



Google

SRE Classroom

Alex Perry @ Google, Los Angeles

Andrew Suffield @ Goldman Sachs, London

Agenda



Basics of Non-Abstract Large System Design

Problem Statement

Group Work Session (with break)

Example Solution

Before we start...

There will be a handout.

We will share a link with further reading at the end.

Relax. Have fun.

Keep that little paper with the number on your seat.

Slides will be shared.

Disclaimer: Ain't Nobody Got Time for That!



Due to time constraints this is not a complete discussion of large system design abstractions/solutions/patterns.

There are other talks that discuss this in more depths.



Basics

Basics - Topics

- Requirements & Scaling
- Dealing with Loss
- Keeping State
- Non-abstract Design



Basics: Requirements & Scaling

Requirements



Identify Service Level Indicators (SLI) and Objectives (SLO)

- Data freshness
- Availability
- Latency

This is a job for an SRE.

Requirements



$$\text{SLI} = \frac{\text{good events}}{\text{total events}} \times 100\%$$

An SLI is the proportion of valid events that are good.

It is a time series which falls between 0 and 1.

Requirements



Sample SLO: 99%ile of queries returns valid result within 100ms

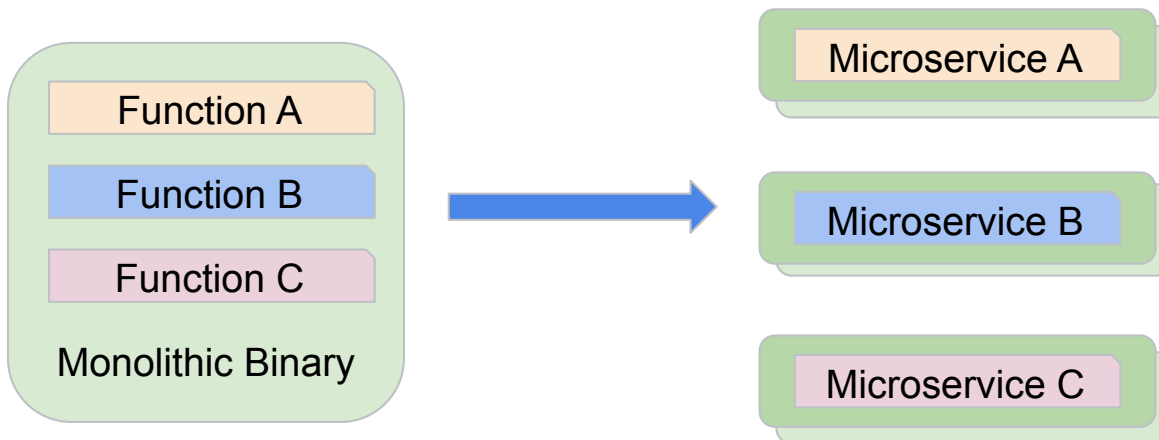
Service Level Agreement (SLA) is the contract containing all relevant SLOs and the punishment if it's violated.

Design for current and consider future scale: What about 10x?

Design with Microservices



Replace a monolithic binary with distinct microservices.



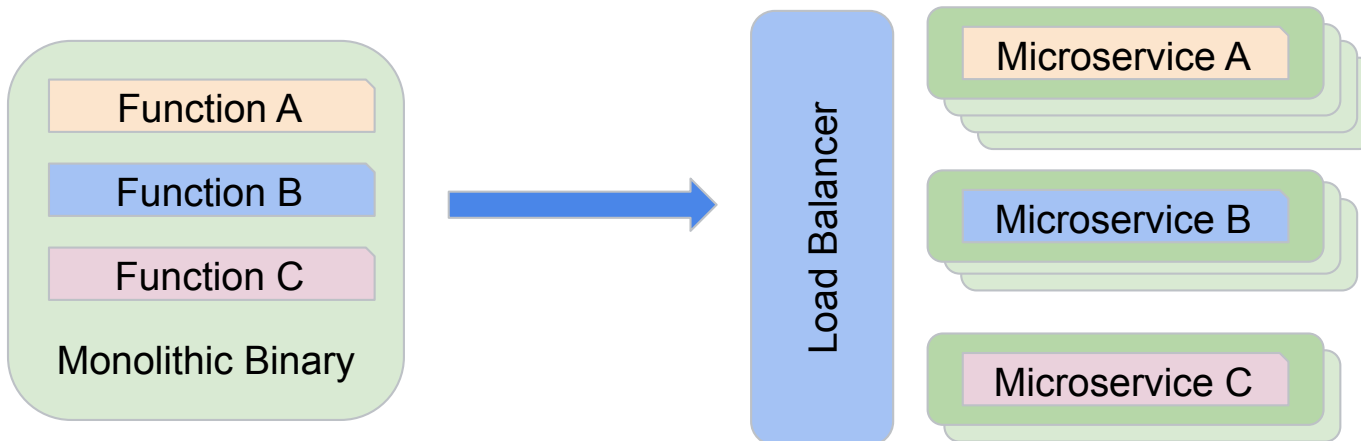
Microservices Add Flexibility



Allows horizontal scaling. Makes functional scaling easier.

Each component can scale independently after deployment.

Add load balancers for work distribution.



Scaling in Distributed Systems



Typically refers to horizontal scaling.

Other relevant dimensions of scaling:

- Geographical: More physical locations
- Functional: Different feature sets

Dealing with Loss and Failure

Failure Domains



Failure is not an option. It is certain to occur.

- Single process
- Machine
- Rack enclosure
- Datacenter / Zone
- Servers running the same binary
- Global configuration system
- Everything used by a single customer



Comic by KC Green - gunshowcomic.com

Understanding and Dealing with Failure



Preferred: Make it somebody else's problem.

Cloud providers offer various ready-made solutions.

Be aware of their limitations.

Designing to Defend against Failure



Decouple: Spread responsibilities across multiple processes.

Avoid global changes: Use a multi-tiered canary rollout.

Spread risk: Don't depend on one backend.

Degrade gracefully: Keep serving if configs are corrupt or fail to push.

Achieving Reliability



Run $N+1+1$ stacks geographically distributed.

N = deployment large enough to deal with standard load

Why +1s?

- +1 or more for planned maintenance (repairs)
- +1 or more for unplanned maintenance (failures)

Keeping State and Data

Keeping State



Useful for consistency, performance and reliability.

Regardless of the amount of data involved, it all comes down to global consensus.

Always prefer stateless. It is easier.

Speaking about Consistency



Global consistency critical for large systems.

Find authoritative instances of other services (leader election).

Quick discovery of sharded data.

Various algorithmic implementations for this: Paxos, Raft

Consistency - CAP



CAP = consistency (C), availability (A), partition resilience (P)

Networks are not reliable, but partitions are rare.

Need to be able to converge after partition.

E.g. Spanner: AP plus eventual consistency.

Hot Data and Hotspotting



Some data is accessed more frequently than others.

Frequent access to the same data can cause servers to overload.

It can also create bottlenecks via I/O.

One solution: Sharding.

Caches



Move data access to faster I/O to handle hot data (e.g. to RAM).

Improves user experience.

Reduces load on backend.

Caching can be for performance or capacity, but caching for capacity is hard.

Think about whether your cache can support your SLO

Non-abstract Design

Suggestions



For each microservice consider:

- Disk I/O
- QPS
- Network Bandwidth

Use back of the envelope numbers! These are easier to calculate in head.

- One day has 25 hours.
- Each month has 30 days.
- Each year has 300 days. Or 400.



Problem Statement

Problem Statement



Design a photo-serving system that will have an initial global user base of 1 million users, each user uploading 50 photos per day.

The initial UI provides two pages: a thumbnail view with 10 thumbnails per page, based on a search result (user, latest, topical, hashtag); and a detail view of a specific image.

An uploaded compressed picture is 4 MB in size.

Each thumbnail is 256 KB.

SLOs



Serve the thumbnail page (10 thumbnails) within 250 ms at the 99.9 pctl (HTTP 200).

Serve the image page (1 image, full-resolution) within 200ms at the 99.9 pctl (HTTP 200).

These latency SLOs do not apply to data older than 30 days.

Available Infrastructure



Network (SLA 99.99 % availability).

Storage system (SLA read 100 ms at 95 pctl, write 200 ms at 95 pctl).

Storage system handles downsampling, provides thumbnails on request.

Globally replicated, eventually consistent in the order of minutes.

Multiple datacenters in 3 regions: Europe, Asia, North America.

Each unavailable for one week each quarter for planned maintenance.

Available Hardware



HDD Server: 24G RAM, 8 cores, 2x2TB hard drives, 10Gb ethernet

SSD Server: 24G RAM, 8 cores, 2x1TB SSD, 10Gb ethernet

As many as you need if you can make a case for it.

Form up in Groups!



The number on your seating card is your group.

Each group has sticky notes, markers, and a facilitator.

Suggestion: Start designing a system diagram for one datacenter.

⇒ Does it meet the SLO?

Then expand the solution for multiple datacenters.

⇒ Does it still meet the SLO?

⇒ Run numbers and create a bill of materials

Form up in Groups!



Perfect solution not required. Use back of the envelope style of reasoning.

Laptop not required.

You'll work in groups, by tables, assisted by facilitators.

Relax, this is not a test ;)

Have fun tackling a technical problem together.

Group Work Time



Ask any facilitator if you have any questions.

We have a break for refreshments from 15:30 to 16:00.

We will present a reference solution at 17:00, and then have Q&A.



Example Solution

Example Solution

A load balancer.

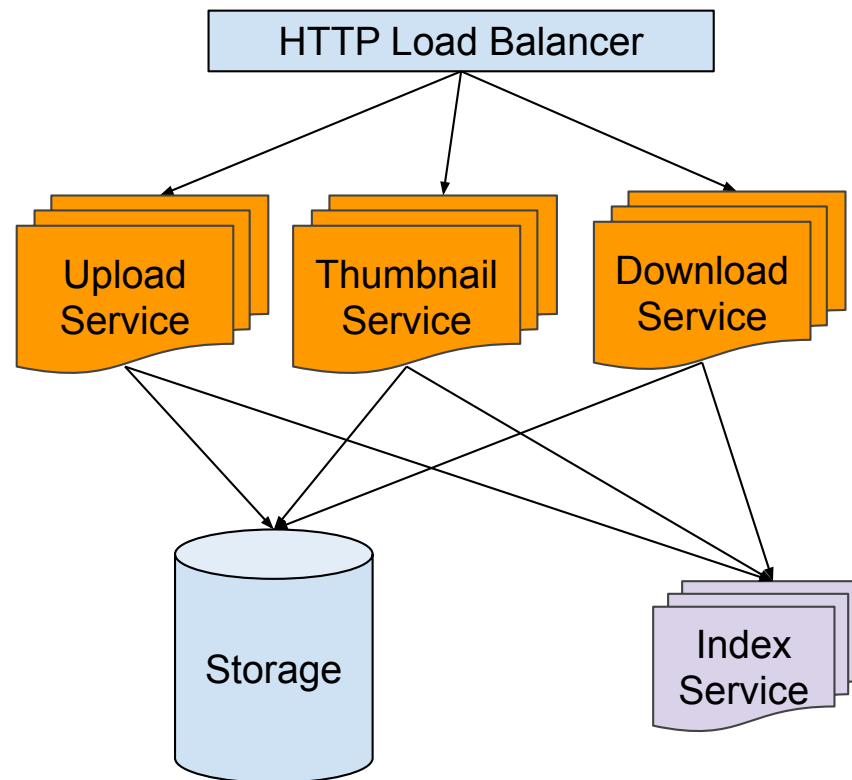
Dedicated web servers.

An index service.

A storage service.



SRE Classroom



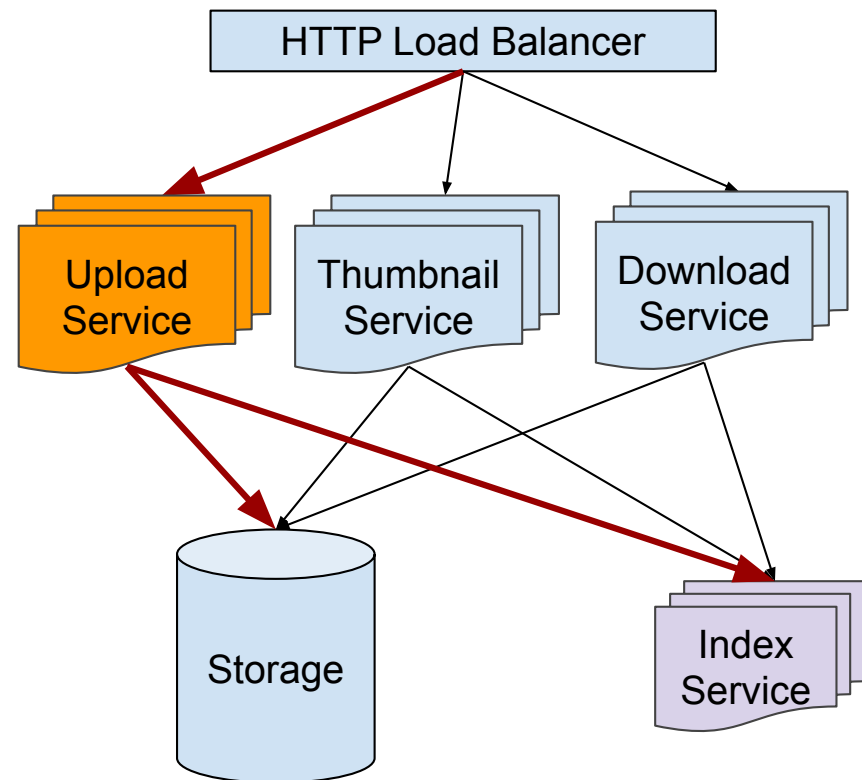
Uploading Pictures



Send query to upload service.

Pipe data to storage.

Update index service.



Searching



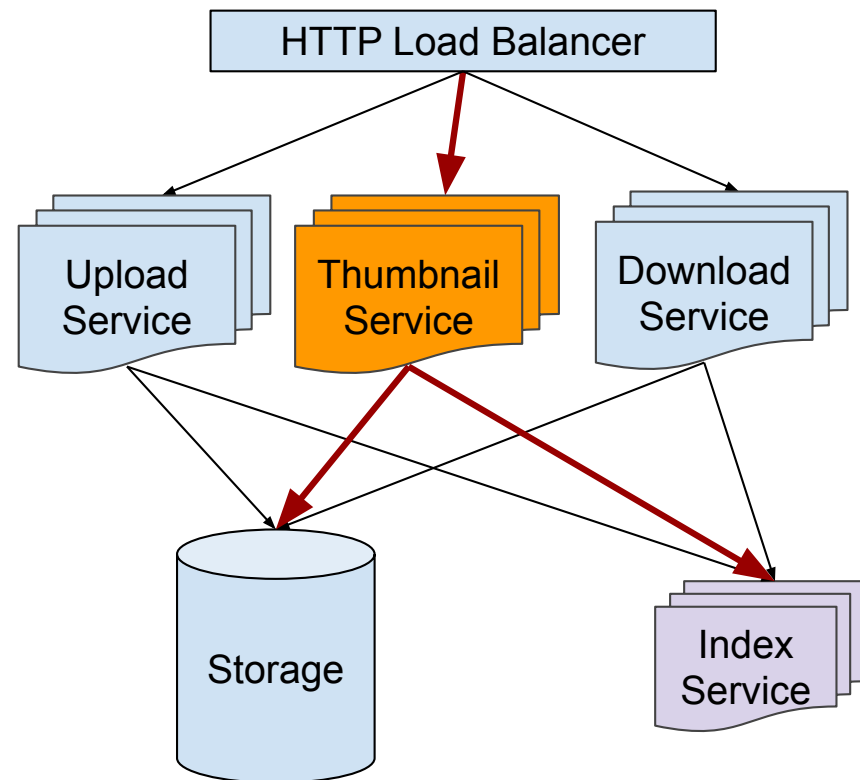
Send query to thumbnail service.

Or query index service.

Retrieve thumbnails from storage.

Return thumbnail page.

Serve cached data if possible.





Downloading a Full Size Picture

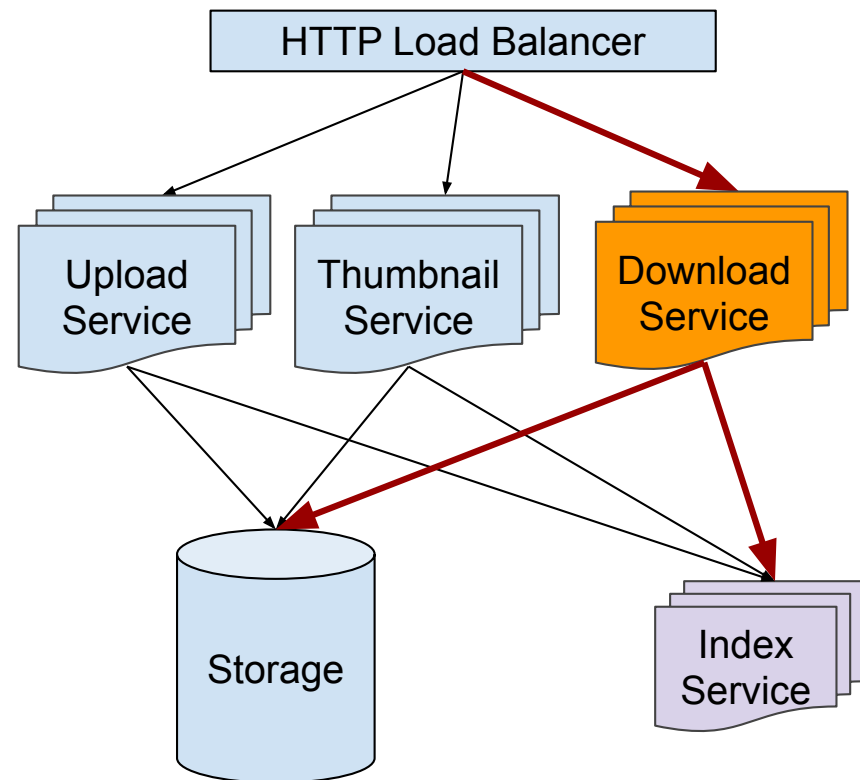
Send query to download service.

Or query index service.

Retrieve picture blob from storage.

Serve picture data.

Serve cached data if possible.



Index Service

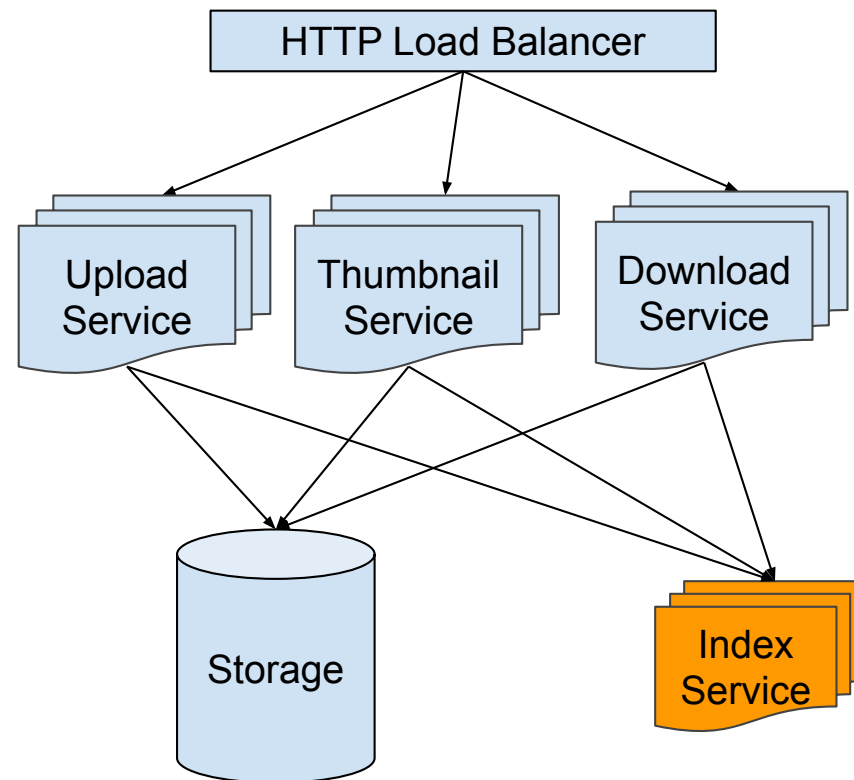
Stores metadata (tags, storage ID).

Assume metadata size of 8 KB per picture.

Globally consistent view.



SRE Classroom



Example Solution - Bill of Materials

General QPS and Storage



1.000.000 daily users with 50 uploads each

$1.000.000 * 50 / 24h / 3600s \approx \mathbf{579 \text{ qps}}$

Data uploaded per day: $1.000.000 * 50 * 4MB = \mathbf{191 \text{ TB}}$

For 30 days: $191 \text{ TB} * 30 = \mathbf{5.59 \text{ PB}}$

Thumbnail pages - assume 5 searches per user organically.

$1.000.000 * 5 / 24h / 3600s \approx \mathbf{58 \text{ qps}}$

Assume each search results in one full size download as well $\approx \mathbf{58 \text{ qps}}$

Upload Service - Bandwidth



Incoming picture data per second:

$579 \text{ QPS} * 4 \text{ MB} = 2314 \text{ MB/s} \rightarrow 2314 / 1000 \text{ MB/s} \rightarrow 2.314 \text{ NICs}$

Theoretical maximum for data streaming to storage including handling incoming:

$1000 \text{ MB/s} / 4\text{MB} = 250 \text{ QPS per NIC}$

That means we need 3 servers worth of NIC.

Upload Service - Timing



RPCs to storage backend are streamed in parallel.

No user-facing SLO for storing.

Storage write starts within 200 ms. Network time per picture 8 ms.

Actual blocking time on NIC is 8 ms -> 125 qps

Need 5 servers if we completely rely on remote storage.

Thumbnail Service - Bandwidth



58 qps incoming for thumbnail pages. 10 pics of 256 KB per result page.

$58 \text{ qps} * 10 \text{ pictures} * 256 \text{ KB} \approx 145 \text{ MB/s}$

One NIC is sufficient.

Thumbnail Service - Timing

Storage reads start within 100 ms of RPC call.

Index service query ~0.5 ms

Network time per results page: 2.5 ms

SLO wait time per query 103 ms at 58 QPS.

Actual blocking time on NIC is 2.5 ms -> 400 qps

Need 1 server if we completely rely on remote storage.



Download Service - Bandwidth

58 qps incoming for full size pics. Each pic is 4 MB.

$58 \text{ qps} * 4 \text{ MB} \approx 232 \text{ MB/s}$

One NIC is sufficient.



Download Service - Timing

Storage reads start within 100 ms of RPC call.

Network time per results page: 4 ms

Index service query ~ 0.5 ms

SLO wait time per query ~105 ms.

Actual blocking time on NIC is 4 ms -> 250 qps

Need 1 server if we completely rely on remote storage.



Index Service - Bandwidth

Metadata per photo is 8 KB.

579 qps from upload service containing single metadata.

$579 \text{ qps} * 8 \text{ KB payload} \approx 4.5 \text{ MB/s}$

Search queries return metadata for 10 entries:

$58 \text{ qps} * 10 * 8 \text{ KB reply} \approx 4.5 \text{ MB/s}$

One NIC can easily do this.



Index Service - Storage



1.000.000 daily users * 50 pics * 30 days = 1.500.000.000 entries

1.500.000.000 entries * 8 KB \approx 11.2 TB

3 HDD servers or 6 SSD servers.

SSD more future proof for I/O, but HDD are cheaper.

Index Service - I/O Timing



~579 write qps result in ~579 IOPS

~58 search qps result in ~580 IOPS

HDD: Seek 10 ms + read/write 8K seq 0.04 ms \approx 10.04 ms

~99 IOPS per disk -> **12 servers needed**

SSD: (random 4K read) 0.02 ms + (4K seq r/w) 0.00003 ms \approx 0.02 ms

~ 50K IOPS per disk -> **1 server needed**

Index Service - Write QPS



Remote storage is only eventually consistent.

Need a globally replicated search index.

Writes bounded by global consensus. 150 ms RTT = ~ 6 qps for writes.

579 write QPS = ~ 100 entries per batch

Since this is only metadata **one NIC can deal with this.**

Index Service - Sizing Needs



	HDD	SSD
Bandwidth	1	1
Storage	3	6
IOPS	12	1
Total	12	6

Load Balancer - Bandwidth



Load balancer is fully network bound and needs to be able to deal with all traffic.

Upload traffic: 579 QPS, ~2314 MB/s

Thumbnail traffic: 58 QPS, ~145 MB/s

Full size traffic: 58 QPS, ~232 MB/s

Total: 695 QPS, ~2691 MB/s

We need 3 servers.

Footprint - Summary



	Bandwidth	Timing	Storage
Load Balancer	3	n/a	n/a
Upload Service	3	5	n/a
Thumbnail Service	1	1	n/a
Download Service	1	1	n/a
Index Service	1	1	6 (SSD)
TOTAL per DC	16		

Footprint - Increased Robustness



Account for non-uniform load (multiply by 1.25).

Get rid of SPOF.

	Bandwidth	Timing	Storage
Load Balancer	4	n/a	n/a
Upload Service	3	6	n/a
Thumbnail Service	2	2	n/a
Download Service	2	2	n/a
Index Service	2	2	8 (SSD)
TOTAL per DC	22		

Global Footprint

22 servers per cluster.

3 DCs * 22 servers = 66 servers global footprint

In reality and given our architecture you can start with less and use auto-scaling for each component.



Improving with caches



We have unused resources on most servers. Use it for caches.
Performance of cache depends on access speed.

RAM: ~24 GB/s

SSD: ~1 GB/s, HDD: ~200 MB/s

Network: ~1 GB/s

Caches can be distributed but then have added bandwidth costs.

So Many Unexplored Angles

Monitoring and alerting

Realistic storage backend behavior (degradation)

Long-term storage (hot and cold storage)

Capacity growth (per year, retention)

Privacy requirements (GDPR anyone?)

Toil (rollout, maintenance) - more a process thing though





Wrap-up



Questions?

Questions?

Thank you!

To our facilitators in the room!

salim@google.com



Further reading

<https://goo.gl/LGea4R>

