Edgeless Systems

# Confidential AI inference in practice: What's required and how to implement it
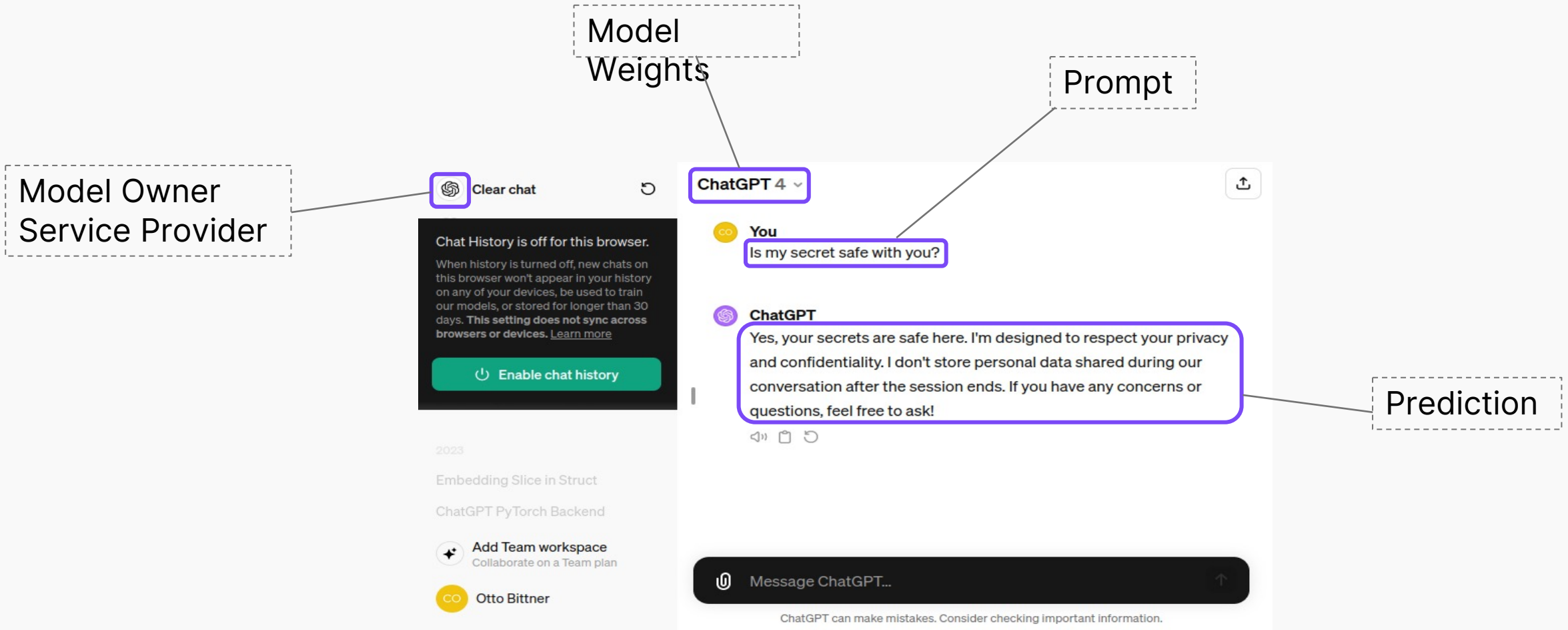
2024-03-17          Otto Bittner                    Open Confidential Computing Conference

# 01 Challenge

Company logos are used for illustrative purposes only.

Edgeless Systems

Prompt

HuggingChat  New Chat

Is my secret safe with you?

Service Provider

Today

🔒 Secure information

Yes, your secret is safe with me. I am a machine learning model and do not have the ability to disclose any information. I will treat all input as confidential and use it only for the purpose of providing a response.

Prediction

Login

Theme

Models  7

Assistants  New

Search web ⓘ

Ask anything

Settings

About & Privacy

Model: mistralai/Mixtral-8x7B-Instruct-v0.1  Generated content may be inaccurate or false.  ⤒ Share this conversation

Model Weights

Model Owner

MISTRAL AI_

# Challenge?

# Challenge



Prompts

Weights

MISTRAL AI_

🤗 Hugging Face

# Challenge

# 02 Goals

# Goal #1

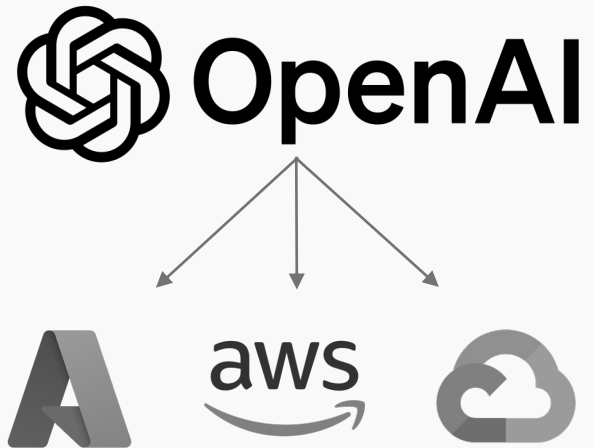# Goal #2

OpenAI

# Goal #3



Company logos are used for illustrative purposes only.
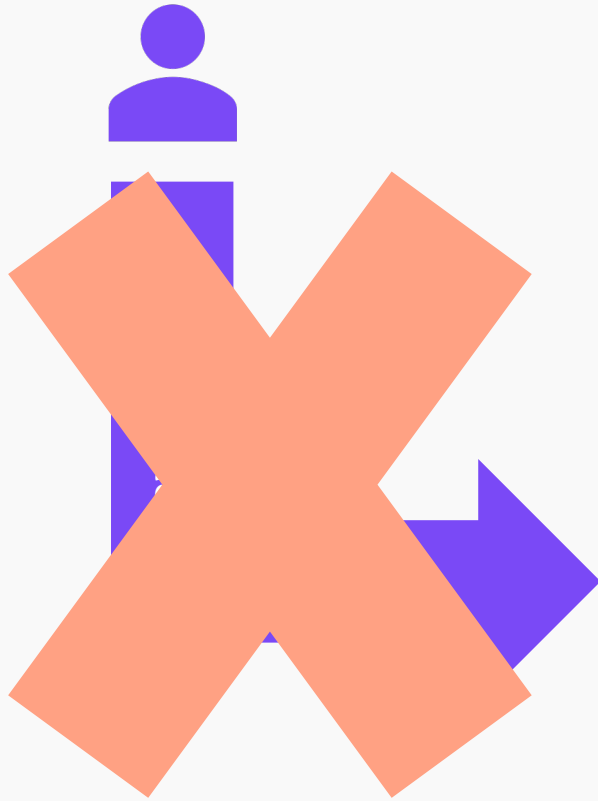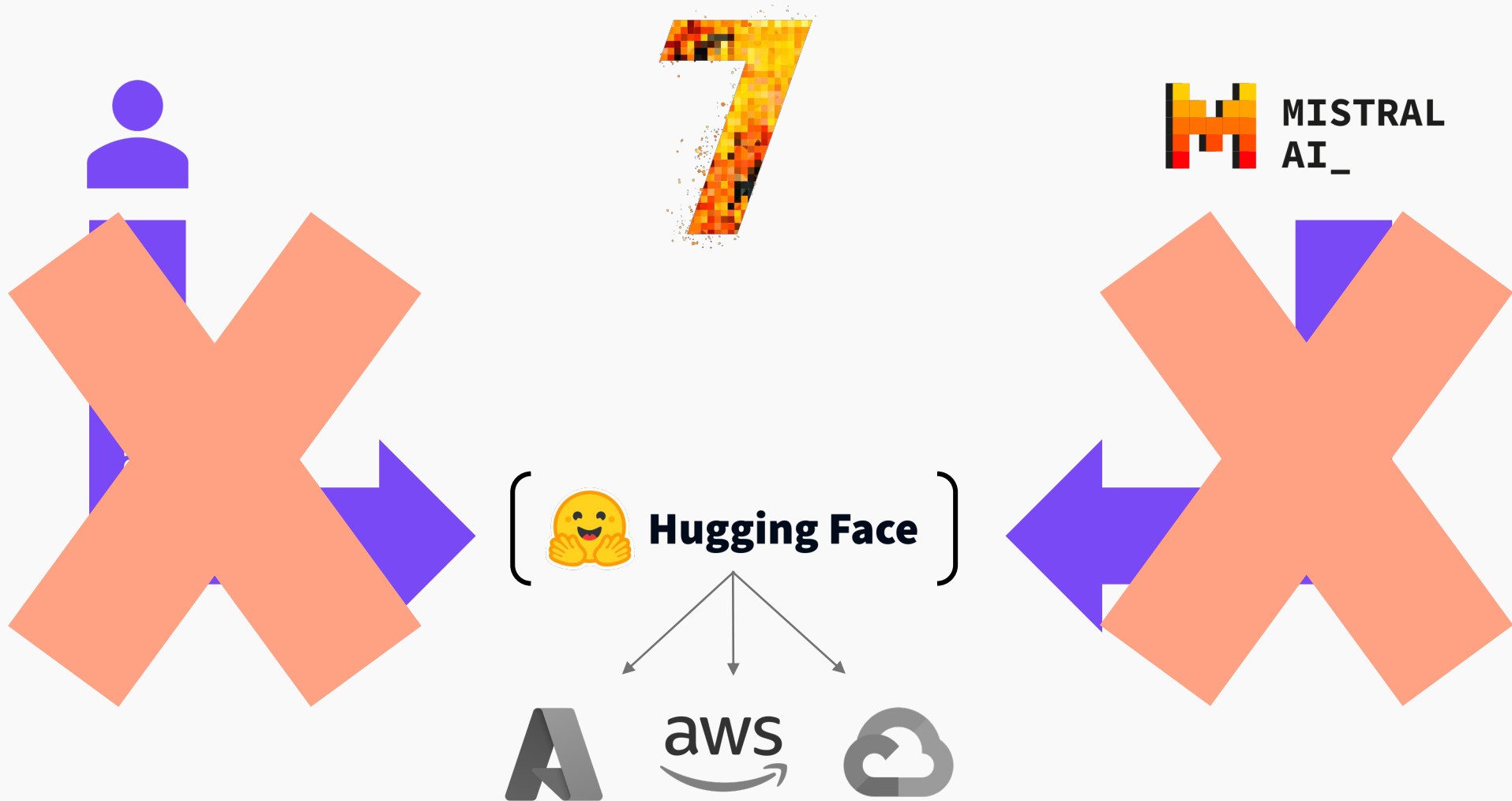
# Goals - specific

- A **service** provider should have neither access to **prompts** nor **weights**


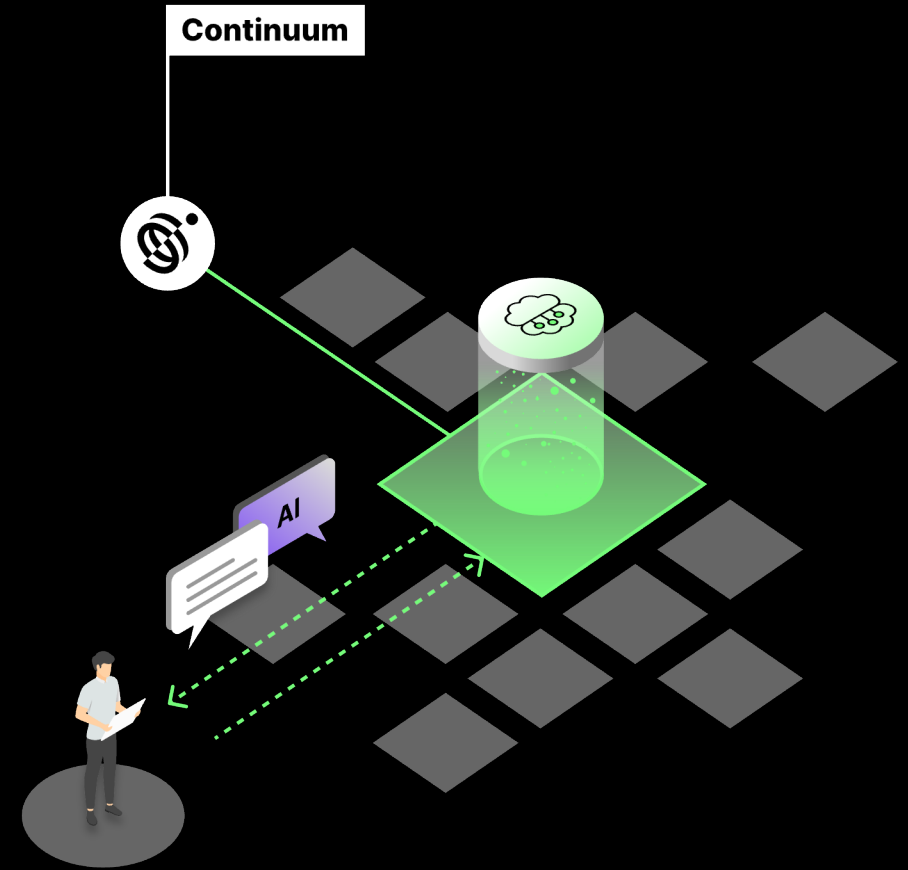- An **infrastructure** provider should have neither access to **prompts** nor **weights**
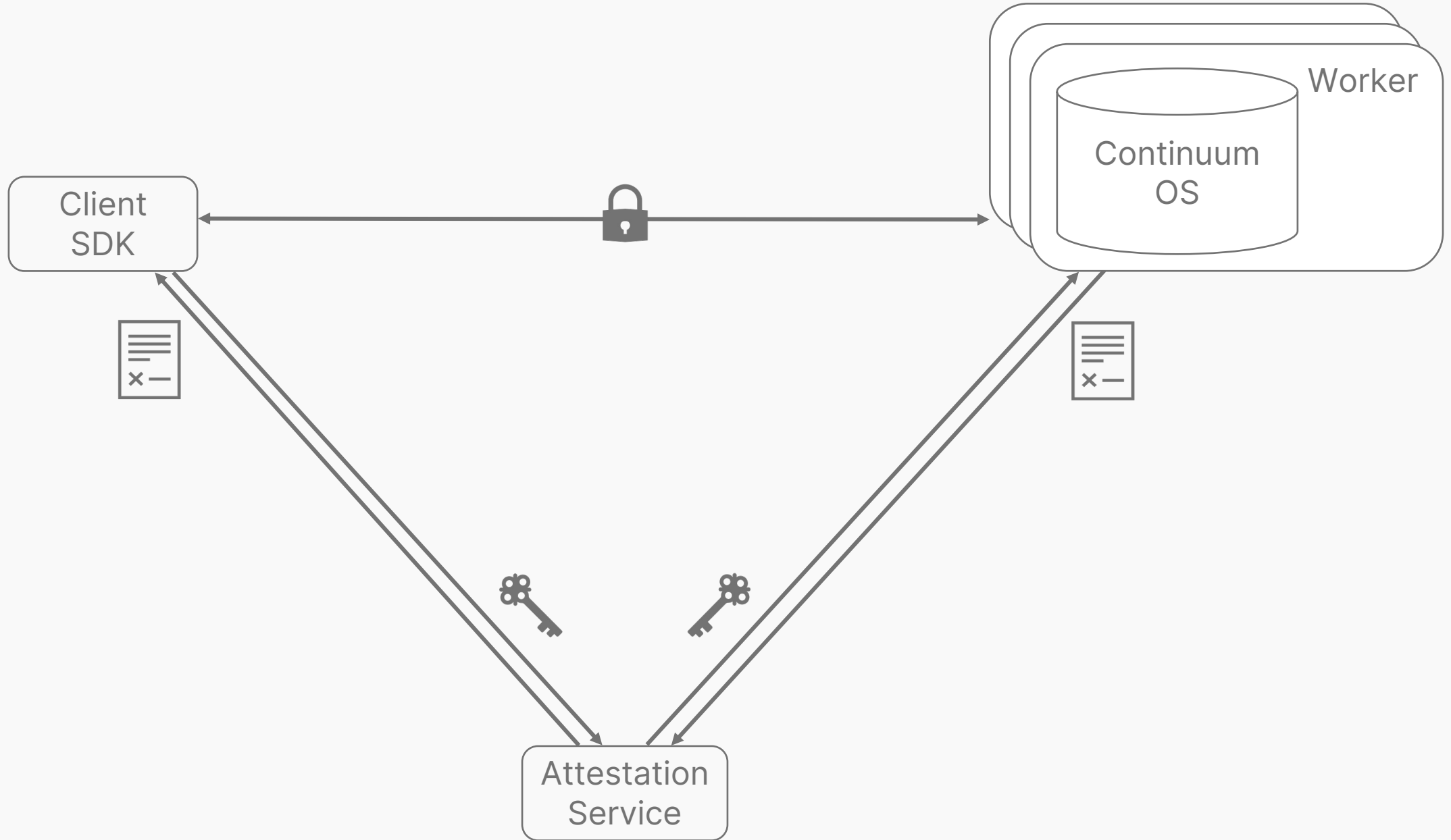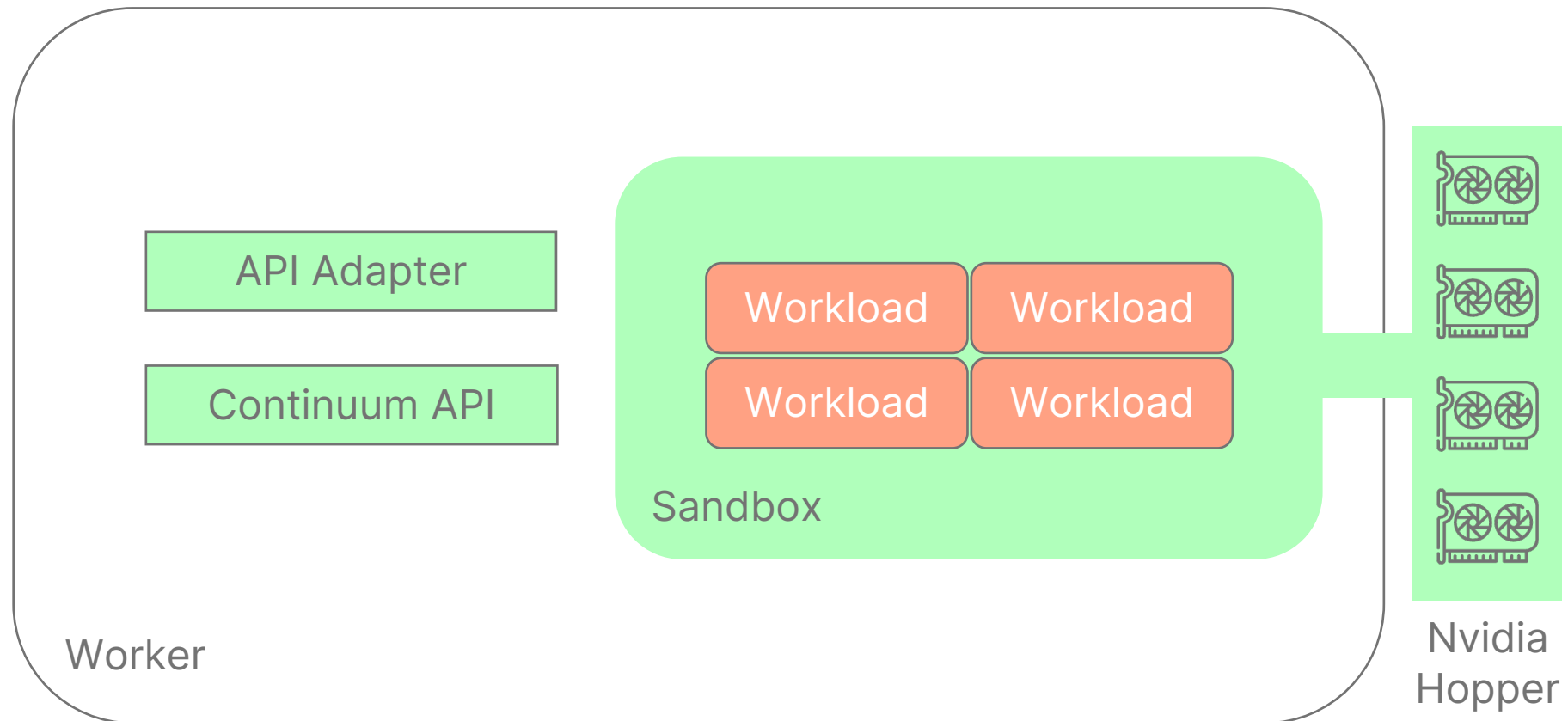
03 Solution

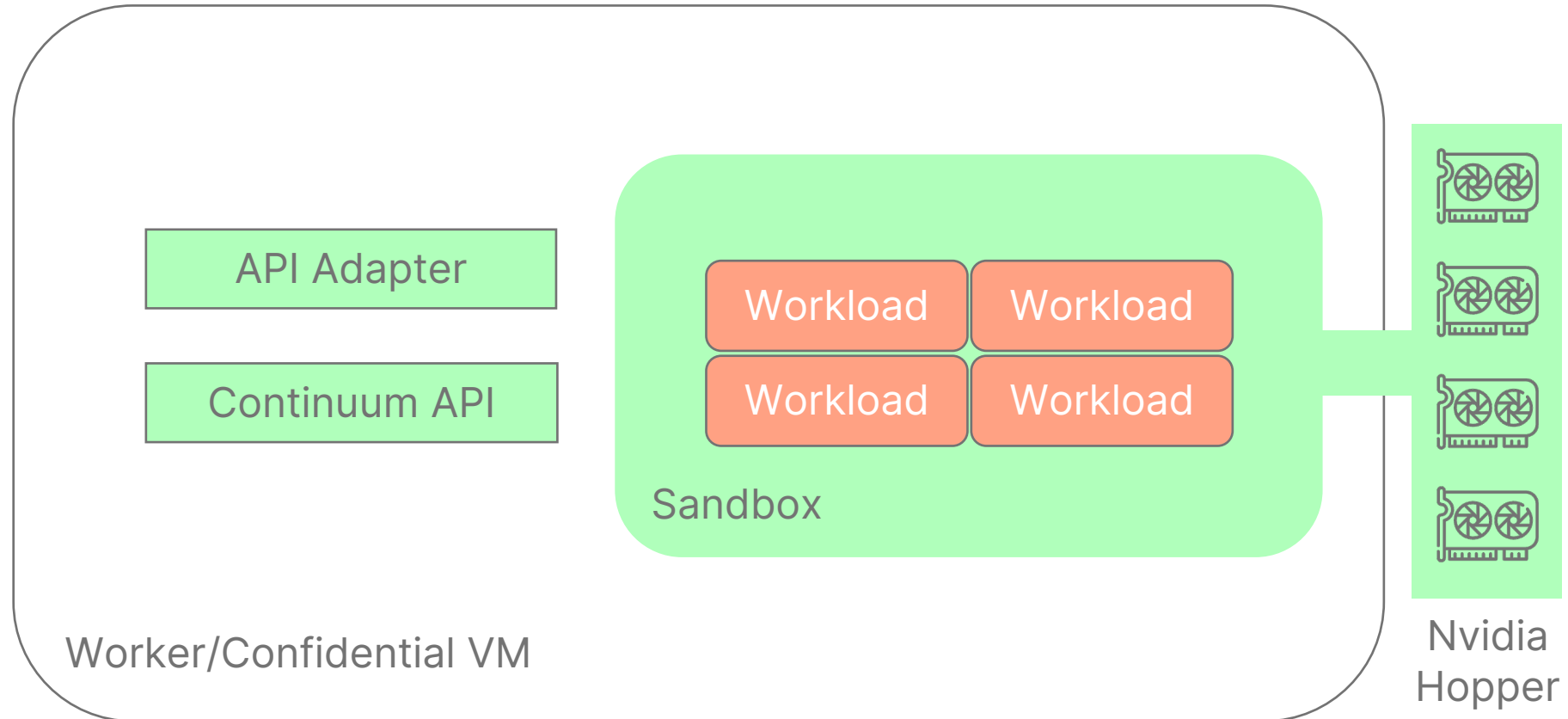# Continuum

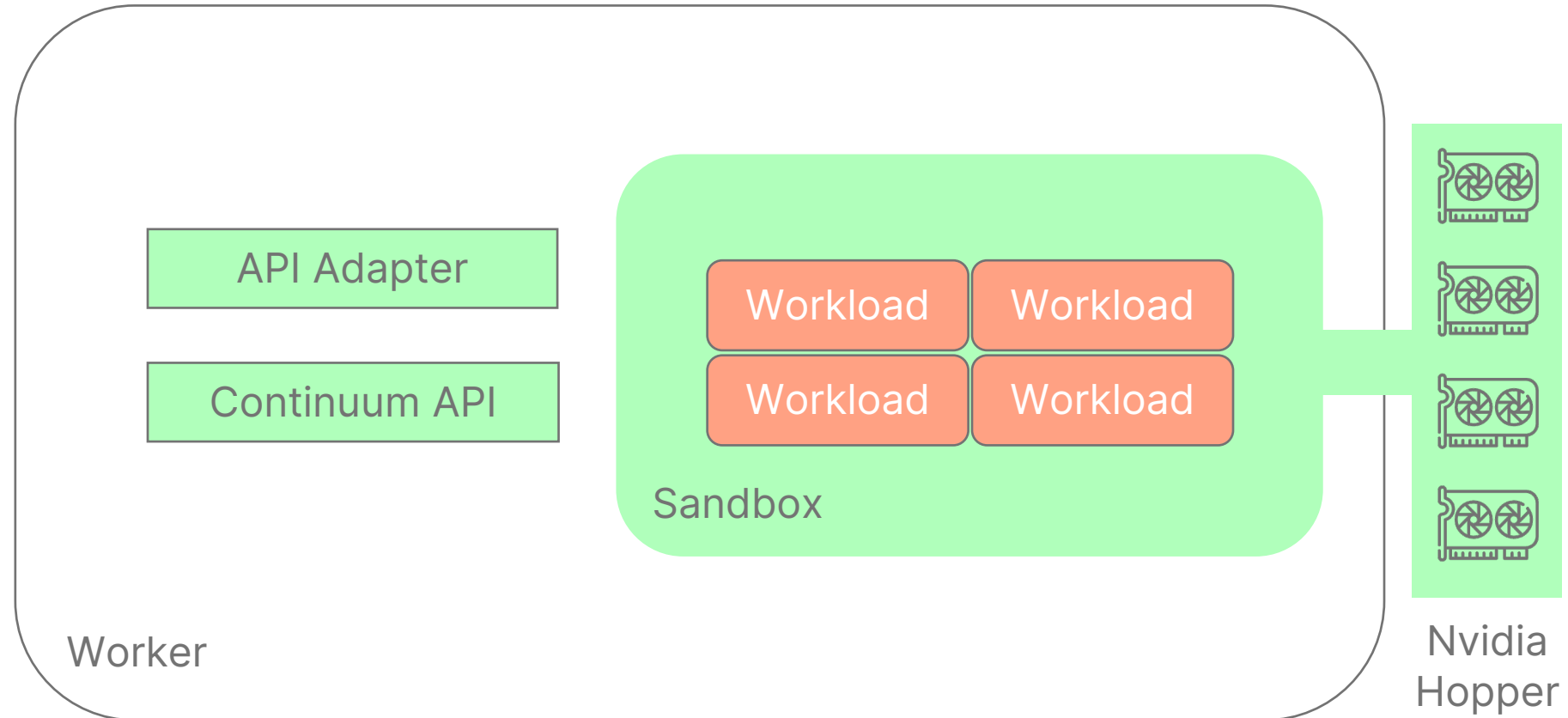# Continuum

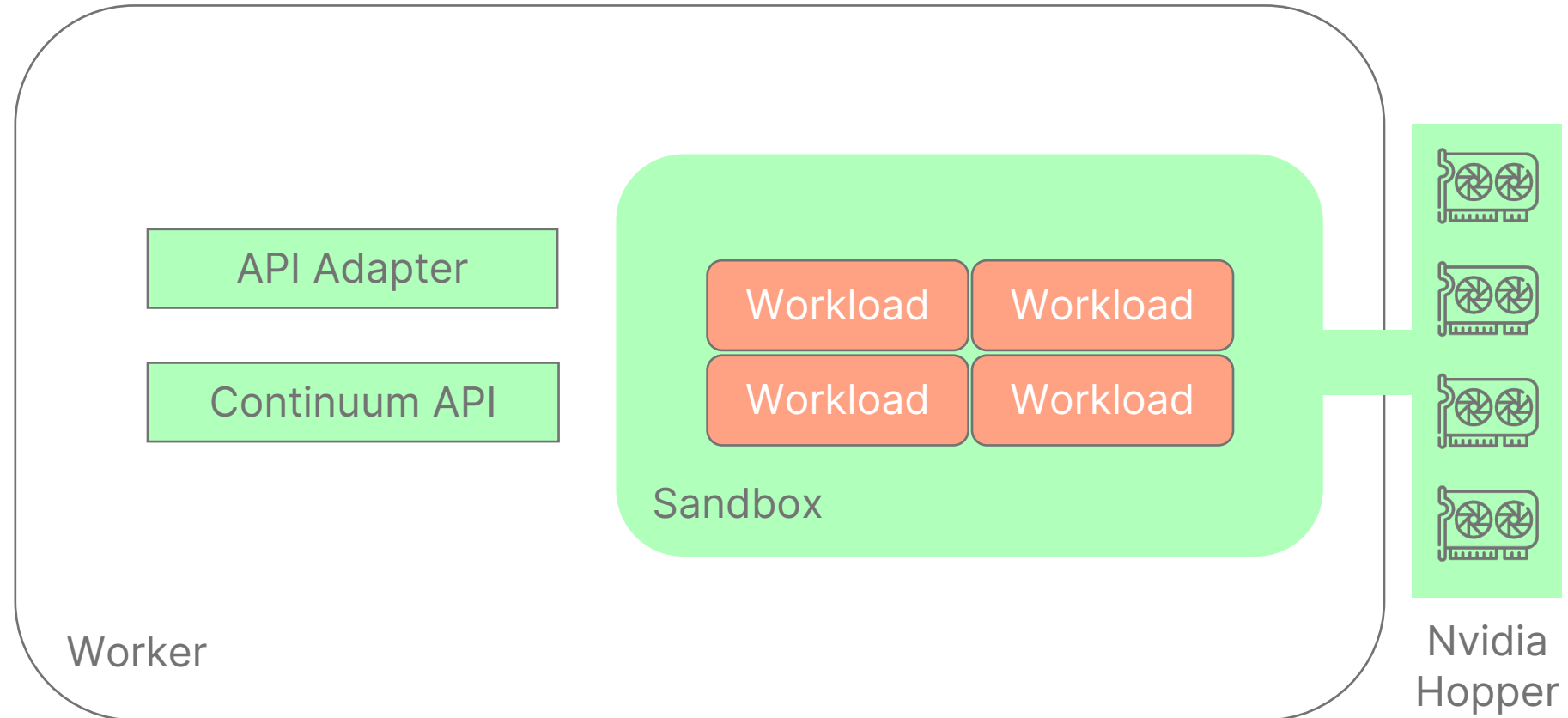# 04 Implementation

# Architecture

# Architecture

# Architecture

Sandbox prevents:
- Network access
- Persistent storage access

# Architecture

API Adapter provides:
- Transparent encryption

API Adapter

Continuum API

Sandbox

Workload  Workload

Workload  Workload

Worker
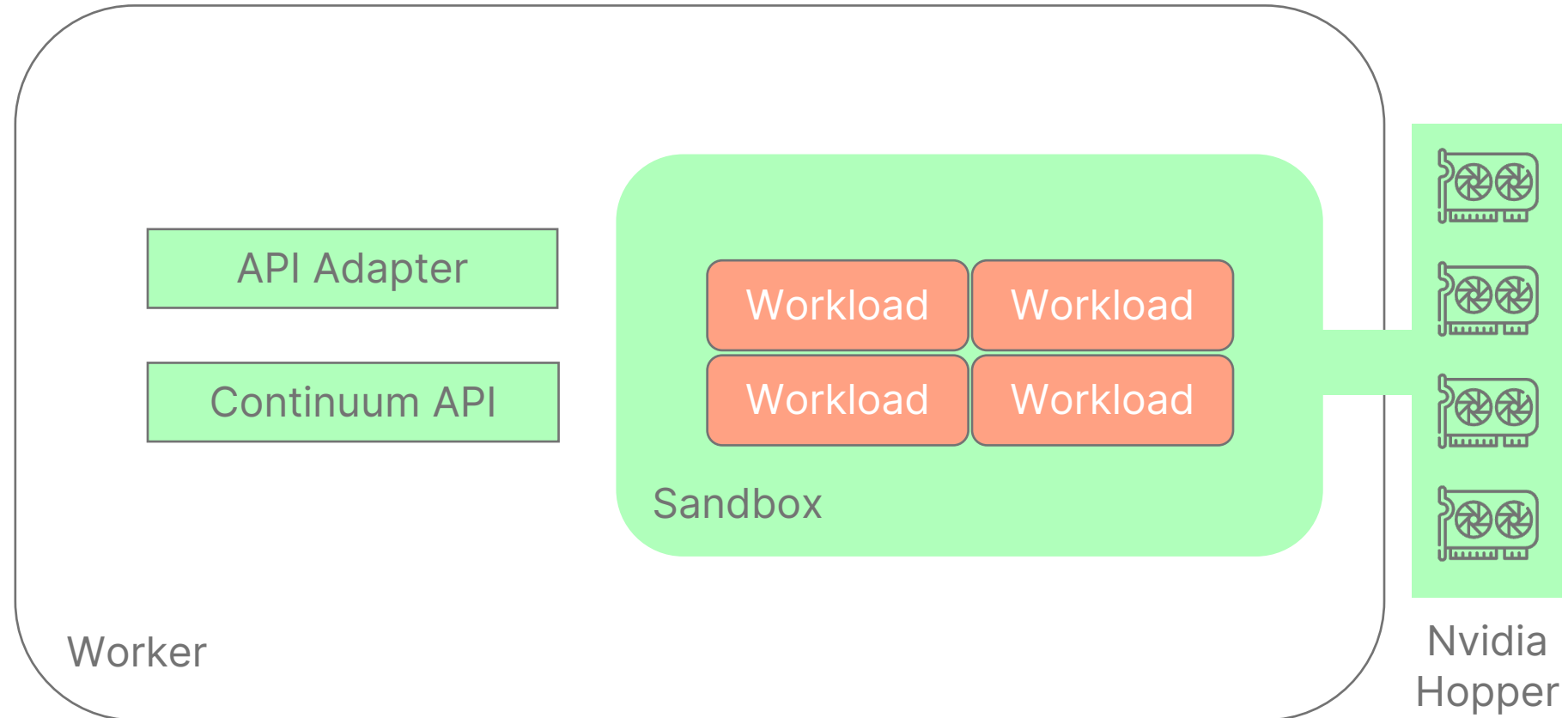
Nvidia Hopper

# Architecture

Continuum API provides:
- Deployment API
- Secret management
- Remote attestation

API Adapter

Continuum API

Workload    Workload

Workload    Workload

Sandbox

Worker

Nvidia Hopper

Edgeless
Systems

05 Demo

Upload Attestation Service Manifest & Secrets

Attestation Service Manifest

Choose File no file selected

Encryption Secrets

Choose File no file selected

# 06 Outlook

# Outlook

- Working with early partners for evaluation
- More Client SDKs
- Open source release in Q2/24.

Edgeless Systems

Thank you!