

LEAP out of GPU quota

Status: Draft

Overview

The LEAP project was seeing heavy GPU usage and was running out of GPU quota, triggering some user server start failures. We asked for more quota and this was granted.

OWNER OF REVIEW PROCESS

Yuvi Panda

Where we got lucky

No comments added

IMPACT TIME

Aug 26 at 11:31 to Aug 26 at 13:45

What Went Well?

- 1. Our new alerts for user server startup failure fired, and investigating it to prevent false positives surfaced a new issue for us

DURATION

2h 13m 22s

What Didn't Go So Well?

- 1. We don't have a specific alert just for quotas being close to done, as that would have *prevented* this issue rather than resolve it after the fact.

*All times listed in this report are in Pacific Time (US & Canada).

Action Items

- 1. Add alerts for any quota being close to 90% full <https://github.com/2i2c-org/infrastructure/issues/2265>

Timeline

Aug 26, 2025

11:21 AM

User tries to spawn a GPU server, but new nodes immediately fail due to lack of quota. Visible in logs at https://console.cloud.google.com/logs/query;query=resource.type%3D%22k8s_cluster%22%0Aresource.labels.location%3D%22us-central1%22%0Aresource.labels.cluster_name%3D%22leap-cluster%22%0AlogName%3D%22projects%2Fleap-pangeo%2Flogs%2Fcontainer.googleapis.com%252Fcluster-autoscaler-visibility%22%0A%2528jsonPayload.resultInfo.results.errorMessage.messageId%3D%22scale.up.error.quota.exceeded%22%2529;cursorTimestamp=2025-08-26T18:24:17.223461691Z;duration=P3D?project=leap-pangeo with a "scale.up.error.quota.exceeded" error.

11:31 AM

Triggered through the API.

Description: [FIRING:1] Server Startup Failed leap prod (take immediate action) (View Message)

INCIDENT #1273

LEAP out of GPU quota

11:31 AM

GPU server spawn fails after 10minutes

1:37 PM

Additional GPU quota requested via the Cloud Console and immediately granted

1:44 PM

Community informed of this action via Freshdesk (<https://2i2c.freshdesk.com/a/tickets/3795>)

1:45 PM

Resolved by Yuvi Panda through the website.

INCIDENT #1273

LEAP out of GPU quota