



BUILDING FOR THE ROAD AHEAD

# DETROIT 2022

**DETROIT 2022**



**KubeCon**



**CloudNativeCon**

North America 2022

BUILDING FOR THE ROAD AHEAD

**DETROIT 2022**

**October 24-28, 2022**



**Andrew Newdigate**

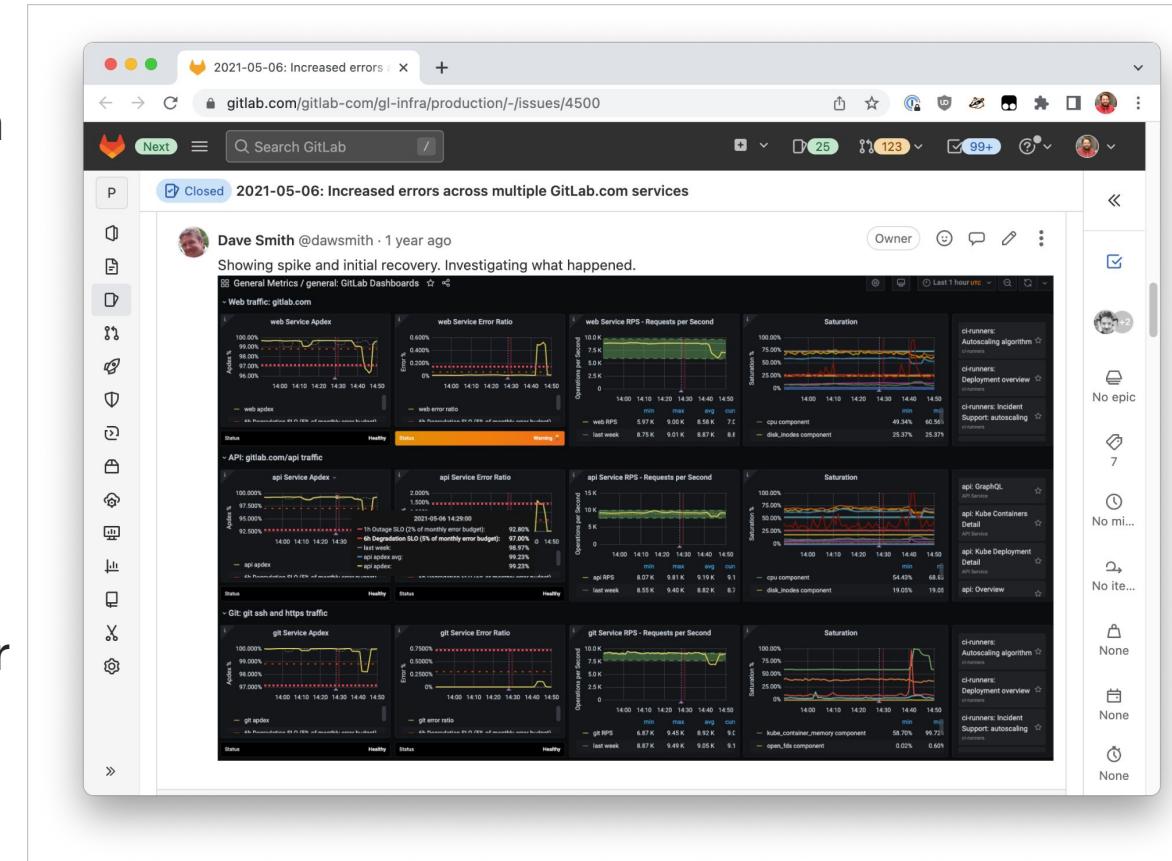
Distinguished Engineer, Infrastructure  
*GitLab, Inc.*



[twitter.com/suprememoocow](https://twitter.com/suprememoocow)

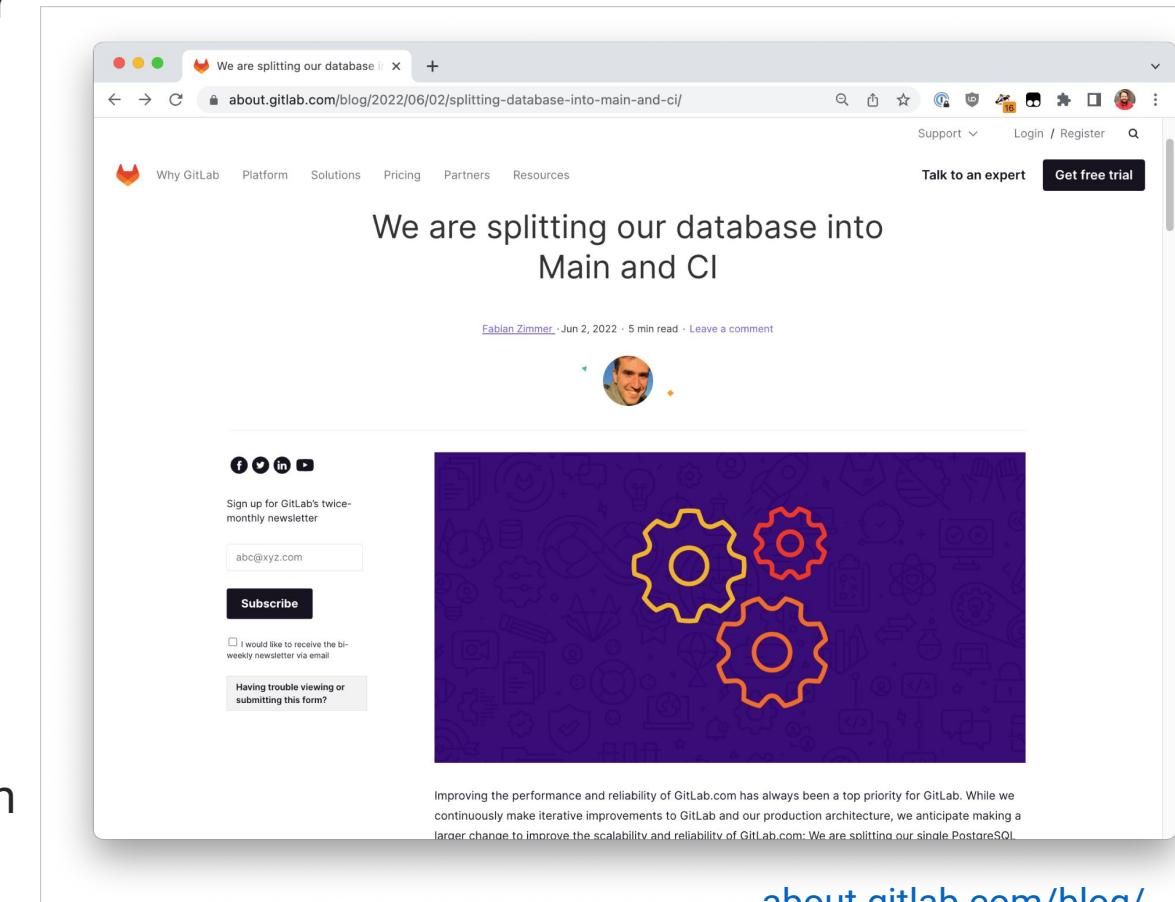
# A series of Postgres saturation events..

- In **mid-2021**, GitLab.com experienced a series of database incidents, related to high resource utilization on our primary Postgres database.
- Inside the company, concerns were raised about growth rate of GitLab.com and **the ability for our primary database to support this growth**.
- There were concerns that **the database would hit maximum capacity**, leading to outages due to further saturation.



# It's probably time we break up, Postgres

- Path forward would be split our main Postgres cluster into two separate clusters, through a process known as **functional decomposition**.
- This would be a **complicated software operation**, which would need to be prepared and performed carefully to avoid data-loss or outages.
- **Rushing the migration would be risky.**
- At the same time, **taking too long would be disastrous** if the Postgres server hit capacity and was no longer able to deal with further growth in traffic.



[about.gitlab.com/blog/](https://about.gitlab.com/blog/)

# Is there enough runway?

- It was important to know whether the **timeframes for the functional decomposition project were within the available capacity timeframes** for the main Postgres cluster.
- To answer this, a **capacity planning review** for Postgres was performed.
- This was done with Tamland, **GitLab.com's capacity planning tool**.
- The review confirmed that we had enough runway left, allowing the team to focus on the migration.

The screenshot shows a Google Document window titled "Postgres Capacity Review". The document contains the following sections:

- Executive Summary**: Existing saturation metrics were analysed, and, after discussion with [redacted], [redacted], a set of additional metrics were added to our Postgres fleet.
- Assumptions regarding the Completion of Current Ongoing Projects**: This report highlights several known high risk areas, all of which are currently being addressed by ongoing projects. The assumption is that these ongoing projects continue at their current pace to completion. Specifically:

Type	Workstream	Reference Date	Estimated Completion
Env Alloc	Adult Primavera Kav Ministrations	2021-06-24T00:00:00Z	August 2021 / Mon

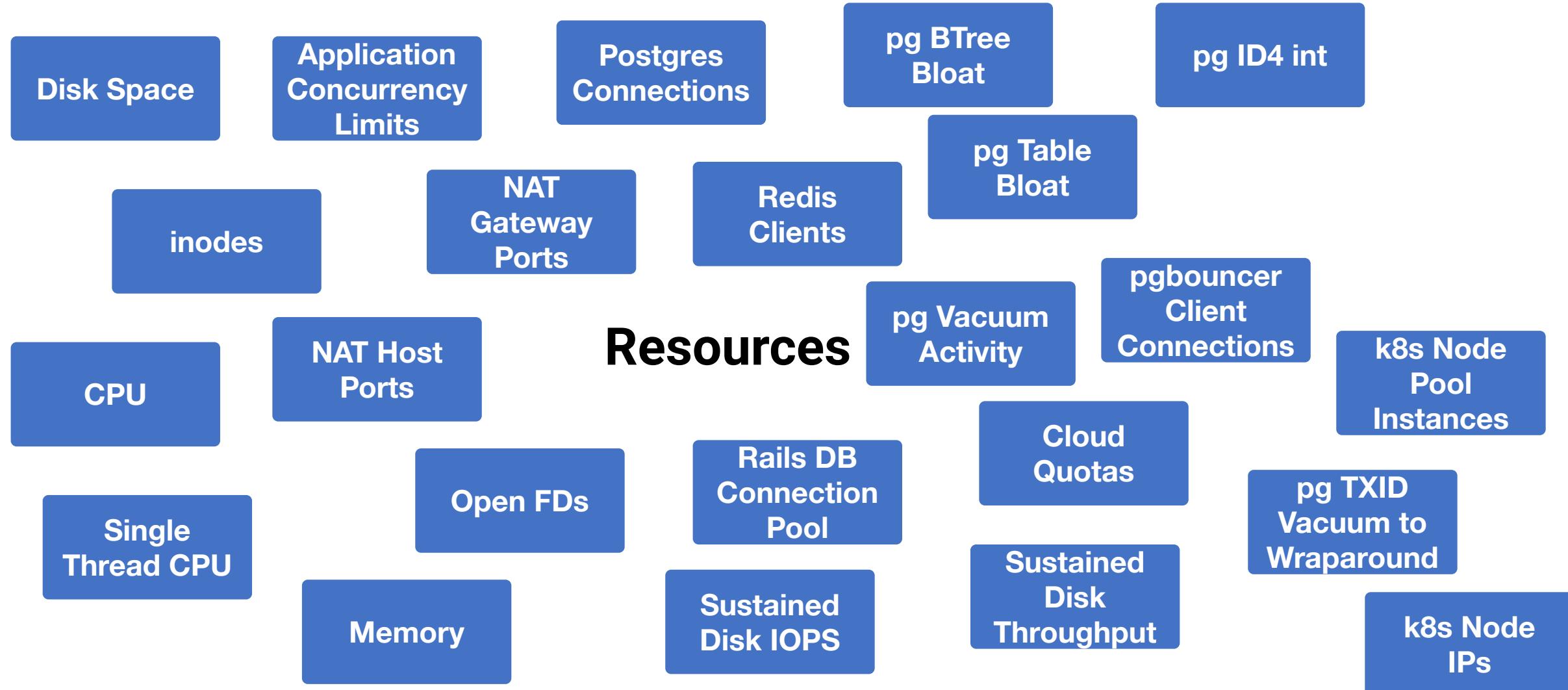
***“The goal of capacity management is controlling uncertainty. In the midst of the unknown, the service must be available now and continue to run in the future. A challenging but rewarding and delicate balance of tradeoffs is in play: efficiency vs. reliability, accuracy vs. complexity, and effort vs. benefit.”***

SRE Best Practices for Capacity Management, Luis Quesada Torres and Doug Colish

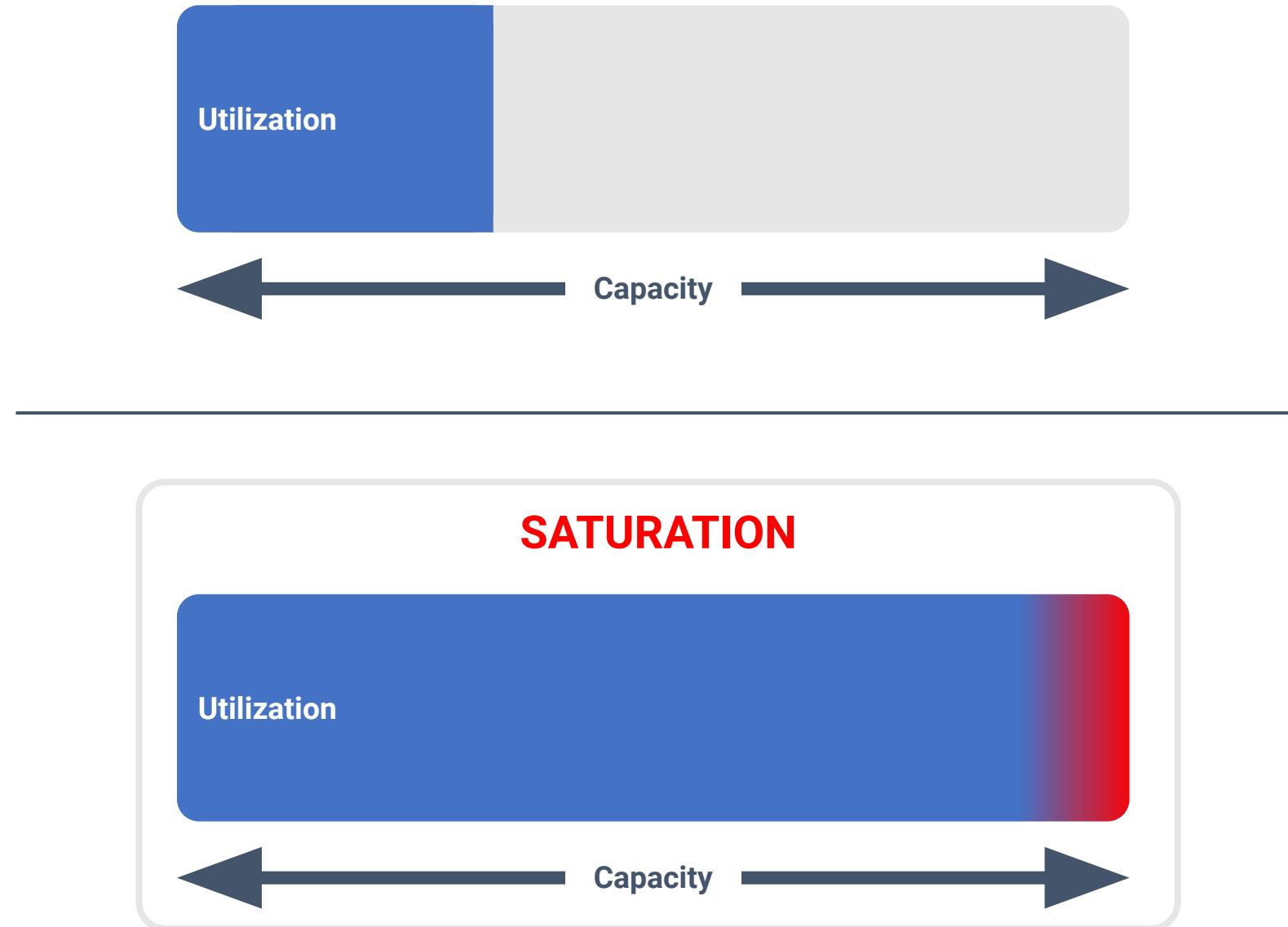
<https://research.google/pubs/pub50108/>

# First, Some Definitions

# Definitions: Resources



# Definitions



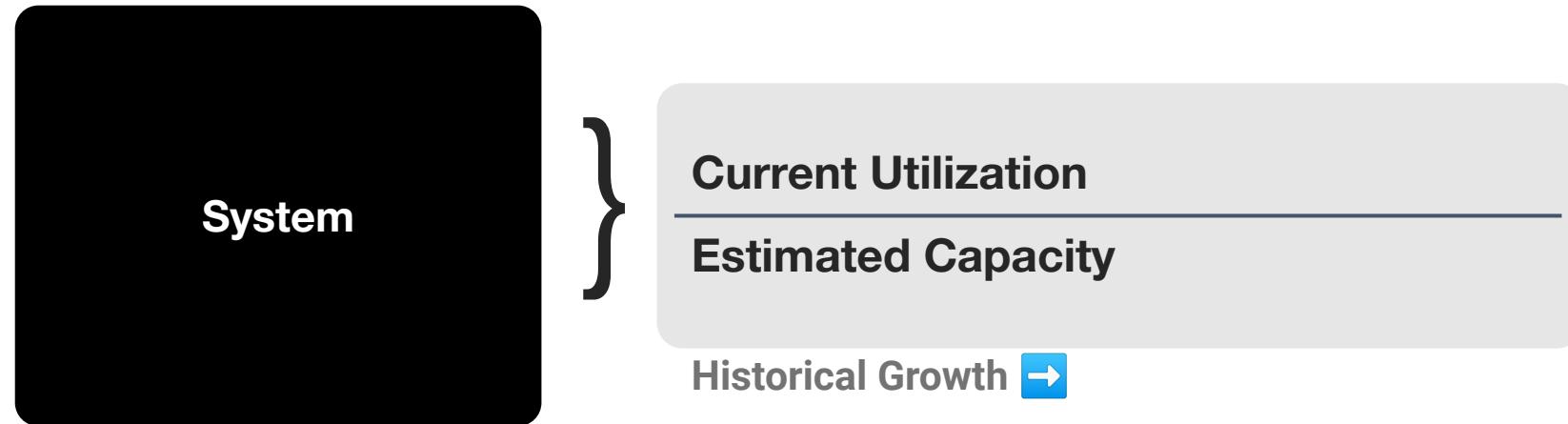
# Saturation Events vs. Mitigation Time

Resource	Saturation Failure Mode	Migitation	Mitigation Time
<b>CPU</b>	Queuing, Increased application latency	Increase maximum CPU capacity in cluster.	Minutes
<b>Memory</b>	OOM kills, transient errors	Increase maximum memory capacity in cluster.	Minutes
<b>Inodes / Disk Space</b>	Write failures, corruption	Increase disk capacity	Potentially hours
<b>Postgres Transaction ID Wraparound</b>	Database shutdown	Postgres Recovery	Hours, <b>days</b>
<b>Redis Single-Threaded CPU</b>	Increased application latency, errors	Refactor application to shard Redis, migrate to Redis cluster.	<b>Weeks, months</b>
<b>Postgres 32 bit ID Wraparound</b>	Database shutdown	Migrate from 32 bit IDs to 64 bit IDs	<b>Months</b>
<b>Postgres Primary CPU, no further vertical scaling possible</b>	Increased application latency, errors	Postgres Functional Decomposition	<b>Months</b>



# GitLab.com's approach to capacity planning...

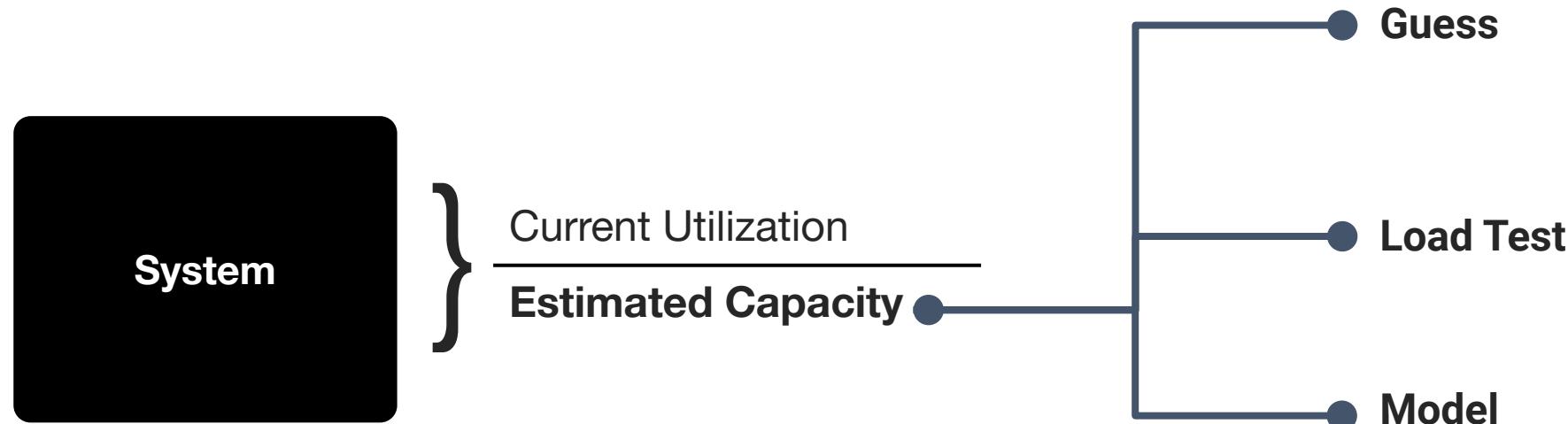
# System Capacity Planning



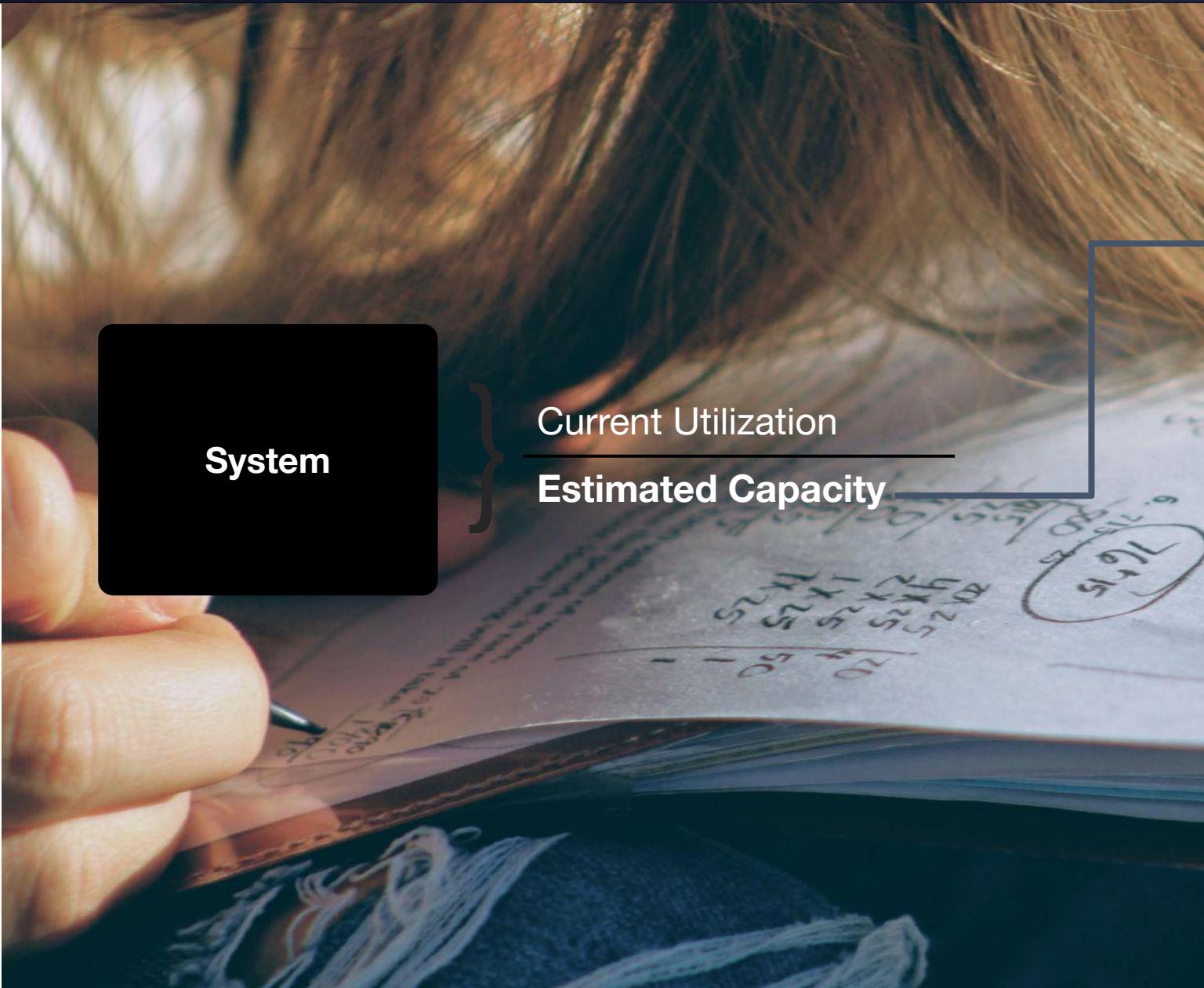
## Demand Signals

- Monthly Active Users
- Requests per Second
- Concurrent Users
- etc

# Estimating System Capacity is Hard



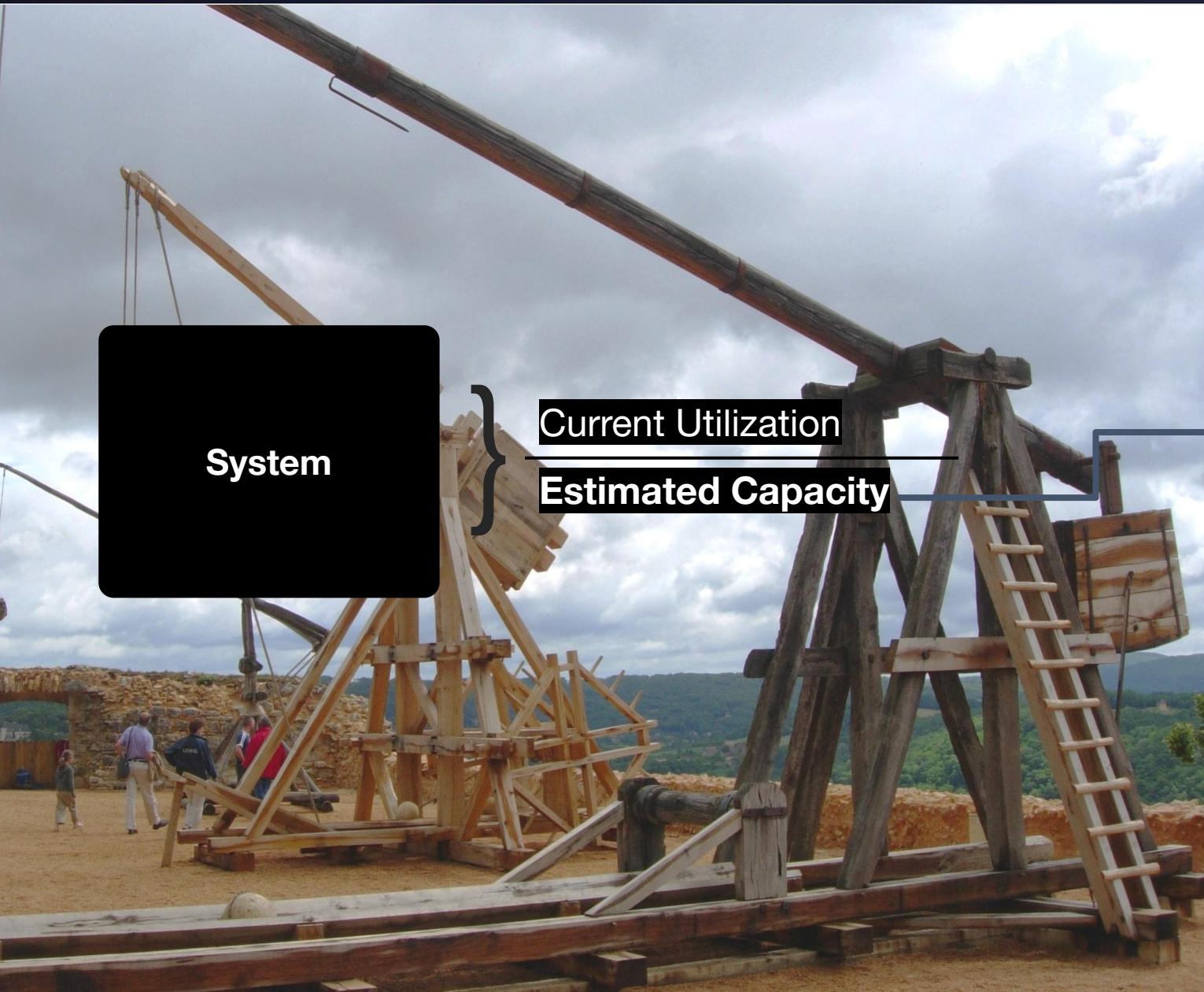
# Guessing Maximum Capacity



## Guess

- aka, the “Back of the Envelope” method
- Likely inaccurate
- Not scalable
- Not automatable

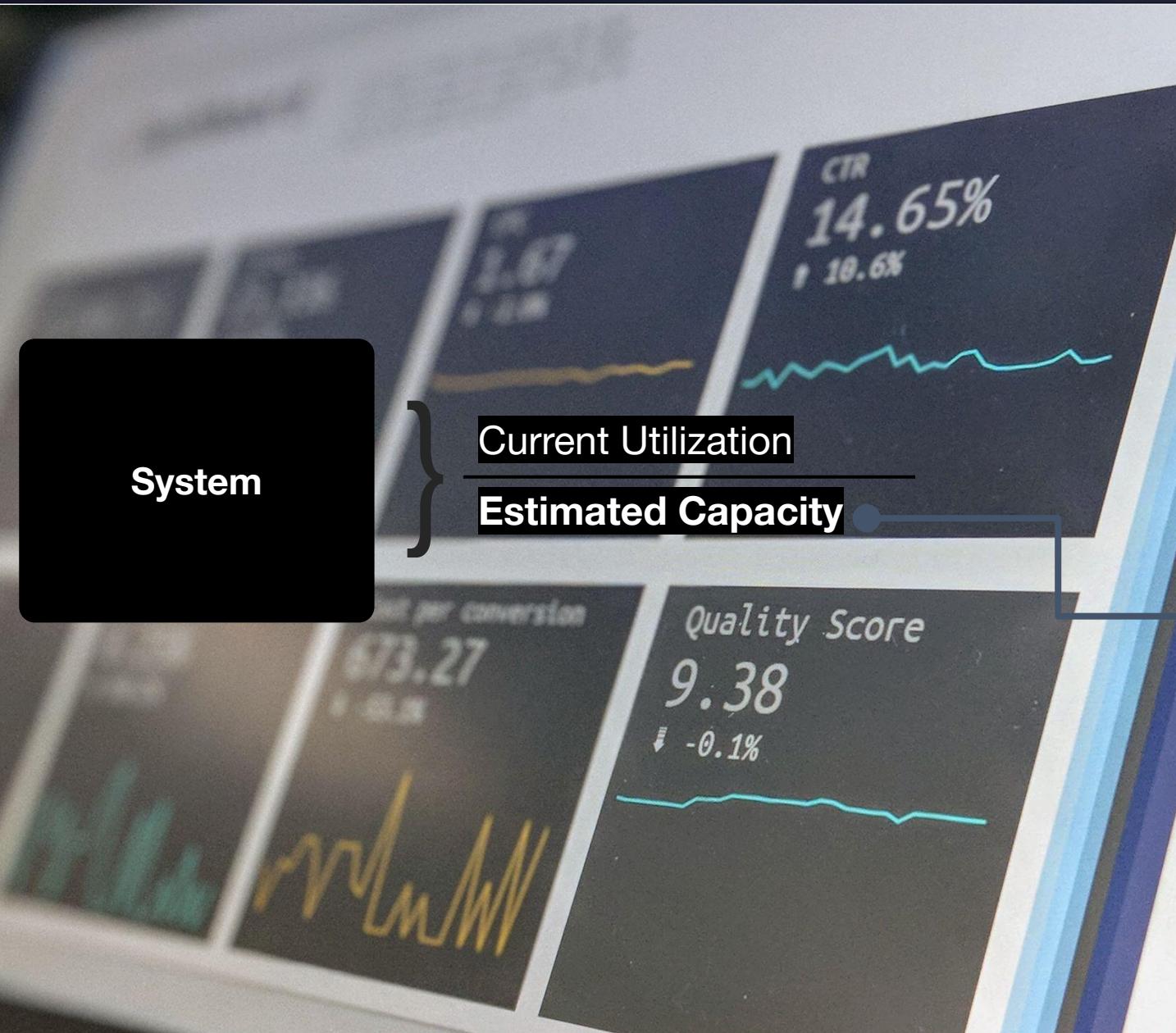
# Load Testing Maximum Capacity



## Load Test

- Setting up an equivalent testbed can be time-consuming and expensive.
- May miss important failure modes
- Difficult to guess correct workloads
- New features require that testing be redone.

# Modelling Maximum Capacity



# Estimating System Capacity is Hard

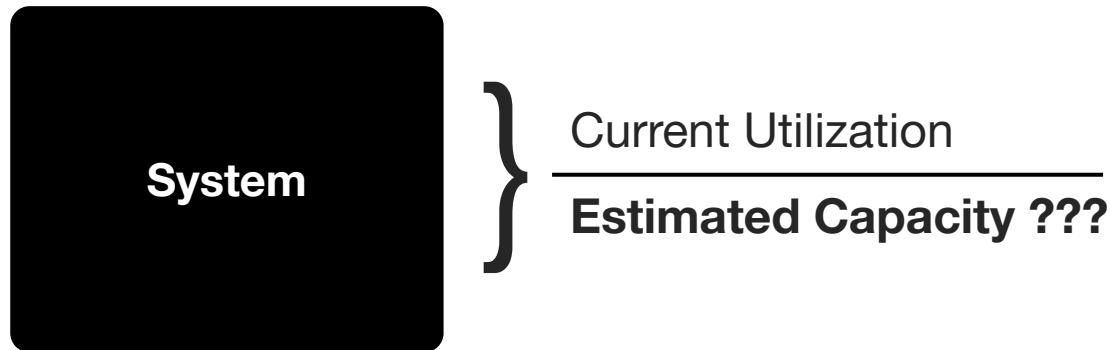
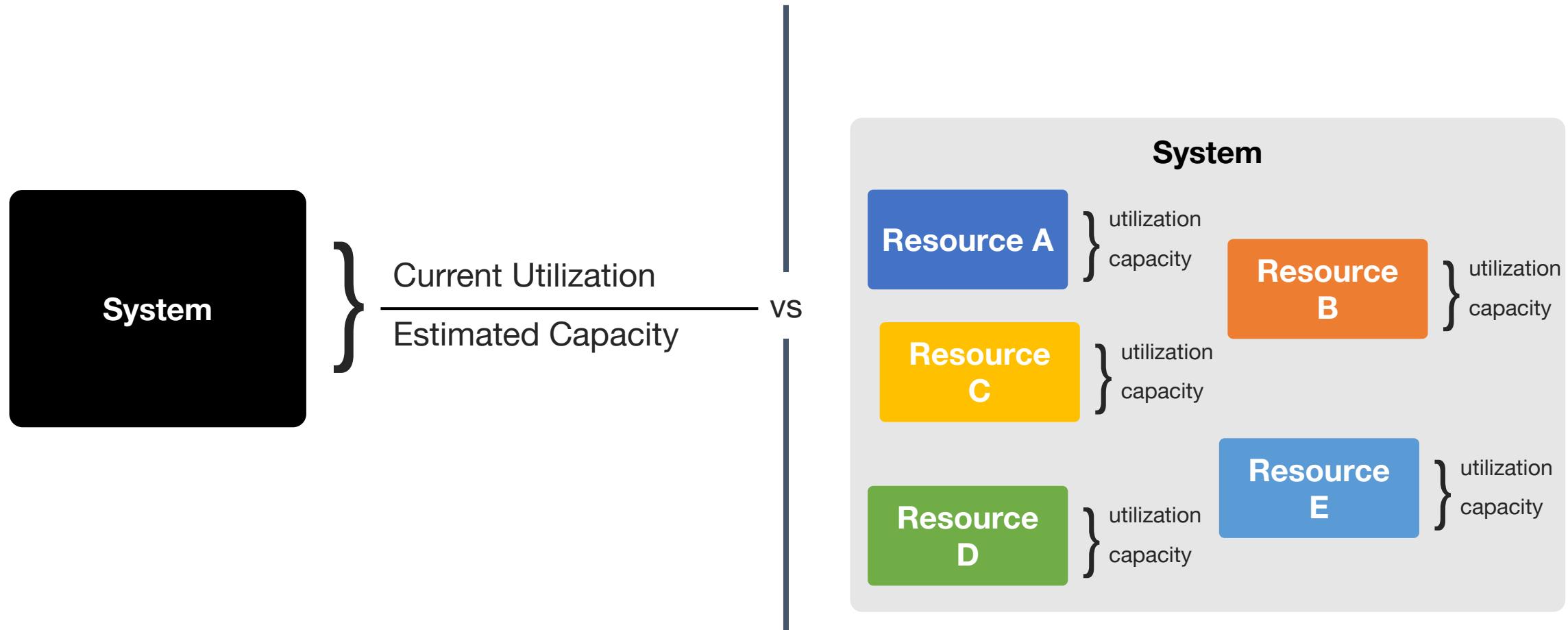


Photo by [Jon Tyson](#) on [Unsplash](#)

# Resource Capacity, not System Capacity

Instead of treating system capacity as a blackbox, break the system down into the individual resources within the system, and measure, track and forecast each one individually (aka whitebox).



# Resource/Service Matrix

	Service A	Service B	Service C	Service D	Service E	Service F
cgroup memory		✓	✓	✓		✓
CPU	✓		✓	✓	✓	✓
Disk Space	✓	✓			✓	✓
GCP Quota Limits	✓					
Disk inodes	✓	✓	✓	✓		✓
PGBouncer Connections Primary		✓				
Single Threaded CPU	✓		✓			
Private Runners			✓			
Public Runners			✓			
Disk IOPS		✓		✓		
Disk Read Throughput		✓		✓		
Disk Write Throughput		✓		✓		
NAT Gateway Port Allocation	✓					
NAT Host Port Allocation	✓					

# Building a Prometheus Model for Monitoring Utilization...

# Measuring Utilization Consistently

## Using Recording Rules

```
gitlab_resource_utilization:ratio{service="postgres", resource="primary_cpu"}  
  
gitlab_resource_utilization:ratio{service="postgres", resource="pgbouncer_connection_pool"}  
  
gitlab_resource_utilization:ratio{service="redis", resource="redis_memory"}  
  
gitlab_resource_utilization:ratio{service="gitaly", resource="total_disk_capacity"}  
  
gitlab_resource_utilization:ratio{service="web", resource="rails_db_connection_pool"}  
  
gitlab_resource_utilization:ratio{service="nat", resource="nat_host_port_allocation"}
```

# Prometheus Recording Rules

```
- record: gitlab_component_saturation:ratio
  labels:
    resource: gcp_quota_limit
  expr: |
    gcp_quota_usage
    /
    gcp_quota_limit
```

1

```
232 - record: gitlab_component_saturation:ratio
233   labels:
234     component: open_fds
235   expr: |
236     process_open_fds
237     /
238     process_max_fds
```

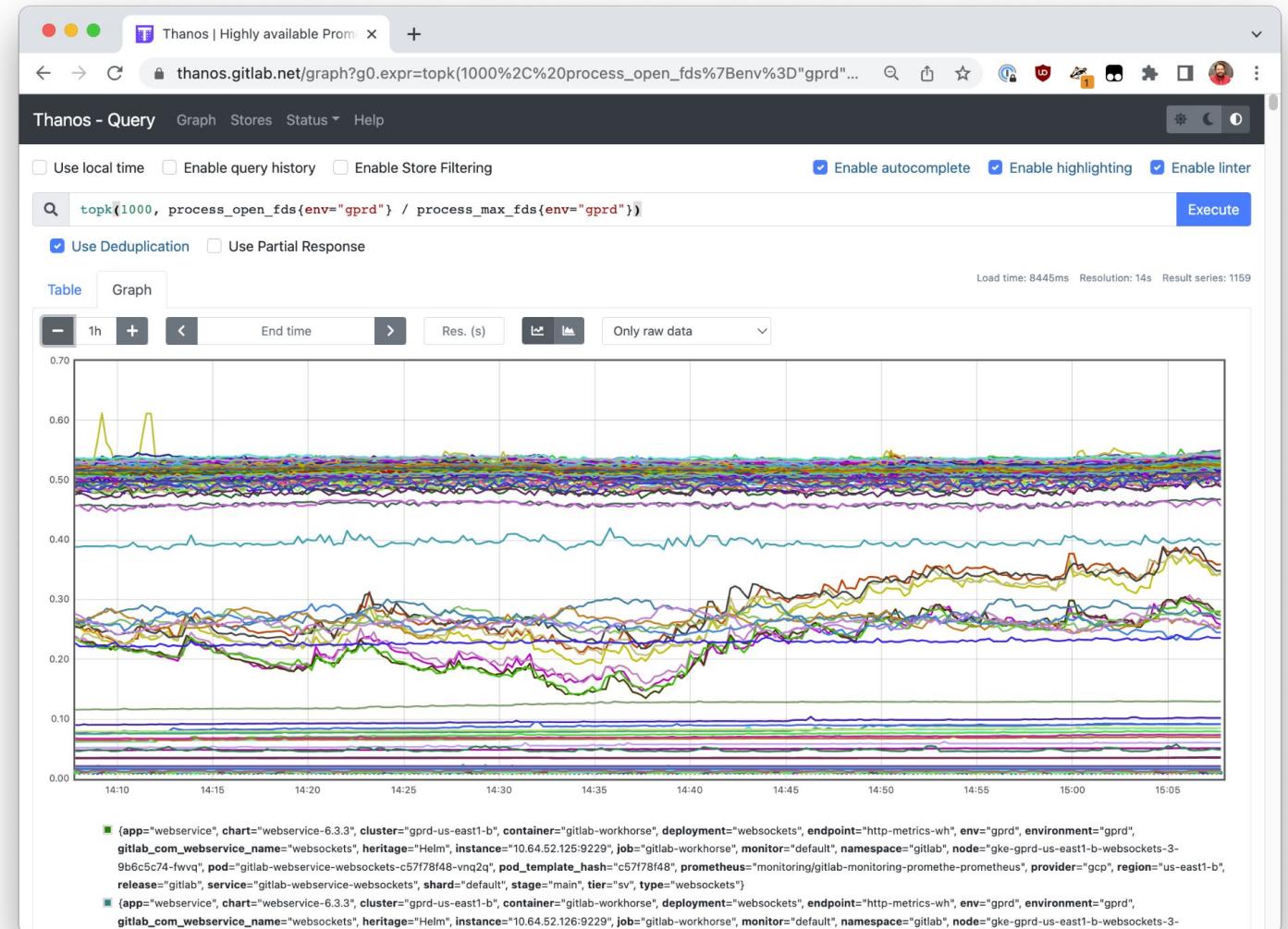
2

```
- record: gitlab_component_saturation:ratio
  labels:
    component: kube_persistent_volume_claim_disk_space
  expr: |
    kubelet_volume_stats_used_bytes
    /
    kubelet_volume_stats_capacity_bytes
```

3

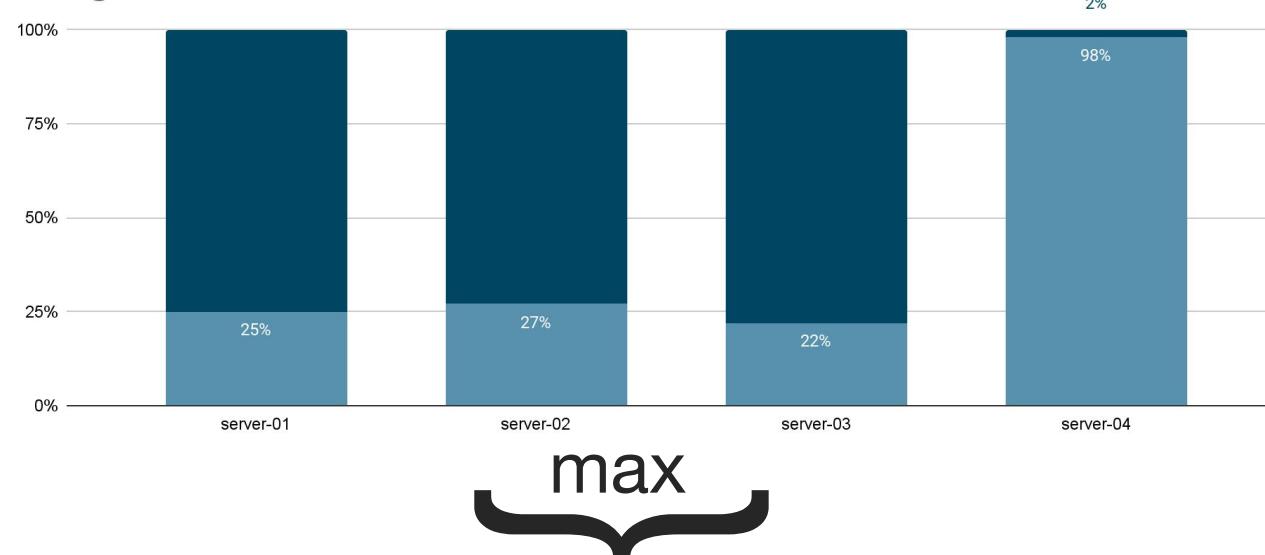
Robust Perception: [Alerting on approaching open file limits](#)

# Aggregating to Reduce Data

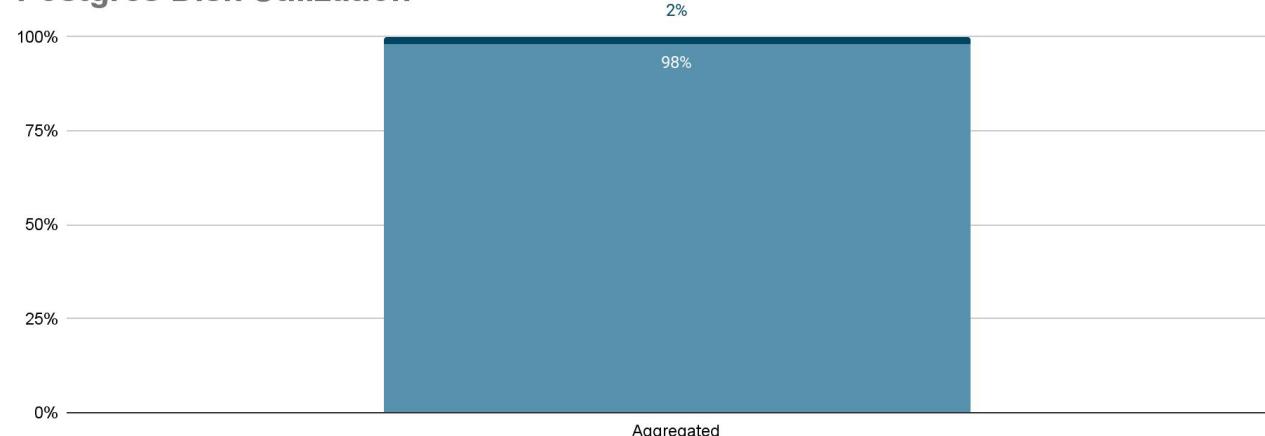


# Aggregating Utilization Across Resources

Postgres Disk Utilization

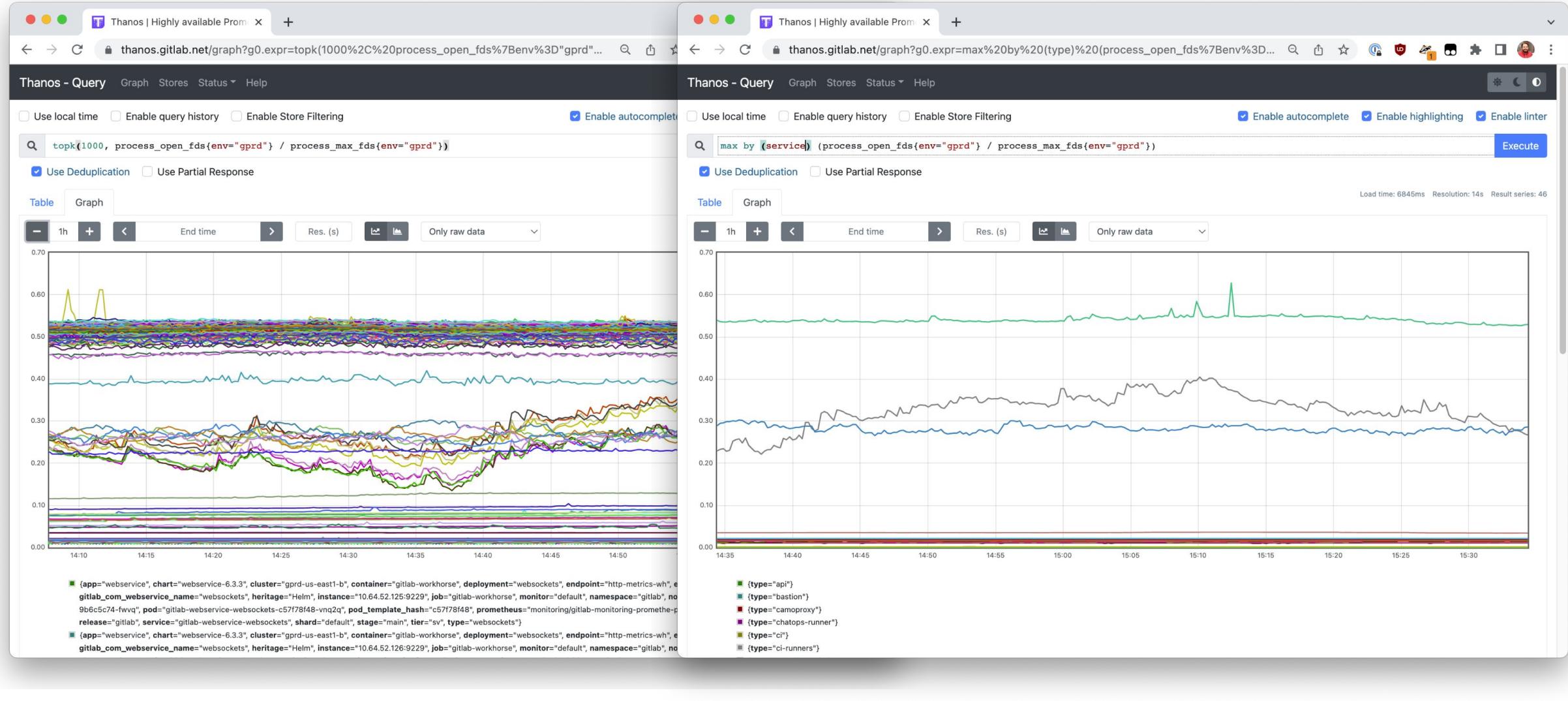


Postgres Disk Utilization



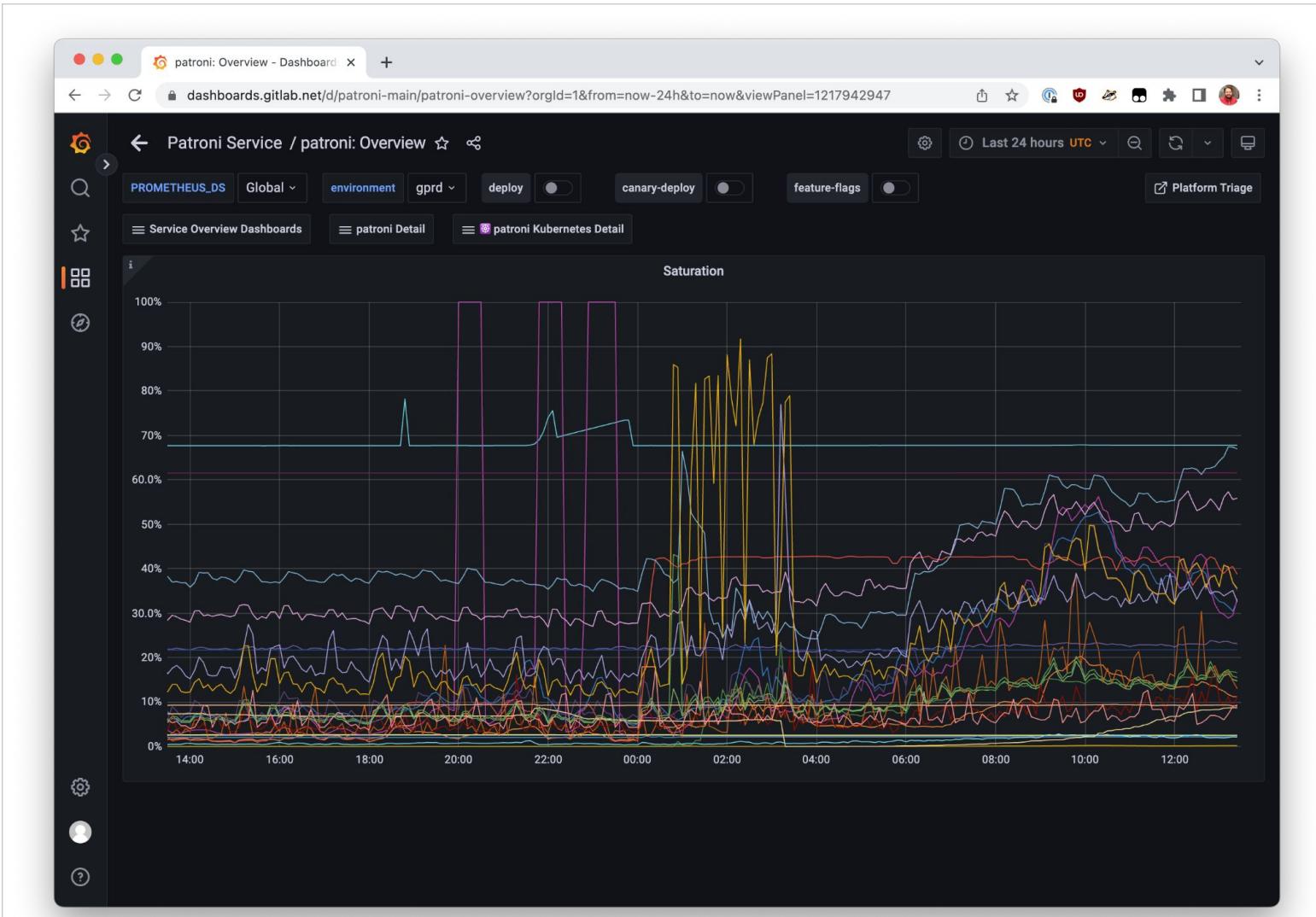
```
max by(service) (  
  1 -  
  node_filesystem_avail_bytes  
  /  
  node_filesystem_size_bytes  
)
```

# Aggregating to Reduce Data



```
5  {
6    disk_space: resourceSaturationPoint({
7      title: 'Disk Space Utilization per Device per Node',
8      horizontallyScalable: true,
9      description: /**
10       Disk space utilization per device per node.
11     */,
12     grafana_dashboard_uid: 'sat_disk_space',
13     resourceLabels: ['device'],
14     // We filter on `fqdn!=" "` to filter out any nameless workers. This is done mostly for the ci-runner fleet
15     query: /**
16       1 - (
17         node_filesystem_avail_bytes{%(selector)s}
18         /
19         node_filesystem_size_bytes{%(selector)s}
20       )
21     */,
22     slos: {
23       soft: 0.85,
24       hard: 0.90,
25       alertTriggerDuration: '15m',
26     },
27   }),
28 }
```

# Real-time Utilization Monitoring



# Real-time Utilization Monitoring (Detail)



# Alerting

Slackline APP 11:05

Alert Firing: The Node Scheduler Waiting Time resource of the frontend service (main stage) has a saturation exceeding SLO and is close to its capacity limit.

This means that this resource is running close to capacity and is at risk of exceeding its current capacity limit.

Details of the Node Scheduler Waiting Time resource:

Measures the amount of scheduler waiting time that processes are waiting to be scheduled, according to [ CPU Scheduling Metrics ](<https://www.robustperception.io/cpu-scheduling-metrics-from-the-node-exporter>).

A high value indicates that a node has more processes to be run than CPU time available to handle them, and may lead to degraded responsiveness and performance from the application.

Additionally, it may indicate that the fleet is under-provisioned.

Grafana Snapshot (86 kB) ▾

node\_schedstat\_waiting component saturation: Node Scheduler Waiting Time

Saturation %

0% 20% 40% 60% 80% 100%

03:00 04:00 05:00 06:00 07:00 08:00 09:00

min max avg current

Node	min	max	avg	current
fe-01-lb-gprd.c.gitlab-production.internal.default	6.11%	17.75%	9.41%	17.75%
fe-02-lb-gprd.c.gitlab-production.internal.default	6.98%	20.04%	10.85%	19.93%
fe-03-lb-gprd.c.gitlab-production.internal.default	5.03%	15.03%	7.90%	15.03%
fe-04-lb-gprd.c.gitlab-production.internal.default	6.79%	18.58%	9.99%	18.26%
fe-05-lb-gprd.c.gitlab-production.internal.default	6.08%	17.27%	9.33%	17.27%
fe-06-lb-gprd.c.gitlab-production.internal.default	7.26%	18.69%	10.90%	18.21%

View in Grafana

🔥 alertname: component\_saturation\_slo\_out\_of\_bounds

⚡ alert\_type: cause

🌐 environment: gprd

component: node\_schedstat\_waiting

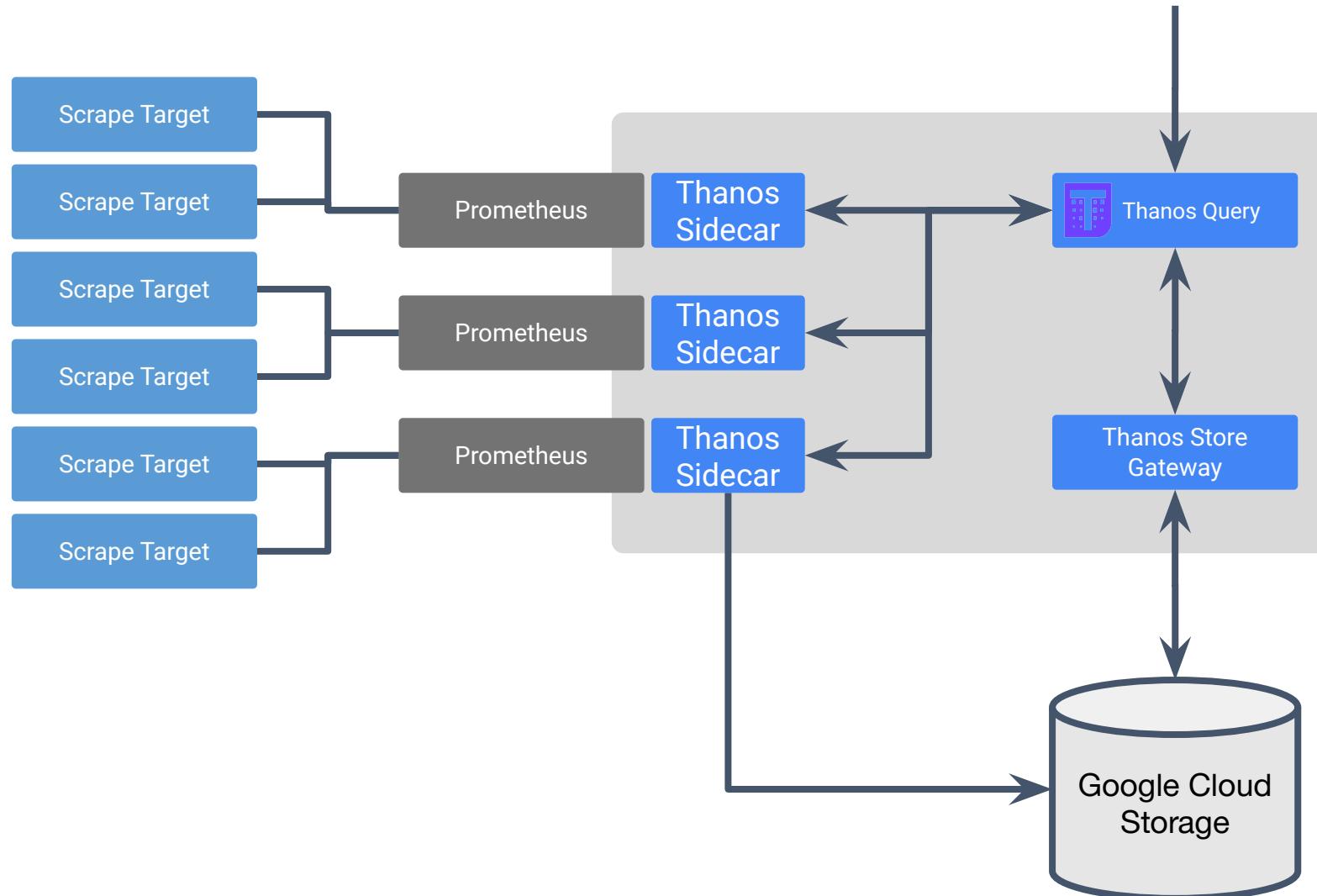
.SEVERITY. severity: s4

- **Alerting thresholds are usually less than 100%.**
- Alerting threshold depends on the nature of the resource being monitored, severity of saturation, etc.
- **Some can be very low** (this example has a threshold of only 20%)
- When an alert fires, **Slackline** tool will generate a snapshot of the **saturation detail graph** and post it to Slack (example left)

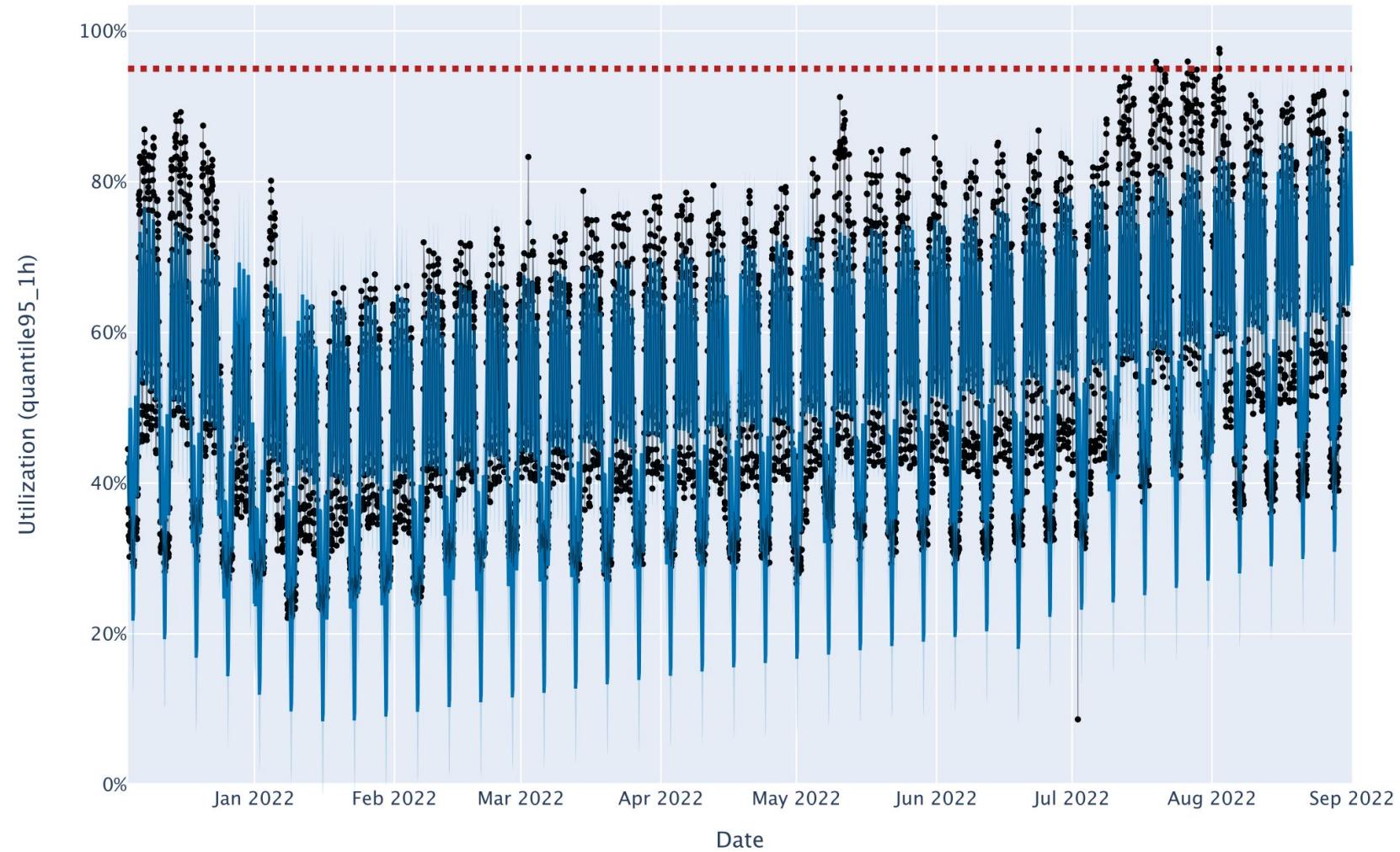
# Tamland, Forecaster



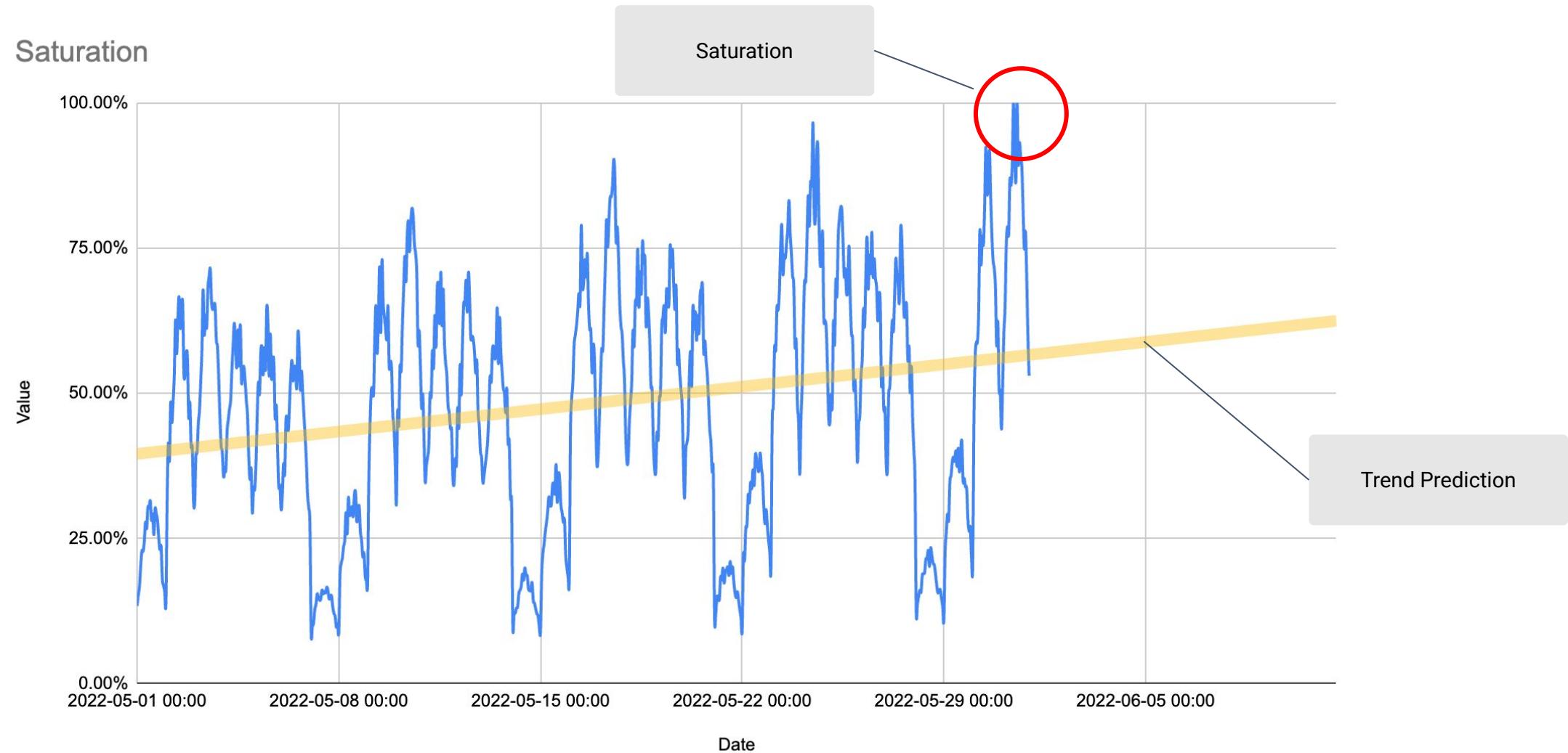
# Using Thanos for Long-Term Metrics Storage



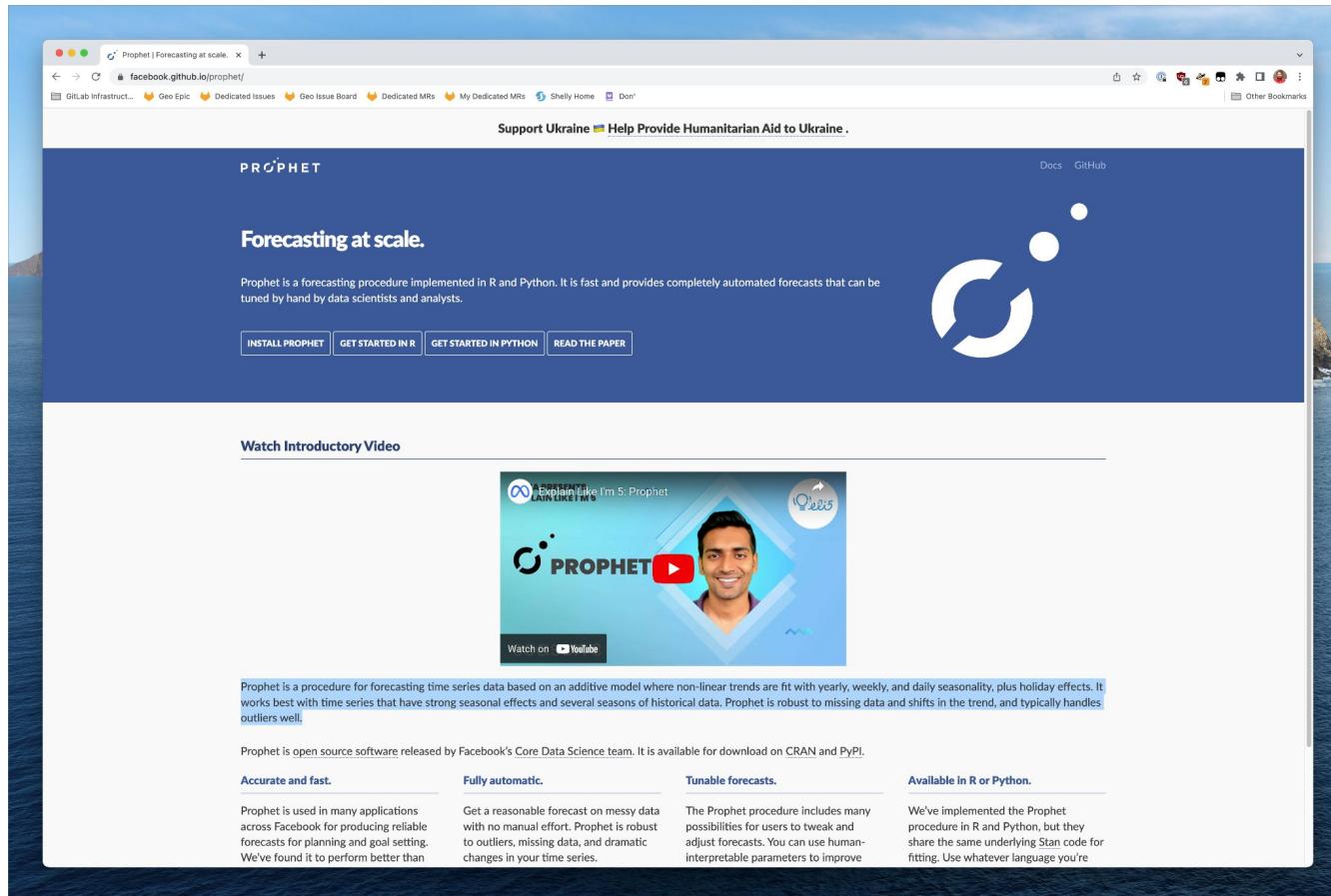
# Forecasting



# First Attempt at Forecasting: Linear Regression



# Second Attempt: Prophet

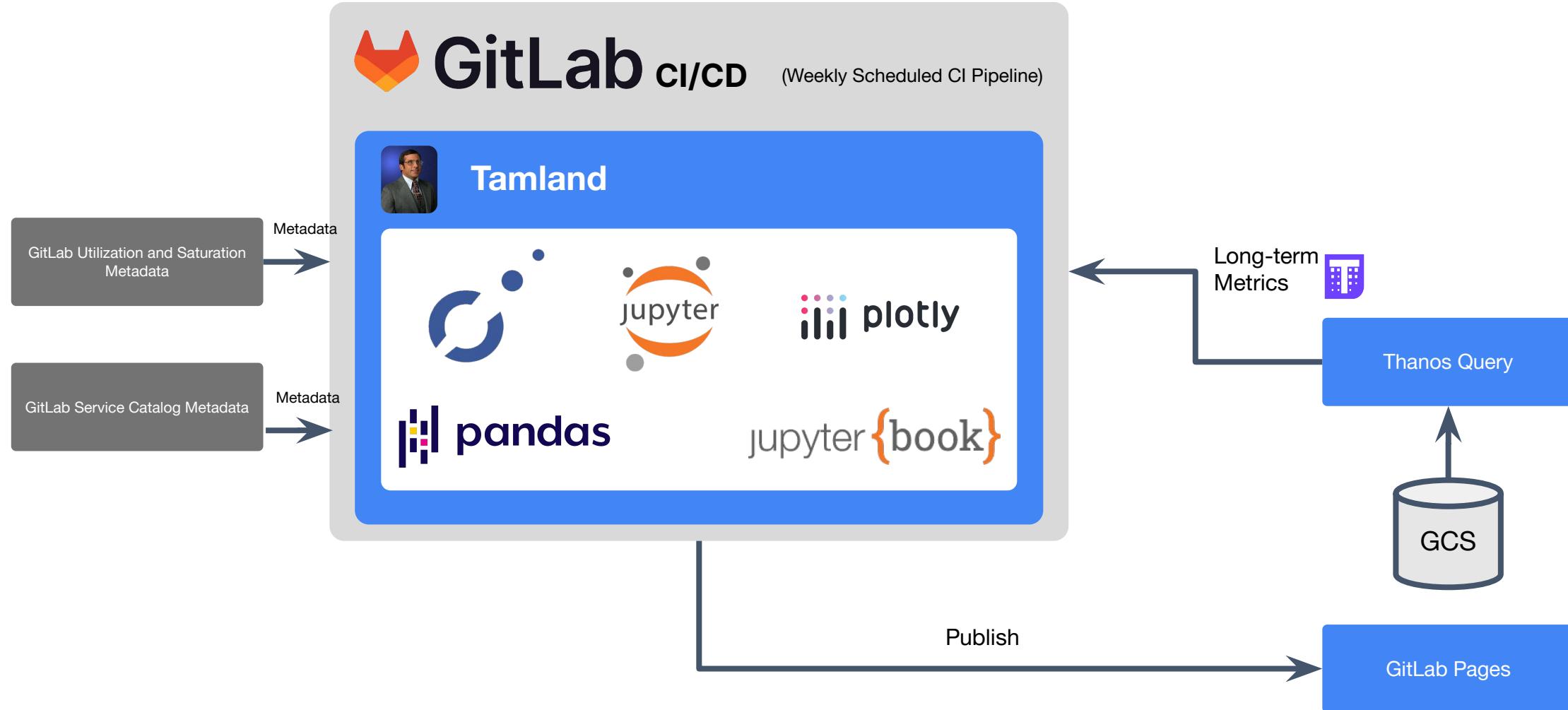


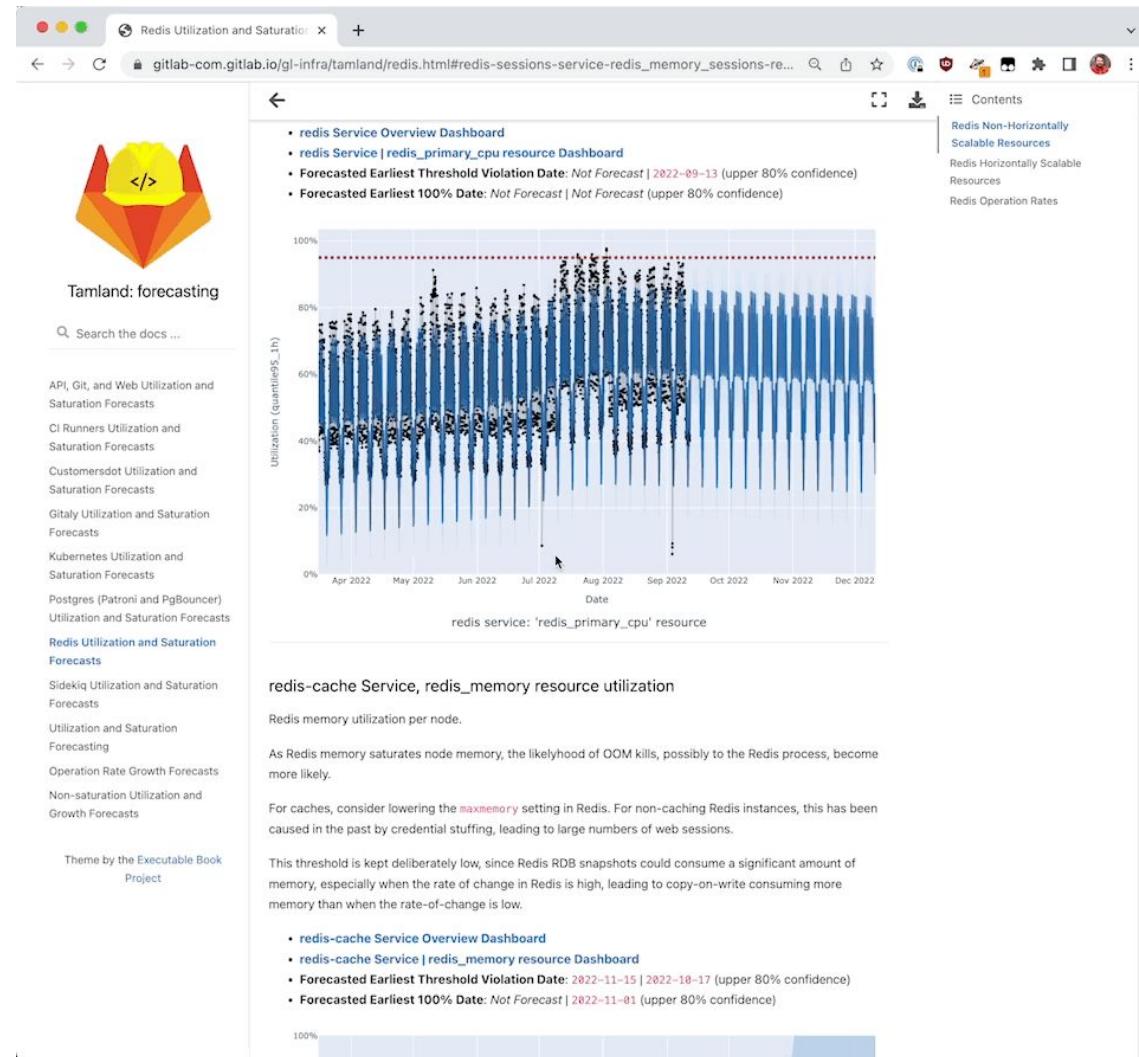
[facebook.github.io/prophet](https://facebook.github.io/prophet)

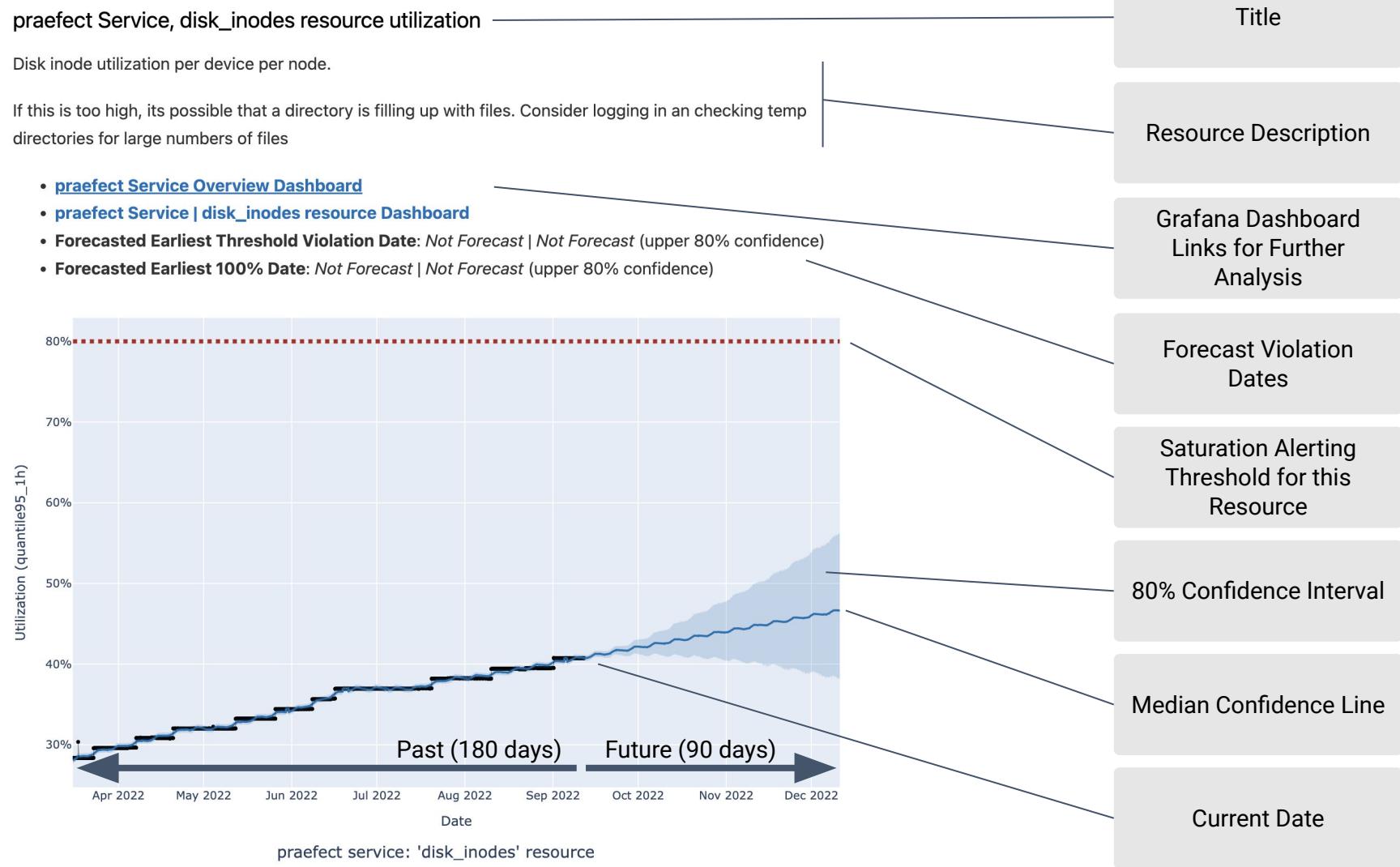
# Why Prophet?

- Building forecasting models is a **specialist data-science skill** and is difficult to do.
- **Prophet is designed for ease-of-use**, for people without a statistics or data-science major.
- **Scalable**: no need to tune each forecast independently. Prophet is able to work well, with some generic tuning, to forecast across 400 different utilization series for GitLab.com.
- Works well with data with multiple seasonalities - in our case, our data is very strongly seasonal on **hourly and weekly seasons**.
- **Handles outliers** (eg, during an incident) well.
- **Handles missing data** (eg, during a monitoring outage) well.
- **Automatic changepoint detection**. Can tell when the trend has changed and adjust the forecast.

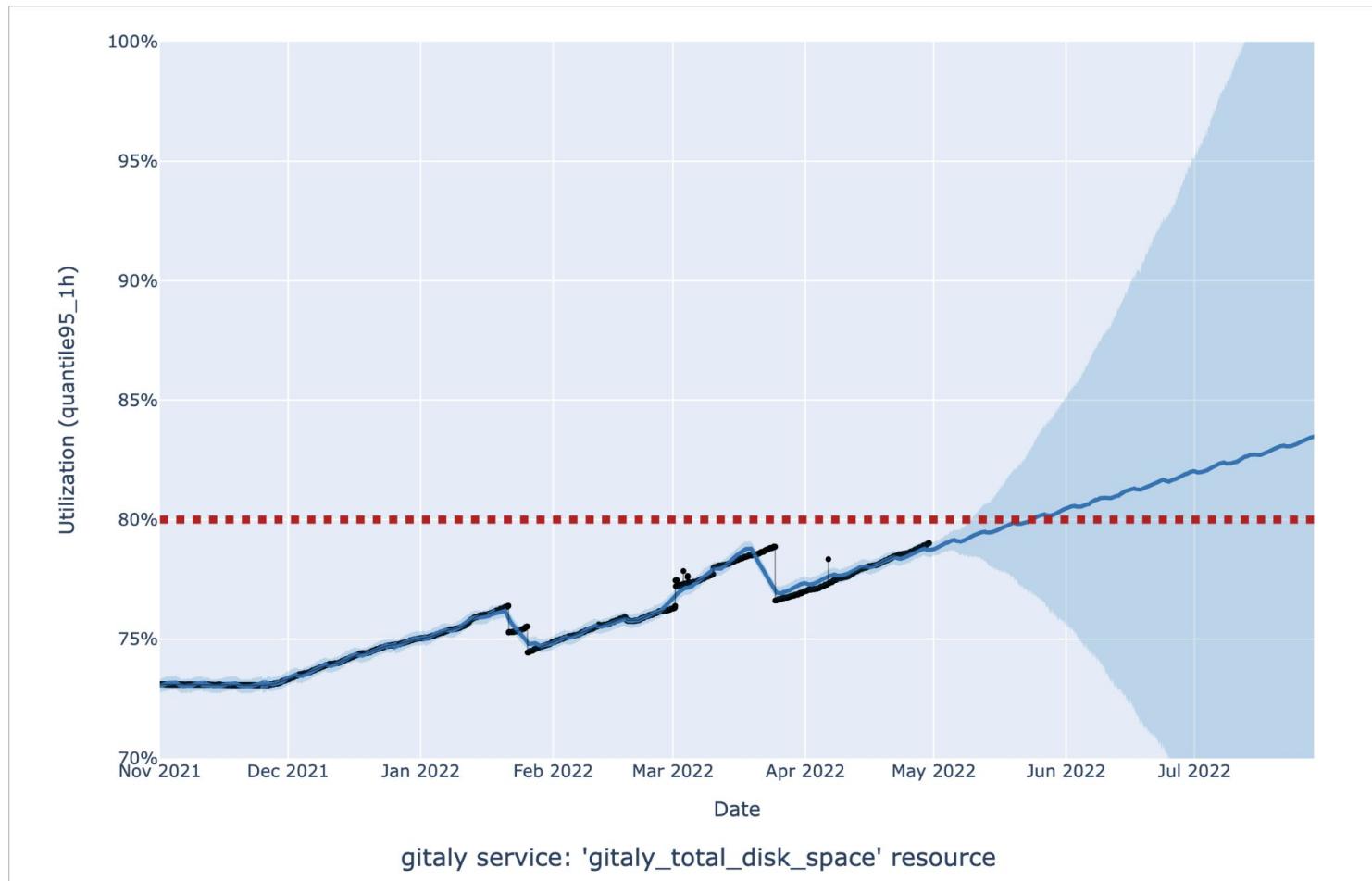
# Thanos Long Term Metrics + Facebook Prophet + Jupyter = Tamland





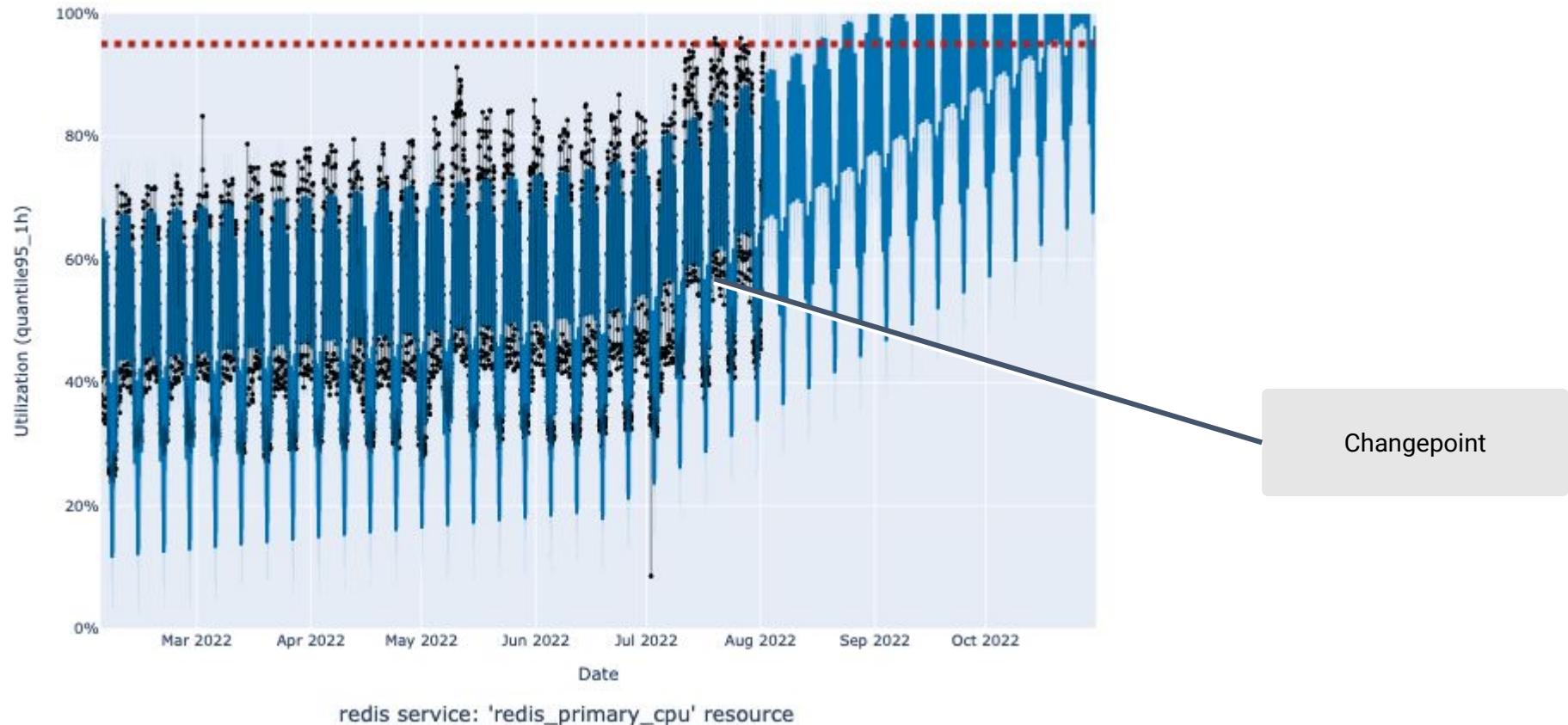


# Planning the Rollout of Extra Gitaly Nodes for Git Storage



# Changepoints

- Forecasted Earliest Threshold Violation Date: 2022-08-16 | 2022-08-02 (upper 80% confidence)
- Forecasted Earliest 100% Date: 2022-08-30 | 2022-08-09 (upper 80% confidence)



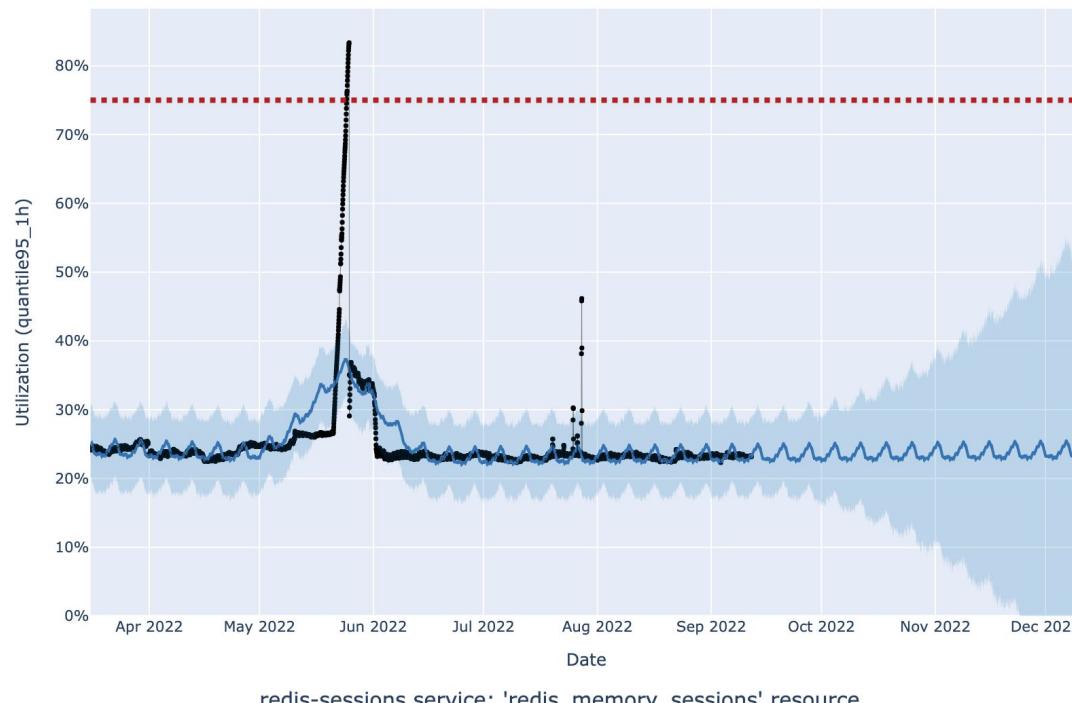
# Runaway Session State Bug

## redis-sessions Service, redis\_memory\_sessions resource utilization

Redis maxmemory utilization per node

On the sessions Redis we have maxmemory and an eviction policy as a safety-valve, but do not want or expect to reach that limit under normal circumstances; if we start evicting we will start logging out users slightly early (although only the longest inactive sessions), so we want to be alerted some time before that happens.

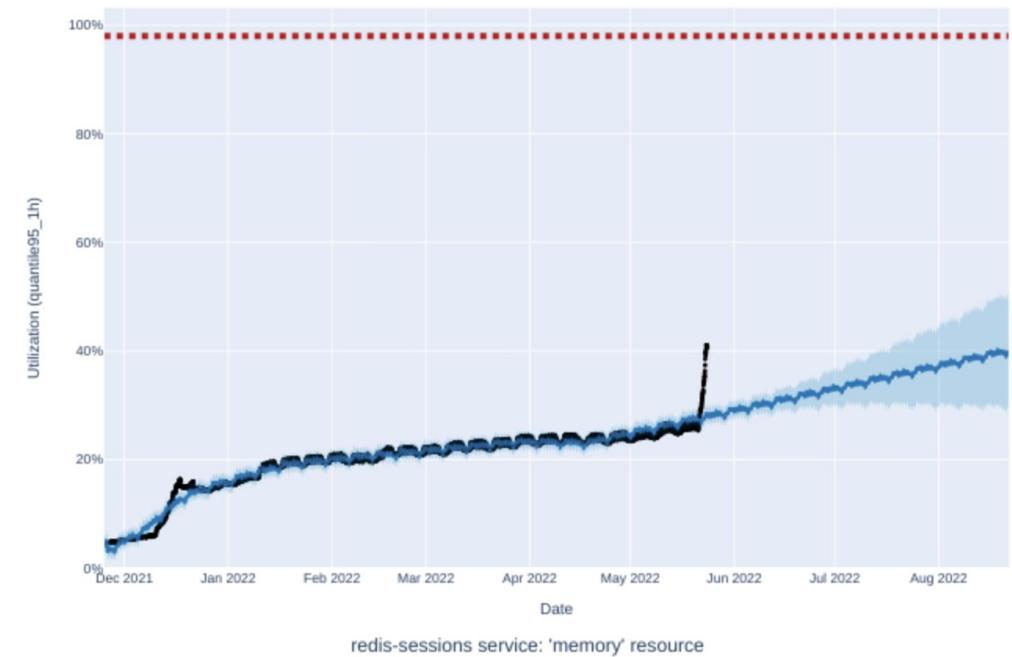
- [redis-sessions Service Overview Dashboard](#)
- [redis-sessions Service | redis\\_memory\\_sessions resource Dashboard](#)
- **Forecasted Earliest Threshold Violation Date:** Not Forecast | Not Forecast (upper 80% confidence)
- **Forecasted Earliest 100% Date:** Not Forecast | Not Forecast (upper 80% confidence)



Bob Van Landuyt @reprazent · 3 months ago

## redis-sessions Service, memory resource utilization

- [redis-sessions Service Overview Dashboard](#)
- [redis-sessions Service | memory resource Dashboard](#)
- **Forecasted Earliest Threshold Violation Date:** Not Forecast | Not Forecast (upper 80% confidence)
- **Forecasted Earliest 100% Date:** Not Forecast | Not Forecast (upper 80% confidence)



This doesn't look like it's close to saturation just yet, though the behaviour is highly uncharacteristic. We should probably figure out what happened here. It looks like something might have changed on 2022-05-22 that made us store more

Though `redis_memory_sessions` looks worse, this takes `maxmemory` into account.

- **Weekly forecasting for 363 service/resource combinations.**
- **366 Capacity Planning issues opened/investigated (of which 262 have been closed).**
- **Tamland report generation job takes around ~90 minutes to run (with caching of metrics) on a dedicated n1-standard-8 runner.**



GIFSOUP.COM

# Future Plans

- **Further decouple** Tamland from GitLab.com Infrastructure, allowing more widespread usage.
- **Improve accuracy:** Allow further customization of forecasting methods and alerting parameters for different resource classes.
- **Better ownership** - allow teams to receive notifications for the services they manage.
- **Experiment with other forecasting libraries**, such as [NeuralProphet](#) or LinkedIn's [Greykite](#).

Neural Prophet



The GitLab.com capacity planning process:

- **Avoid estimating total system capacity**
- Use a **consistent resource utilization monitoring** strategy across all your services.
- Use a **scalable forecasting library** - Prophet
- **Setup solid engineering workflows** to review and escalate potential issues.



# You Stay Classy KubeCon NA!

**Tamland Source Code:**

<https://gitlab.com/gitlab-com/gl-infra/tamland>

**GitLab Runbooks Saturation Definitions:**

<https://gitlab.com/gitlab-com/runbooks/-/tree/master/libsonnet/saturation-monitoring>

**Interesting in this sort of thing? Join us!**

<https://about.gitlab.com/jobs/all-jobs/#Engineering>



Please scan the QR Code above to  
leave feedback on this session