# System Design Interview Cheat Sheet

## Interview Framework

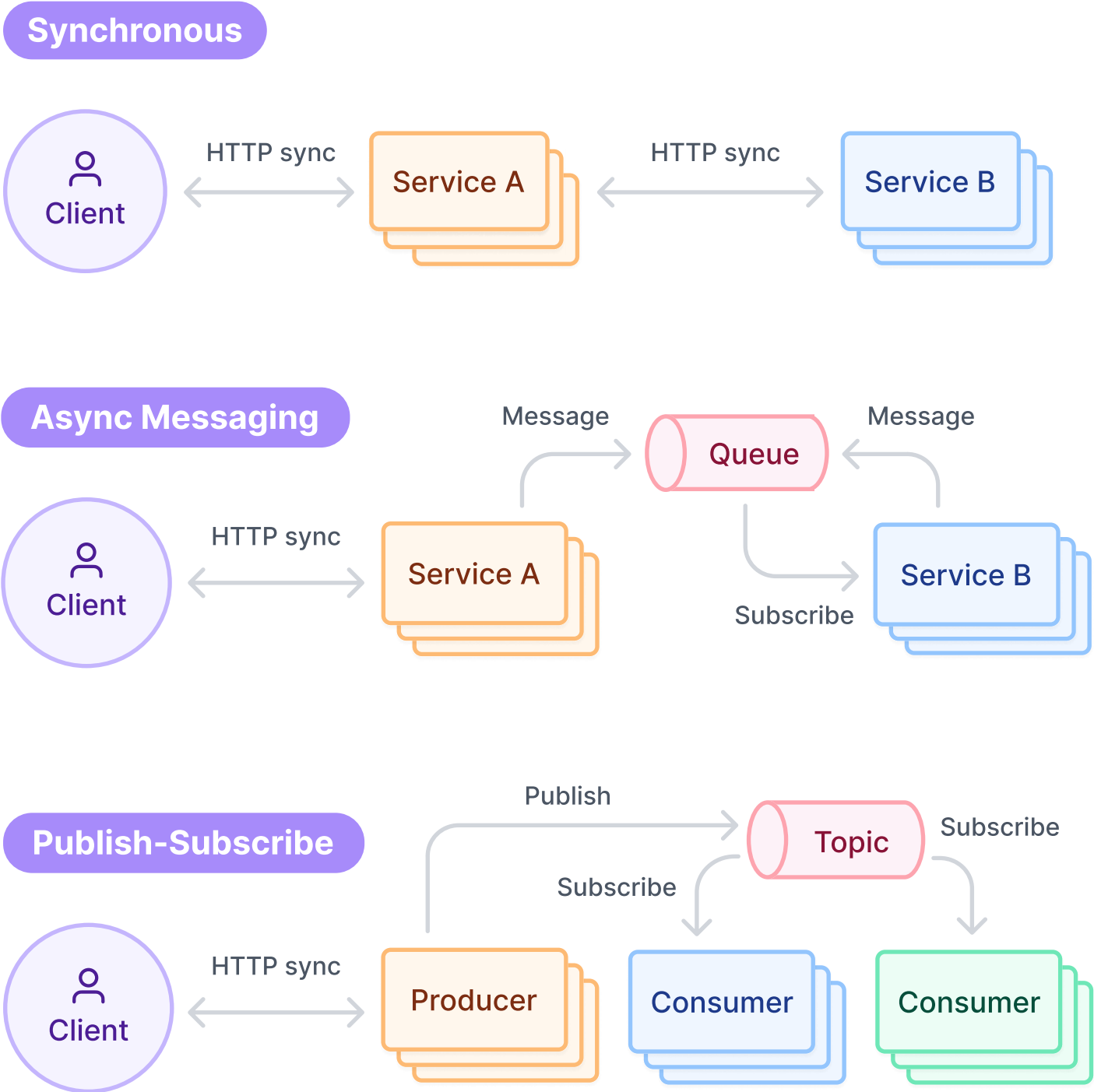| Step 1 — 10 min. | Step 2 — 10 min. | Step 3 — 10 min. | Step 4 — 10 min. | Step 5 — 5 min. |
|---|---|---|---|---|
| **Understand the Problem** Gather more information about the system requirements and constraints. | **High-Level Design** Explain how each part of the system works together. Start by defining APIs. They are the foundation of the architecture. | **Deep-Dive** Examine system components in detail. Your interviewer may pick a specific area or ask you what you'd like to explore. | **Improve the Design** Take a step back. What are the bottlenecks? How does it scale? | **Wrap Up** Summarize the requirements, justify your decisions, suggest alternatives, and answer any questions. |

## API Design Choices

Explain how each part of the system works together. Start by defining APIs and the overall design patterns that your application will use.

| | REST | RPC | GraphQL |
|---|---|---|---|
| **Properties** | ✓ resource-oriented<br>✓ data-driven<br>✓ flexible | ✓ action-oriented<br>✓ high performance | ✓ single endpoint<br>✓ strongly-typed requests<br>✓ no data overfetching<br>✓ self-documenting |
| **Data** | JSON, XML, YAML, HTML, plain text | JSON, XML, Thrift, Protobuf, FlatButters | JSON |
| **Use cases** | ✓ web-based apps<br>✓ cloud apps<br>✓ client-server apps<br>✓ cloud computing services<br>✓ developer APIs | ✓ complex microservices system<br>✓ IoT applications | ✓ high-performance mobile apps<br>✓ complex systems and microservice-based architectures |

**Synchronous**

Client ↔ HTTP sync ↔ Service A ↔ HTTP sync ↔ Service B

**Async Messaging**

Client ↔ HTTP sync → Service A → Message → Queue → Message → Service B (Subscribe)

**Publish-Subscribe**

Client ↔ HTTP sync → Producer → Publish → Topic → Subscribe → Consumer / Consumer
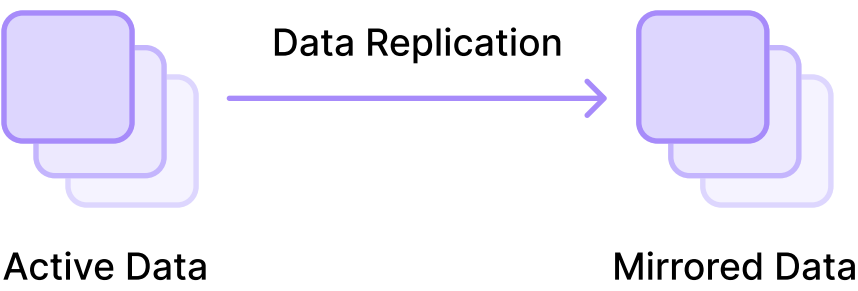
## Scalability

Consider the scale of your system. How many users and requests will the server support? What happens with increased demand?
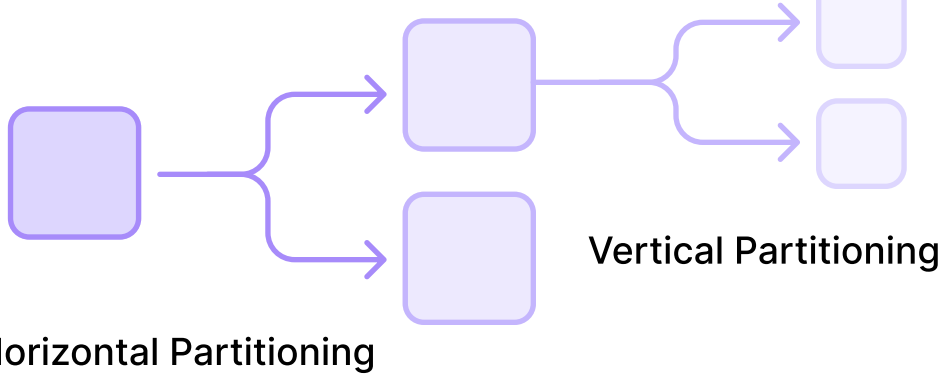
### Replication

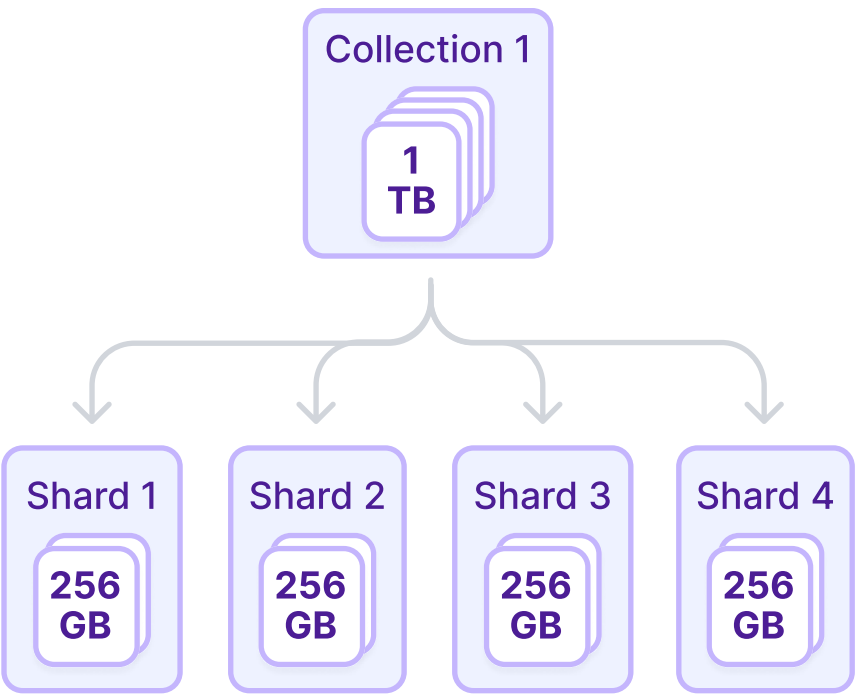Is the data important enough to make copies? How important is it to keep all copies the same?

Active Data → Data Replication → Mirrored Data

### Partitioning

Partitions contain a subset of the whole table. Each partition is stored on a separate server.

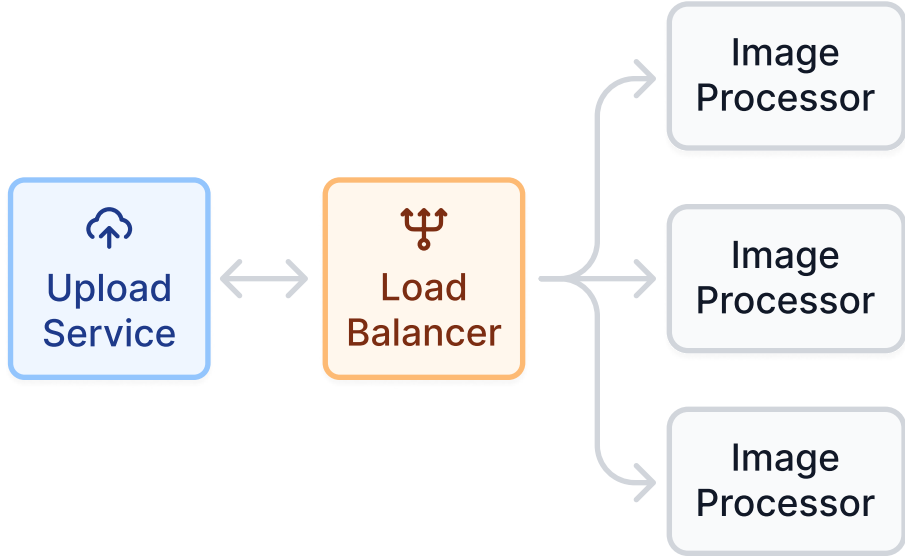Horizontal Partitioning / Vertical Partitioning

### Sharding

Sharding allows a system to scale as data increases, but not all data is suitable for sharding.

Collection 1 — 1 TB
- Shard 1 — 256 GB
- Shard 2 — 256 GB
- Shard 3 — 256 GB
- Shard 4 — 256 GB

### Load Balancing

Load balancing distributes incoming traffic across multiple servers or resources.

Upload Service → Load Balancer → Image Processor / Image Processor / Image Processor

## Caching

| In-memory Cache | Distributed Cache |
|---|---|
| ✓ **Latency** - in-memory cache is faster because it doesn't require a network request like distributed. | ✓ **Sharing data / Consistency** - data can be shared across machines with a distributed cache.<br>✓ **Availability** - distributed cache is not affected by individual server failures. |

- No. items
- Cache Miss & Hit
- Disk & Memory Usage

- Write-Through
- Read-Through
- Write-Around
- Write-Back

**Popular caches:**
- In-memory
- Redis
- Memcached
- AWS Elasticache
- GCP Memorystore

**Eviction:**
- LRU (Least Recently used)
- LFU (Least Freq. used)
- FIFO
- MRU
- Random Eviction
- Least Used
- On-Demand Expiration
- Garbage Collection

- Storing user sessions
- Communication between microservices
- Caching frequent database lookups