# Cloud Cascade

cloudcascade.io

Shardool Pathak - Founder
contact@cloudcascade.io
https://www.linkedin.com/in/shardool-p-86a066125/

# How to
# Seamlessly Automate and Abstract away HPC Infrastructure

# Problem Context

**Current HPC infra and cluster setup involves...**

**Steep Learning Curve:** Weeks of training needed to master cloud-specific HPC deployment(AWS, GCP)

# Problem Context

**Current HPC infra and cluster setup involves...**

**Steep Learning Curve:** Weeks of training needed to master cloud-specific HPC deployment(AWS, GCP)

**Custom Effort Required:** Each cloud provider demands a unique setup and configuration

# Problem Context

**Current HPC infra and cluster setup involves...**

**Steep Learning Curve:** Weeks of training needed to master cloud-specific HPC deployment(AWS, GCP)

**Custom Effort Required:** Each cloud provider demands a unique setup and configuration

**Non-Transferable Skills:** Knowledge does not easily transfer between different cloud platforms

# Problem Context

**Current HPC infra and cluster setup involves...**

**Steep Learning Curve:** Weeks of training needed to master cloud-specific HPC deployment(AWS, GCP)

**Custom Effort Required:** Each cloud provider demands a unique setup and configuration

**Non-Transferable Skills:** Knowledge does not easily transfer between different cloud platforms

**High Costs:** User are often surprised with high hourly costs for clusters post-deployment

# Central Idea

**Effortless HPC cluster deployment for scientific computing, research, ML/AI model training etc.**

**Simplicity:** intuitive, natural language interactions

# Central Idea

**Effortless HPC cluster deployment for scientific computing, research, ML/AI model training etc.**

**Simplicity:** intuitive, natural language interactions

**Freedom:** Work seamlessly across any cloud (AWS, Google Cloud, Azure, Oracle etc.)

# Central Idea

**Effortless HPC cluster deployment for scientific computing, research, ML/AI model training etc.**

**Simplicity:** intuitive, natural language interactions

**Freedom:** Work seamlessly across any cloud (AWS, Google Cloud, Azure, Oracle etc.)

**Speed:** let AI take care of infra complexity, zero learning curve

# Features

- Automated Config Generation w/ Validation

- Pre-Deployment Cost Breakdown

- Seamless Cluster Management via UI

# Automated Config Generation w/ Validation

# Pre-Deployment Cost Breakdown

# Demo

https://youtu.be/G9WWXuJiGeE?si=IW9eYZhfKwY9BNQd

# Architecture

**Cloud Cascade**
(hosted in GCP Cloud Run)

**Frontend**
[angular, typescript]

**Backend**
scalable
[Go]

AWS Parallel Cluster

GCP HPC
Toolkit

Postgres

LLM APIs
Claude Sonnet

## Customer owned projects

**AWS**

Stack

HPC
cluster

**GCP**

Blueprint

HPC
cluster

# Next Features

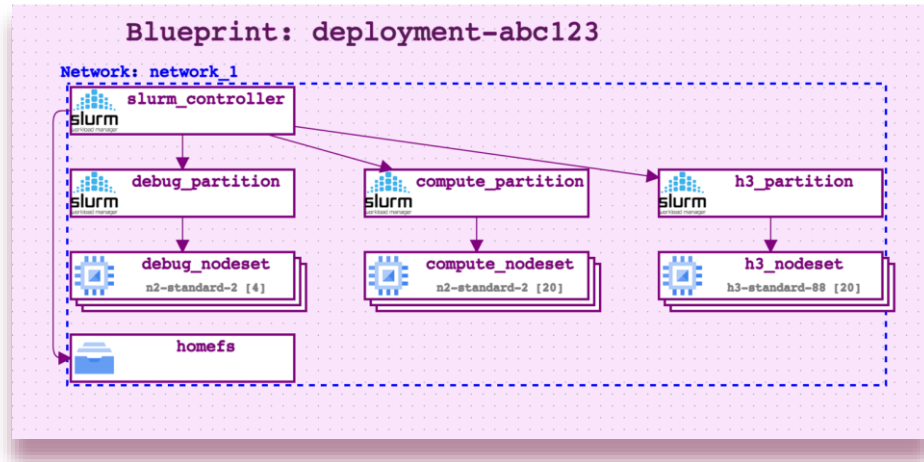**More Cloud Support (Azure, NVidia, Oracle)**

**Cloud to Cloud Conversion (AI agent)**

    **Machine / core equivalence**

    **Price equivalence**

**Cluster Diagrams & Visualization**

**Cluster Diagrams & Visualization - GCP Blueprint**

# Questions?