

The background features a dark navy blue field with abstract, overlapping shapes in vibrant magenta and deep red. Thin, light blue lines intersect diagonally across the composition. The text is positioned on the left side.

AWS re:Invent

DECEMBER 2 – 6, 2024 | LAS VEGAS, NV

FSI320

Bloomberg: Lessons learned from building and training LLMs on AWS

Phil Vachon

(he/him)

Head of Infrastructure and Security
Office of the CTO
Bloomberg

Vadim Dabravolski

(he/him)

Team Lead, AI Engineering
Bloomberg

Vasu Chari

(he/him)

Principal Solutions Architect
AWS



Agenda

01 In the beginning . . .

02 The training process

03 Managing training
infrastructure

04 The results

05 Applying our research in
gen AI-enhanced products

06 If we had to do it again today?

In the beginning . . .



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Imagine a time before ChatGPT . . .

May 2020:

“Language models are
few-shot learners”
paper (GPT-3)

**The scaling
laws are
established:**

The more parameters,
the better

**We start to see
the potential
of LLMs:**

They transform everything

LLMs' potential for financial services

GPT-3 showed amazing results but required a huge number of parameters – 175 billion for GPT-3

Financial systems have their own terminology

Our customers expect our insights to be correct, relevant, and timely – could a general-purpose model “as-is” meet these requirements?

We decided to build and train a GPT-3-like model ourselves, tailored to financial services use cases

The training process

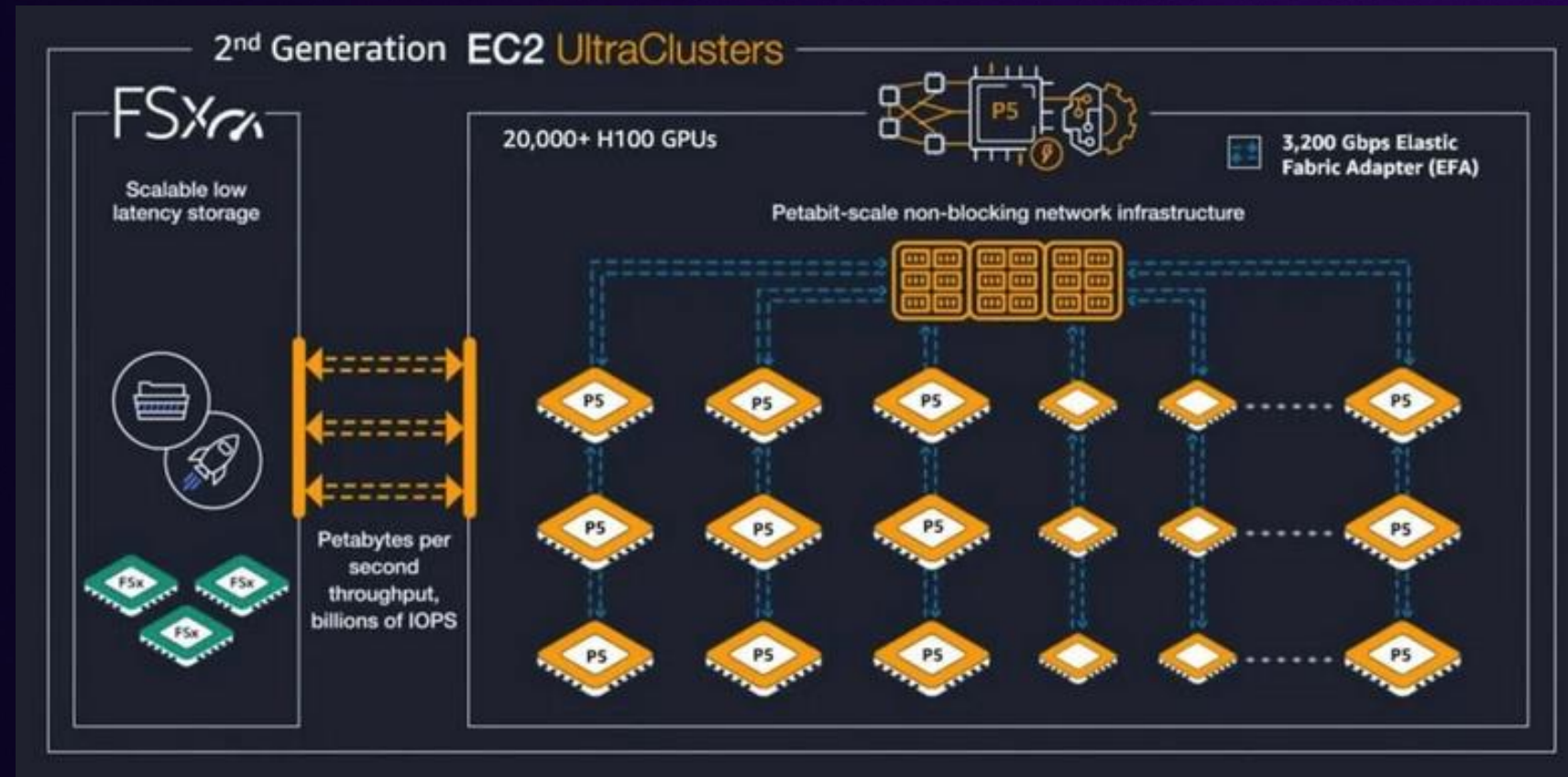


Scaling laws and compute

Our training dataset:
~700B tokens

Compute resources:
1.3M GPU hours (NVIDIA
A100 GPUs, 40 GB RAM)

Our model size:
50B parameters



The power of AWS



Accessible

Amazon SageMaker is a fully managed service, with a native model parallel training framework SMP



Available

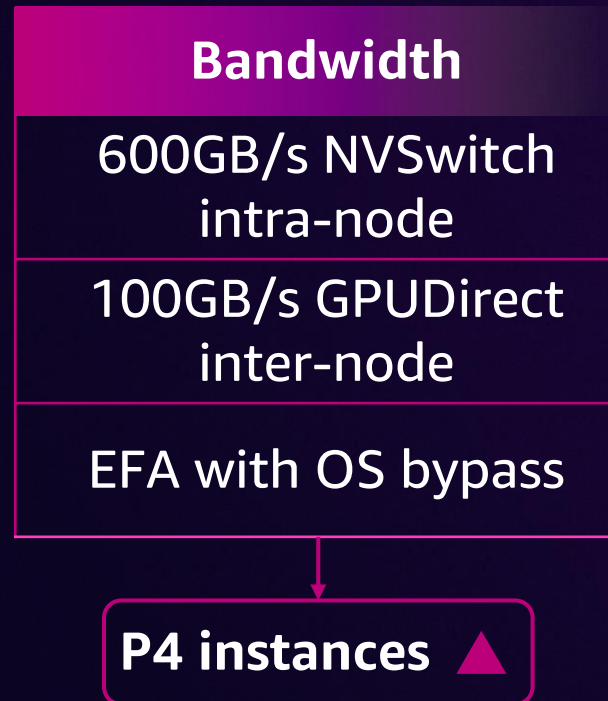
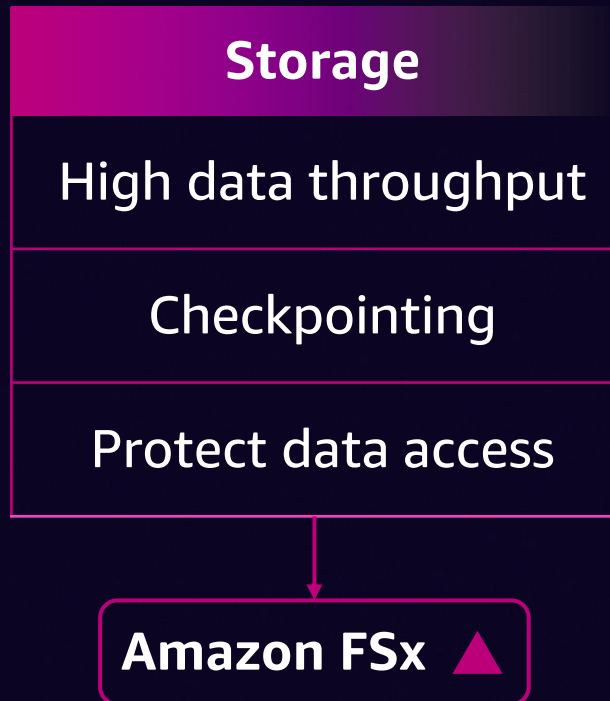
p4d.24xlarge instances (8 NVIDIA A100 GPUs) were available to schedule



Collaborative

Bloomberg's AI researchers and Amazon SageMaker engineers shared knowledge, accelerating the process and enabling just 9 Bloomberg AI researchers to achieve amazing results

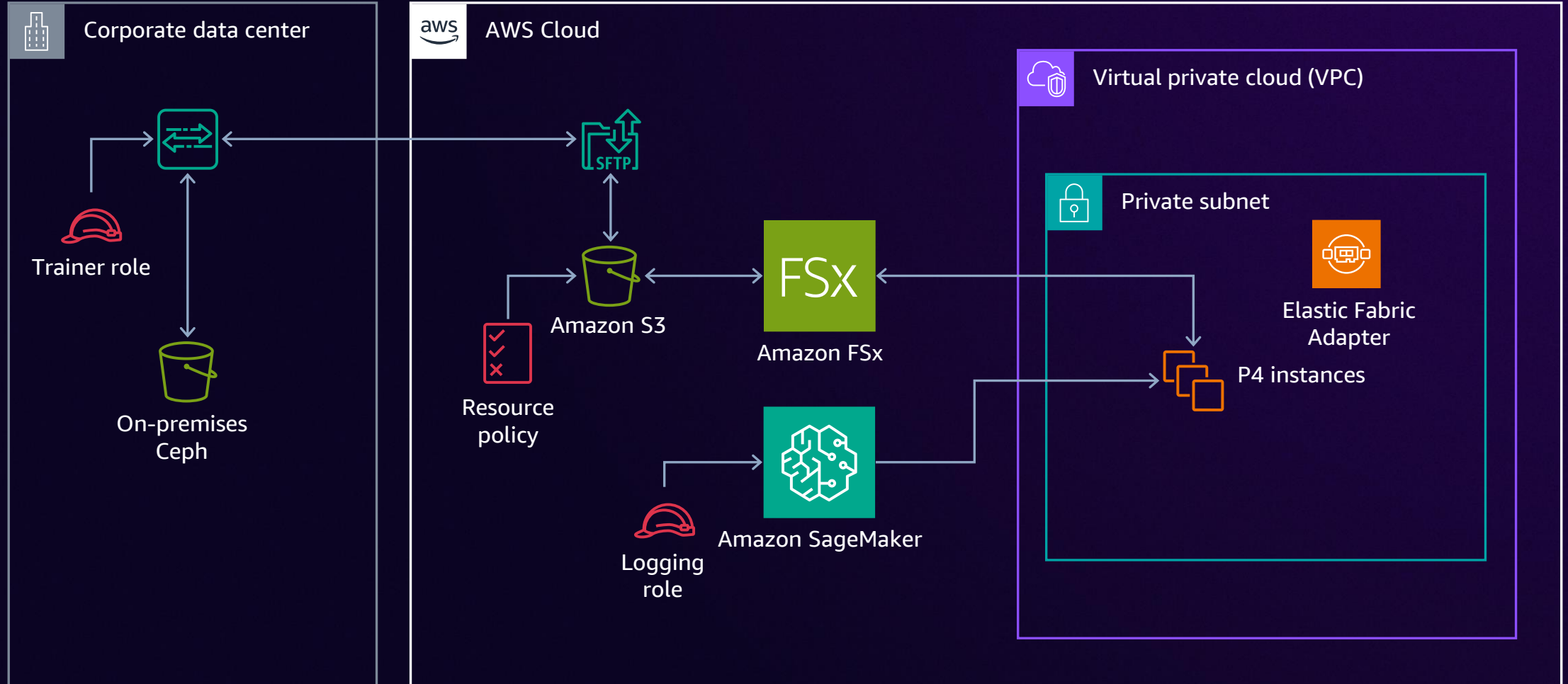
Model design impacts architecture



Managing training infrastructure



ML training environment architecture



The power of AWS



Lots of bytes

Our training dataset was
3.5 TB

Checkpointing generated
~600 GB of data every two
hours

0111010
0101110
0111010

Throughput

GPUs need lots of data to
keep them busy

Checkpointing costs precious
computation time



Protecting our data

Training datasets,
checkpoints, and model
weights are high value

Robust, resilient, and reliable
storage, with the right
security controls, was
essential

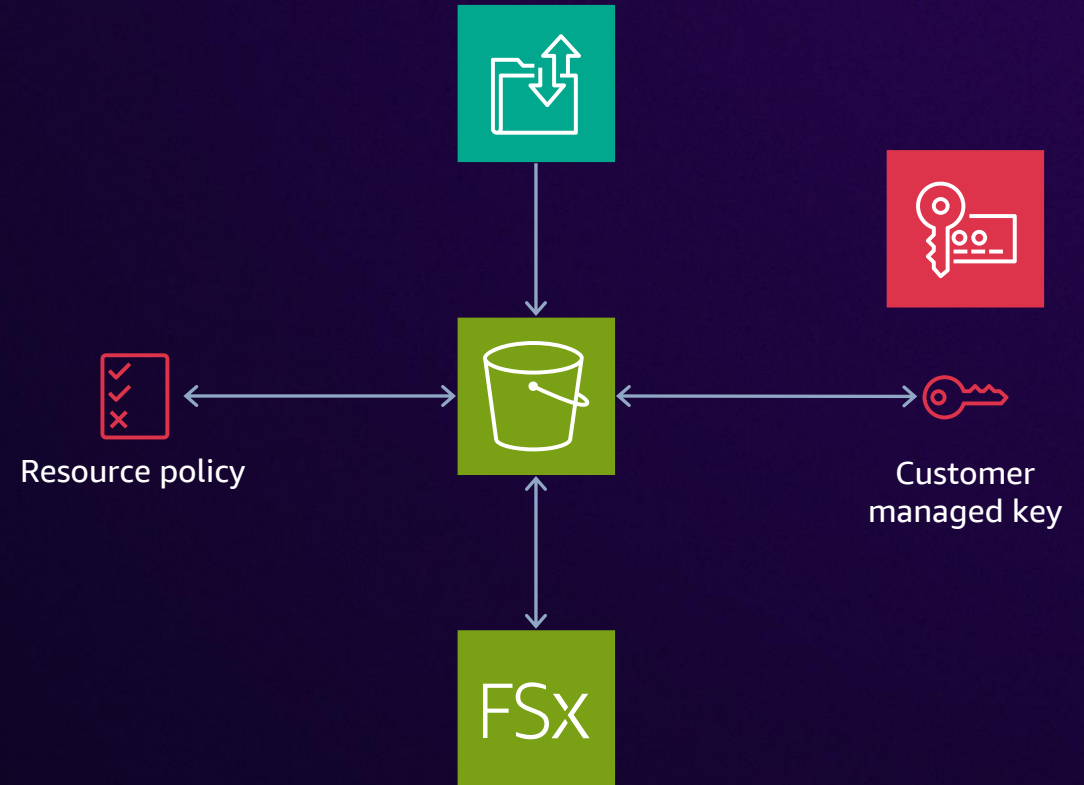
By leveraging Amazon FSx for Lustre for access, Amazon S3 for near-line and backup storage, and synchronizing between the two, we had a high level of assurance our data was safe and secure

Protect the perimeter

Training data is valuable; model weights are even more valuable

Accidents can happen if you don't put guardrails in place

Training on permitted data must be easy, reliable, and secure



Access control challenges

Safe experimentation

The training process was as much a discovery process as a training one
Need to enable rapid iteration while keeping an eye on things

Enterprise controls

Ensure our standard observability, internet connectivity, and data protection rules are enforced



AWS Management
Console



Universal 2FA
(and SSO)



Access policies



Model
weights



Training
data

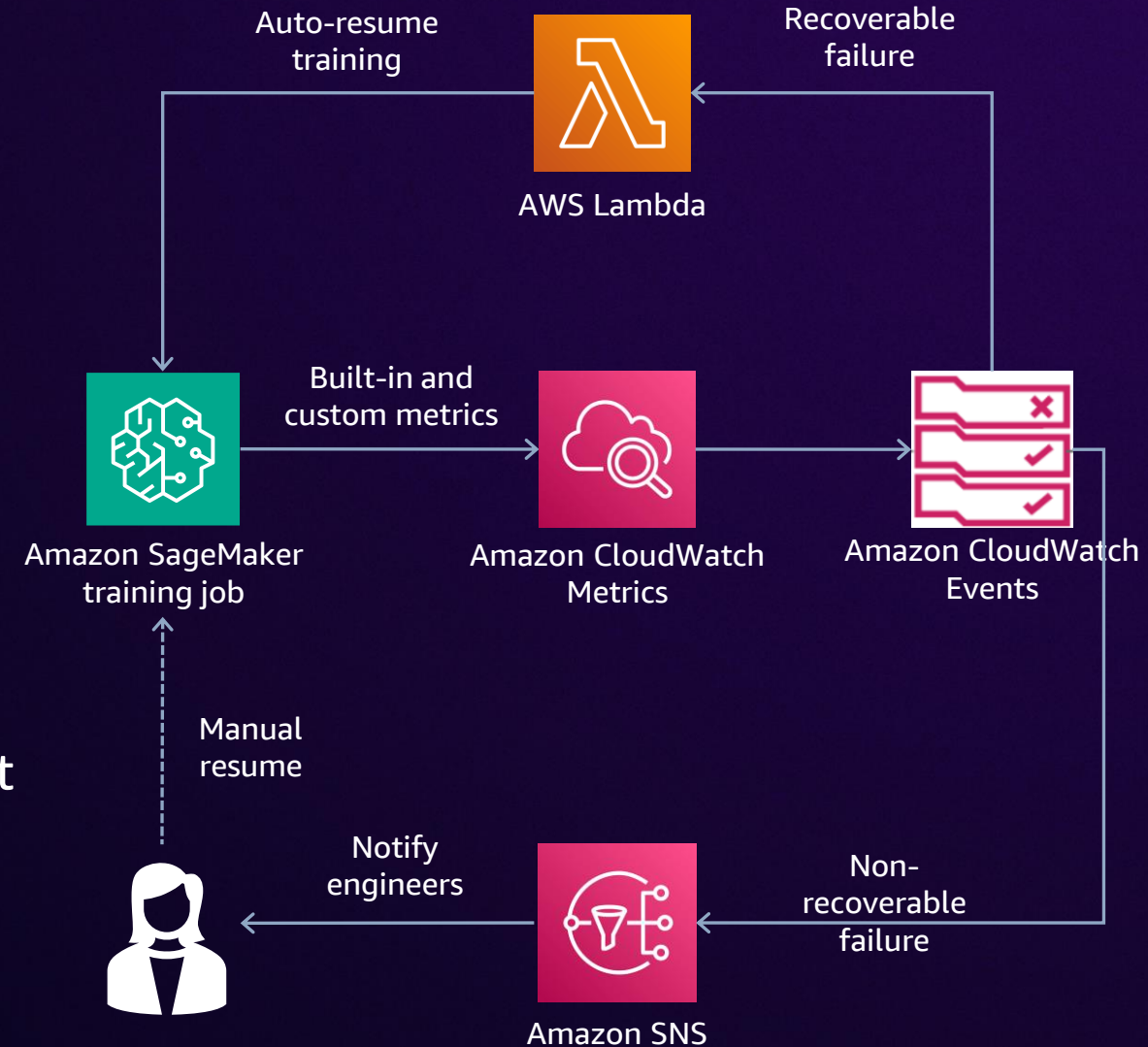
Operational management of training jobs

LLM training presents unique operational challenges:

- Long-running, high-cost, stateful
- Large number of failure modes
- High volume of diverse data signals

Remediation depends on type(s) of failure, such as:

- Automatic resumption from latest checkpoint
- Manual inspection and update by personnel

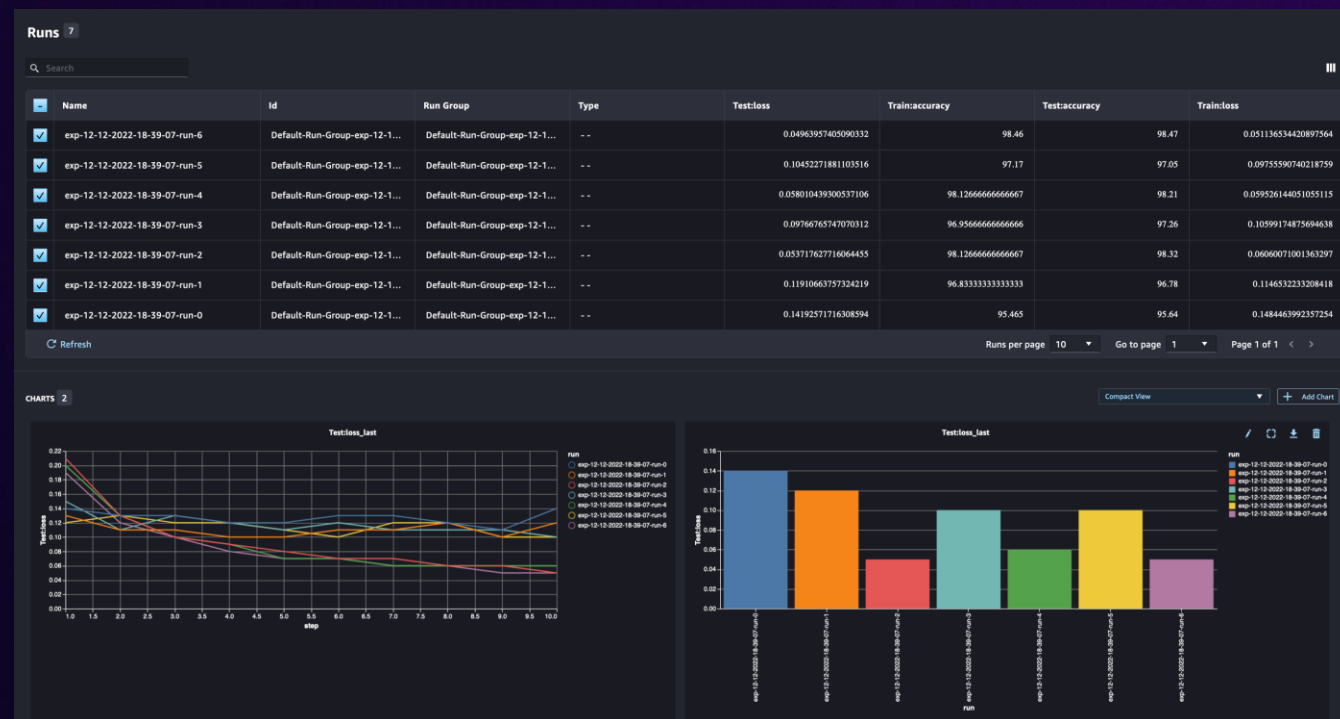


Failure mode: Underutilization of GPUs

Common problem in LLM training, which leads to unnecessary waste and longer delivery time

Detection and remediation:

1. Calculate effective FLOPS per GPU per training step
2. Stream effective FLOPS to CloudWatch
3. Configure metric to stop training if effective FLOPS < 50% of theoretical FLOPS

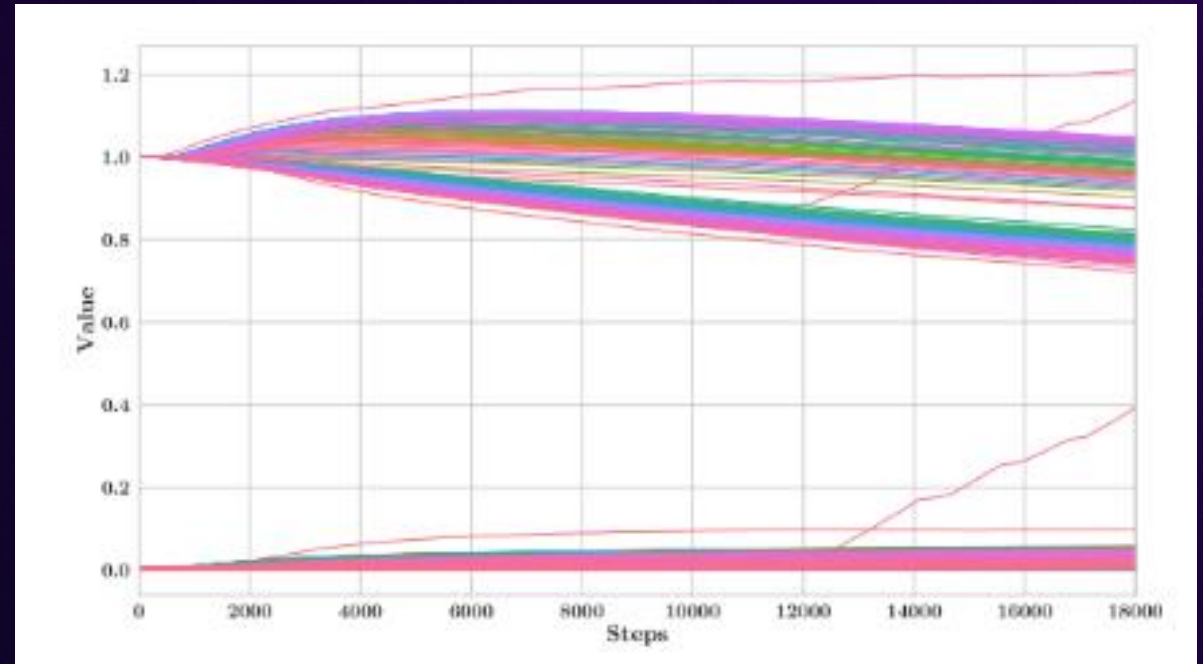


Failure mode: Training convergence

Be prepared for your loss not to decrease smoothly

Detection and remediation:

1. Instrumented model state and introduced custom metrics (such as gradient norm per layer)
2. Stream metrics on each forward and backward pass
3. Aggregate metrics via CloudWatch dashboards



Monitoring gradient norm per layer during training; one layer stands out

The results



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

BloombergGPT research model

BloombergGPT: A Large Language Model for Finance

Shijie Wu^{1,*}, Ozan İrsoy^{1,*}, Steven Lu^{1,*}, Vadim Dabravolski¹, Mark Dredze^{1,3}, Sebastian Gehrmann¹, Prabhanjan Kambadur¹, David Rosenberg², Gideon Mann¹

¹ Bloomberg, New York, NY USA

² Bloomberg, Toronto, ON Canada

³ Computer Science, Johns Hopkins University, Baltimore, MD USA

Abstract

The use of NLP in the realm of financial technology is broad and complex, with applications ranging from sentiment analysis and named entity recognition to question answering. Large Language Models (LLMs) have been shown to be effective on a variety of tasks; however, no LLM specialized for the financial domain has been reported in literature. In this work, we present BLOOMBERGGPT, a 50 billion parameter language model that is trained on a wide range of financial data. We construct a 363 billion token dataset based on Bloomberg's extensive data sources, perhaps the largest domain-specific dataset yet, augmented with 345 billion tokens from general purpose datasets. We validate BLOOMBERGGPT on standard LLM benchmarks, open financial benchmarks, and a suite of internal benchmarks that most accurately reflect our intended usage. Our mixed dataset training leads to a model that outperforms existing models on financial tasks by significant margins without sacrificing performance on general LLM benchmarks. Additionally, we explain our modeling choices, training process, and evaluation methodology. We release Training Chronicles (Appendix C) detailing our experience in training BLOOMBERGGPT.

Applying our research in gen AI-enhanced products



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Guiding principles: Using LLMs in our products

Derive answers from our trusted sources of information

Build system features to assist users

Provide transparency and attribution

Analyzing and visualizing data

Learning to code to analyze data

Walt Disney closing prices for the last month

```
for('DIS US EQUITY'])  
get(px_last(dates=range(-1m,0d)))
```

Tesla volume curve for last 5 days

```
for(['TSLA US EQUITY'])  
get(cumsum(px_volume(dates=range(-4d,0d))))
```

Median gross margin for members of SPX

```
for(members(['SPX INDEX']))  
get(median(group(gross_margin)))
```

30 day rolling sharpe over the last month for Walmart

```
for(['WMT US EQUITY'])  
get(rolling(sharpe_ratio(calc_interval=range  
(-30d,0d)), iterationdates=range(-1m,0d)))
```

Return the ratio of buy analyst recommendations
to total recommendations for holdings of Starbucks

```
for(holdings('SBUX US EQUITY'))  
get(tot_buy_rec / tot_analyst_rec)
```

Natural language simplifies data analysis



BQL Assistant

BQL <GO>

Prompt

Revenue growth next quarter for Ford and its peers

Editable BQL Query

```
for(peers('F US EQUITY'))
get(name,pct_chg(dropna(sales_rev_turn(
fpo=range(0,1),fpt=q)),negative_values=
abs_base))
```

Table

Ticker	Name	Revenue Growth
F US Equity	Ford Motor Co	-13.83
MBG GR Equity	Mercedes-Benz Group	10.34
7203 JP Equity	Toyota Motor Corp	-2.99
7201 JP Equity	Nissan Motor Co Ltd	1.63
GM US Equity	General Motors Co	-11.48
VOW GR Equity	Volkswagen AG	-8.01
7269 JP Equity	Suzuki Motor Corp	-2.76

Gathering and synthesizing insights

Transcript summarization



Document Viewer

DOCV <GO>

Apple Inc Earnings Call



Automated Summary

Guidance

- Company gross margin was near the high end of the guidance range.
- Operating expenses were at the low end of the guidance range.
- Gross margin is expected to be between 45.5% and 46.5%.
- The company provided guidance for the September quarter, assuming the macroeconomic outlook doesn't worsen from current projections.
- Total company revenue is expected to grow year-over-year at a rate similar to the June quarter.
- Services revenue is expected to grow double digits at a rate similar to the first three quarters of this fiscal year.

FINAL TRANSCRIPT

2024-08-01

Apple Inc (AAPL US Equity)

Q3 2024 Earnings Call

Company Participants

- [Luca Maestri](#), Senior Vice President and Chief Financial Officer
- [Suhasini Chandramouli](#), Director of Investor Relations
- [Tim Cook](#), Chief Executive Officer

Other Participants

- [Amit Daryanani](#), Evercore ISI
- [Atif Malik](#), Citi

Transcript summarization



Document Viewer

DOCV <GO>

Apple Inc Earnings Call



Automated Summary

Guidance

- Company gross margin was near the high end of the guidance range.
- Operating expenses were at the low end of the guidance range.
- Gross margin is expected to be between 45.5% and 46.5%.
- The company provided guidance for the September quarter, assuming the macroeconomic outlook doesn't worsen from current projections.
- Total company revenue is expected to grow year-over-year at a rate similar to the June quarter.
- Services revenue is expected to grow double digits at a rate similar to the first three quarters of this fiscal year.

During the quarter, we returned over \$32 billion to shareholders, including \$3.9 billion in dividends and equivalents, and \$26 billion through open market repurchases of 139 million Apple shares.

As we move ahead into the September quarter, I'd like to review our outlook, which includes the types of forward-looking information that Suhasini referred to at the beginning of the call. The color we're providing today assumes that the macroeconomic outlook doesn't worsen from what we are projecting today for the current quarter.

We expect foreign exchange to continue to be a headwind, to have a negative impact on revenue of about 1.5 percentage points on a year-over-year basis. We expect our September quarter total company revenue to grow year-over-year at a rate similar to the June quarter.

We expect Services revenue to grow double digits at a rate similar to what we reported in the first three quarters of this fiscal year. We expect gross margin to be between 45.5% and 46.5%. We expect OpEx to be between \$14.2 billion and \$14.4 billion. We expect OI&E to be around negative \$50 million, excluding any potential impact from the mark-to-market of minority investments, and our tax rate to be around 16.5%.

If we had to do it again today?



Now, it's just out of the box . . .

EC2 P5 UltraClusters

- NVIDIA H100 GPUs increase effective TFLOPS

SageMaker HyperPod

- Turnkey training infrastructure, tools, and libraries
- Automatic checkpointing and healing from hard faults
- Maximize utilization of hardware, without manual intervention

FSDP support in SMPv2

- Some use cases we run on-premises due to data sensitivity
- Create and maintain a single code base for SageMaker or on-premises
- Use the same training loops on-premises and on AWS as we scale up

One model does not rule them all

Model selection criteria: Evaluate strengths and weaknesses of each model against specific use cases

- Accuracy on the downstream task
- Latency/throughput/cost
- Tool support/agents
- RAG versus fine-tuning when adding knowledge to the model
- Multimodality
- Continuous evaluation

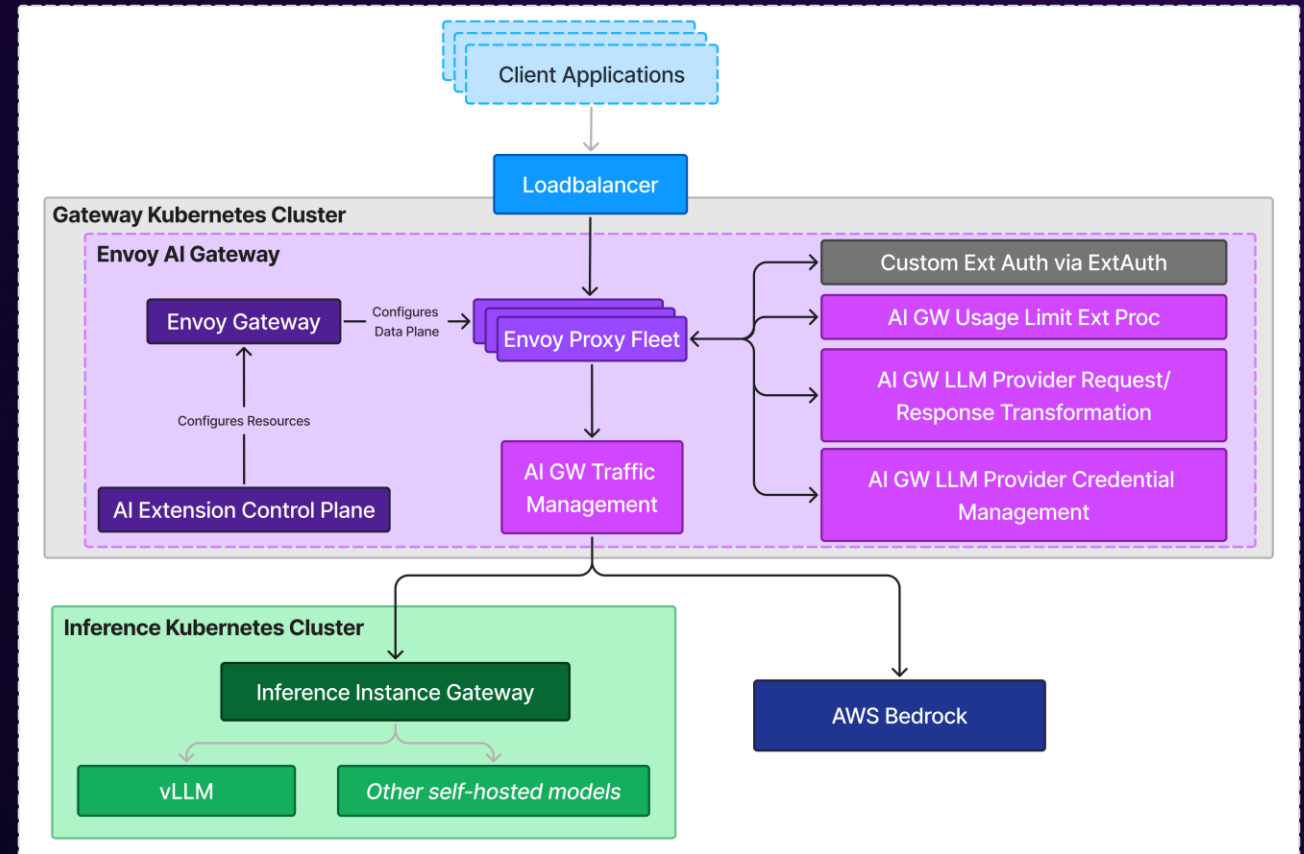
Using multiple LLMs: Envoy AI Gateway



Centralized access: Standard APIs, controlled and auditable way to access self-trained, open source, and/or commercial models

Enforce access policies: Support various methods to authenticate and manage requests to a wide range of AI models

Monitor costs: Integrated observability and cost controls (such as token-based limits) optimize expenses and mitigate risks



So, do you **actually** need a
custom foundation model?

Thank you!

Phil Vachon

pvachon1@bloomberg.net
[@pvachonnyc](https://twitter.com/pvachonnyc)

Vadim Dabravolski

vdabravolski@bloomberg.net

Vasu Chari

vxchari@amazon.com



Please complete the session survey in the mobile app