

Stop Guessing, GET Requesting

Back of the envelope estimations to get started with.

Scheduled Sync:



-> Request -> Transform -> Stash in DB



-> Request

-> Transform

-> Stash in DB

Questions to answer:

- * Concurrent requests?
- * Single process?
- * Bottlenecks?

Info to collect and estimates to make

- * How many requests?
- * avg response time?
- * origin limits?
- * Sync frequency?
- * DB limitations?

How?

- * Service Docs
- * Existing metrics
- * Comparables
- * Simulation
- * “Guesstimates”

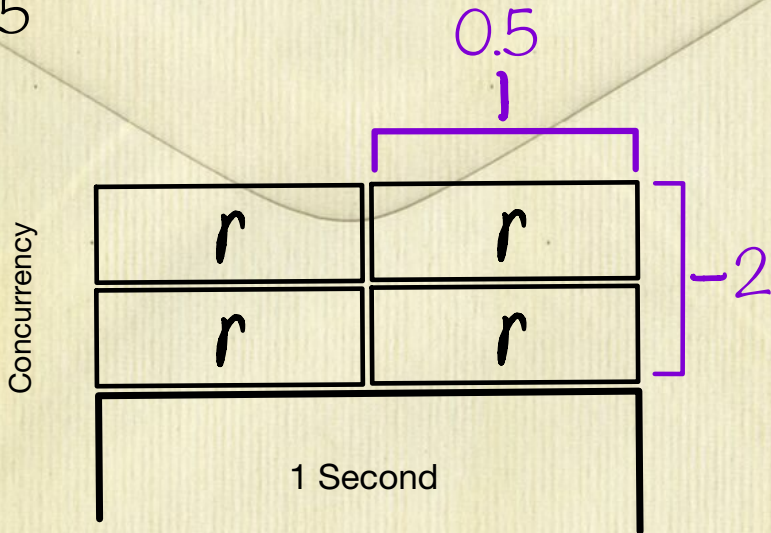
Request Considerations

- * Stay under rate limit
- * Don't kill the target service

Little's Law

$$C = RL$$

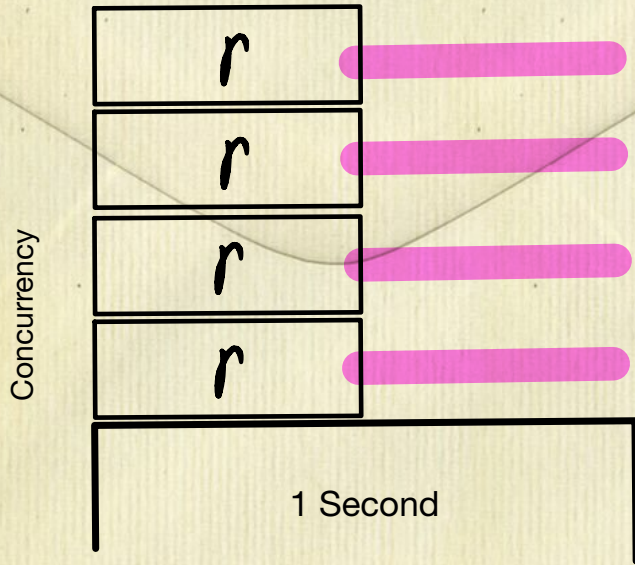
$$2 = 4 * 0.5$$



Little's Law

$$C = RL$$

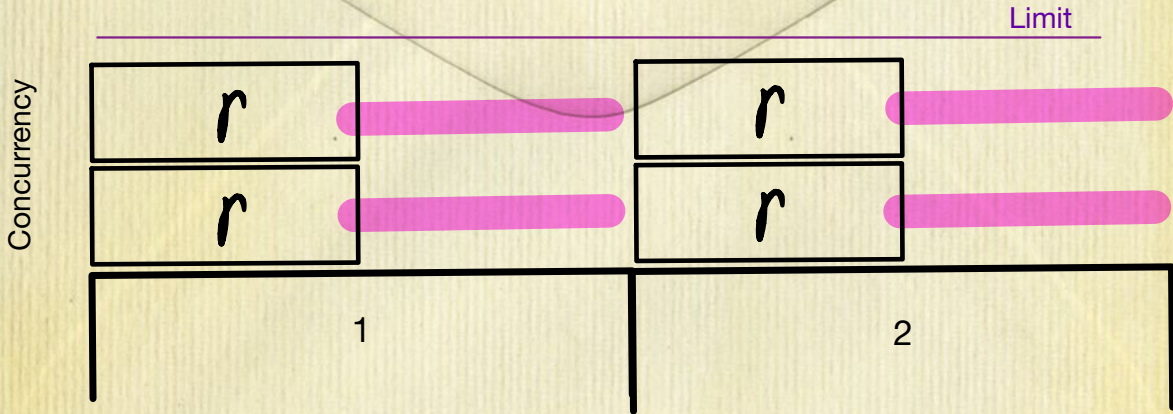
$$4 \uparrow = 4 * 1 \uparrow$$



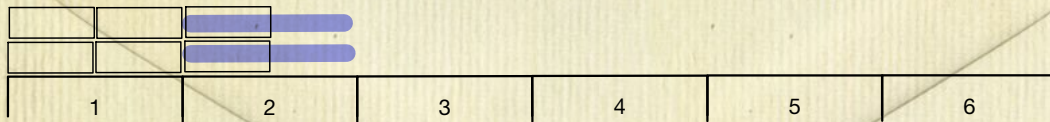
Little's Law

$$C = RL$$

$$2 = 2 \downarrow^* 1 \uparrow$$



Client

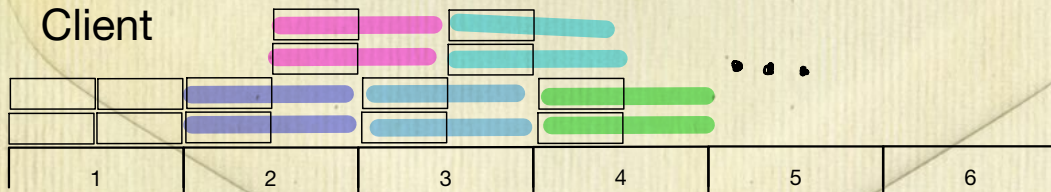


Origin

Hiccup
↓



Client



Origin

Hiccup
↓

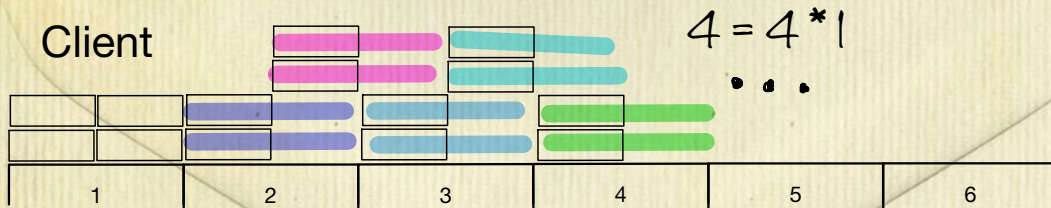
Max Workers



Origin Backlog



Client



Origin

Hiccup
↓

$$2 = 4 * 0.5$$

Max Workers



Origin Backlog

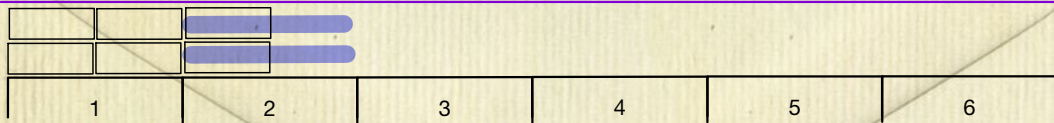
$$2 = 4 * 0.5$$



Client



Limit



Origin

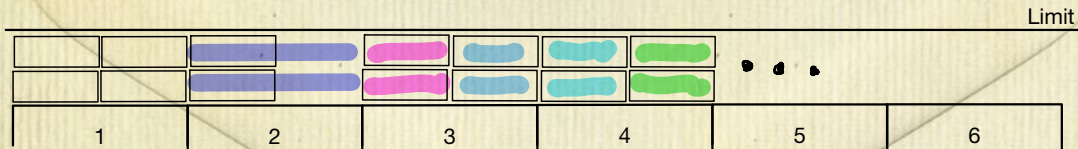
Hiccup



Max Workers

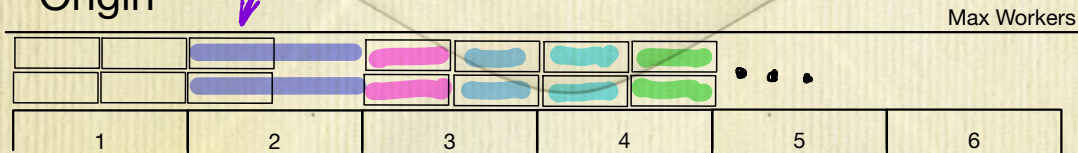


Client

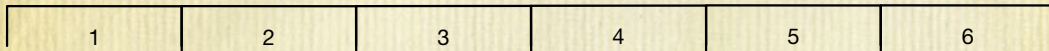


Origin

Hiccup
↓



Origin Backlog



Client

$$2 = 4 * 0.5$$

Limit



Origin

Hiccup
↓

$$2 = 4 * 0.5$$

Max Workers

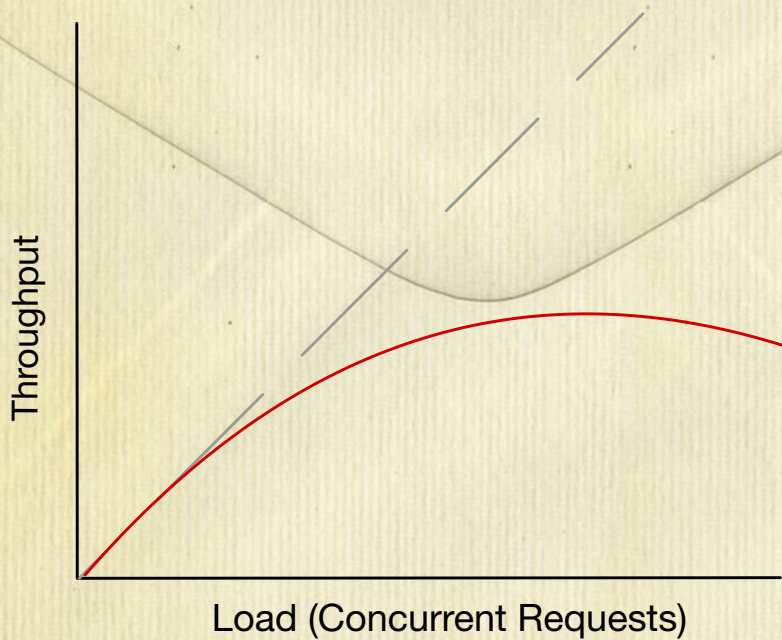


Origin Backlog

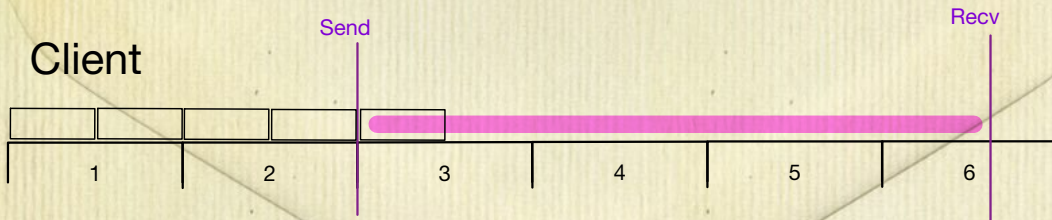
$$0 = 0 * 0$$



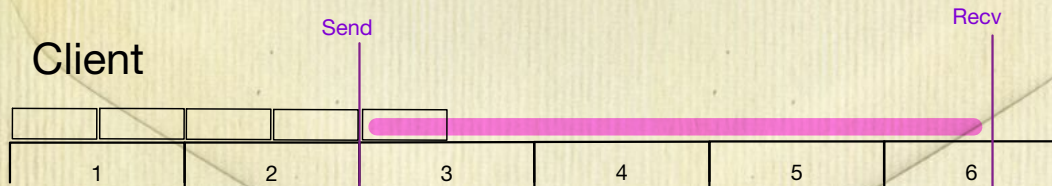
Why we care



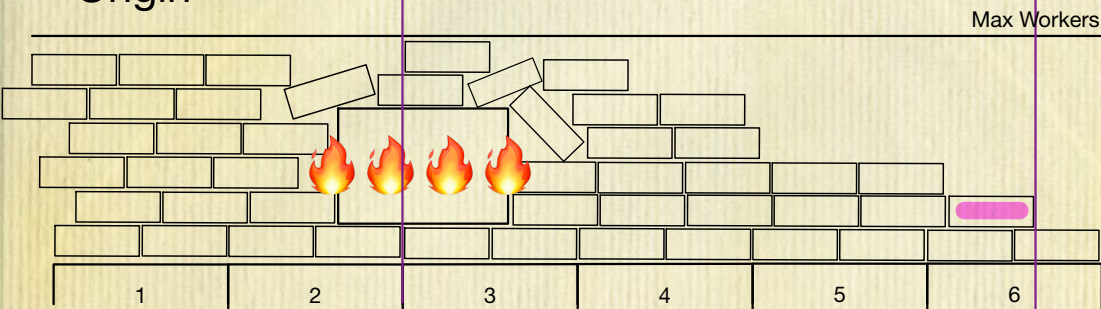
Client



Client



Origin

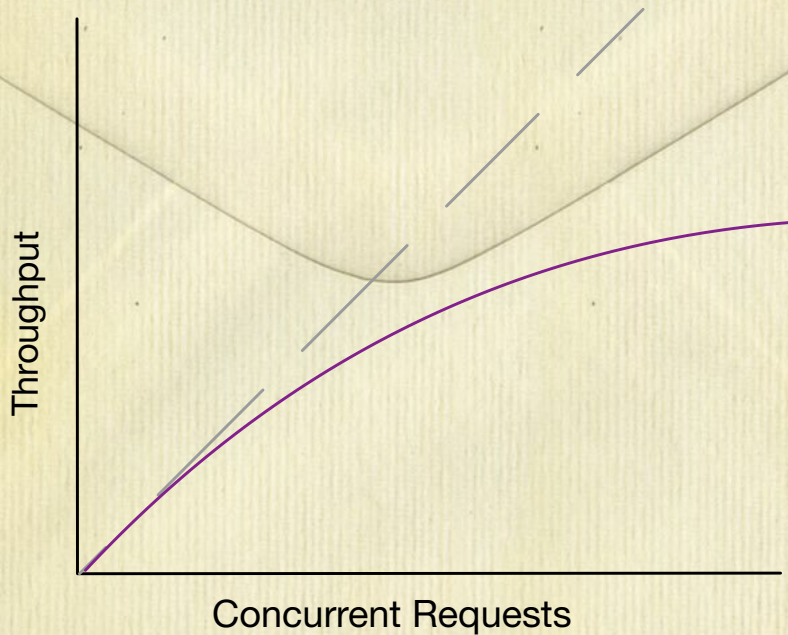


Back to our Scheduled Sync Task

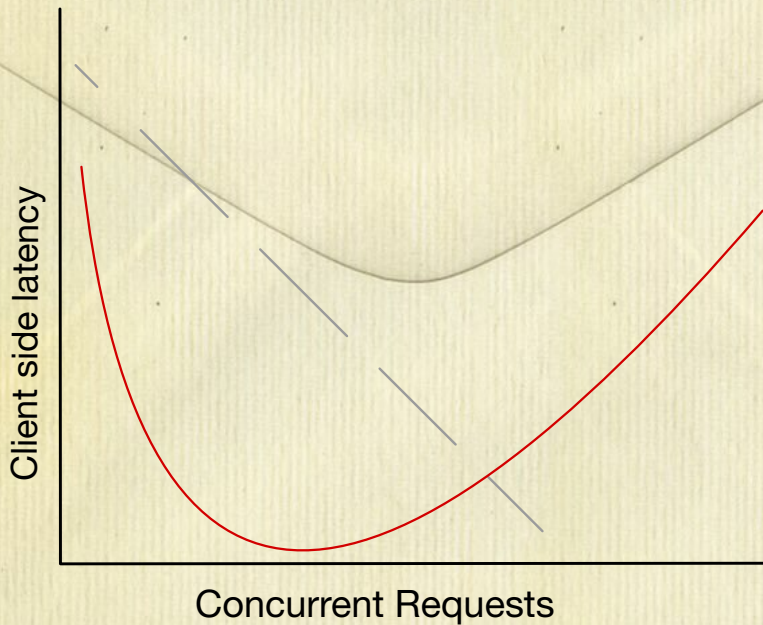
1,800,000 requests to be made.

$$1,000 * 0.5 = 500 \quad \leftarrow \text{Little's Law}$$

...but, diminishing returns



...because



We need more GIL's

$$1 / 0.01 = 100 \text{ p/s}$$

$$100 < 1,000$$

Amdahl's Law

G: processing time averaged per request

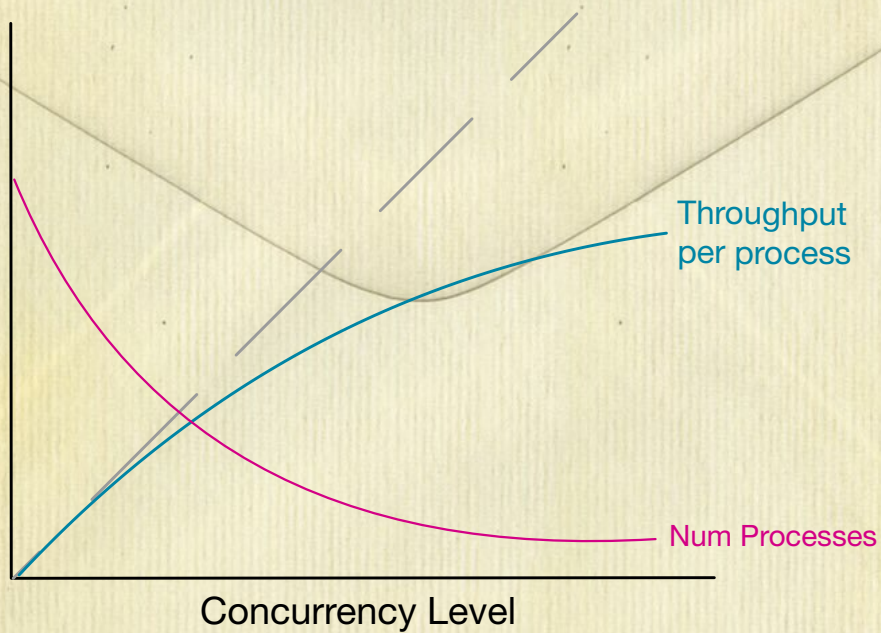
T: target rate

$$R(C) = 1 / (G + (L / C))$$

Calculating number of processes needed.

$$P(C, T) = T / R(C)$$

Amdahl's Law



$R(50) \rightarrow 50 \text{ requests/s}$

$$1,000 / 50 = 20$$

What about the DB?

100 rows per request
0.5 response time.

Client A



Client B



Client C



Client A



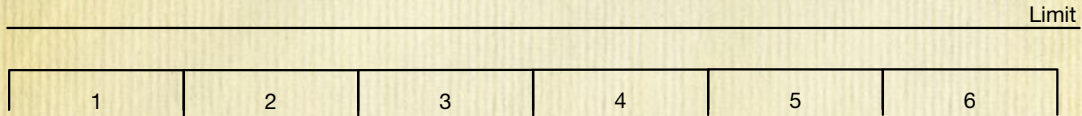
Client B



Client C



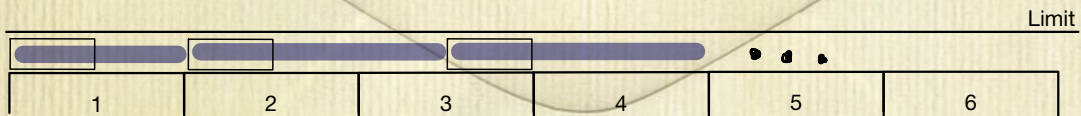
DB



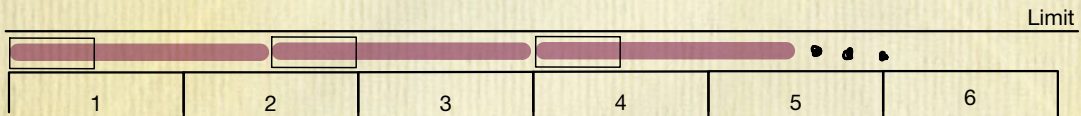
Client A



Client B



Client C



DB



Client A

$$l = 2/3 * 1.5$$

Limit



Client B

$$l = 2/3 * 1.5$$

Limit



Client C

$$l = 2/3 * 1.5$$

Limit



DB

$$l = 2 * 0.5$$

Limit



$$1 = 2 * 0.5$$

$$2 * 100 = 200 \text{ rows/second}$$

$$200 < 1,000$$

The DB is a bottleneck

Do we care? Enough?

$$1,800,000 / 200 / 60 / 60 = \sim 2.5 \text{ hrs}$$

Other considerations

- * DB connection limits (e.g. PostgreSQL sans PGBouncer)

What if tasks overlap?

Things to note

- * Only as good as the inputs
- * Model a distribution of possible inputs
- * Just one example.

Key Takeaway

1. Resource contention -> queuing -> Queuing theory
2. Metrics and models

GitHub: steven-cutting

References:

- * "Stop Rate Limiting! Capacity Management Done Right" by Jon Moore, StrangeLoop 2017
- * Amdahl's Law, Jakob Jenkov, jencov.com
- * Probability and Statistics with Reliability, Queuing, and Computer Science Applications, Kishor S. Trivedi