# AWS re:Invent

DECEMBER 2 – 6, 2024 | LAS VEGAS, NV

AIM380

# High performance distributed model training with Amazon SageMaker

**Anirudh Viswanathan**

Sr. Product Manager, Technical,
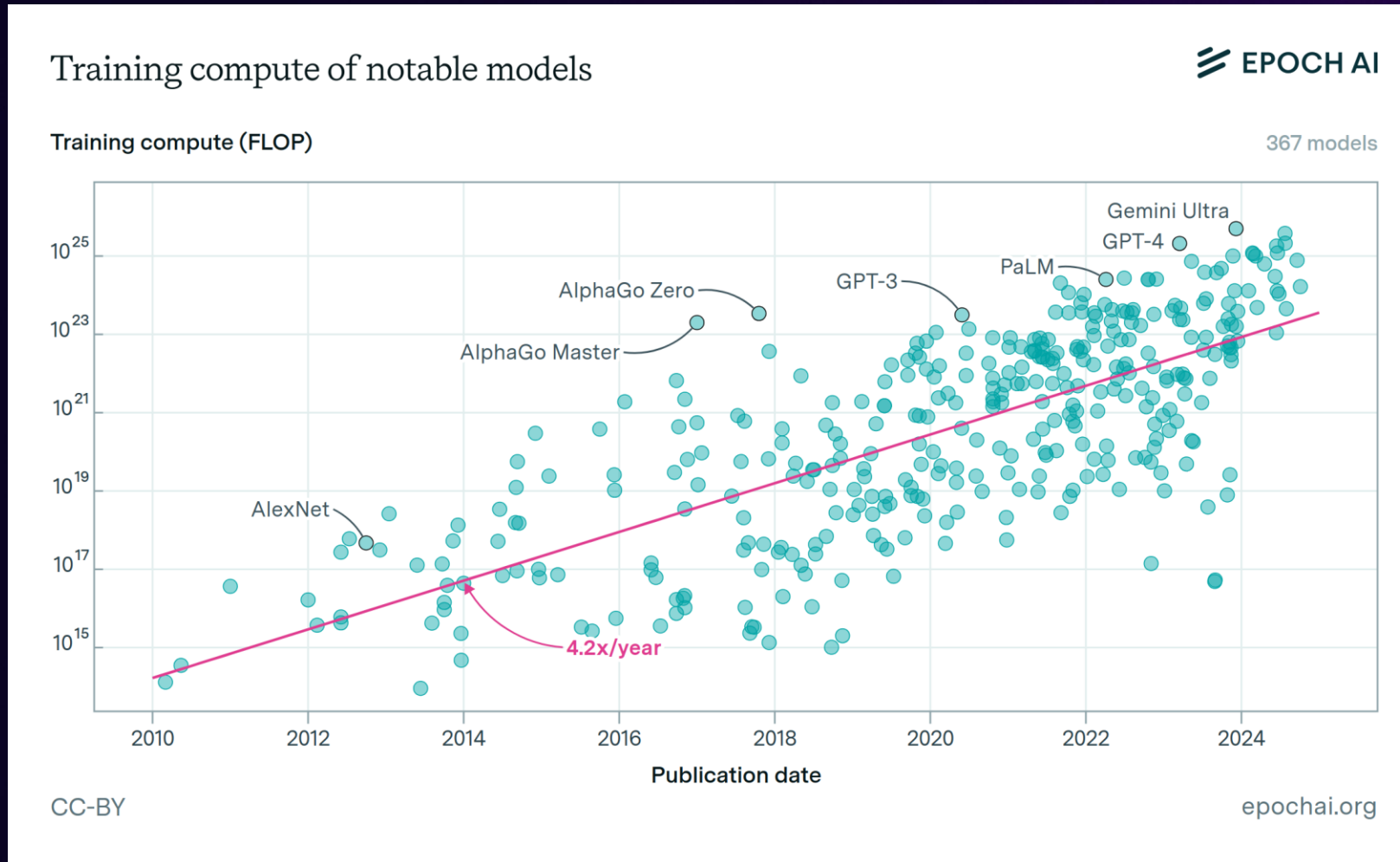Amazon SageMaker
AWS

**Sanjay Dorairaj**

Software Development Manager,
Amazon SageMaker
AWS
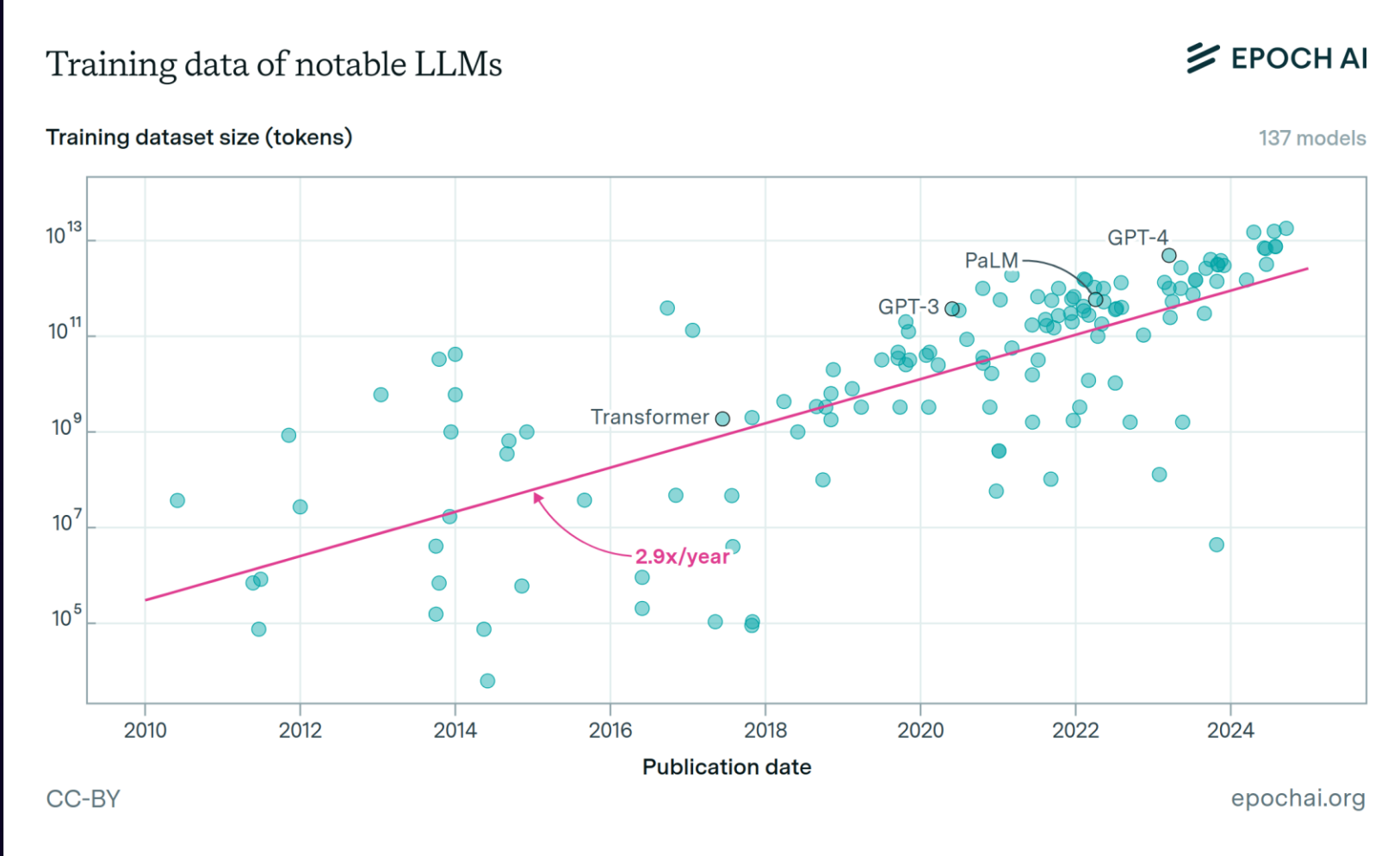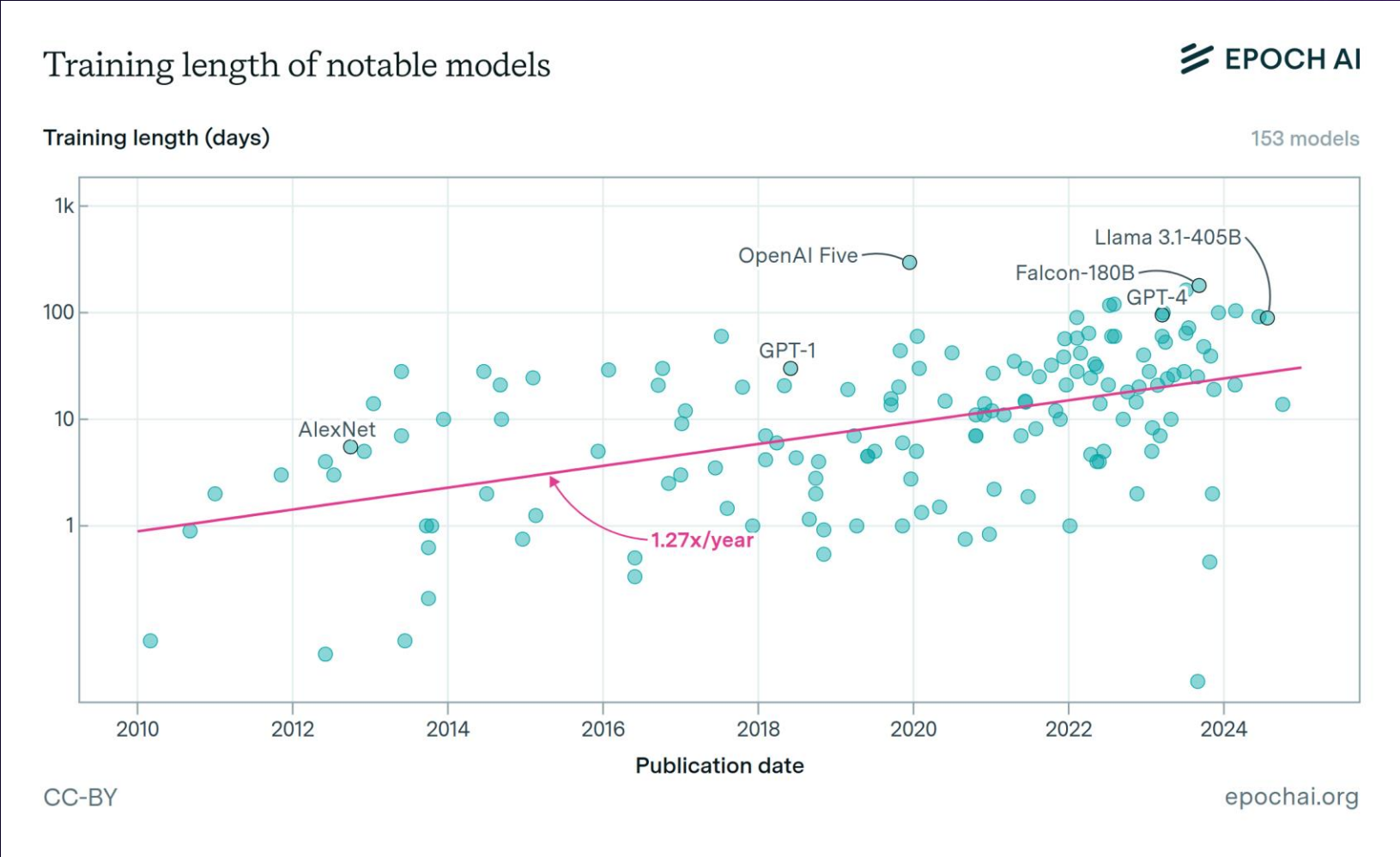
**Antonio Ginart, Ph. D.**

Lead Scientist,
AI Research
Salesforce

# Training compute of foundational models is growing by ~5x per year, doubling ~6 months
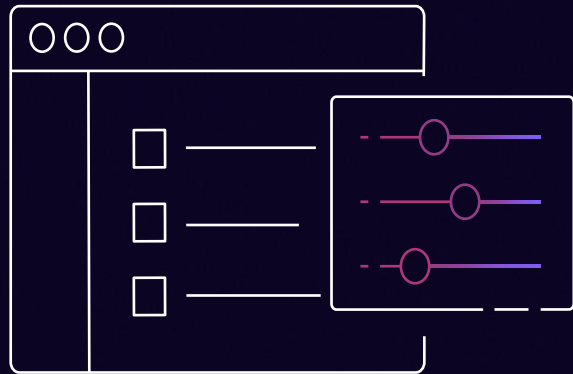
# Dataset sizes are doubling every 8 months



Training data of notable LLMs — EPOCH AI

Training dataset size (tokens) — 137 models

2.9x/year

Labeled points: Transformer, GPT-3, PaLM, GPT-4

Publication date

epochai.org

# Time to market spans months of investment and continues to increase



Training length of notable models — EPOCH AI

Training length (days) — 153 models

Models labeled: AlexNet, GPT-1, OpenAI Five, GPT-4, Falcon-180B, Llama 3.1-405B

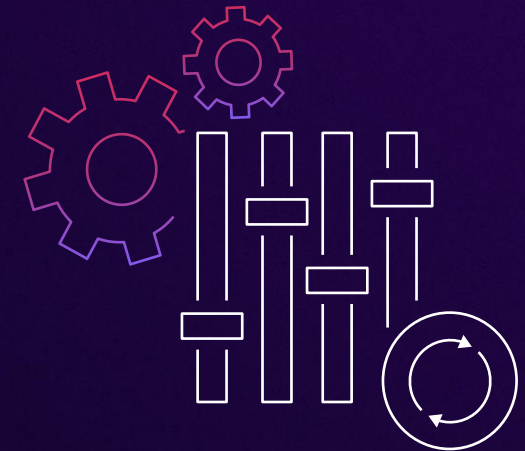1.27x/year

Publication date

epochai.org

# Unique challenges to manage hardware resources efficiently for large-scale FM training

Cluster provisioning and management

Infrastructure stability

Distributed training performance

# Amazon SageMaker HyperPod

Scale and accelerate generative AI model development across thousands of AI accelerators

Designed for scale

Reduce model training time by 40%

Resilient environment

Advance observability and control

# Top AI companies use HyperPod to train and deploy models

# Optimizing the distributed training software stack

ILLUSTRATIVE OSS STACK

**Framework**

| PyTorch Distributed (FSDP, DTensor, DCP) | NeMo & Megatron-Core |
| --- | --- |
| | PyTorch Lightning |

PyTorch Core

**Network and hardware programming**

| NCCL | |
| --- | --- |
| AWS OFI NCCL | Accelerators SDK & libs (CUDA, Neuron) |
| libfabric | |

**Device**

| Accelerator driver (GPU, TRN) | EFA device & kernel driver |

# Optimizing the distributed training software stack

**Recipes**

SageMaker HyperPod recipes

**Framework**

SageMaker model parallelism (SMP)

PyTorch Lightning

PyTorch Core

NCCL

AWS OFI NCCL

libfabric

Accelerators SDK & libs (CUDA, Neuron)

**Device**

Accelerator driver (GPU, TRN)

EFA device & kernel driver

# Optimizing the distributed training software stack

| | |
|---|---|
| **Recipes** | SageMaker HyperPod recipes |

| | | |
|---|---|---|
| **Framework** | SageMaker model parallelism (SMP) | PyTorch Lightning |

PyTorch Core

| | |
|---|---|
| NCCL | |
| AWS OFI NCCL | Accelerators SDK & libs (CUDA, Neuron) |
| libfabric | |

| | |
|---|---|
| **Device** | Accelerator driver (GPU, TRN) | EFA device & kernel driver |

# PyTorch fork with best-in-class, composable training techniques that are mutually compatible

Megatron

Transformer engine

Tensor parallelism

Context parallelism

Mixture of experts & expert parallelism

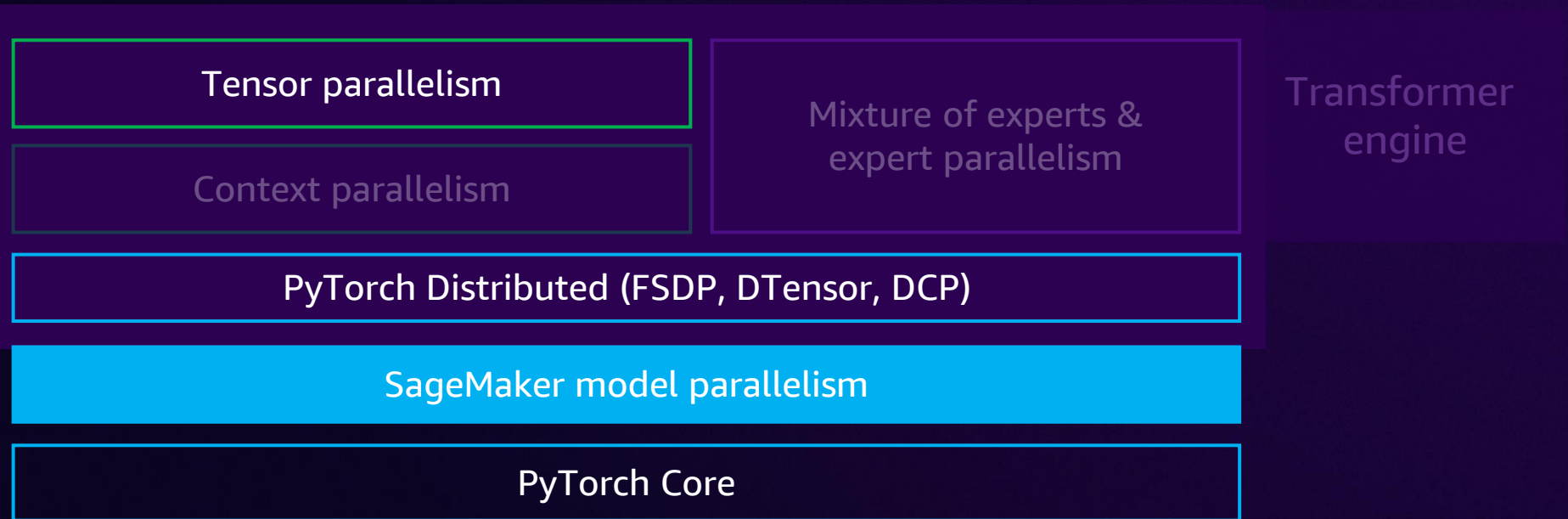PyTorch Distributed (FSDP, DTensor, DCP)

SageMaker model parallelism

PyTorch Core

**Also includes: Delayed parameter initialization, activation checkpointing, activation offloading, etc.**

# PyTorch fork with best-in-class, composable training techniques that are mutually compatible

Megatron

Tensor parallelism

Context parallelism

Mixture of experts & expert parallelism

Transformer engine

PyTorch Distributed (FSDP, DTensor, DCP)

SageMaker model parallelism

PyTorch Core

**Supports training with TP + FSDP**

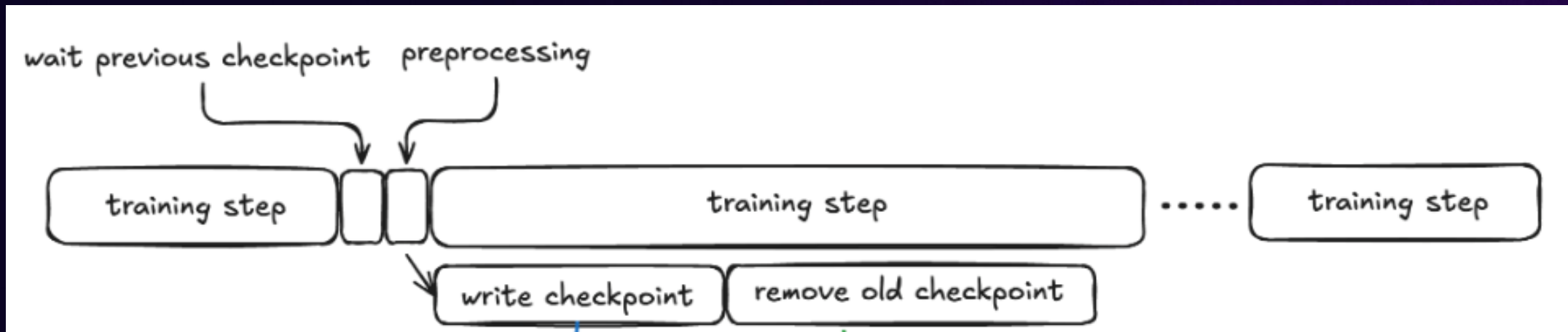# PyTorch fork with best-in-class, composable training techniques that are mutually compatible

Megatron

| Tensor parallelism | Mixture of experts & expert parallelism | Transformer engine |
|---|---|---|

Context parallelism

PyTorch Distributed (FSDP, DTensor, DCP)

SageMaker model parallelism

PyTorch Core

**Supports training with CP + FSDP, FSDP + FP8 etc.**

# Async checkpointing for accelerated training

- Compatible with native PyTorch Distributed Checkpoint (DCP)



**Optimizations**

- Smart metadata caching system for faster checkpoint saving

- Native Amazon S3 support to use S3 links as the checkpoint saving/loading destination

# Optimizing the distributed training software stack

**Recipes**

SageMaker HyperPod recipes

**Framework**

| SageMaker model parallelism (SMP) | PyTorch Lightning |
|---|---|

PyTorch Core

| NCCL | |
|---|---|
| AWS OFI NCCL | Accelerators SDK & libs (CUDA, Neuron) |
| libfabric | |

**Device**

| Accelerator driver (GPU, TRN) | EFA device & kernel driver |
|---|---|

# Audience poll

How many of you are currently pre-training and fine-tuning LLMs?

How many of you are using SageMaker for LLM training?

How many of you plan to start LLM training in the next 6 months?

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

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



~minutes

**Select a model**

~hours

**Configure framework**

~weeks

**Optimize model training**

~months

Compute Nodes

**Production pre-training & fine-tuning**

# Why FM training falters: The technical bottleneck

- 70+ pre-training parameters and 100+ fine-tuning parameters
- Each parameter choice triggers cascading updates to other parameters
- Default choices may be suboptimal
- Leads to cost overruns, missed deadlines, and reduced productivity

# Amazon SageMaker HyperPod recipes

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

Tested and validated by AWS for foundational models such as Llama & 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

**Select**

a model training and
fine-tuning recipe
on GitHub

**Set up
prerequisites**

Resource limits,
AWS credentials,
a training cluster

**Run the recipe**

on Amazon SageMaker
HyperPod
*or*
SageMaker
training jobs

# Getting started

https://github.com/aws/
sagemaker-hyperpod-recipes

# How it works

## AMAZON SAGEMAKER

## HYPERPOD TRAINING RECIPES REPOSITORY

Open source implementation that includes launcher and recipes collection

Built on NeMo foundations (launcher, configuration hierarchy)

Over 30 recipes with different configurations

| SageMaker optimized models (GPU) | AWS Neuron optimized models (Trainium) |
|---|---|
| Native NeMo models | Custom models |

---

📖 README    💬 Code of conduct    ⚖️ Apache-2.0 license    ⚖️ Security

# Amazon SageMaker HyperPod recipes

## Overview

Amazon SageMaker HyperPod recipes help customers get started with training and fine-tuning popular publicly available foundation models in just minutes, with state-of-the-art performance. The recipes provide a pre-configured training stack that is tested and validated on Amazon SageMaker.

Please see Amazon SageMaker HyperPod recipes for documentation.

The recipes support Amazon SageMaker HyperPod (with Slurm or Amazon EKS for workload orchestration) and Amazon SageMaker training jobs.

Amazon SageMaker HyperPod recipes include built-in support for:

- Model parallelism - tensor parallelism and context parallel
- Automated distributed checkpointing
- Distributed optimizer
- Accelerators: NVIDIA H100 (ml.p5), NVIDIA A100 (ml.p4), and AWS Trainium (ml.trn1)
- Fine-tuning: Full, QLoRA, LoRA
- AWS Instances: ml.p5.48xlarge, ml.p4d.24xlarge, and ml.trn1.32xlarge instance families
- Supported Models: Llama, Mistral, Mixtral models
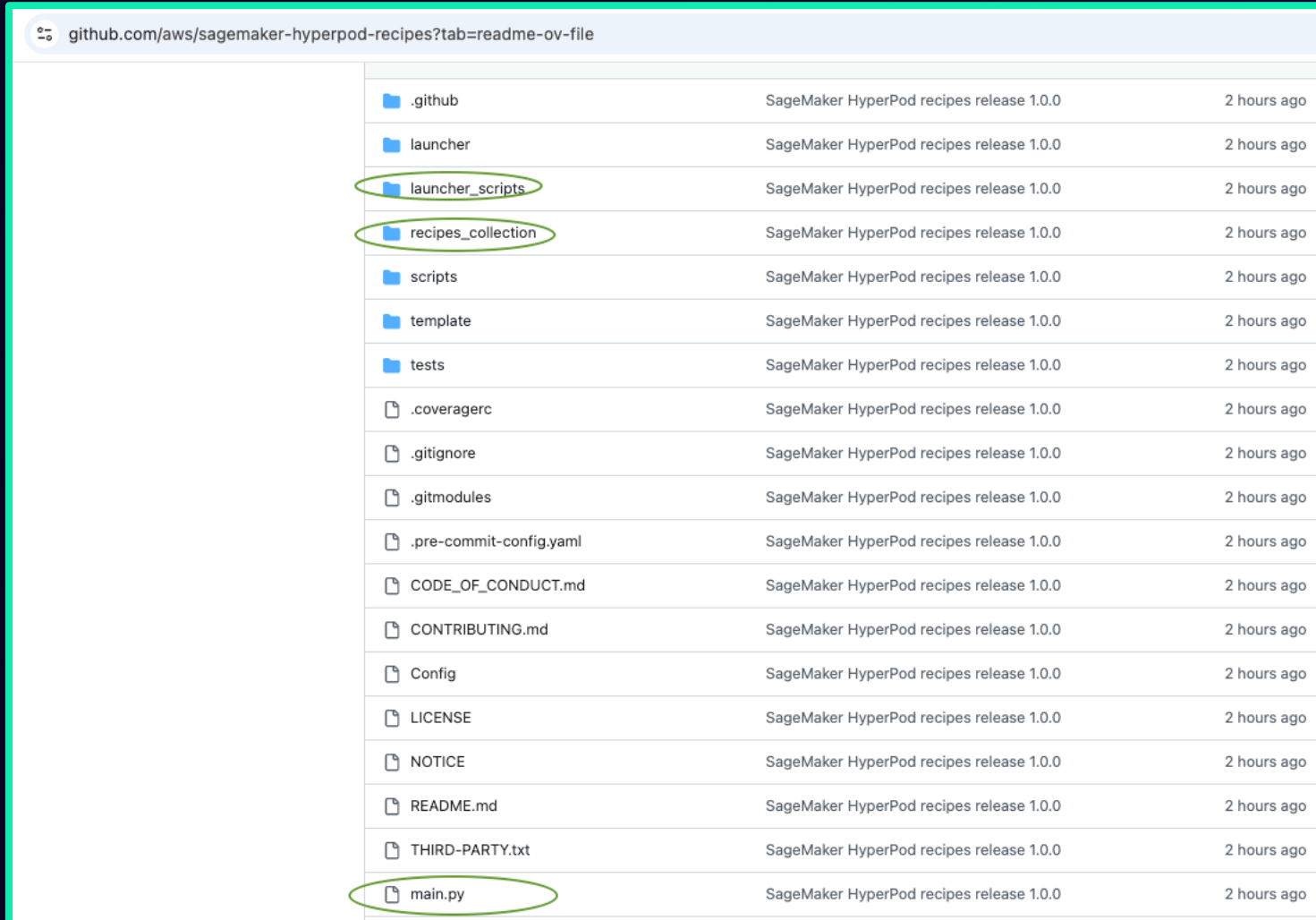- Model Evaluation: Tensorboard

## Model Support

### Pre-Training

List of specific pre-training recipes used by the launch scripts.

| Source | Model | Size | Sequence length | Nodes | Instance | Accelerator | Recipe | Script |
|---|---|---|---|---|---|---|---|---|
| Hugging Face | Llama 3.2 | 11b | 8192 | 4 | ml.p5.48xlarge | GPU H100 | link | link |

# Contents of the repository

github.com/aws/sagemaker-hyperpod-recipes?tab=readme-ov-file

| | | |
|---|---|---|
| .github | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| launcher | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| launcher_scripts | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| recipes_collection | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| scripts | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| template | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| tests | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| .coveragerc | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| .gitignore | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| .gitmodules | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| .pre-commit-config.yaml | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| CODE_OF_CONDUCT.md | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| CONTRIBUTING.md | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| Config | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| LICENSE | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| NOTICE | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| README.md | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| THIRD-PARTY.txt | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |
| main.py | SageMaker HyperPod recipes release 1.0.0 | 2 hours ago |

**launcher_scripts** contains pre-configured bash scripts for model training

**recipes_collection** contains Hydra-based YAML training and fine-tuning recipes

**main.py** is a Nemo-style launcher

All dependencies for HyperPod recipes are present in a docker container and an enroot-based filesystem

# Getting started: Step 1 – Pick a recipe

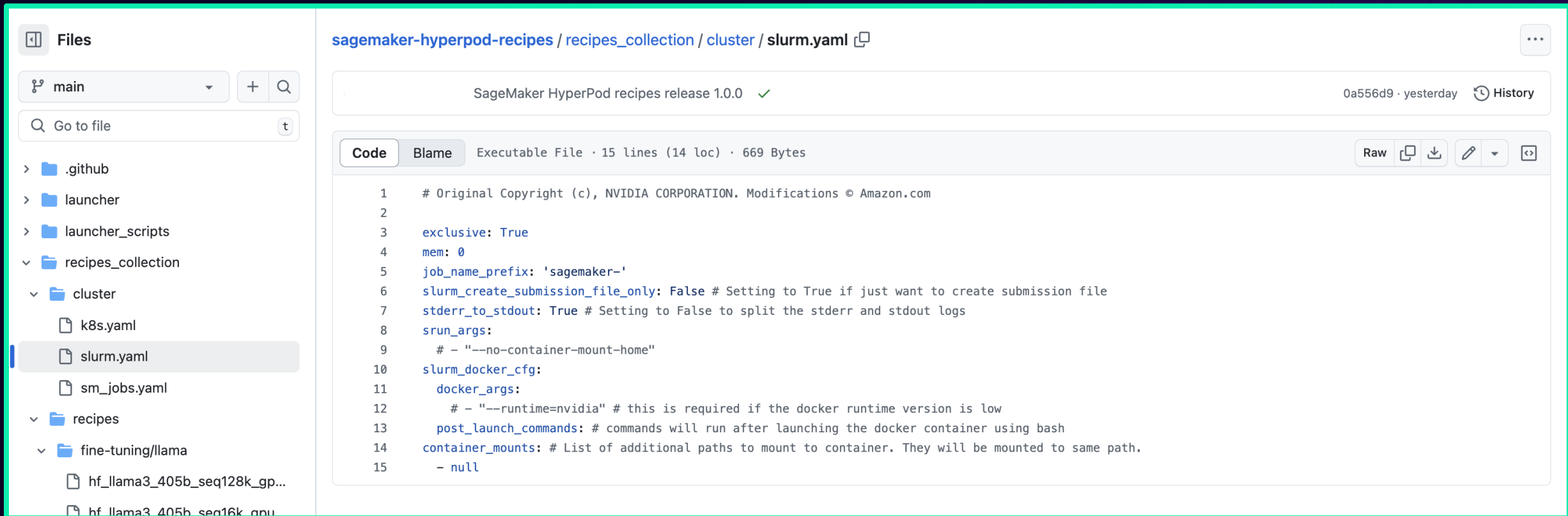# Getting started: Step 2 – Specify the cluster config

**Files**

main

Go to file

> .github
> launcher
> launcher_scripts
∨ recipes_collection
  ∨ cluster
    k8s.yaml
    slurm.yaml
    sm_jobs.yaml
  ∨ recipes
    ∨ fine-tuning/llama
      hf_llama3_405b_seq128k_gp...
      hf_llama3_405b_seq16k_gpu

sagemaker-hyperpod-recipes / recipes_collection / cluster / **slurm.yaml**

SageMaker HyperPod recipes release 1.0.0  ✓          0a556d9 · yesterday  🕙 History

**Code**   Blame          Executable File · 15 lines (14 loc) · 669 Bytes          Raw  ⧉  ⬇  ✏ ∨  <>

```
1    # Original Copyright (c), NVIDIA CORPORATION. Modifications © Amazon.com
2
3    exclusive: True
4    mem: 0
5    job_name_prefix: 'sagemaker-'
6    slurm_create_submission_file_only: False # Setting to True if just want to create submission file
7    stderr_to_stdout: True # Setting to False to split the stderr and stdout logs
8    srun_args:
9      # - "--no-container-mount-home"
10   slurm_docker_cfg:
11     docker_args:
12       # - "--runtime=nvidia" # this is required if the docker runtime version is low
13     post_launch_commands: # commands will run after launching the docker container using bash
14   container_mounts: # List of additional paths to mount to container. They will be mounted to same path.
15       - null
```

https://github.com/aws/sagemaker-hyperpod-recipes

# Getting started: Step 3 – Update root-level config

```
 3   defaults:
 4     - _self_
 5     - cluster: slurm  # set to `slurm`, `k8s` or `sm_jobs`, depending on the desired cluster
 6     - recipes: training/llama/hf_llama3_8b_seq16k_gpu_p5x16_pretrain # select desired config inside the training directory
 7     - override hydra/job_logging: stdout
 8
 9   cluster_type: slurm  # bcm, bcp, k8s or sm_jobs. If bcm, k8s or sm_jobs, it must match - cluster above.
10   # If using sm_jobs cluster_type, set sm_jobs_config. See cluster/sm_jobs.yaml for example.
11
12   hydra:
13     run:
14       dir: .
15     output_subdir: null
16
17   debug: False
18
19   instance_type: p5.48xlarge
20   base_results_dir: null  # Location to store the results, checkpoints and logs.
21
22   container: null
23
24   git:
25     repo_url_or_path: null
26     branch: null
27     commit: null
28     entry_script: null
29     token: null
30
31   env_vars:
32     NCCL_DEBUG: WARN  # Logging level for NCCL. Set to "INFO" for debug information
33
34   # Do not modify below, use the values above instead.
35   training_config: ${hydra:runtime.choices.recipes}
```

# Getting started – Slurm/k8s

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

Launcher scripts

```
bash launcher_scripts/<model>/<launcher-script>
```

HyperPod Slurm/k8s

Recipes on SageMaker HyperPod (Amazon EKS)

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

HyperPod CLI

Recipes on SageMaker HyperPod (Slurm./k8s)

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

NeMo-style launcher

# Getting started – SageMaker training jobs

Recipes on SageMaker training jobs

```
estimator = PyTorch(

    training_recipe=recipe-name
    ...,
)

estimator.fit(...)
```

SageMaker Python SDK

# Putting it all together

https://github.com/aws/sagemaker-hyperpod-recipes

# Putting it all together

https://github.com/aws/sagemaker-hyperpod-recipes

# Putting it all together

https://github.com/aws/sagemaker-hyperpod-recipes

# Putting it all together

User

- SageMaker HyperPod CLI
- SageMaker Python SDK

NeMo-style launcher

**SageMaker HyperPod Training Recipes**
- Parse model configuration from recipe
- Get training script path
- Parse cluster configuration (inc. container to use)
- Create launch script for SageMaker HyperPod (Slurm/EKS) and SageMaker Training jobs
- Launch job

SageMaker optimized model (GPU)

**SageMaker HyperPod Training Adapter for NeMo**
- SageMaker Distributed Training Libraries

Neuron optimized model (Tranium)

**NeuronX Training Toolkit**
- NeuronX Distributed

Native NeMo models (GPU)

**Nvidia NeMo**
- Megatron Core

Custom models

**Custom**

- SageMaker HyperPod Slurm
- SageMaker HyperPod EKS
- SageMaker Training Jobs

# Putting it all together

https://github.com/aws/sagemaker-hyperpod-recipes

# Putting it all together

https://github.com/aws/sagemaker-hyperpod-recipes

# Recap

Over 30 recipes with different configurations

Run training with a single line of code

Recipes support SageMaker HyperPod and training jobs

| SageMaker optimized models (GPU) | AWS Neuron optimized models (Trainium) |
| :---: | :---: |
| Native NeMo models | Custom models |

---

📖 README  🛡 Code of conduct  ⚖ Apache-2.0 license  ⚖ Security

# Amazon SageMaker HyperPod recipes

## Overview

Amazon SageMaker HyperPod recipes help customers get started with training and fine-tuning popular publicly available foundation models in just minutes, with state-of-the-art performance. The recipes provide a pre-configured training stack that is tested and validated on Amazon SageMaker.

Please see Amazon SageMaker HyperPod recipes for documentation.

The recipes support Amazon SageMaker HyperPod (with Slurm or Amazon EKS for workload orchestration) and Amazon SageMaker training jobs.

Amazon SageMaker HyperPod recipes include built-in support for:

- Model parallelism - tensor parallelism and context parallel
- Automated distributed checkpointing
- Distributed optimizer
- Accelerators: NVIDIA H100 (ml.p5), NVIDIA A100 (ml.p4), and AWS Trainium (ml.trn1)
- Fine-tuning: Full, QLoRA, LoRA
- AWS Instances: ml.p5.48xlarge, ml.p4d.24xlarge, and ml.trn1.32xlarge instance families
- Supported Models: Llama, Mistral, Mixtral models
- Model Evaluation: Tensorboard

## Model Support

### Pre-Training

List of specific pre-training recipes used by the launch scripts.

| Source | Model | Size | Sequence length | Nodes | Instance | Accelerator | Recipe | Script |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Hugging Face | Llama 3.2 | 11b | 8192 | 4 | ml.p5.48xlarge | GPU H100 | link | link |

# Salesforce AI research

# A little bit about Salesforce AI Research …

**We are Salesforce's AI Research org**

- **Foundational R&D**
  Pushing forward state-of-the-art models for enterprise AI

- **Customer incubation**
  Customer-centric pathfinding with cutting-edge AI for high-value use cases

- **Product innovation**
  Incorporating AI models and technologies into generally available features

# Salesforce AI & Amazon SageMaker HyperPod
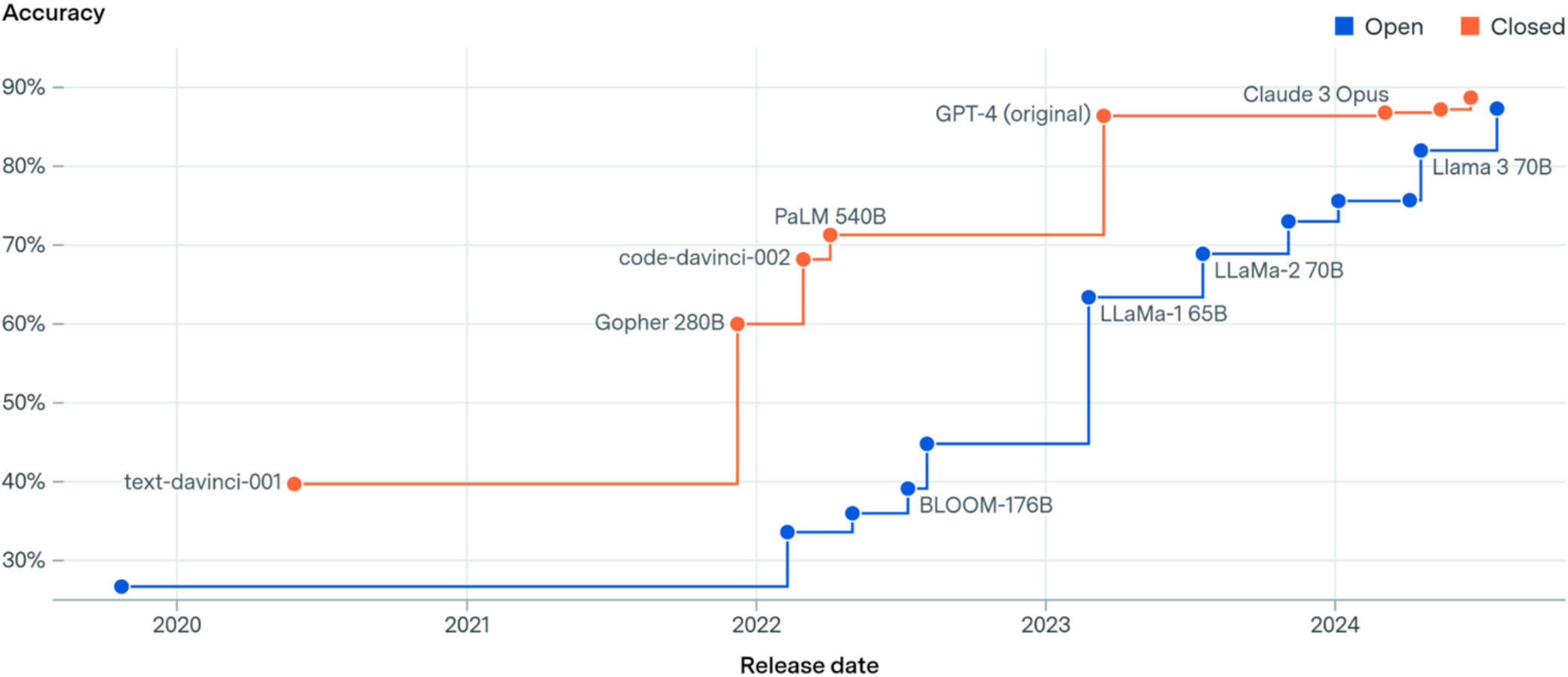
We're big users of SageMaker for years

This year, we've leveraged HyperPod

*"Supercomputer at our fingertips"*

# Why do model training and fine-tuning matter?



Top-performing open and closed AI models on MMLU benchmark

≩ EPOCH AI

Accuracy

■ Open   ■ Closed

Models labeled: Claude 3 Opus, GPT-4 (original), PaLM 540B, code-davinci-002, Gopher 280B, text-davinci-001, BLOOM-176B, LLaMa-1 65B, LLaMa-2 70B, Llama 3 70B

Release date

# Salesforce AI
# Fine-tuning FMs with HyperPod recipes

SageMaker HyperPod recipes have great features that saved us _**significant**_ infra work

**FSDP + Context Parallelism + Low-Rank Adaptation (LoRA)**
**Unique Feature Set**

**Fine-tune Llama-3.1-405B at full 128K context with 4 nodes**
1 node with Q-LoRA

# Production models at Salesforce trained on HyperPod

# Agentforce: Powering next-gen AI for Sales with xGen-Sales LLM

First proprietary sales-focused large language model (LLM)

Key features: call summarization, customer profiling, contact enrichment, and pipeline tracking

Users can interactively engage with AI-generated call summary

AI-generated summaries favored >50% of the time against human-generated



Contributors:
Semih Yavuz, Xinyi Yang, Srijan Bansal, Donna Tran, John Emmons, Jason Lee , Erik Nijkamp, Bo Pang, Egor Pakhomov, Akash Gokul, Antonio Ginart, Yingbo Zho

# Agentforce: Powering next-gen AI for Sales with xGen-Sales LLM

# Agentforce: Powering next-gen AI for Sales with xGen-Sales LLM

Developed using human-in-the-loop reinforcement learning and diverse sales datasets

Planned integration with Agentforce, live customer pilots are already underway!



Our Approach to Building xGen-Sales

Proprietary Data → Filtering → Relevant Transcripts

High-Quality SFT Data → Labeling

XGen-instruct → Training → XGen-Sales

# Ragforce: Building *contextually faithful* LLMs with SFR-RAG

**Retrieval Augmented Generation:** Combine semantic search + in-context learning + generative LLM

- Augments LLMs with enterprise data retrieval, but quality depends heavily on document hygiene
  - Conflicting versions, outdated info, and inconsistent formatting can confuse the model

- How should models handle prioritization of information given rich metadata?
  - Prioritization semantics based on sources, tags, and dates
  - "Official" tag takes priority over "Draft" tag unless date is <6 months old

# Ragforce: Building *contextually faithful* LLMs with SFR-RAG

- Trained a contextual LLM (9B) for RAG generation

  - Faithful: Trained to minimize hallucinations, even in adverse settings (conflicting info) & precisely follow complex prioritization rules

  - Lightweight and high-performing: State-of-the-art aggregate performance with 9B params, beating larger models (104B+)

  - Trained on HyperPod using 2x p5.48xlarge nodes



**SFR-RAG: Towards Contextually Faithful LLMs**

Xuan-Phi Nguyen*    Shrey Pandit    Senthil Purushwalkam    Austin Xu

Hailin Chen    Yifei Ming    Zixuan Ke    Silvio Savarese    Caiming Xong

Shafiq Joty

Salesforce AI Research

# Rankforce: State-of-the-art enterprise reranking with SFR-LlamaRank



**Role of a Reranker in RAG Systems**

Knowledge Base → 1. Fetch → Retriever
Query → Retriever
Retriever → 2. Initial Results → Reranker
Reranker → 3. Refined Results → Generator

**Reranker's Role:**
- Refines initial retrieval results for improved relevance
- Enhances quality of information for generation
- Uses language models to capture nuanced relevance beyond semantic similarity

Lead contributors: A. Ginart, N. Kodali, J. Emmons

# **Rankforce:** State-of-the-art enterprise reranking with SFR-LlamaRank

- **Core technical specs:**
  - Built on Llama3-8B-Instruct with human-guided training (4x p5.48xlarge)
  - Supports 8K document chunks with fast inference (<200ms/4 docs on H100)
- Superior code search performance
- Linear, calibrated scoring
  (0.9+ highly relevant to ~0.0 irrelevant)
- Horizontally scalable for enterprise deployment

### Results: Hit Rate @ K = 8

| Model | Avg | SQuAD | TriviaQA | NCS | TrailheadQA |
|---|---|---|---|---|---|
| SFR LlamaRank | 92.9% | 99.3% | 92.0% | 81.8% | 98.6% |
| Cohere Rerank V3 | 91.2% | 98.6% | 92.6% | 74.9% | 98.6% |
| Mistral-7B QLM | 83.3% | 87.3% | 88.0% | 60.1% | 97.7% |
| Embeddings Only | 73.2% | 93.2% | 88.3% | 18.2% | 93.2% |

Lead contributors: A. Ginart, N. Kodali, J. Emmons

# Judgeforce
# What is a reward model?

- Learn a model that predicts human preferences based on pairwise comparisons ("A is better than B")

- Used in training models (RLHF) as well as automatic evaluations and even inference

# Judgeforce: Automating LLM eval with SFR-Judge

- Trained 3 LLM-as-judge models (8B, 12B, 70B) for *automatic evaluation*
  - Multifaceted: Trained to evaluate via pairwise comparisons, single rating, and classification

  - State-of-the-art aggregate performance across 13 evaluation benchmarks

  - Trained 70B on HyperPod using 4x p5.48xlarge nodes

### DIRECT JUDGEMENT PREFERENCE OPTIMIZATION

Peifeng Wang*, Austin Xu*, Yilun Zhou, Caiming Xiong, Shafiq Joty
Salesforce AI Research

(a) **Single Rating**: Assign a score between 1 and 5, according to the scoring rubric.
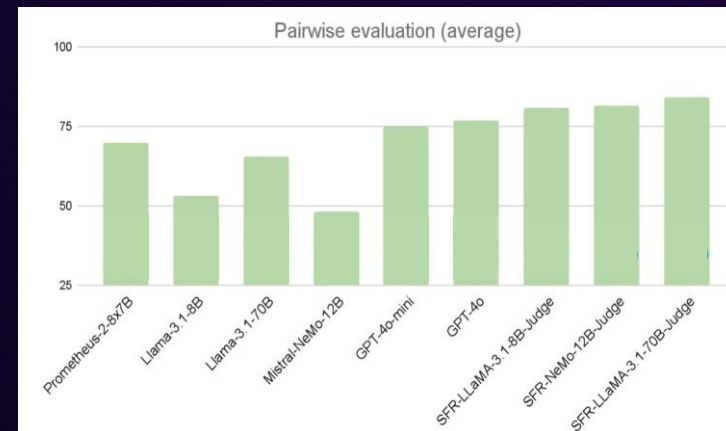
(b) **Pairwise Comparison**: Select Response A or B, that is better for the given instruction.

(c) **Classification**: Does the response meet the requirement of ...?

Pairwise evaluation (average)

# What's next at Salesforce AI?



Are we hitting entropy rate of internet corpus?
*Can we scale up post-training dataset sizes?*

What will inference-time compute scaling be?
*2025 may be the year of inference-time scaling!*

**Multi-Agent Systems:**
- **Components:** Specialized models for tasks –
  Re-rankers, reward models
  xGen-Sales

- **Assistants:** Dispatch directly with humans – real-time, multi-modal, low-latency
  SFR-RAG

- **Agents:** Long-running AI with tools - Minutes, hours, or even days
  - Powered by reasoning-style LLMs

# Get started with
## HyperPod recipes

# AMAZON SAGEMAKER HYPERPOD RECIPES

https://github.com/aws/
sagemaker-hyperpod-recipes



📖 README    Code of conduct    ⚖ Apache-2.0 license    ⚖ Security

## Amazon SageMaker HyperPod recipes

### Overview

Amazon SageMaker HyperPod recipes help customers get started with training and fine-tuning popular publicly available foundation models in just minutes, with state-of-the-art performance. The recipes provide a pre-configured training stack that is tested and validated on Amazon SageMaker.

Please see Amazon SageMaker HyperPod recipes for documentation.

The recipes support Amazon SageMaker HyperPod (with Slurm or Amazon EKS for workload orchestration) and Amazon SageMaker training jobs.

Amazon SageMaker HyperPod recipes include built-in support for:

- Model parallelism - tensor parallelism and context parallel
- Automated distributed checkpointing
- Distributed optimizer
- Accelerators: NVIDIA H100 (ml.p5), NVIDIA A100 (ml.p4), and AWS Trainium (ml.trn1)
- Fine-tuning: Full, QLoRA, LoRA
- AWS Instances: ml.p5.48xlarge, ml.p4d.24xlarge, and ml.trn1.32xlarge instance families
- Supported Models: Llama, Mistral, Mixtral models
- Model Evaluation: Tensorboard

## Model Support

### Pre-Training

List of specific pre-training recipes used by the launch scripts.

| Source | Model | Size | Sequence length | Nodes | Instance | Accelerator | Recipe | Script |
|--------|-------|------|-----------------|-------|----------|-------------|--------|--------|
| Hugging Face | Llama 3.2 | 11b | 8192 | 4 | ml.p5.48xlarge | GPU H100 | link | link |

# Thank you!

Please complete the session survey in the mobile app