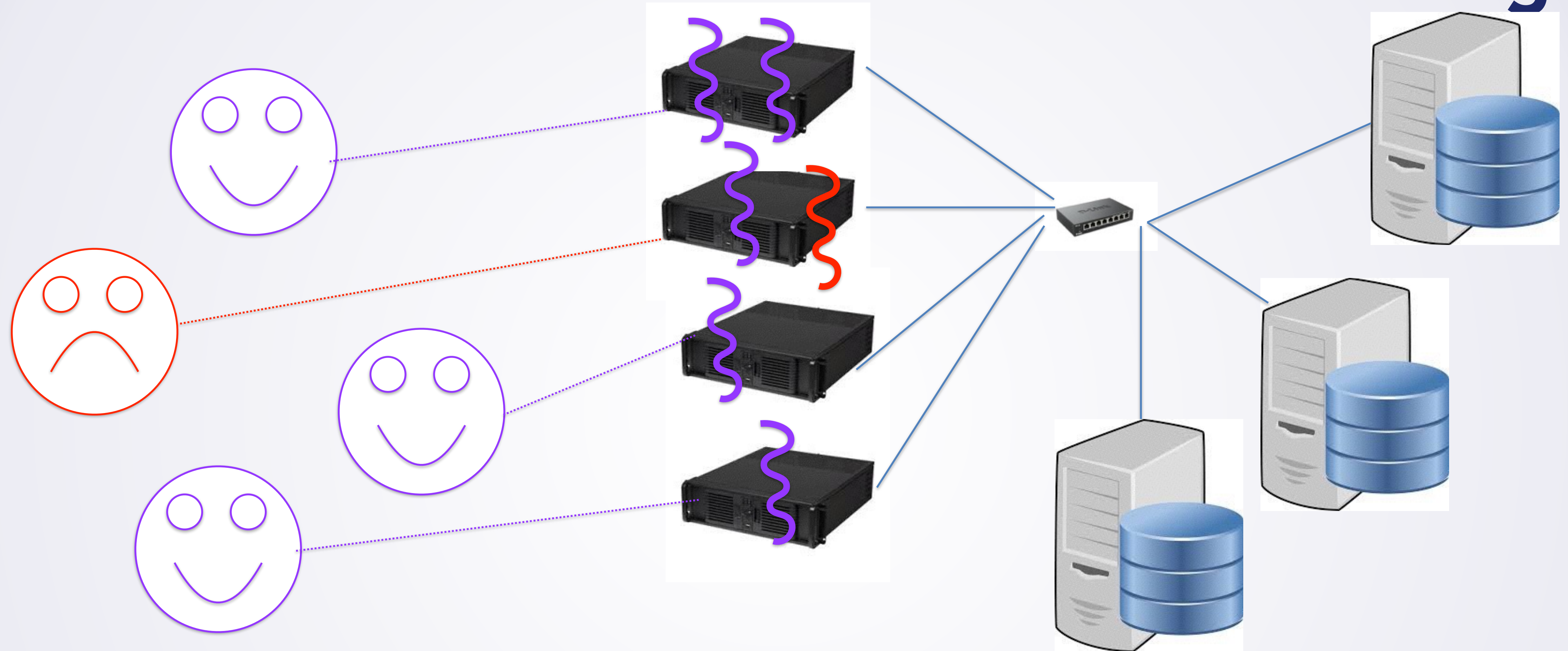
What if we just shut everything down?



Hypothetical picture of what would happen if Google or Facebook were down for 1 minute

- Couldn't we just shut the whole thing down, and upgrade?

Another Reason for Slow Rollout: Testing



- Suppose v2.0 has some bug we didn't catch in testing

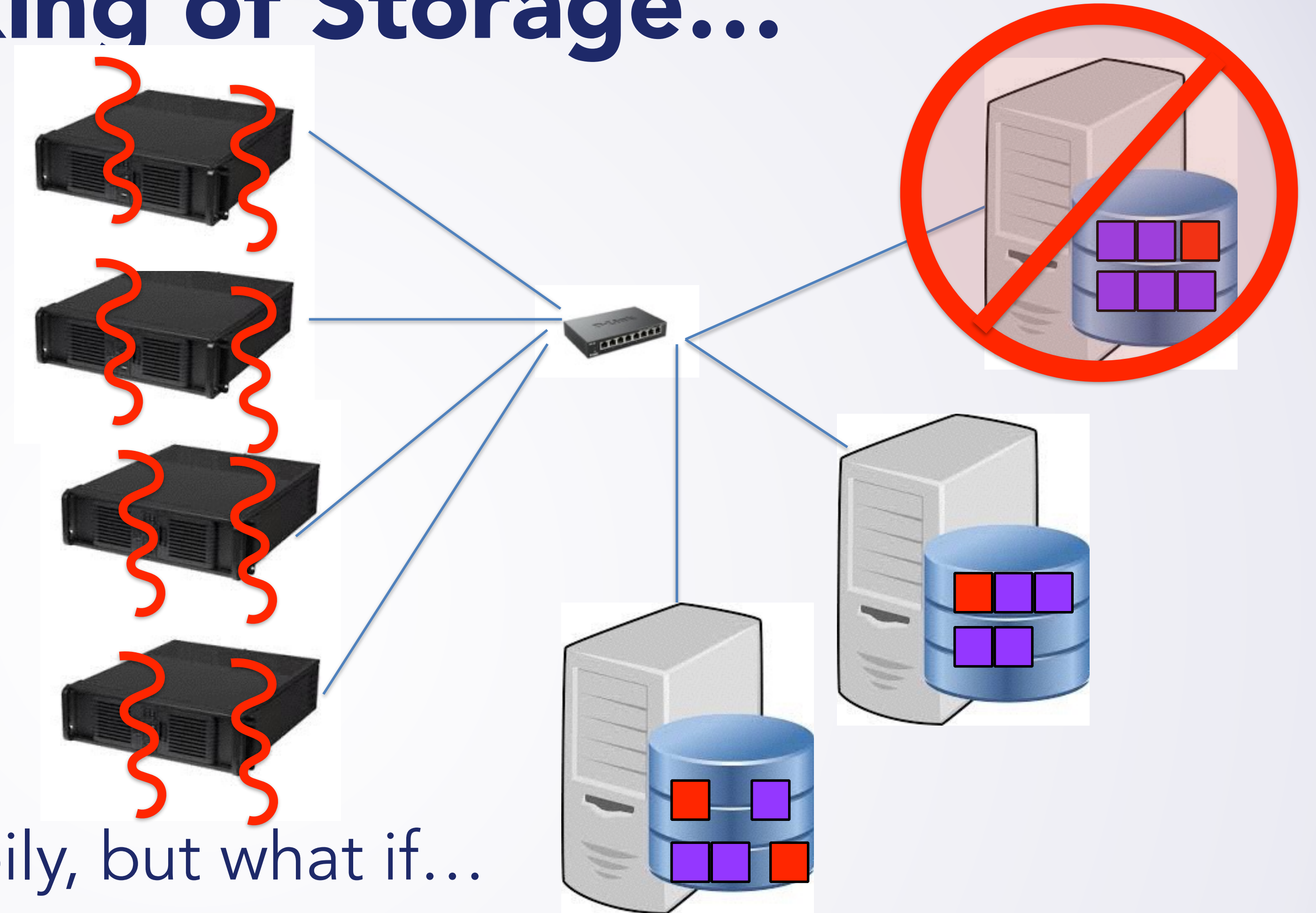
Upgrading Software Versions..

- v2 and v1 work on same data layout/constraints? Great
- v1 can handle v2's layout?
 - v2 [strictly] adds fields or tightens constraints?
 - Design v2 to accept v1 data -> no problems
- v1 can't handle v2's layout?
 - v2 removes/rename fields, relaxes constraints,...
 - Make a v1.9 which: Writes v1 compatible data + accepts v2's data
 - Spin up v1.9s until all v1.0s are gone
 - Then spin up v2.0s

Migrating Data?

- Migrating Data is tricky
 - E.g., change storage tier itself itself?
- Reading:
 - <http://onstartups.com/tabid/3339/bid/97052/How-To-Survive-a-Ground-Up-Rewrite-Without-Losing-Your-Sanity.aspx>

Speaking of Storage...

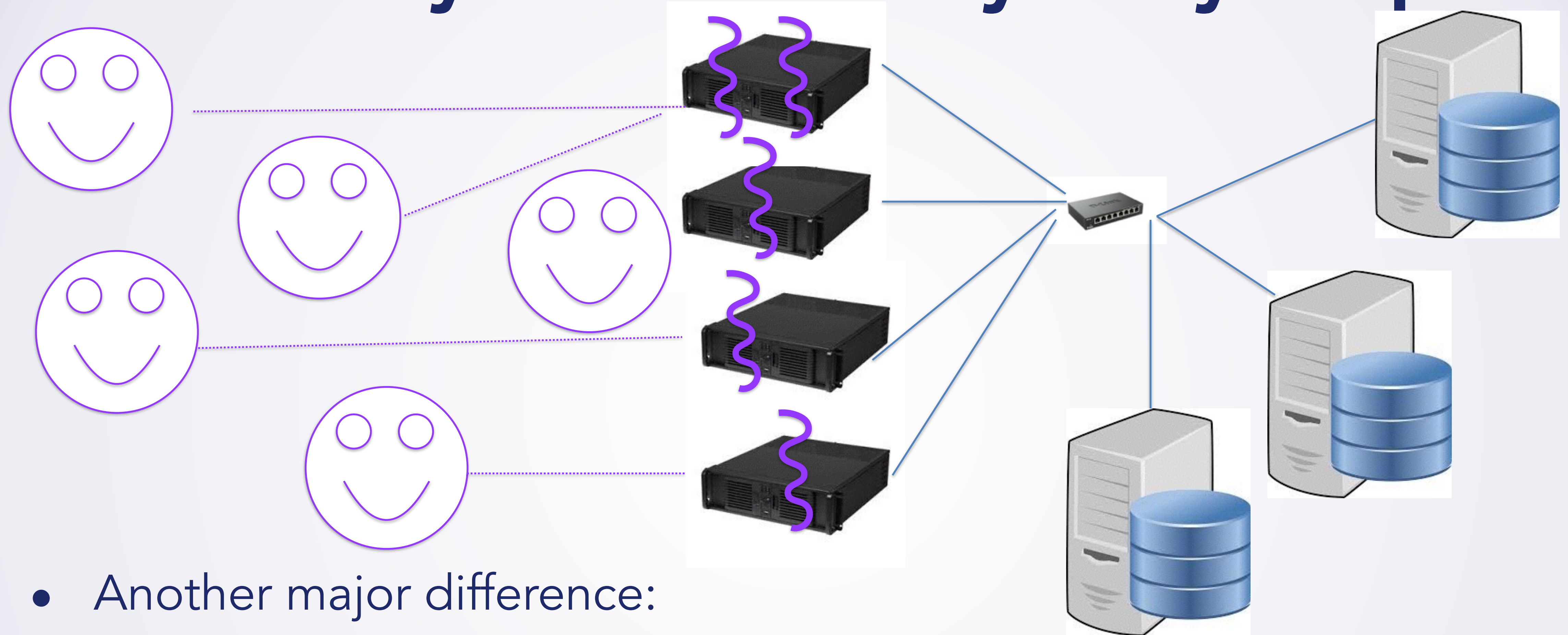


- Our code is running happily, but what if...
 - A storage server fails? Temporarily or Permanently
 - This is part of what we'll discuss right before spring break

Another Major Issue: Configuration!

- Code you have written:
 - Minimal, if any configuration. Likely read at startup
- Servers:
 - Much more configuration: see `/etc/ssh/sshd_config`, `/etc/apache2/*`, etc..
 - Re-read/change while running?
- Warning: changing config as dangerous as changing code!
 - Reading 2:
 - <https://status.cloud.google.com/incident/compute/16007>

Used By You vs Used By Many People



- Another major difference:
 - Things you have written: used by you
 - Server Software: used by (many?) other people...
 - **Complexities?**

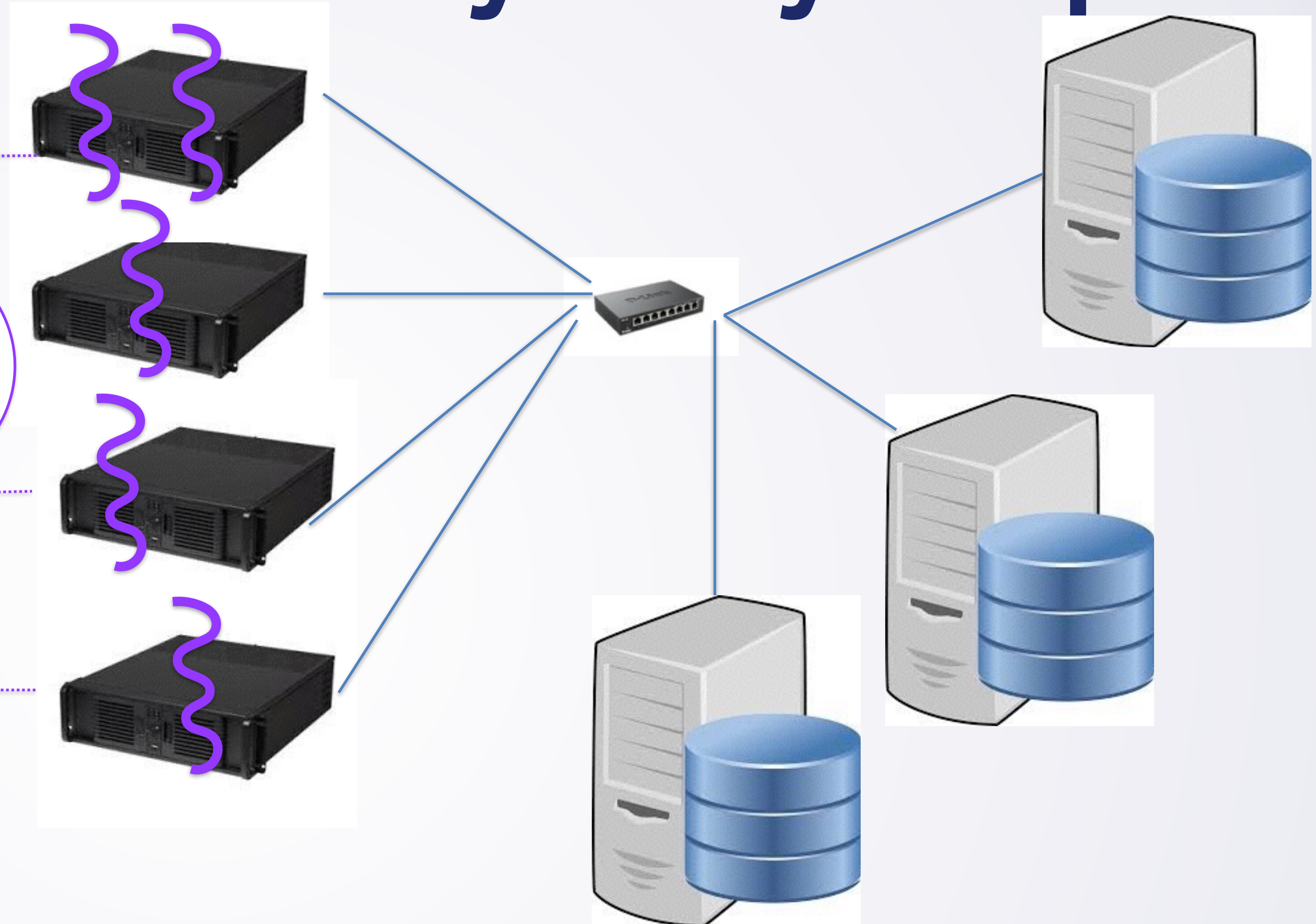
Used By You vs Used By Many People



Book Seat 2A on
flight 1234



Book Seat 2A on
flight 1234



- Concurrency/Scalability
 - Many things going on at once in system
 - Need to handle many requests efficiently

Performance: I feel the need for speed

- Performance: Users care about **speed**
 - Want system to be fast!
- From system perspective:
 - Many users
 - Want to be fast for all of them at once...
- Performance comes in two metrics:
 - Latency: time to complete one request
 - Throughput: requests/second
- Not the same, but they do interact...
- Let us look at non-software example...

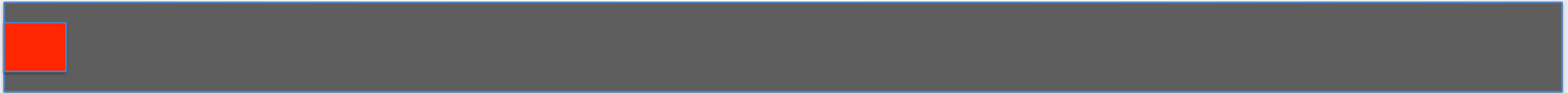


Latency vs Throughput



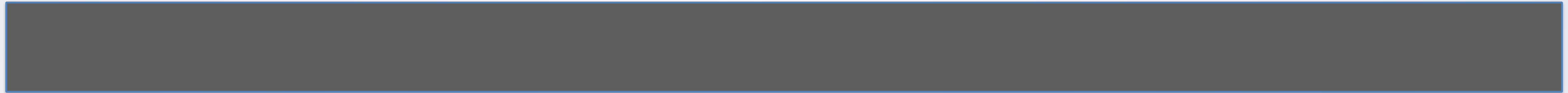
- Here is a "road".
 - 1 lane
 - 70 mph
 - 700 miles long

Latency vs Throughput



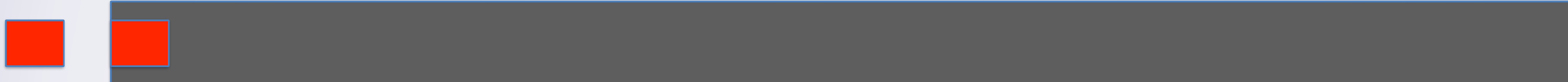
- Latency: 700 miles @ 70 mph = 10 hours to travel

Latency vs Throughput

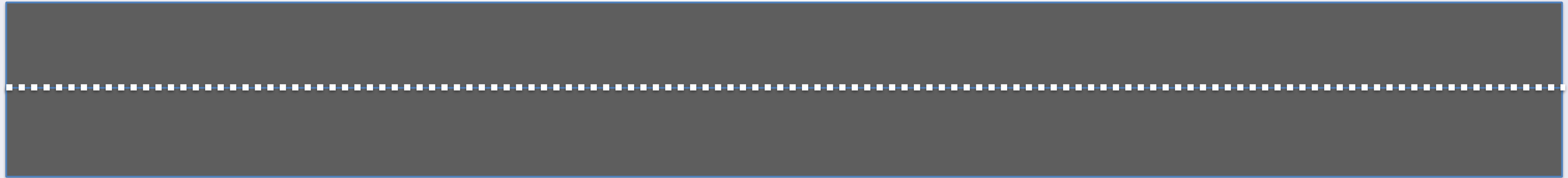


- Latency: 700 miles @ 70 mph = 10 hours to travel
- Throughput: 1 car/ 10 hours = 0.000028 cars/second ?

Latency vs Throughput

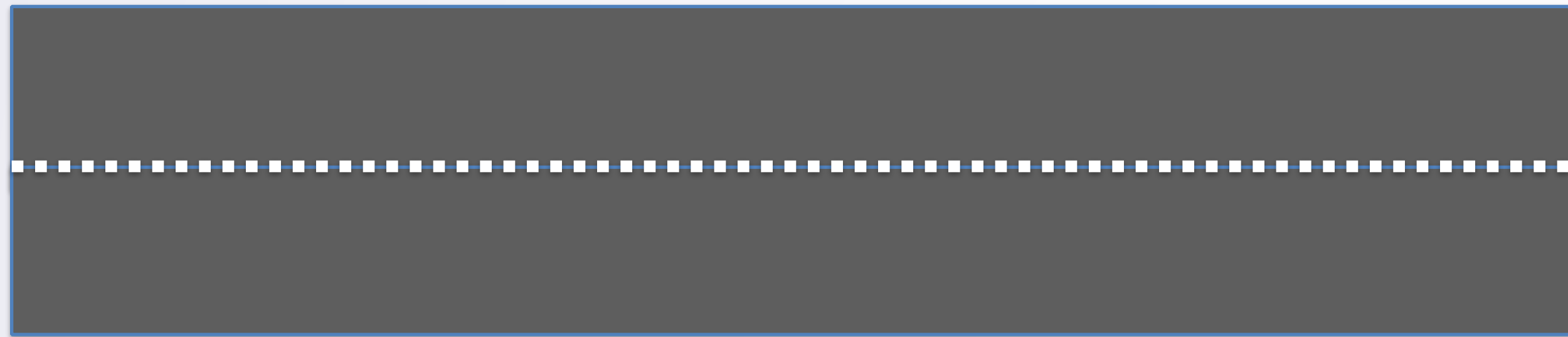
- 
- Latency: 700 miles @ 70 mph = 10 hours to travel
 - Throughput: ~~1 car / 10 hours = 0.000028 cars/second ?~~
 - Throughput: 0.3 cars / second

Latency vs Throughput



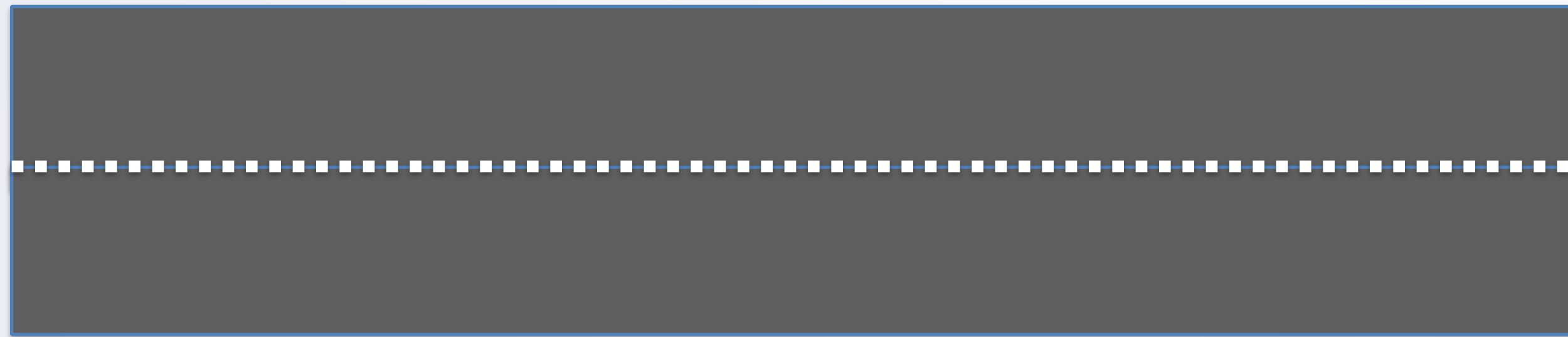
- Different things: can affect one without changing other
 - Another lane? Throughput **improves**, latency **unchanged**

Latency vs Throughput



- Different things: can affect one without changing other
 - Another lane? Throughput **improves**, latency **unchanged**
 - Shorter road? Throughput **unchanged**, latency **improves**

Latency vs Throughput

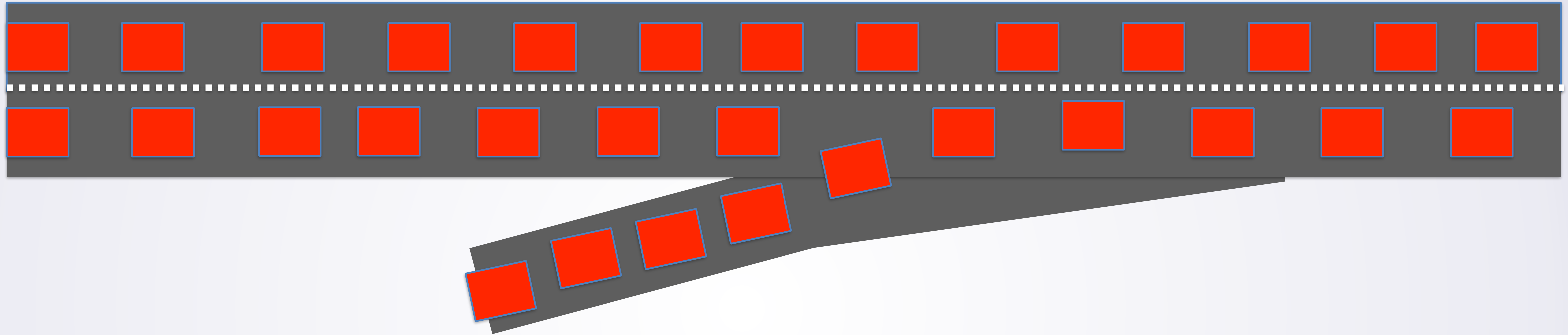


- Different things: can affect one without changing other
 - Another lane? Throughput **improves**, latency **unchanged**
 - Shorter road? Throughput **unchanged**, latency **improves**
 - Cars drive faster? **Both improve (*)**
 - (*) Except that you need more space for safety...

So Which Do We Care About?

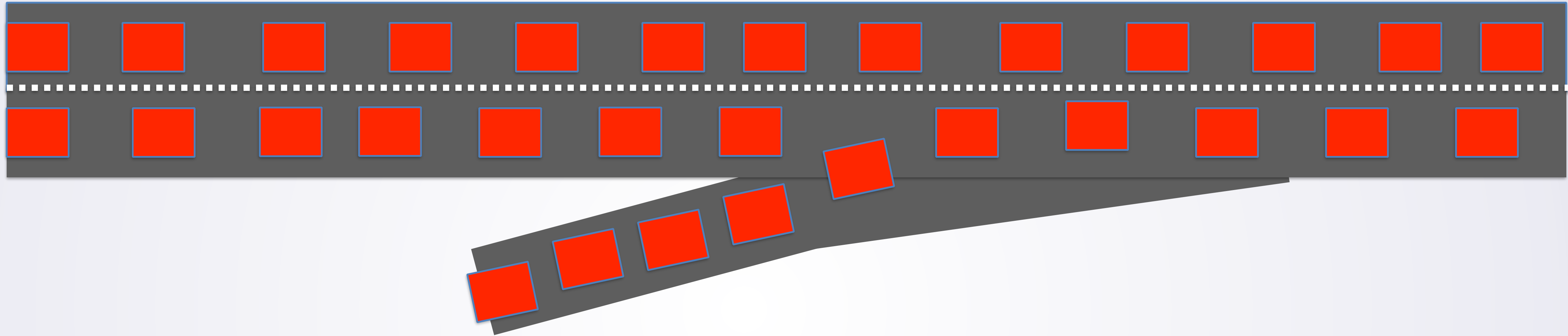
- What matters? Latency or throughput?
 - From a **user's** perspective: **latency**
- From a **system** perspective, **both** matter
 - Need high throughput to get low latency for many users
 - Latency goes up with resource contention and queueing delays
- Back to our road example...

Latency vs Throughput



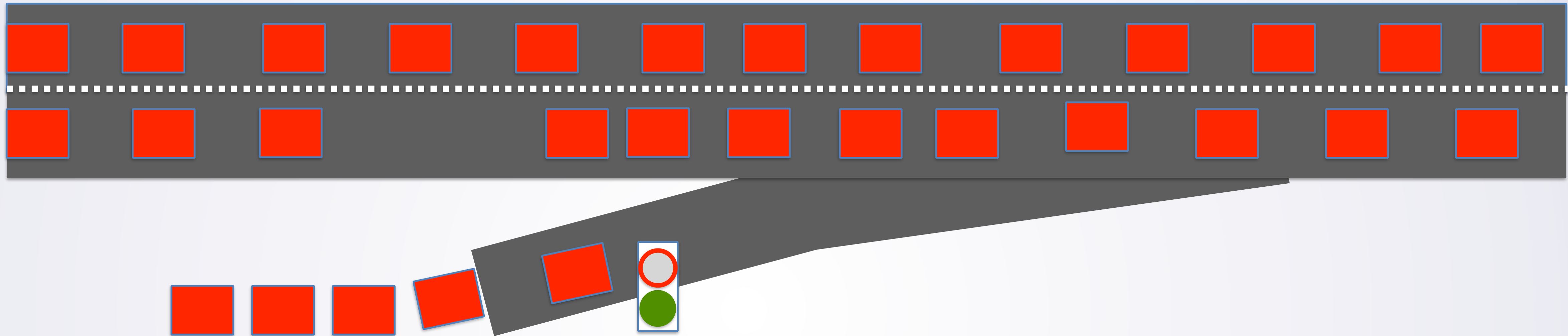
- Heavy traffic, more cars merging in.. What happens?

Latency vs Throughput



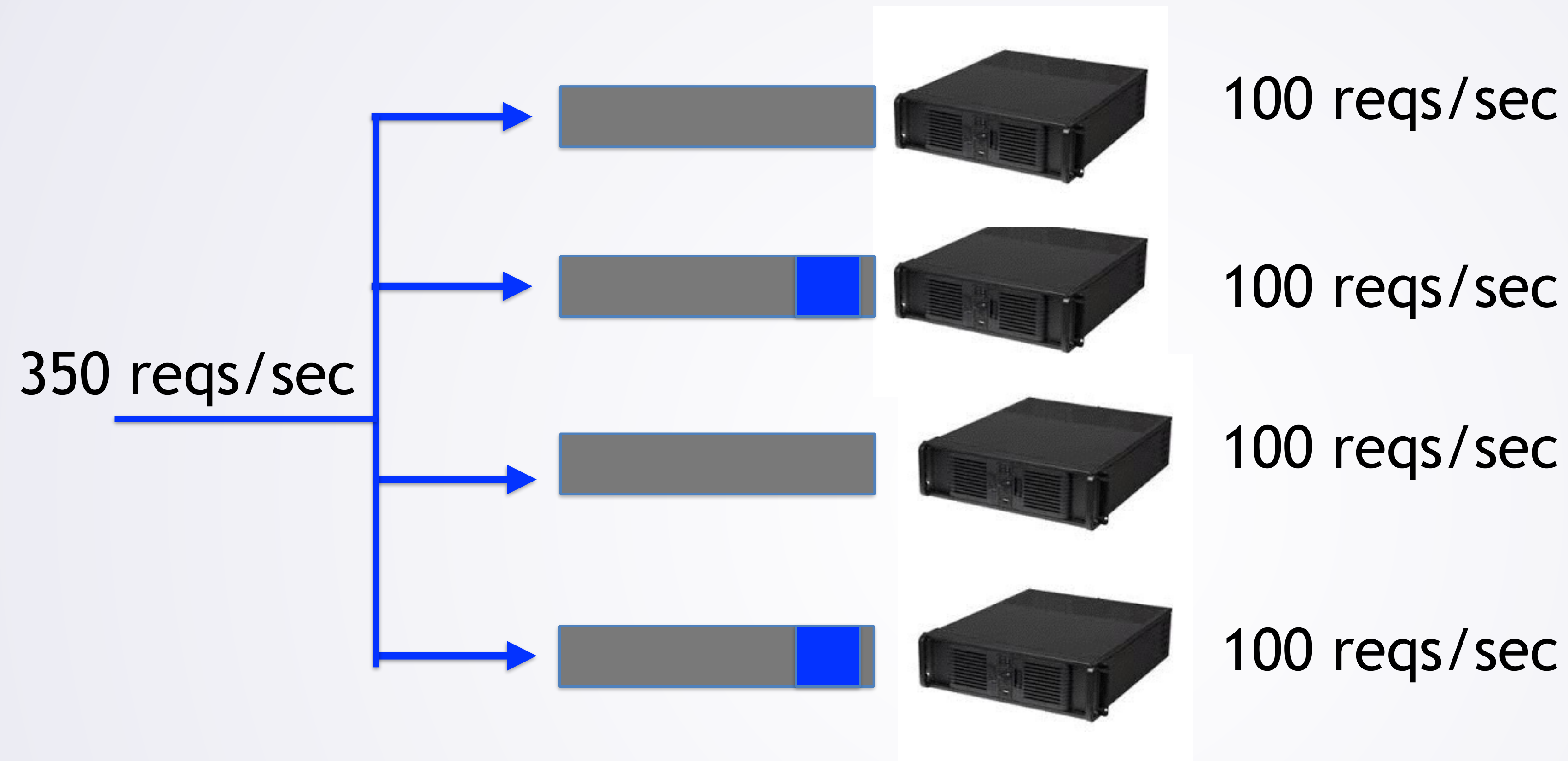
- Heavy traffic, more cars merging in.. What happens?
 - Latency goes up
 - Cars **slow down** due to resource (road space) contention

Latency vs Throughput



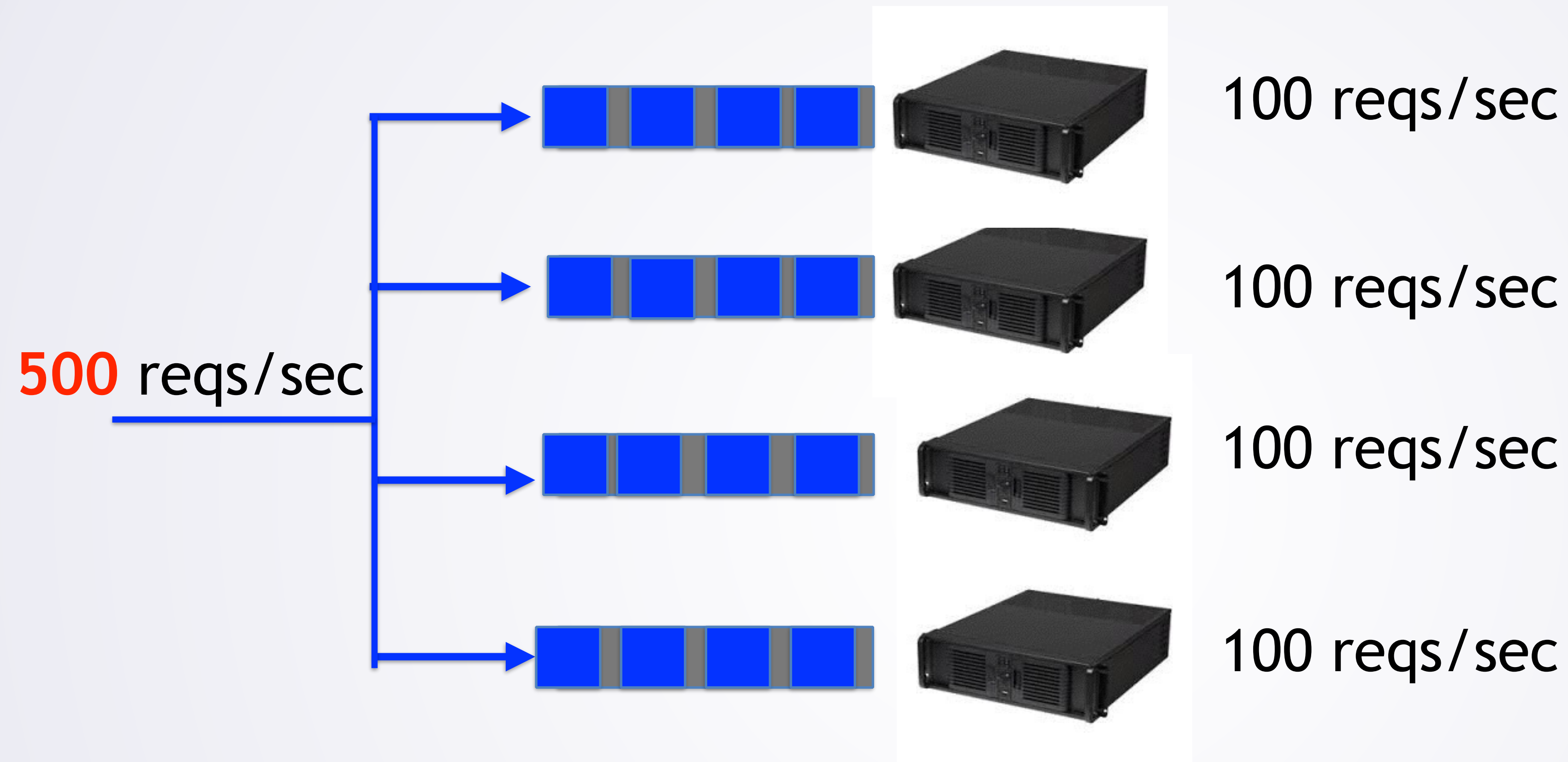
- Alternative: merge traffic lights
 - Traffic queues up (at on ramp)
 - Reduce resource contention (keep speeds higher)
 - Ideally: maintain speed, extra latency comes in queue

Latency vs Throughput



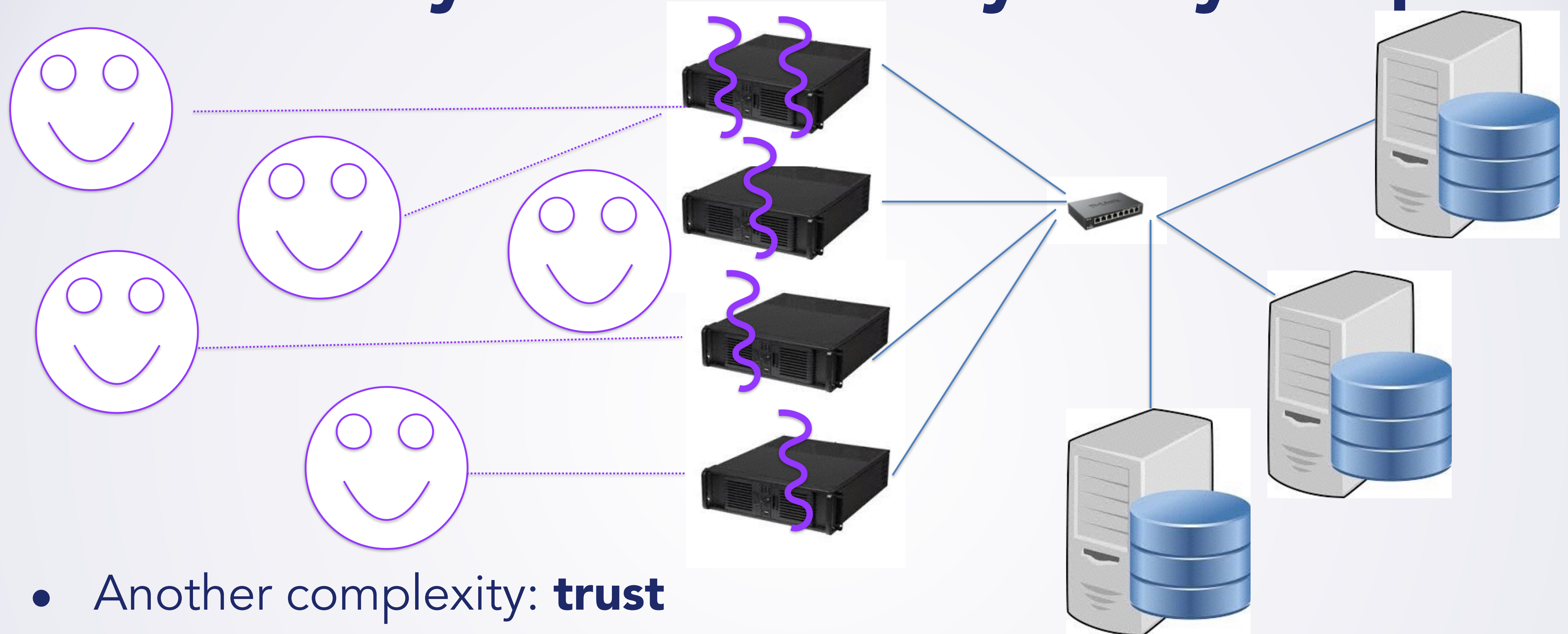
- Adding more systems won't help latency (probably)
 - May experience resource contention (cache, locks, etc...)

Latency vs Throughput



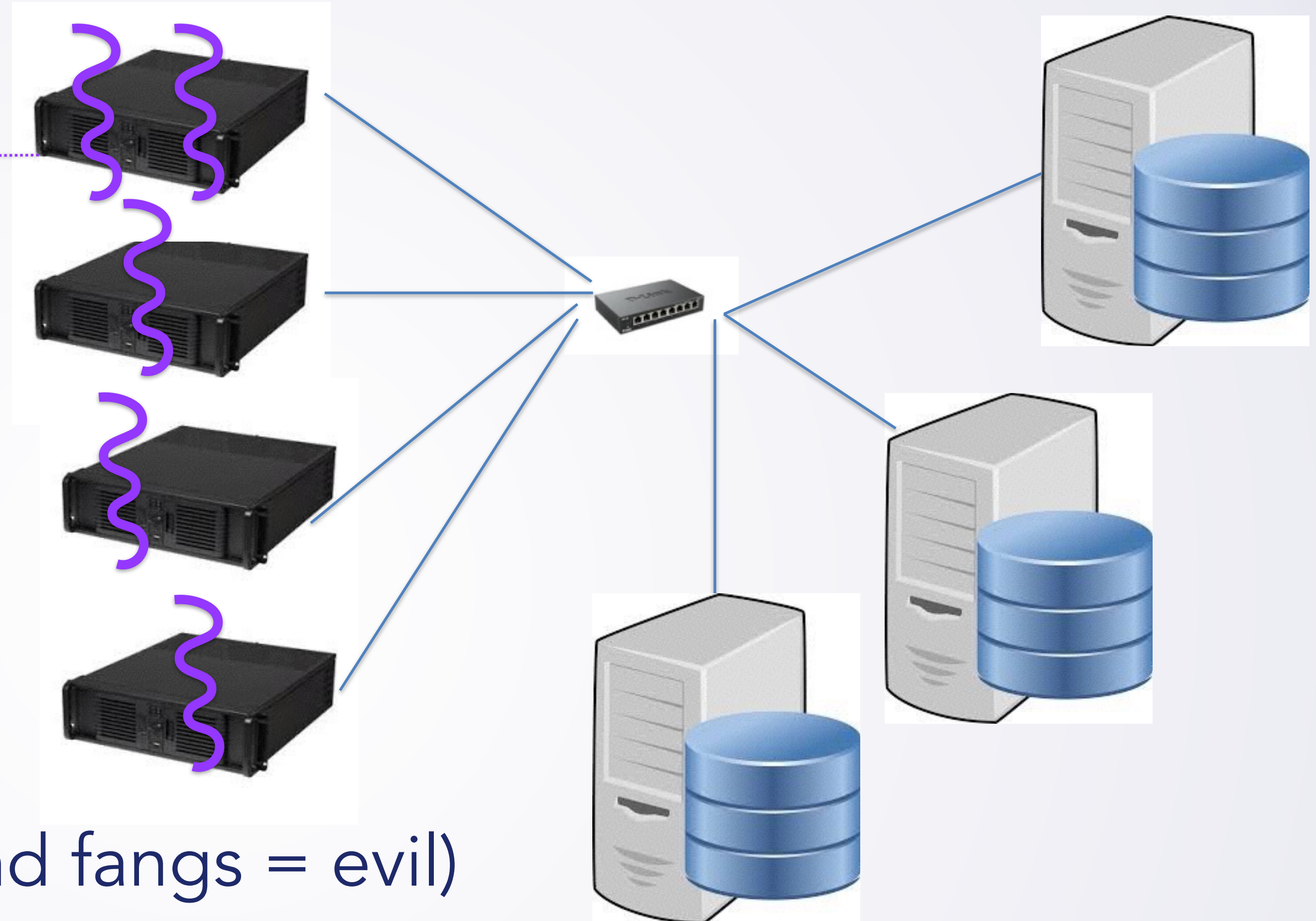
- System is oversubscribed: **queuing delays** add to latency
 - Adding more throughput would reduce latency!

Used By You vs Used By Many People



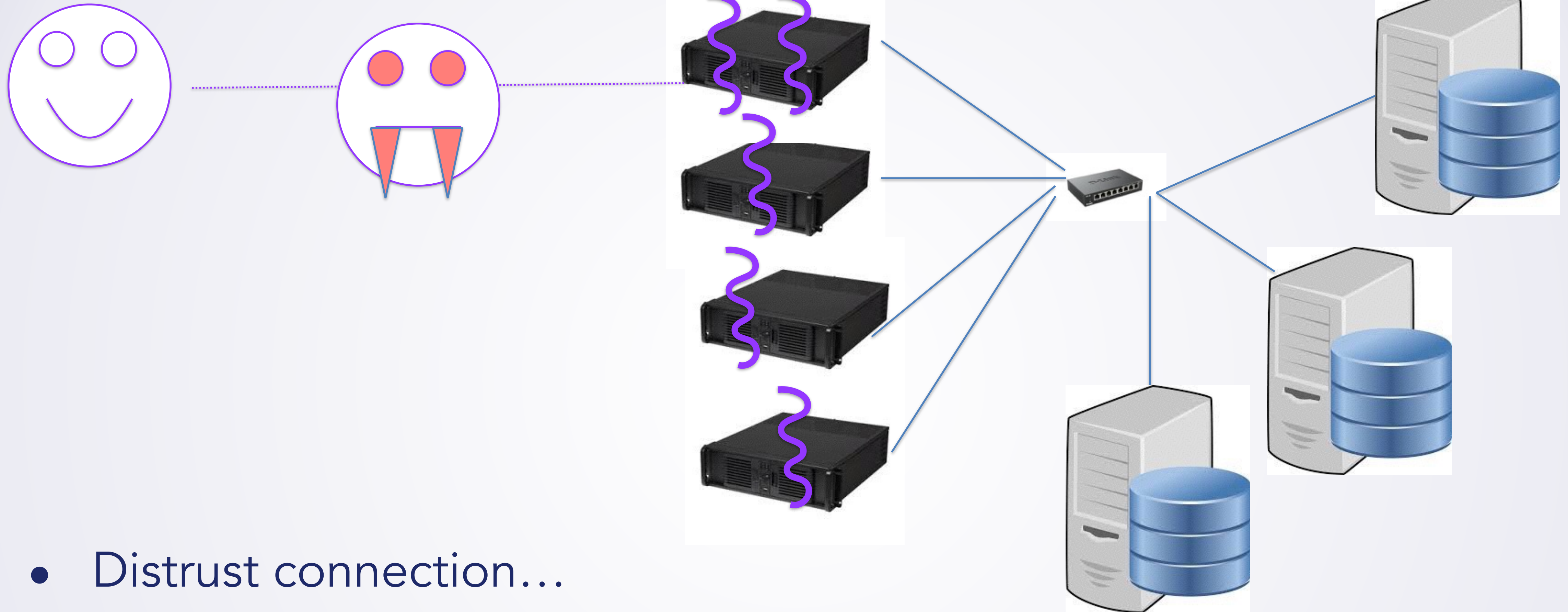
- Another complexity: **trust**
 - Are all those users out there good?

Trust?



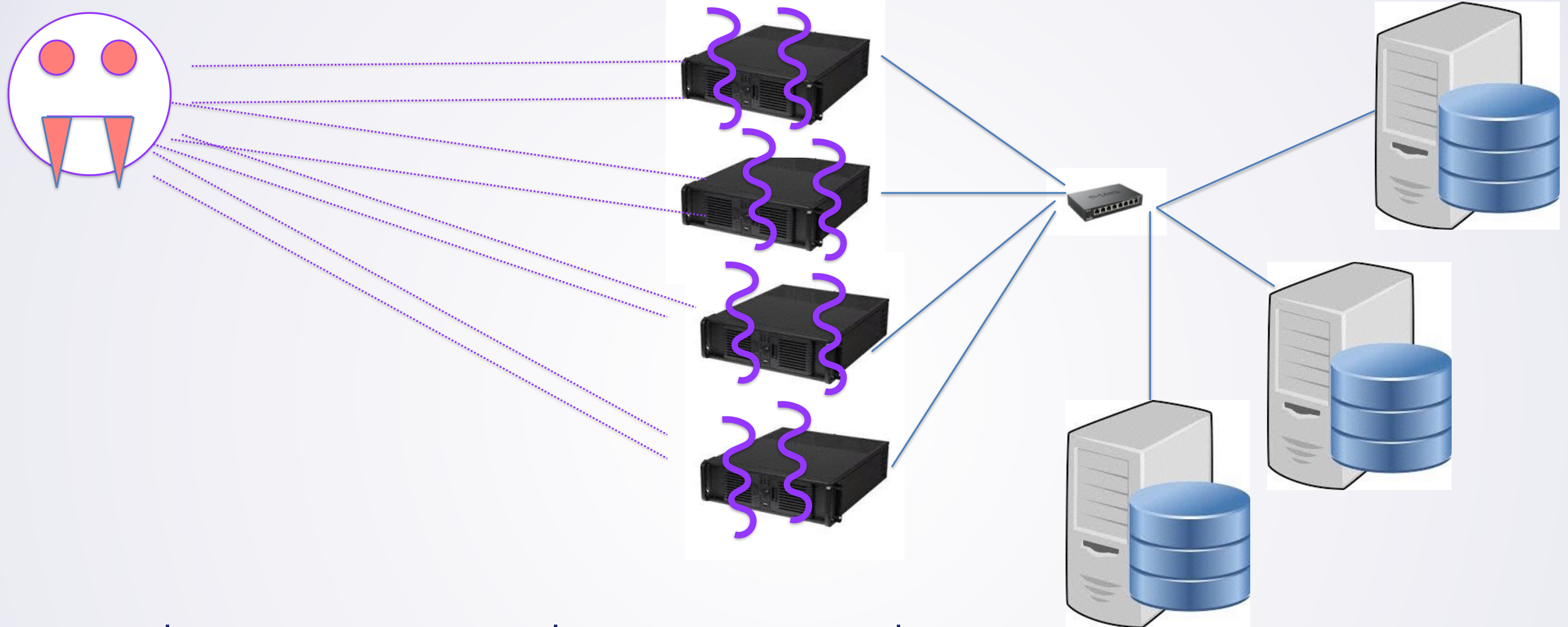
- Might be evil (red eyes and fangs = evil)
 - Steal information
 - Modify information
 - Use server for nefarious purposes (spam,...)

Trust?



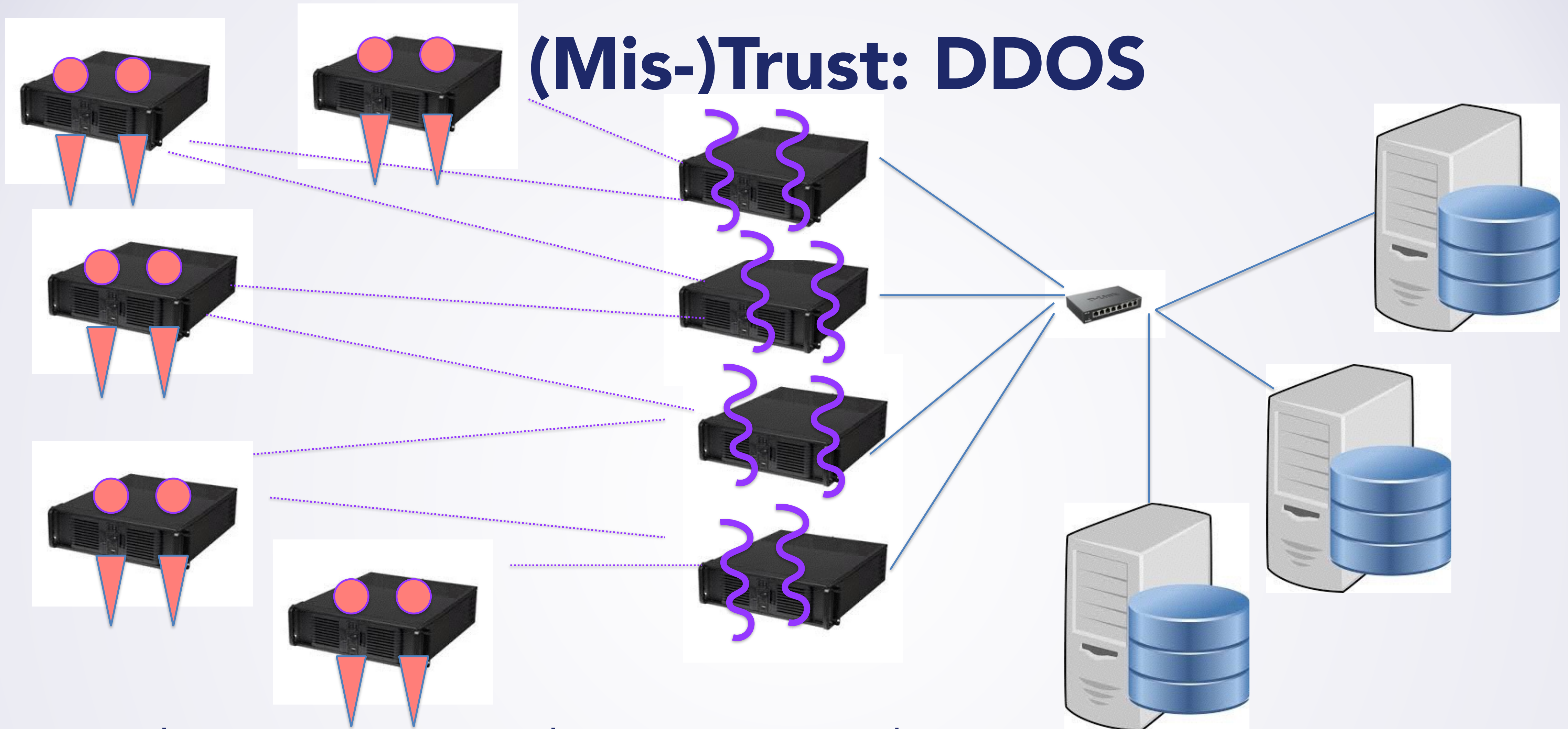
- Distrust connection...
 - Adversary might eavesdrop (passively gather information)
 - Or tamper with connection (actively change what is sent)

(Mis-)Trust: DOS



- Malicious user may also attempt to deny service
 - DOS = Denial of Service

(Mis-)Trust: DDOS

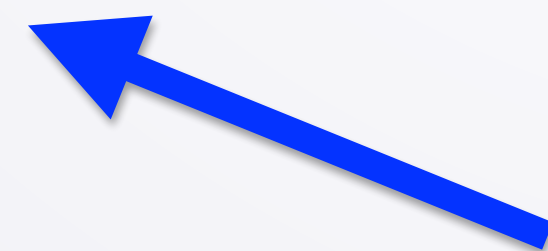


- Malicious user may also attempt to deny service
 - DOS = Denial of Service
 - DDOS = Distributed Denial of Service

What Does The Server Look Like?

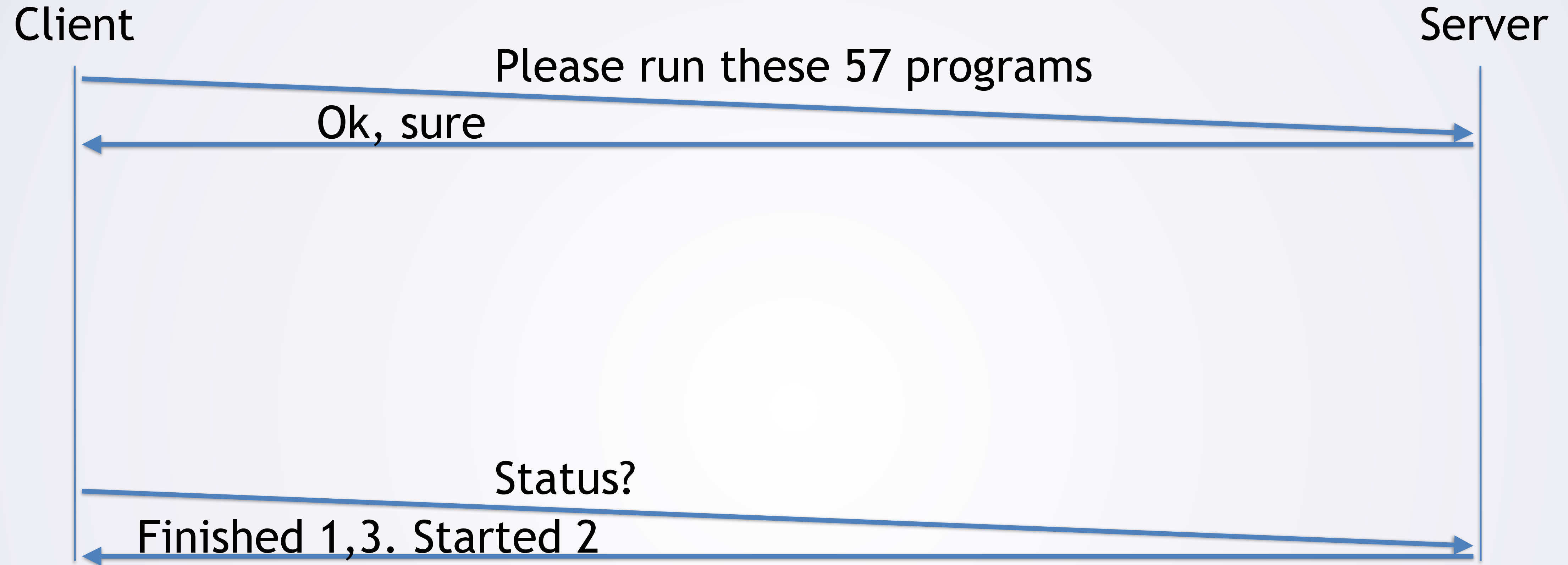


- Now, we've seen a bunch of differences in constraints/requirements
- But what does the server itself look like?
 - ...it depends...



Always the answer in CE

Batch Servers

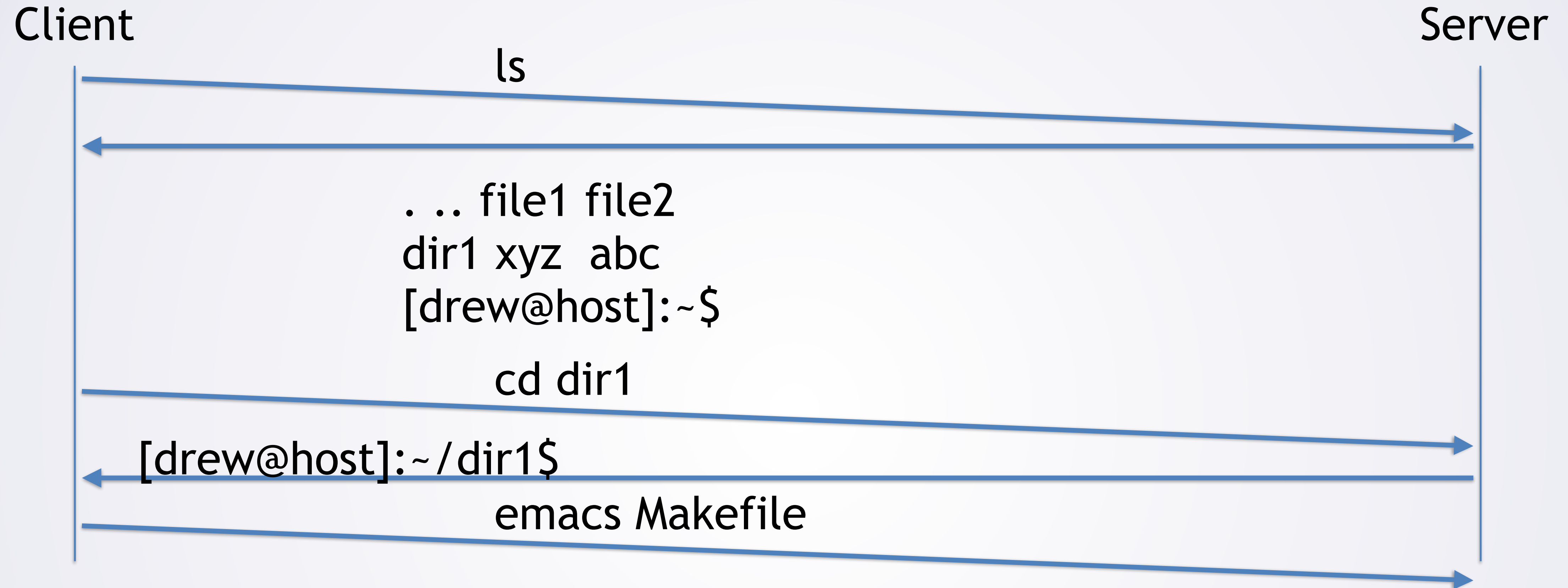


- Submit jobs (possibly in bulk)
- Server will do them later (when it can)

Batch Servers

- Examples:
 - Sun Grid Engine, Condor,...
- Mostly queue requests
 - Possibly with priorities
- Most concerned with throughput
 - Overhead latency \ll job latency
- Running code for user?
 - Generally more trust than most systems

Interactive Servers



- (Many ?) requests, sent/handled frequently

Interactive Servers

- Examples:
 - sshd
 - Game servers (WoW)
- Latency is critical
- Web-servers similar,
 - Just flurry of requests, then close connection

Database Servers / DBMS

- Process queries from clients
- Often must efficiently process many tuples to satisfy query
 - High tuple throughput -> low response latency
- Often have special IO needs, require much RAM
- Quite a complex beast (topic of advanced database classes)
- Examples: Postgres, MySQL, Oracle,....



File Servers

- Put filesystem on remote server
- Why?
 - Use same files on many systems
 - E.g., login to any lab computer, have same home directory
- Compute requirements \ll IO requirements
 - IO slower than compute anyways
- Examples: NFS, AFS,...

Proxy Servers

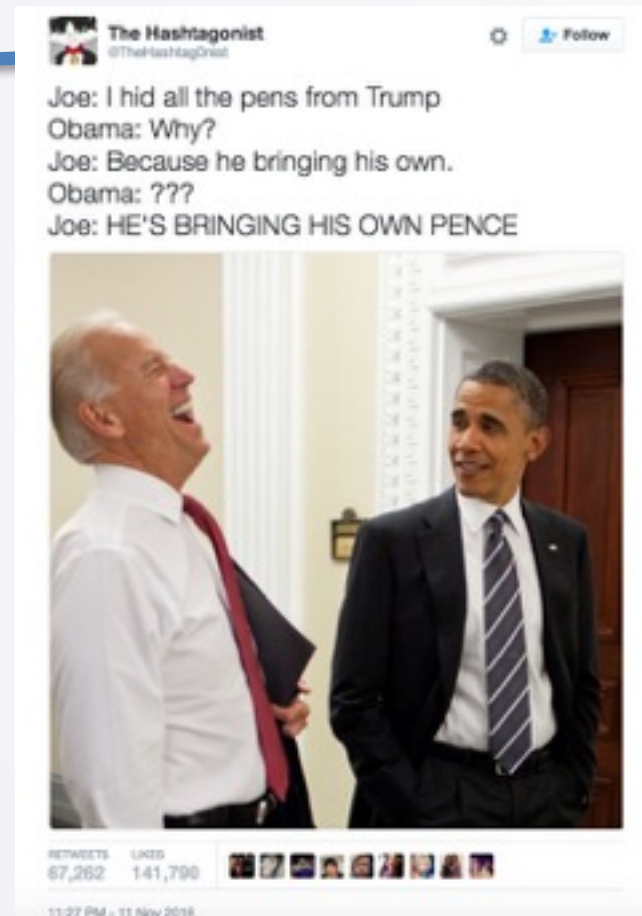
Client

Proxy

Server

GET obama_biden.png

GET obama_biden.png



- Pass requests to "actual" server

...but really...all the same

```
while (true) {  
    req = accept_incoming_request();  
    resp = process_request(req);  
    send_response(req, resp);  
}
```

Note: really need some parallelism

- Pretty much all of these have a unix daemon that
 - Accepts requests
 - Processes them
 - Sends responses

Coming soon: Unix Daemons

```
while (true) {  
    req = accept_incoming_request();  
    resp = process_request(req);  
    send_response(req, resp);  
}
```

- Soon: all the details of how to make this work
 - You'll write a web proxy server
- 650: concurrency + socket programming

Coming soon: Unix Daemons

```
while (true) {  
    req = accept_incoming_request();  
    resp = process_request(req);  
    send_response(req, resp);  
}
```

- Server side web development
 - How to process the request
 - Web-servers (Apache,...) have ways to "hook up" to code to generate content

Next Time:

- Next Time:
 - Slight reordering of schedule (to match 650)
 - Web protocols and technologies
 - HTTP
 - HTML
 - CSS
 - Javascript
 - XML
 - JSON
 - ...