

The background features a dark blue gradient with abstract, glowing shapes in shades of purple and pink. Two thin, light blue lines intersect to form a large 'A' shape. The text is positioned on the left side of the image.

# AWS re:Invent

DECEMBER 2 - 6, 2024 | LAS VEGAS, NV

AIM301 - NEW

# Accelerate gen AI: Amazon SageMaker HyperPod training plans & recipes

**Gal Oshri**

Prin. Product Manager,  
Amazon SageMaker  
AWS

**Giuseppe A. Porcelli**

Sr. Manager,  
Amazon SageMaker  
AWS

**Babak Pahlavan**

Founder & CEO  
NinjaTech AI

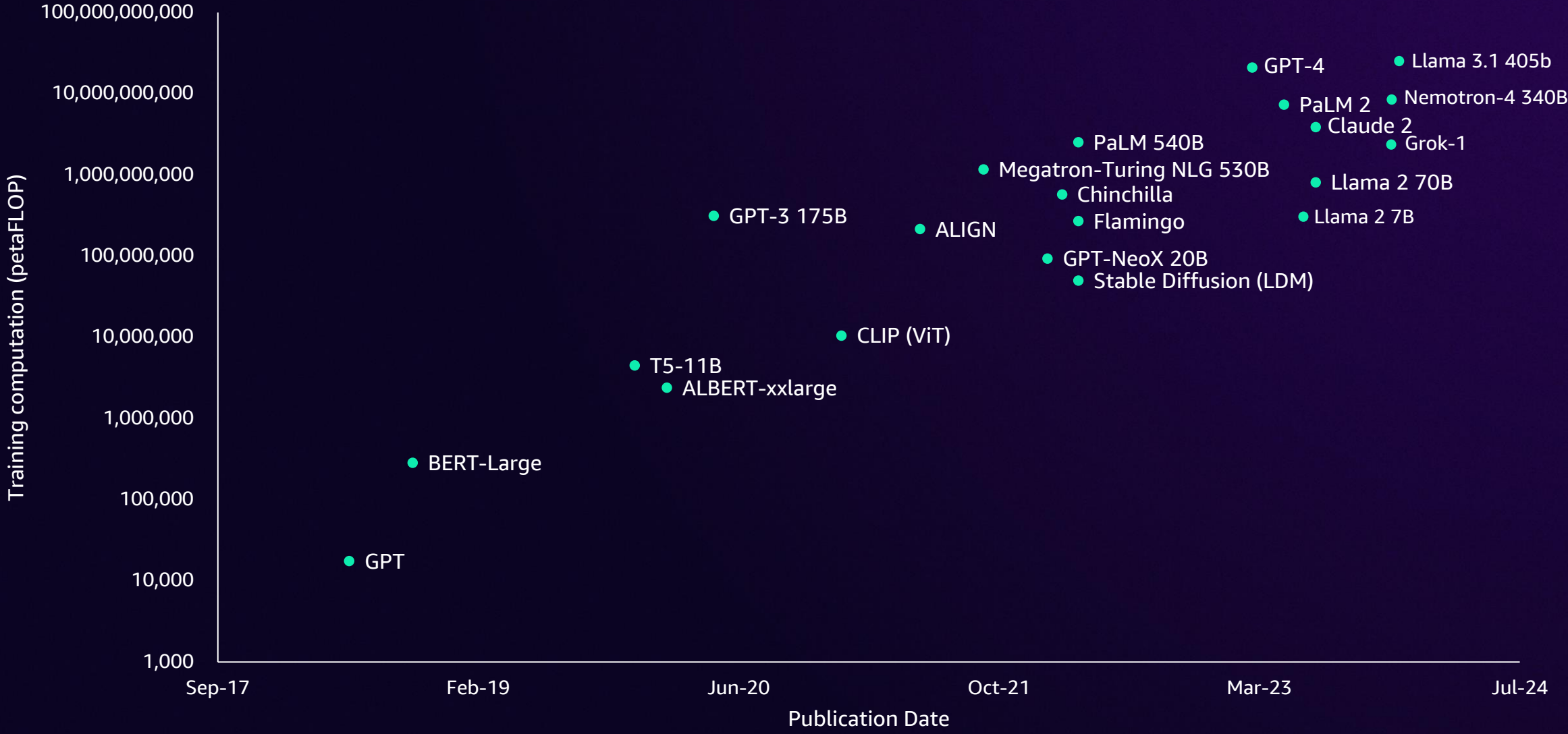


# Agenda

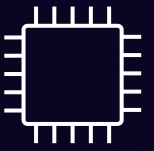
- 01 Challenges with training large-scale models
- 02 Amazon SageMaker HyperPod training plans
- 03 Amazon SageMaker HyperPod recipes
- 04 Demo
- 05 How NinjaTech AI has used Amazon SageMaker HyperPod training plans



# Generative AI model computational demand is growing



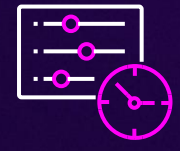
# Challenges with training large-scale models



Hardware



Faults



Timelines



Performance



Cost

# Why Amazon SageMaker HyperPod

REDUCE TRAINING TIME BY UP TO 40% THROUGH RESILIENCY AND PERFORMANCE OPTIMIZATIONS



## Resilient environment

---

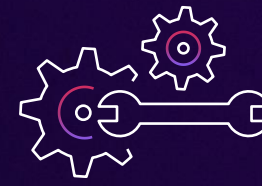
Self-healing clusters reduce training time



## Streamline distributed training

---

SageMaker distributed training libraries improve performance



## Customizable UX and persistent cluster

---

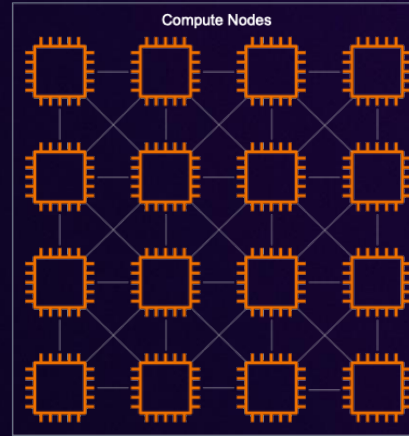
Control over computing environment and workload scheduling

# Setting up compute resources

## On-demand capacity

---

- Flexible
- Variable availability
- Higher-cost option



## Long-term reservations

---

- Utilization
- Predictable availability
- Lower cost

**Challenging to plan training workloads within timelines and budgets**

NEW

# Amazon SageMaker HyperPod flexible training plans

Save weeks of training time and help meet timelines and budgets

- Remove the uncertainty and manual process of capacity reservation
- Automatically set up training infrastructure and clusters
- Ensure FM training meets budgets and timelines





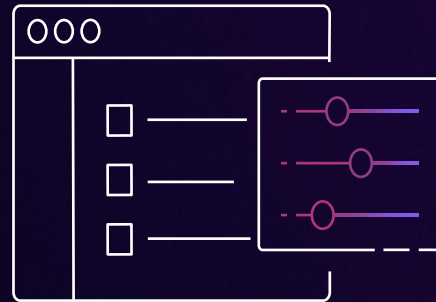
# Get plans that meet your requirements

## Specify your needs



- Instance type
- Number of instances
- Duration of plan
- Start date

## Create a training plan



- Pick recommended offering
- Pay upfront

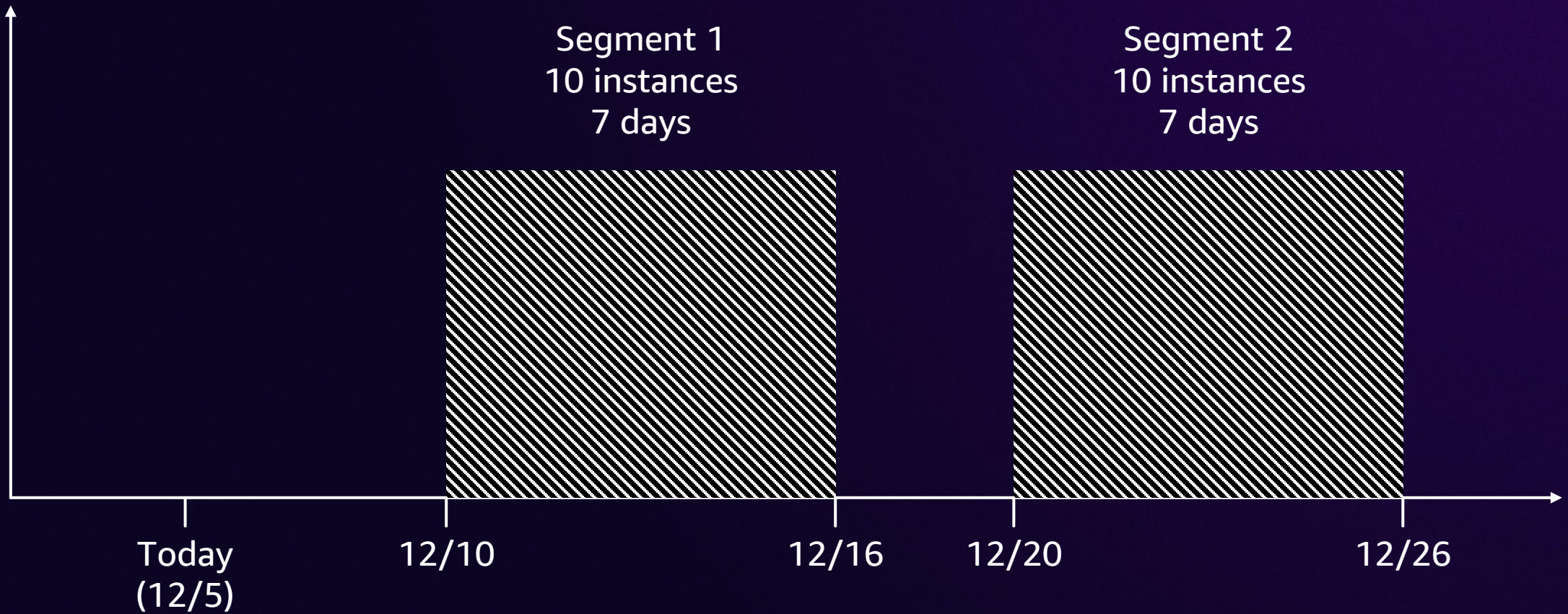
## Use the training plan



- Use the plan in Amazon SageMaker HyperPod to create clusters
- Use the plan in Amazon SageMaker training jobs

# Training plans

**“Create a training plan with 10 instances of ml.p5.48xlarge for 14 days starting 12/10”**



# Training plans with HyperPod

- 1 Create clusters with multiple instance groups
- 2 Specify instance groups to use training plans for compute
- 3 Instance groups scale up when training plan is active
  - ✓ Scale to requested count
  - ✓ Lifecycle scripts executed
  - ✓ Deep health checks
  - ✓ Resiliency

### Create an instance group ✕

**Instance group name**  
Specify a name for this instance group.

**Instance type**

**Quantity**  
Specify the number of instances that you want to allocate to the new instance group.

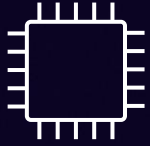
**Instance capacity**

On-demand
On-demand capacity (default)
<b>Training plan</b>
cp-001
cp-002

**Directory path to the on-create lifecycle script**  
Enter the path to the lifecycle configuration script that you want to run on each instance in the instance group after cluster creation. This path should be relative to the S3 path of the lifecycle configuration files you specified under 'S3 path to lifecycle script files'.

Cancel Save

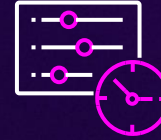
# Benefits of using training plans



Compute access



Resiliency



Predictable timelines



High performance



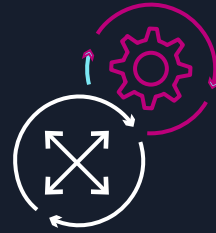
Upfront cost



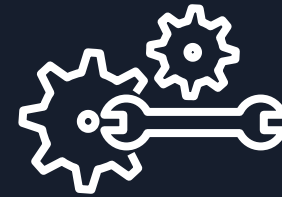
# What generative AI customers are asking for



**Which model  
should I use?**

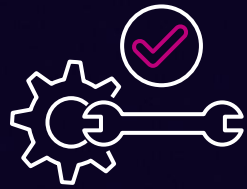


**How do I  
customize  
my model?**



**How can I  
optimize  
training  
performance?**

# Customizing FMs for your business



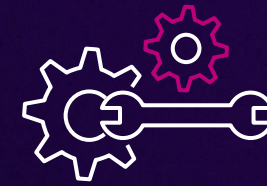
## Fine-tuning

### PURPOSE

Maximizing accuracy  
for specific tasks

### DATA NEED

Small number of  
labeled examples



## Pre-training

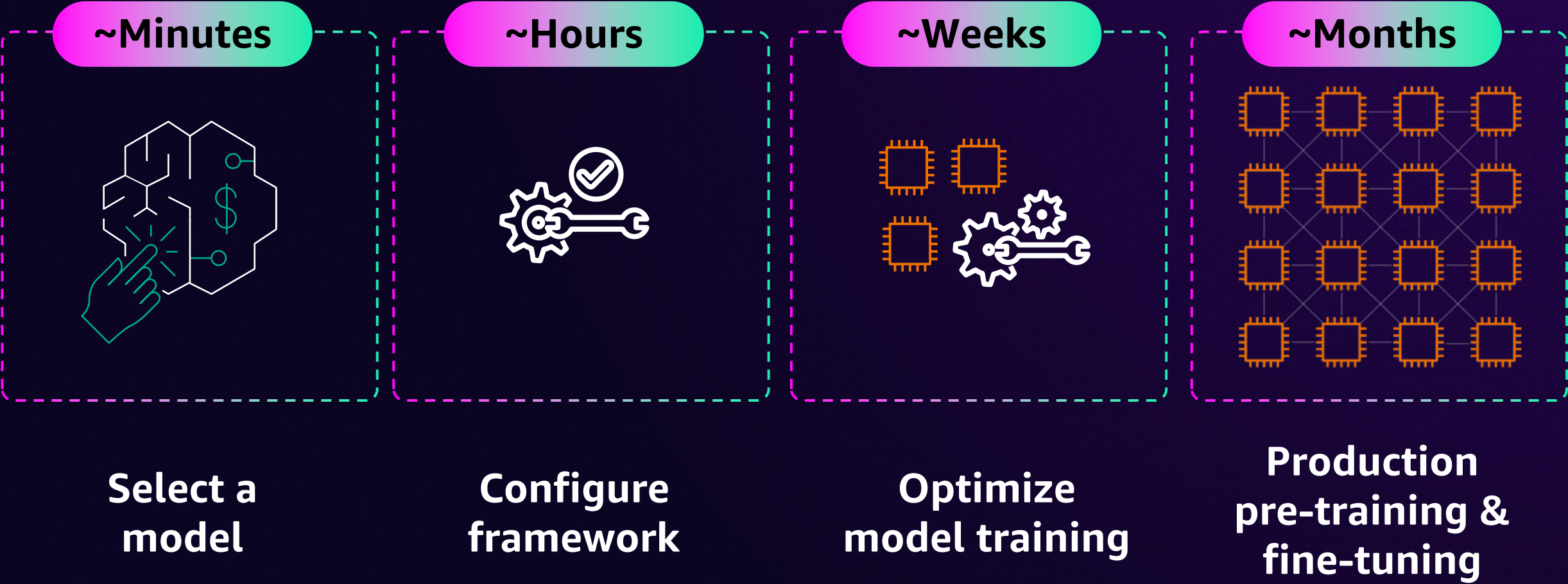
### PURPOSE

Maintaining model  
accuracy for your domain

### DATA NEED

Large number of  
unlabeled datasets

# Optimizing FM pre-training and fine-tuning can take weeks of effort



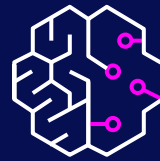
# Training FMs with billions of parameters spans thousands of different training stack configs



3-5 popular training frameworks



Dozens of different tunable parameters



Hundreds of configurable settings



Other performance optimizations like SMP, SMDDP



NOW AVAILABLE

# Amazon SageMaker HyperPod recipes



# Amazon SageMaker HyperPod recipes



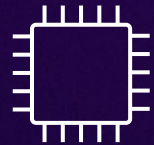
Curated, ready-to-use recipes for pre-training and fine-tuning popular publicly available FMs



Tested and validated for foundational models such as Llama and Mistral



Automatic checkpoints for faster fault recovery and managed end-to-end training loop



Easily switch between GPU-based or Trainium-based instances

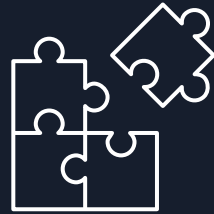
# Getting started in minutes

RUN FM PRE-TRAINING AND FINE-TUNING WITH A SINGLE LINE OF CODE



## Select

Model training and fine-tuning recipe on GitHub



## Setup Prerequisites

Resource limits, AWS credentials, training cluster



## Run the recipe

on Amazon SageMaker HyperPod  
**or**  
Amazon SageMaker training jobs



# Getting started

RUN FM PRE-TRAINING AND FINE-TUNING WITH A SINGLE LINE OF CODE

## Recipes on Amazon SageMaker HyperPod (Slurm)

```
python3 main.py recipes=recipe-name
```

```
run_<model_name>.sh
```

**NeMo-style launcher**  
or  
**Launcher script**

## Recipes on Amazon SageMaker HyperPod (EKS)

```
hyperpod start-job --recipe recipe-name
```

**SageMaker HyperPod CLI**

## Recipes on Amazon SageMaker training jobs

```
model_trainer = ModelTrainer.from_recipe(  
    training_recipe="<recipe_path>",  
    ...  
)  
  
model_trainer.train(wait=False)
```

**SageMaker Python SDK**





# How it works

## AMAZON SAGEMAKER HYPERPOD RECIPES

The screenshot shows the GitHub repository page for 'aws/sagemaker-hyperpod-recipes'. The repository is owned by 'aws' and has 71 branches and 0 tags. The main branch is selected. The repository description is 'Remove setuptools upgrade from installation instructions (#130)'. The repository was last updated 3 days ago by user 'dc5deef' and has 134 commits. The repository contains several folders and files:

Item	Description	Last Updated
.github	Add repolinter to scan source code (#117)	last week
launcher	Only kill Kandinsky docker containers (#127)	4 days ago
launcher_scripts	Mixtral pretrain recipes (#124)	3 days ago
recipes_collection	Mixtral pretrain recipes (#124)	3 days ago
scripts	Add repolinter to scan source code (#117)	last week
template	Merge from bugbash. (#108)	2 weeks ago
tests	Only kill Kandinsky docker containers (#127)	4 days ago
.coveragerc	Increase coverage threshold from 80 to 90 (#105)	2 weeks ago

Open source implementation

Launcher scripts and recipes collection

Built on **NVIDIA NeMo foundations**  
(launcher, configuration hierarchy)

Over **30 recipes** to get started

SageMaker-optimized  
models (GPU)

Neuron-optimized  
models (Trainium)

Native NeMo models

Custom models



# How it works

## EXAMPLE RECIPE

sagemaker-hyperpod-recipes / recipes\_collection / recipes / training / llama / hf\_llama3\_8b\_seq16384\_gpu\_p5x16\_pretrain.yaml

Code Blame 142 lines (128 loc) · 3.54 KB

```
49
50 # Model training configs
51 model:
52   model_type: llama_v3
53 # Base configs
54 train_batch_size: 4
55 val_batch_size: 1
56 seed: 12345
57 grad_clip: 1.0
58 log_reduced_training_loss: True
59
60 # Memory saving /distributed training configs
61 tensor_model_parallel_degree: 4
62 expert_model_parallel_degree: 1
63 context_parallel_degree: 2
64 moe: False
65 activation_checkpointing: False
66 activation_loading_horizon: 1
67 delayed_param: True
68 offload_activations: False
```

## Hydra-based configuration

sagemaker-hyperpod-recipes / recipes\_collection / config.yaml

Code Blame Executable File · 35 lines (27 loc) · 1003 Bytes

```
1 defaults:
2   - _self_
3   - cluster: slurm # set to `slurm`, `k8s` or `sm_jobs`, depending on the desired cluster
4   - recipes: training/llama/hf_llama3_8b_seq8192_gpu # select desired config inside the training directory
5   - override hydra/job_logging: stdout
6
7 cluster_type: slurm # bcm, bcp, k8s or sm_jobs. If bcm, k8s or sm_jobs, it must match - cluster above.
8 # If using sm_jobs cluster_type, set sm_jobs_config. See cluster/sm_jobs.yaml for example.
9
10 hydra:
11   run:
12     dir: .
13     output_subdir: null
14
15 debug: False
16
17 instance_type: p5.48xlarge
18 # TODO: remove
19 data_dir: null # Location to store and read the data.
20 base_results_dir: null # Location to store the results, checkpoints and logs.
21
22 container: null
```

# How it works

## AMAZON SAGEMAKER HYPERPOD TRAINING ADAPTER FOR NEMO

The screenshot shows the GitHub repository page for 'aws/sagemaker-hyperpod-training-adapter-for-nemo'. The repository is owned by 'aws' and has 17 watchers. It is currently on the 'main' branch, with 99 other branches and 0 tags. The repository has 4861594 commits and 191 commits. The repository is organized into several folders and files, each with a corresponding commit message and date:

File/Folder	Commit Message	Date
.github	Rename Adapter (#192)	3 days ago
examples	Fixing adaptor import in custom_pretrain (#195)	2 days ago
requirements	Add profiling dependencies file, install all dependencies (i...	last week
scripts	Cleanups for Adaptor (#185)	5 days ago
src/hyperpod_nemo_adapter	Add option to save the final full checkpoint directly (#193)	3 days ago
tests	fix test_patched_LFA2__init__ (#194)	3 days ago
.gitignore	unit tests for sagemaker_base_model	4 months ago
.pre-commit-config.yaml	Moe support (#57)	2 months ago

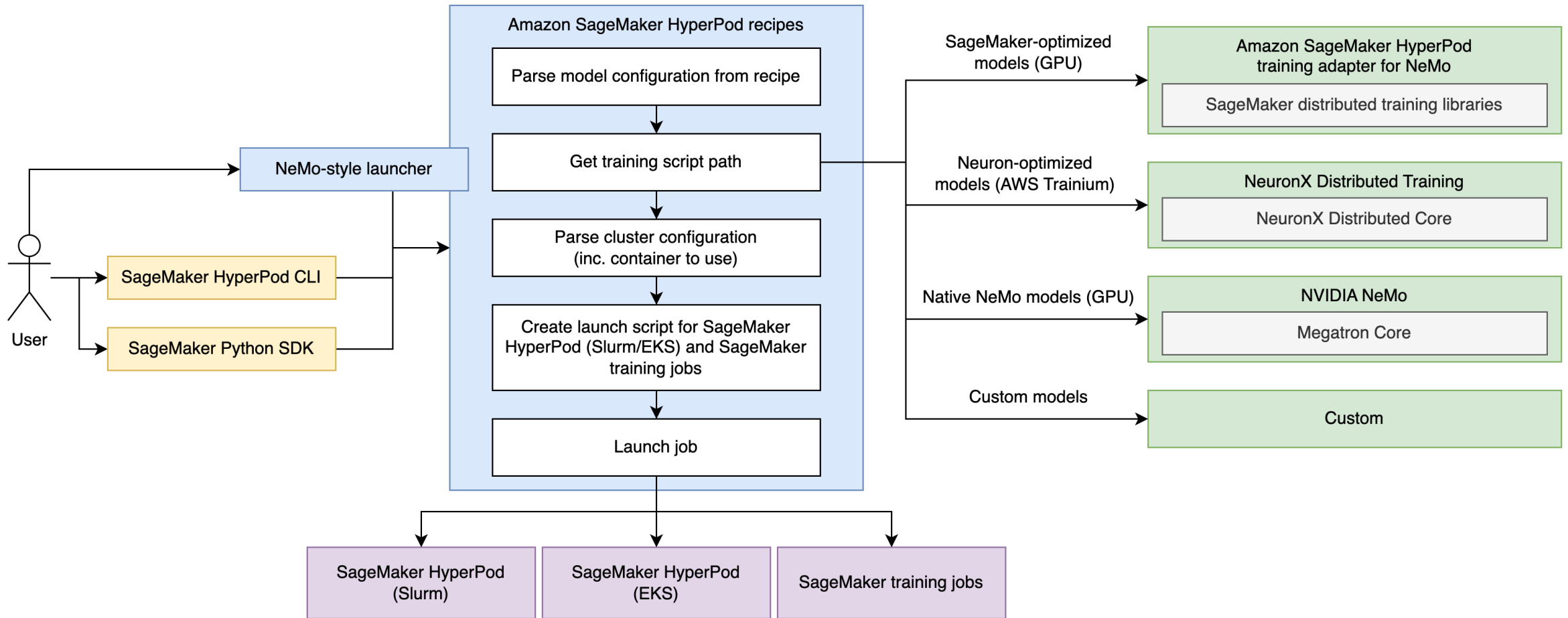
Defines training loop, data loading, and **automatic checkpointing code** for SageMaker-optimized models (GPU)

Implements optimized **SageMaker distributed training strategies**

Built on NVIDIA NeMo backend

# How it works

## JOB SUBMISSION WORKFLOW





# Amazon SageMaker HyperPod training plans and recipes DEMO





# How NinjaTech AI has used SageMaker HyperPod training plans





# NinjaTech AI

Your All-In-One AI Agent for  
Unlimited Productivity

“Netflix of gen AI”: Democratizing access to world’s  
best AI models & AI skills

Babak Pahlavan  
Founder, CEO & CPO  
December 2024



AI



S.



# Ninja

## All-in-one unlimited AI skills

Market Today

One subscription

Multiple subscriptions & models


Unlimited skills & model access



Code generation



Image creation



Online research

And more




Advance reasoning



Image creation



Code generation



Online research



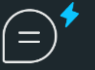
File analysis




Vision analysis




Multilingual processing










Fast inference



More skills  
Scheduling,  
Travel and etc



Access to external models



+more



# SuperAgent AI Assistant

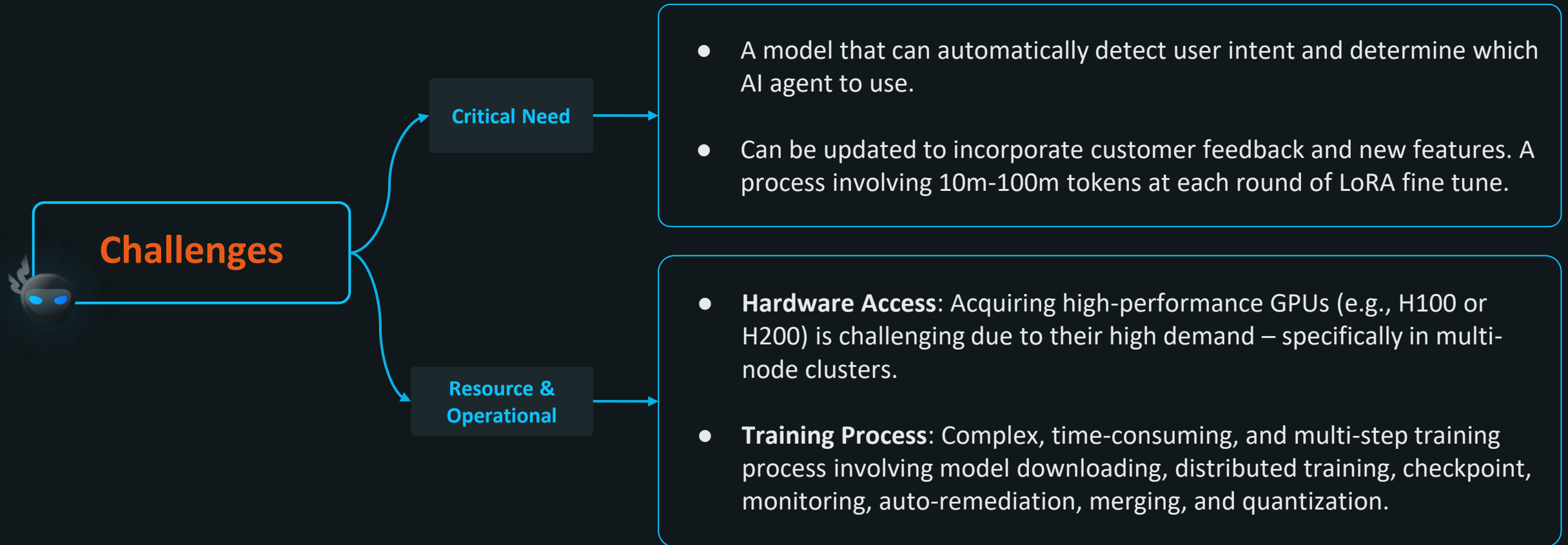
## Platform & Model agnostic all-in-one solution

Get tasks done using the state of art AI agents and world's best LLMs

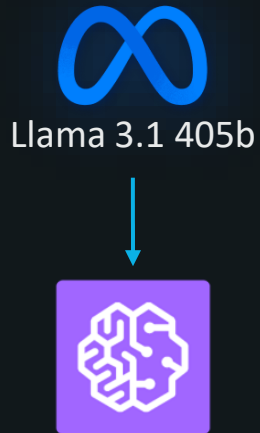


# Our challenges for training: Fine-tuning our 405B model, reducing costs, and automating

**Mission:** To be the one-stop shop for all gen AI solutions, providing a seamless experience for users to access multiple AI agents.



## Planning reliable + affordable compute in advance – HyperPod training plans delivered on this challenge



### Fine-Tuning on HyperPod –

Utilized HyperPod training plans to fine-tune our 405B model on HyperPod, leveraging the power of predictable access to Nvidia H100 GPU capacity.

Efficient

Cost-Effective

High-  
Performance

**Customized Training Plans:** Tailored training plan that matched our specific compute and timeline needs.

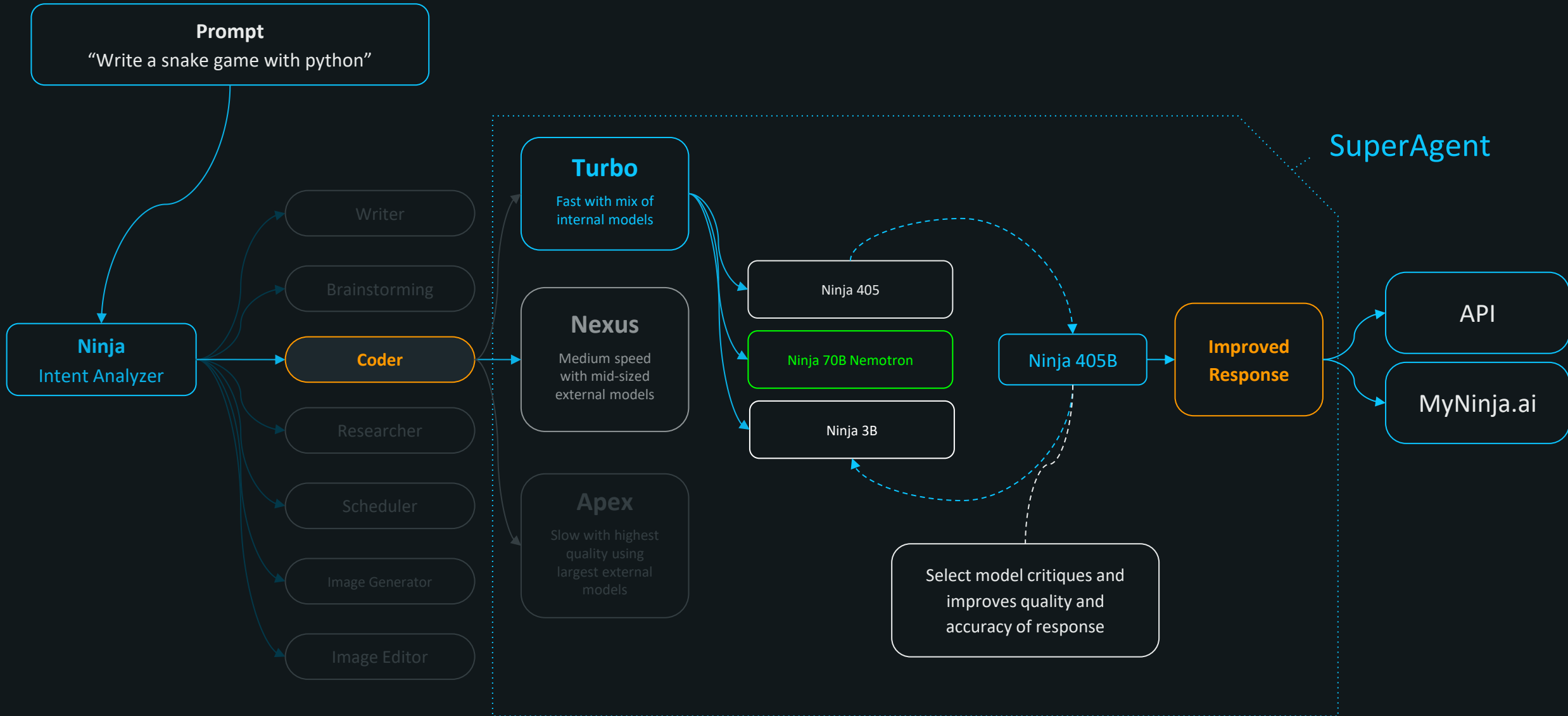
**Automated Execution + Resource Management:** SageMaker HyperPod training plans automated the execution of our training plan, provisioning required resources, setting up of infrastructure, and recovering from failures.

**Flexible Reservation Options:** Flexible reservation options – duration and instance quantity that suited our needs.

**Effortless Integration With HyperPod Clusters:** Seamlessly connected SageMaker HyperPod training plans with Amazon SageMaker HyperPod clusters, which offers expertly managed GPU clusters with advanced monitoring and automatic recovery capabilities.

# SuperAgent

## Real-time Inference Level Optimization (ILO) technology

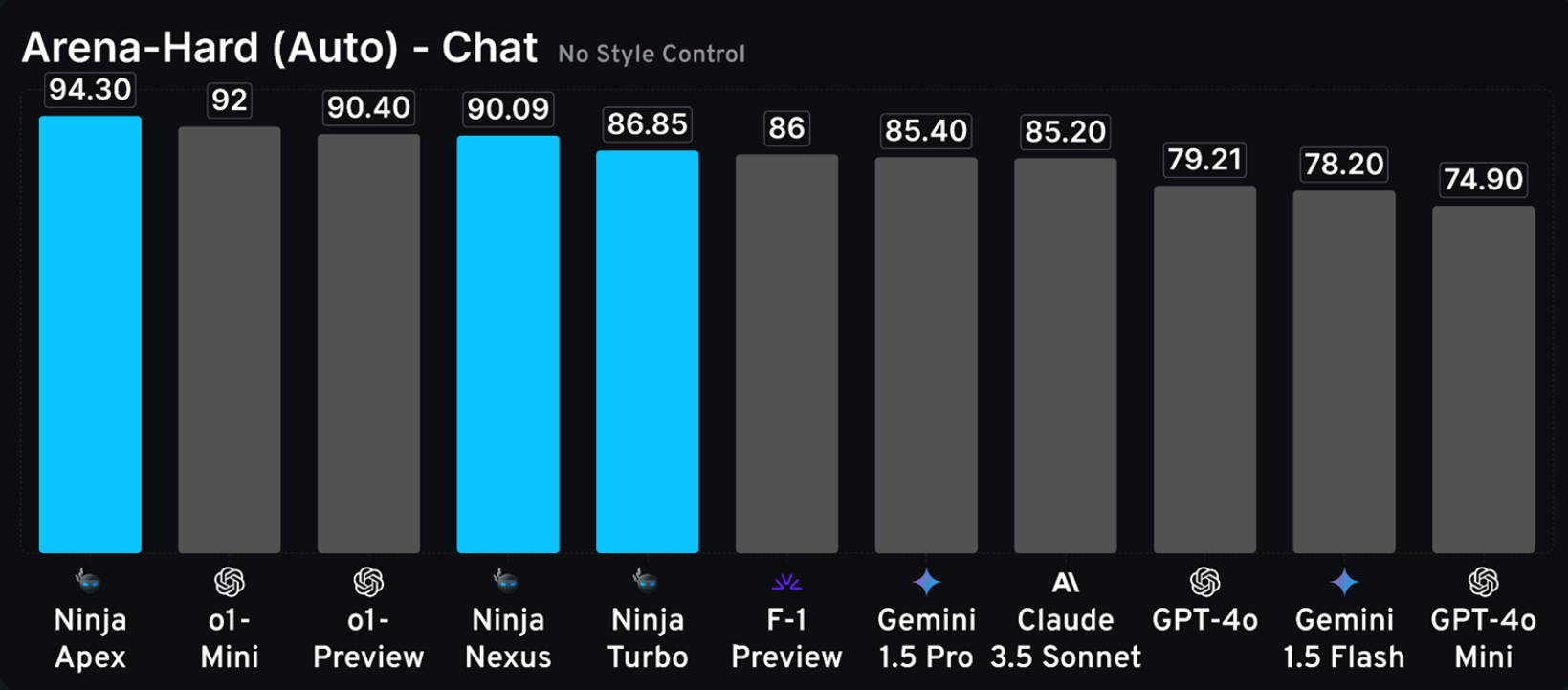




SuperAgent: Agentic Inference level optimization (More tokens → higher accuracy and better results )

# Ninja achieves SOTA (State-of-the-art)

- 1. SuperAgent variations
  - a. Turbo (Fast with internal models)
  - b. Nexus (Medium speed with mid-sized external models)
  - c. Apex (Slow with largest external models)
  
- 2. Accuracy optimizations: Skill based accuracy improvement algorithms (Example: Coding questions will use a different algorithms compared to writing or Research questions).



Demo

# Multi-node trained: Making Llama models talkative

Hi, Aiden!

Ask anything or use / to choose a specific agent



Prompt Library



Improve Prompt



Ninja Agent



IMAGE GENERATOR

Create a flash snap  
image

RESEARCHER

Most popular  
programming  
languages

CODER

Write a Python  
script

WRITER

Write a sincere  
apology

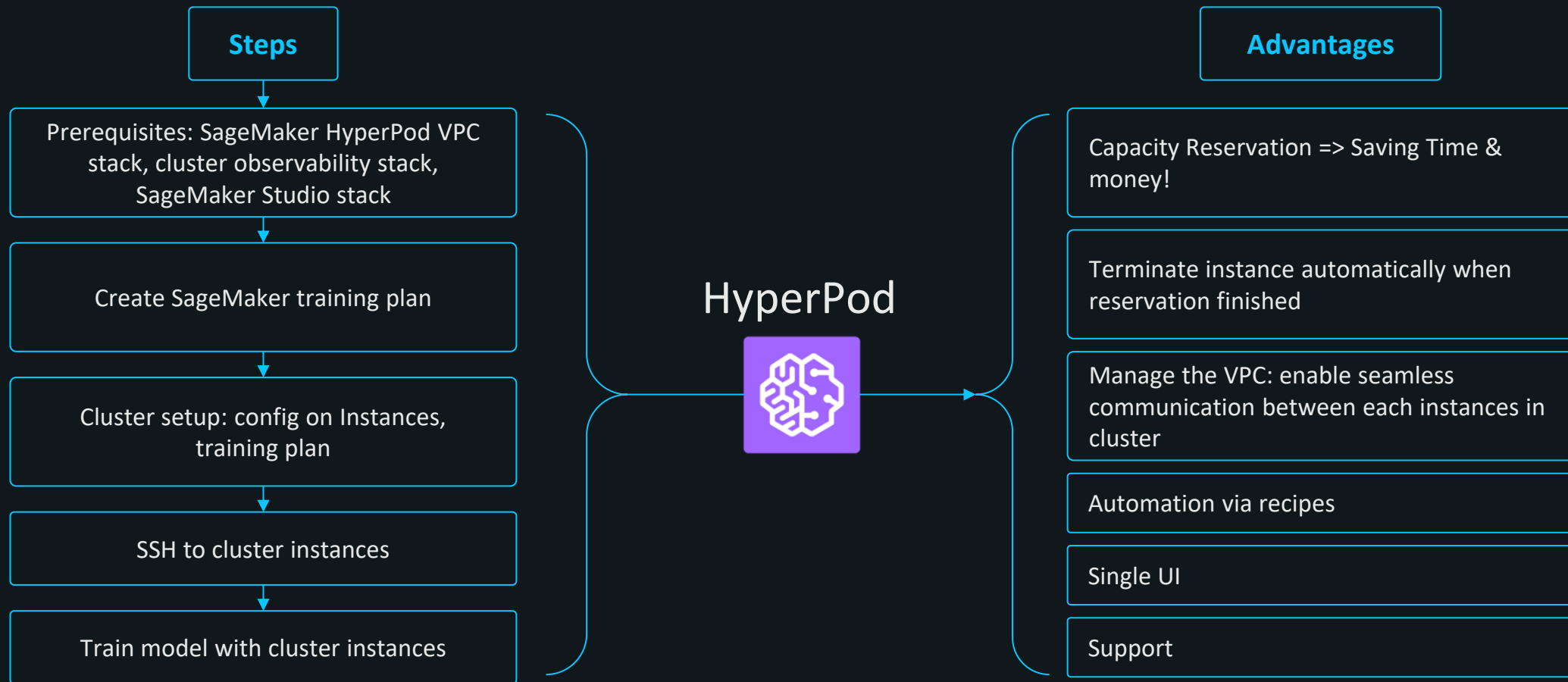
SCHEDULER

What's my  
schedule next  
week?



## Recap

# HyperPod is a game changer, especially for start-ups: It allows seamless multi-node large-scale training.



Try Ninja today!



[MyNinja.ai](https://MyNinja.ai)





# Resources & references

## Service page

<https://aws.amazon.com/sagemaker/hyperpod/>

## Documentation

<https://docs.aws.amazon.com/sagemaker/latest/dg/reserve-capacity-with-training-plans.html>

<https://docs.aws.amazon.com/sagemaker/latest/dg/sagemaker-hyperpod-recipes.html>

## Blogs

<https://aws.amazon.com/blogs/aws/meet-your-training-timelines-and-budgets-with-new-amazon-sagemaker-hyperpod-flexible-training-plans/>

<https://aws.amazon.com/blogs/aws/accelerate-foundation-model-training-and-fine-tuning-with-new-amazon-sagemaker-hyperpod-recipes/>

## Announcements

<https://aws.amazon.com/about-aws/whats-new/2024/12/amazon-sagemaker-hyperpod-flexible-training-plans/>

<https://aws.amazon.com/about-aws/whats-new/2024/12/amazon-sagemaker-hyperpod-recipes/>



# Thank you!

**Gal Oshri**

galoshr@amazon.com

**Giuseppe Angelo Porcelli**

gianpo@amazon.com

**Babak Pahlavan**

babak@ninjatech.ai



Please complete the session survey in the mobile app