

System Design @ Scale

Timeout, Retries And Backoff

"Everything fails, all the time."

Motivation → It will not be possible to avoid failures by 100%, strive for systems to have mechanism in place to handle issues and minimise impact.

Increasing System Availability

① Timeouts

↳ add timeout to any remote call.
connection & request timeout.

what should be timeout value?

not too small not too large

relative terms & totally dependent on system.

① resource consumption ↑

② Backend Traffic & latency ↑

② Too many retries → complete outage.

How to choose?

① p99 latency of dependent services.

↳ consider network latency as well.

↳ latency is closer to p50 → add small padding.

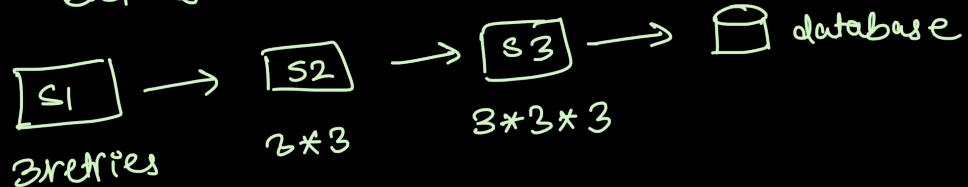
{ In short, careful analysis of system is necessary to come up with timeout values >

② Retries and Backoff

To many retries can cause system overload. though retries are helpful in resolving Transient issues.

Potential problems

- ① Retries at multiple layers in distributed system can cause extensive load on end system.



- ② Circuit Breakers → Cells to downstream is stopped if threshold exceeds.

↳ difficult to test & adds time to recovery
↳ better solution could be token buckets.

- ③ When to Retry? Idempotent API are safe to retry.

API with any side-effects aren't.

↳ which failures to retry for? → net client errors
→ server errors maybe

Backoff → Retries should be executed with Backoff
wait for some time then retry

exponential backoff → wait time increased exponentially after every attempt.
till some maximum attempt.

