KCD
GROWING CLOUD NATIVE TOGETHER
**BEIJING**

# Part 01

## Address challenges of Agentic workload

# What is a Gen AI Agent

Intelligent, autonomous systems

Plan, reason, and act

Access to enterprise data

Ability to use tools

# . . . it leads to challenges

## Complexity
Coding gets complicated

- Complex prompts to limit hallucinations
- Fragile, hard to maintain

## Accuracy
Agent gets confused

- Calling wrong tools
- Passing wrong arguments
- Inconsistent responses

## Cost
Agent gets slower and more expensive

- Frontier models needed
- Prompt sizes grow
- Agents retry steps

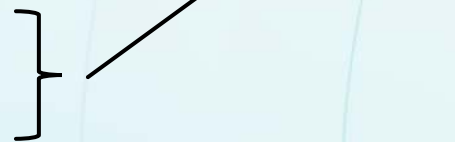# Hints to improve cost, performance and accuracy

- Shorter prompts
- Smaller LLMs

**Task Decomposition Technique**

- Control over workflow
- Concise prompts / context

- Cost effective chipsets
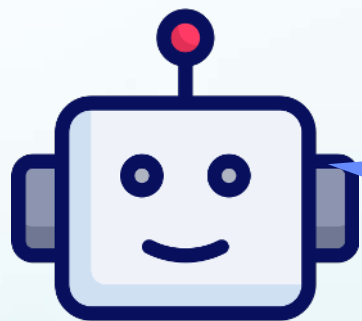- Fast Chipsets

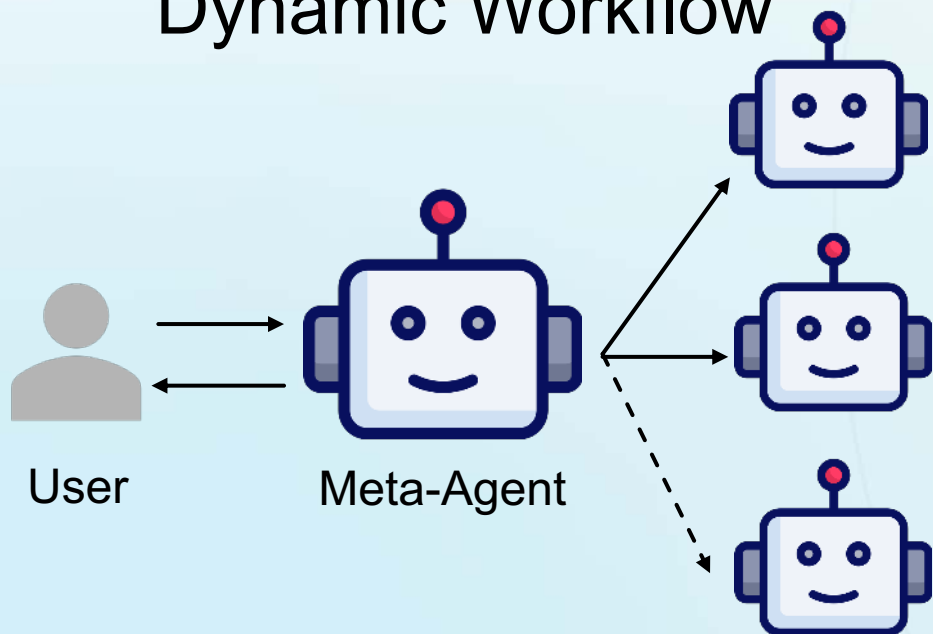**AI Chips Choice**
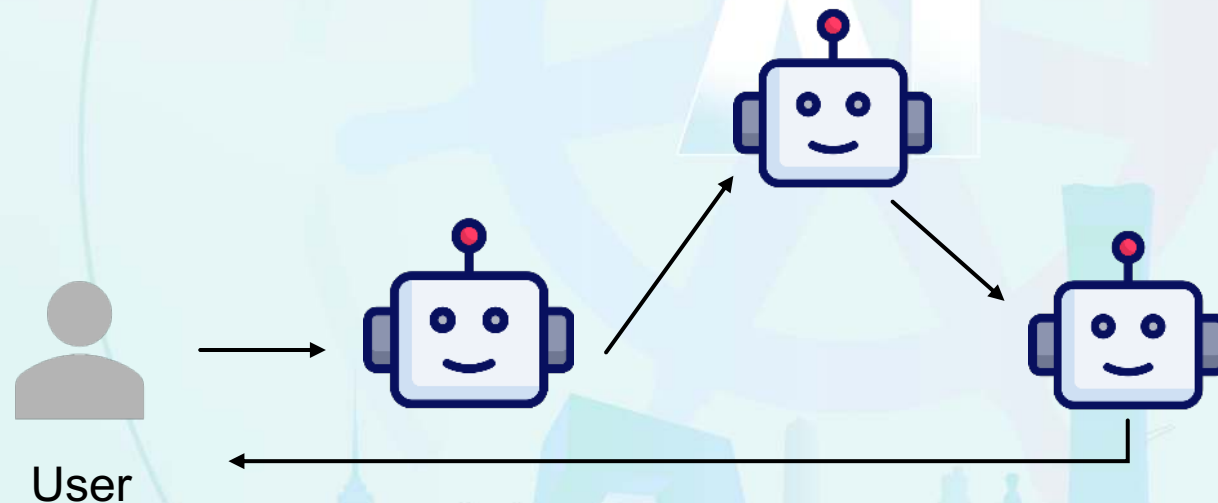
# Part **02**

## Multi-agent workflows

# Agent & workflows



I am an "agent" specializing in a task or just to coordinate (again specialized!)

Dynamic Workflow

User    Meta-Agent

Static Workflow

User

# Why Small Language Model?



Resource

Speed

Customize

# Frugal Architecture Design Patterns

## Task Decomposition

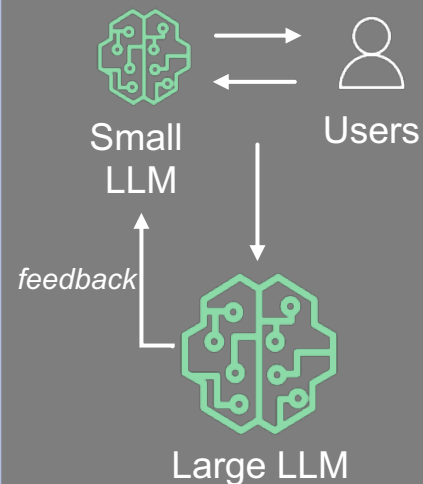**One Large LLM**

Total Task

Small LLM — Task-A
Small LLM — Task-B
Small LLM — Task-C

Fine-tuned / Specialized

## Using Small and Large LLMs Frugally

### Cascaded LLM

*when needed*

*when needed*

### RLAIF with Large LLMs

Small LLM ⇄ Users

*feedback*

Large LLM

### Teacher-Student LLM

Large LLM

Supervised Learning | Reinforcement Learning

Small LLM

## Pre-post processing for Task Simplification

Context

Classical ML → Small LLM

{JSON}

Natural Language

Small LLM

# Workflow : Routing

# Example – Routing

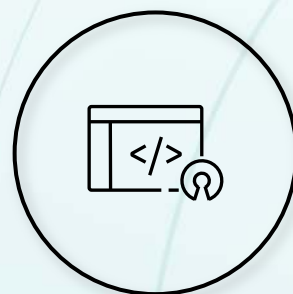# Part 03
## Implementation on K8S

# Why self-host Language Model on K8S

Data privacy
and security

Connecting
to
data sources

Customizing

Accessing
multiple
models
and newer
versions

# Running Agentic workload on K8S

# Workflow: Evaluator-optimizer

# LiteLLM

# Observability - Langfuse

# Implementation on K8S
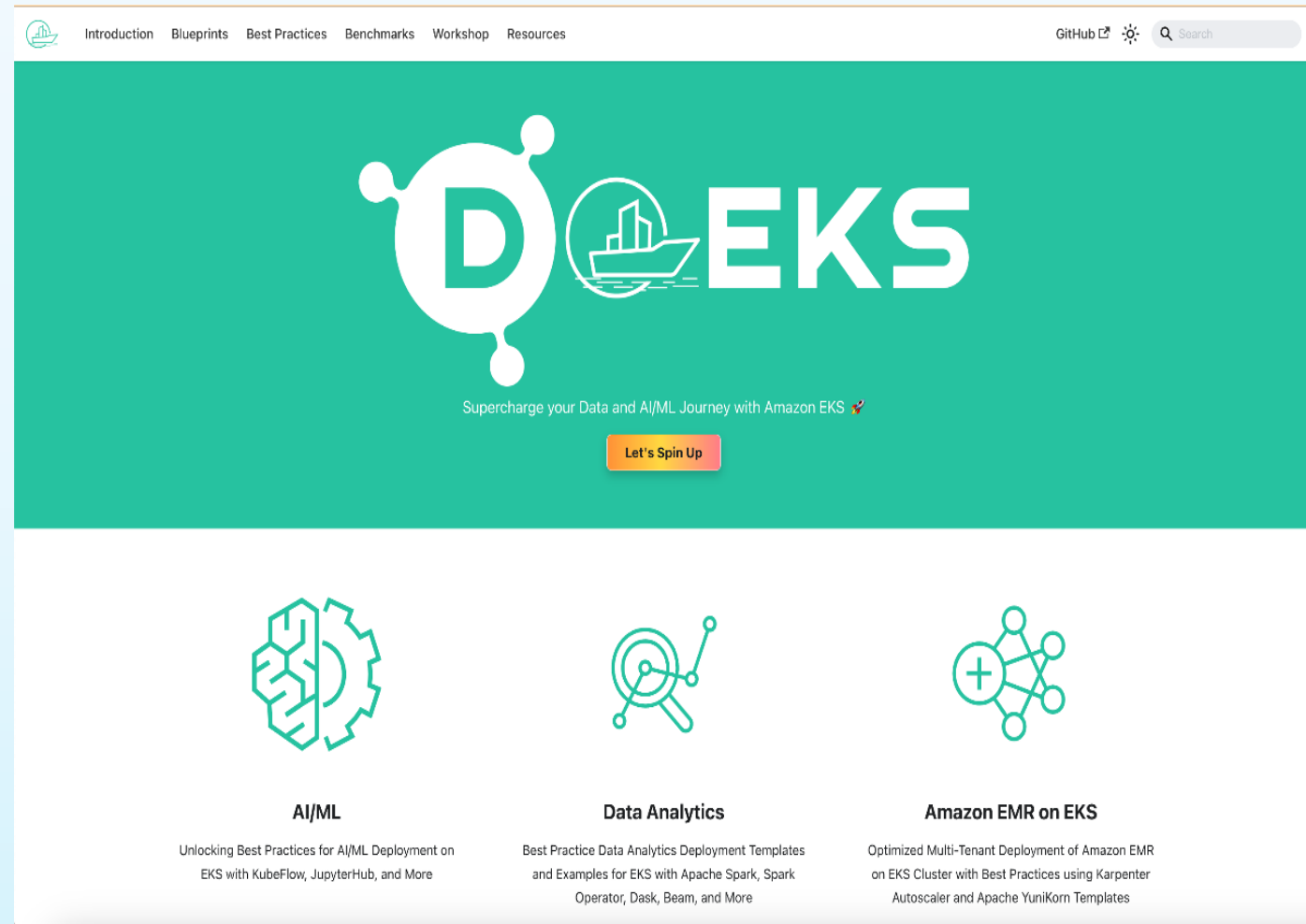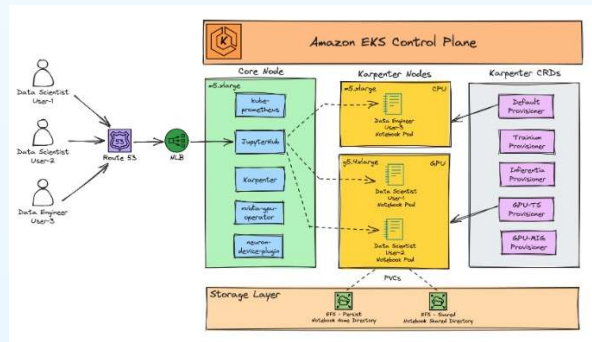
# Part 04
**Key takeaways**

# Key Takeaways

- Choice of Language Models

- Deploy models on AWS Trainium , Inferentia and Graviton

- Multi-agent workflow patterns

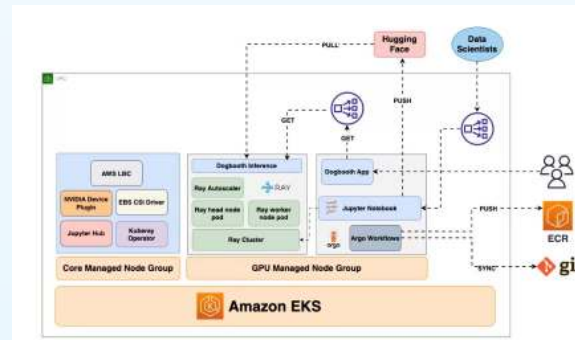- Improve user experience – Centralized Security and observability

# Data on EKS project

➢ Gen AI patterns

➢ IaC templates

➢ Best practices

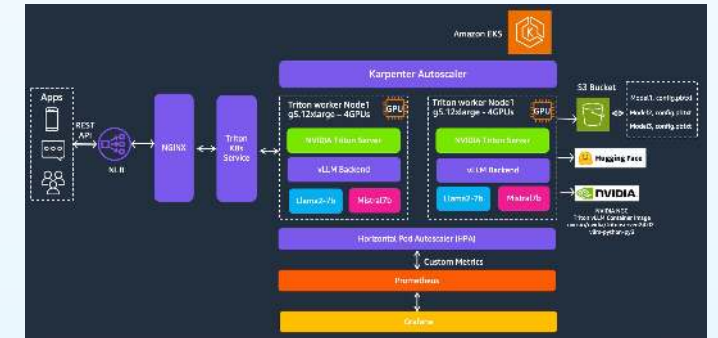➢ Engage with us on GitHub

# Reference patterns for gen AI on EKS
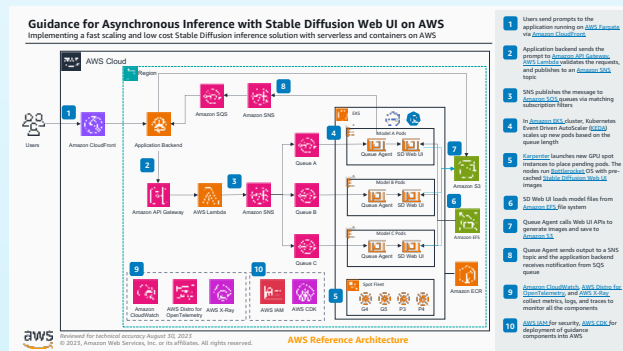


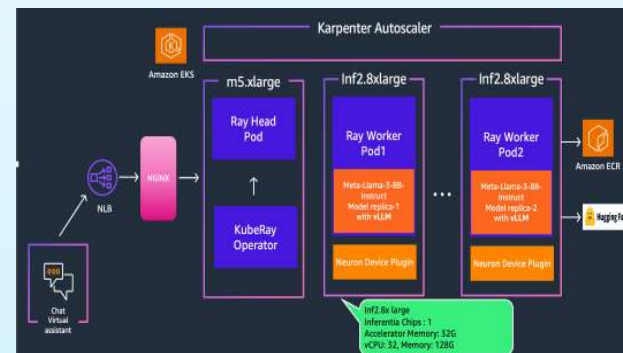Multi-tenant JupyterHub platform



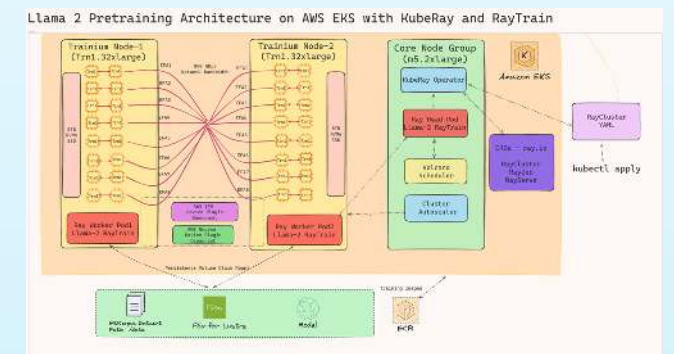End-end generative AI orchestration platform



High performance inference platform on EKS
(NVIDIA Triton with vLLM)



Async inference with Stable Diffusion on EKS



Inference of Llama-3-8B with RayServe/vLLM on EKS



Llama2 distributed pretraining on Trn1 with RayTrain