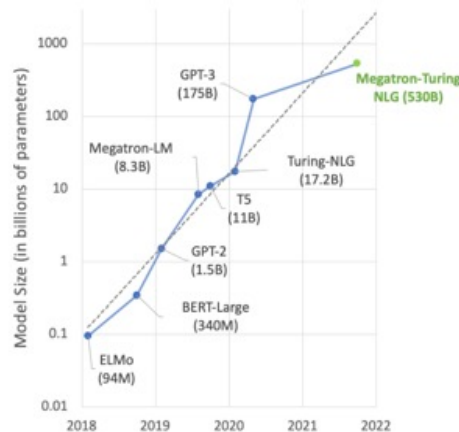At HPE, we believe in changing the way people live and work.

# The World Before November 2022…



Using DeepSpeed and Megatron to Train Megatron-Turing NLG 530B, the World's Largest and Most Powerful Generative Language Model
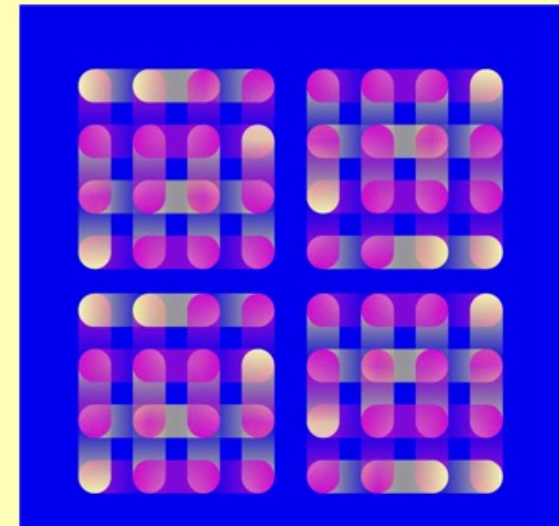
Oct 11, 2021

By Paresh Kharya and Ali Alvi

Top 9 Use Cases of Computer Vision in Manufacturing

GPT-3 powers the next generation of apps

Over 300 applications are delivering GPT-3–powered search, conversation, text completion, and other advanced AI features through our API.

# …and Then This Happened…

OpenAI's chatGPT

OpenAI's GPT-4

Google's Bard

Anthropic's Claude 2

Google's Gemini

AWS's Q
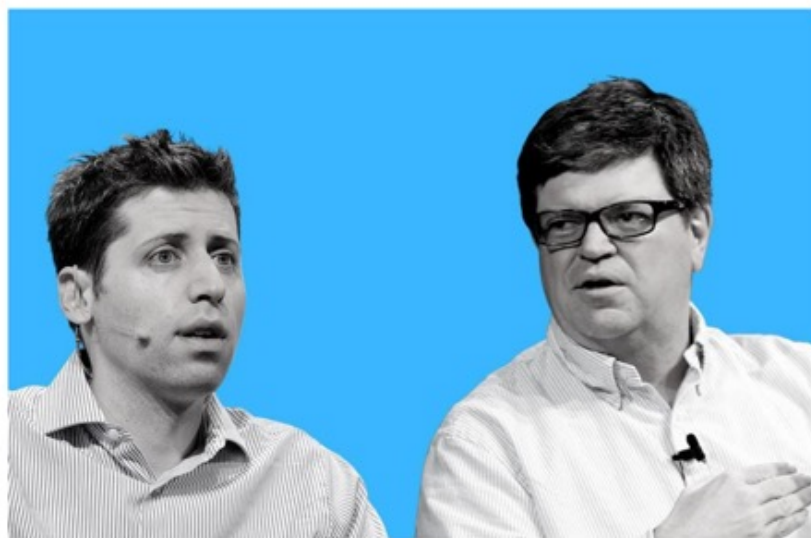
…many more



## The AI Arms Race

Siddharth Sharma · Follow
9 min read · Feb 9

"It's a new day in search. The race starts today… We're going to move fast."
— Satya Nadella, Microsoft CEO

Altman and Lecun (Forbes)

In recent years, the field of Artificial Intelligence has seen a rapid rise in the development of large language models. These models, based on deep

# …and Everything Changed



How Generative AI Will Change Jobs In Financial Services

By Bernard Marr, Contributor.

Follow Author

Jun 10, 2024, 01:43am EDT

Save Article    Comment 0

How Generative AI Will Change Jobs In Financial Services    ADOBE STO

Generative AI Chatbots Dominate Financial Services, But Data Privacy and Security Concerns Loom

Sep 5, 2024    Hubbis

Deepfakes Are Coming for the Financial Sector

…panies using photos or audio to verify customers' identities are …aring for bad actors gaming the system with generative AI

…belle Bousquette  Follow

…3, 2024 7:00 am ET

Share    AA Resize    Listen (2 min)

…penAI Made Me Crazy Videos—Then the …TO Answered (Most of) My Questions    WSJ

AI WOMEN

0:17/10:38

…penAI's new text-to-video AI model, can create realistic scenes. In an exclusive interview, WSJ's Joanna …sat down with the company's CTO, Mira Murati, who explained how it works but ducked questions about …e model was trained. Photo illustration: Preston Jessee for The Wall Street Journal
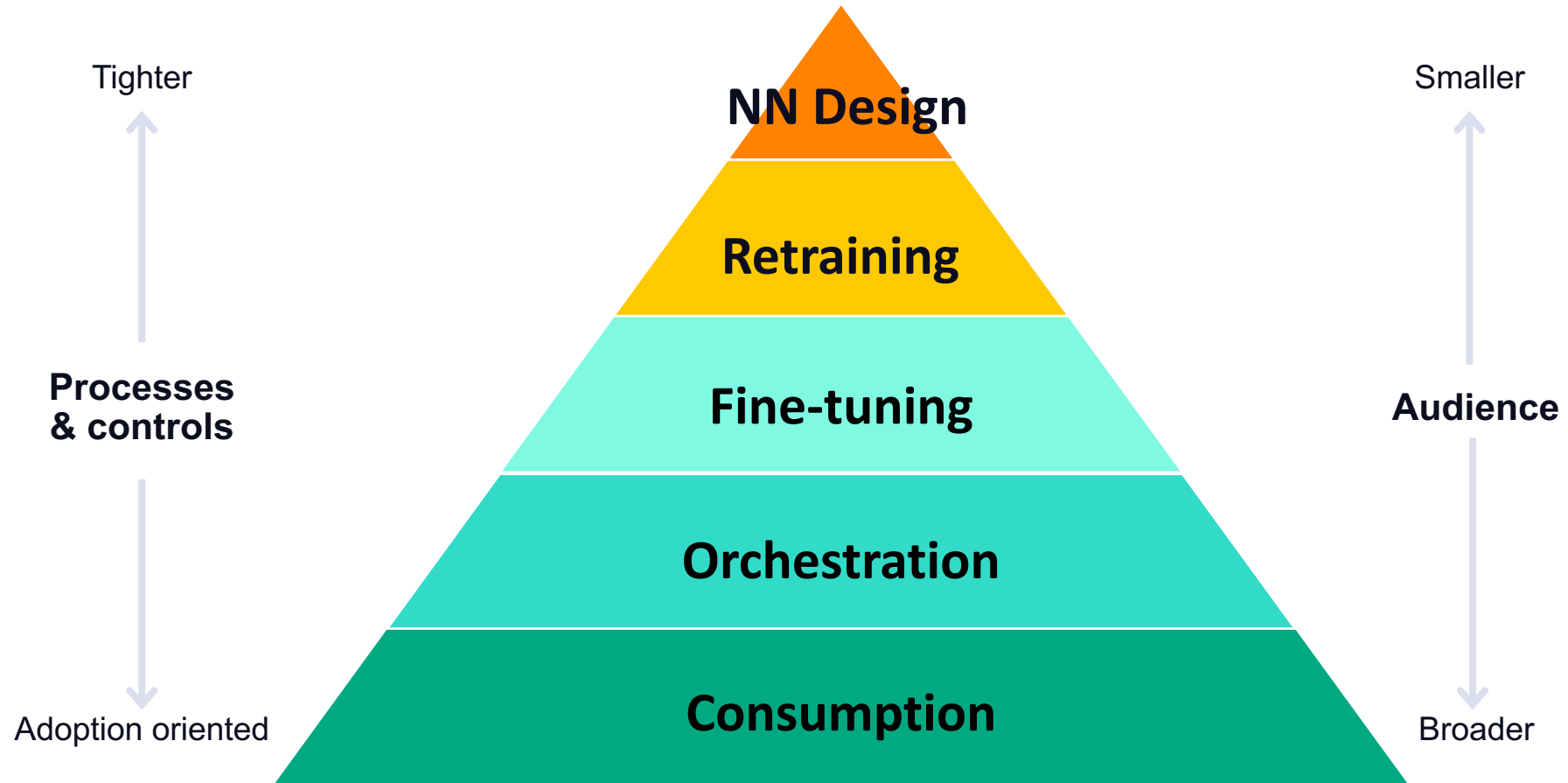
4

# Where Do You Go From Here?



**1** What is needed to **get started ?**

**2** Which **one** (company, model, open source)?

**3** Is it **enterprise ready**?

**4** What about **governance**?

**5** Which **use cases**?

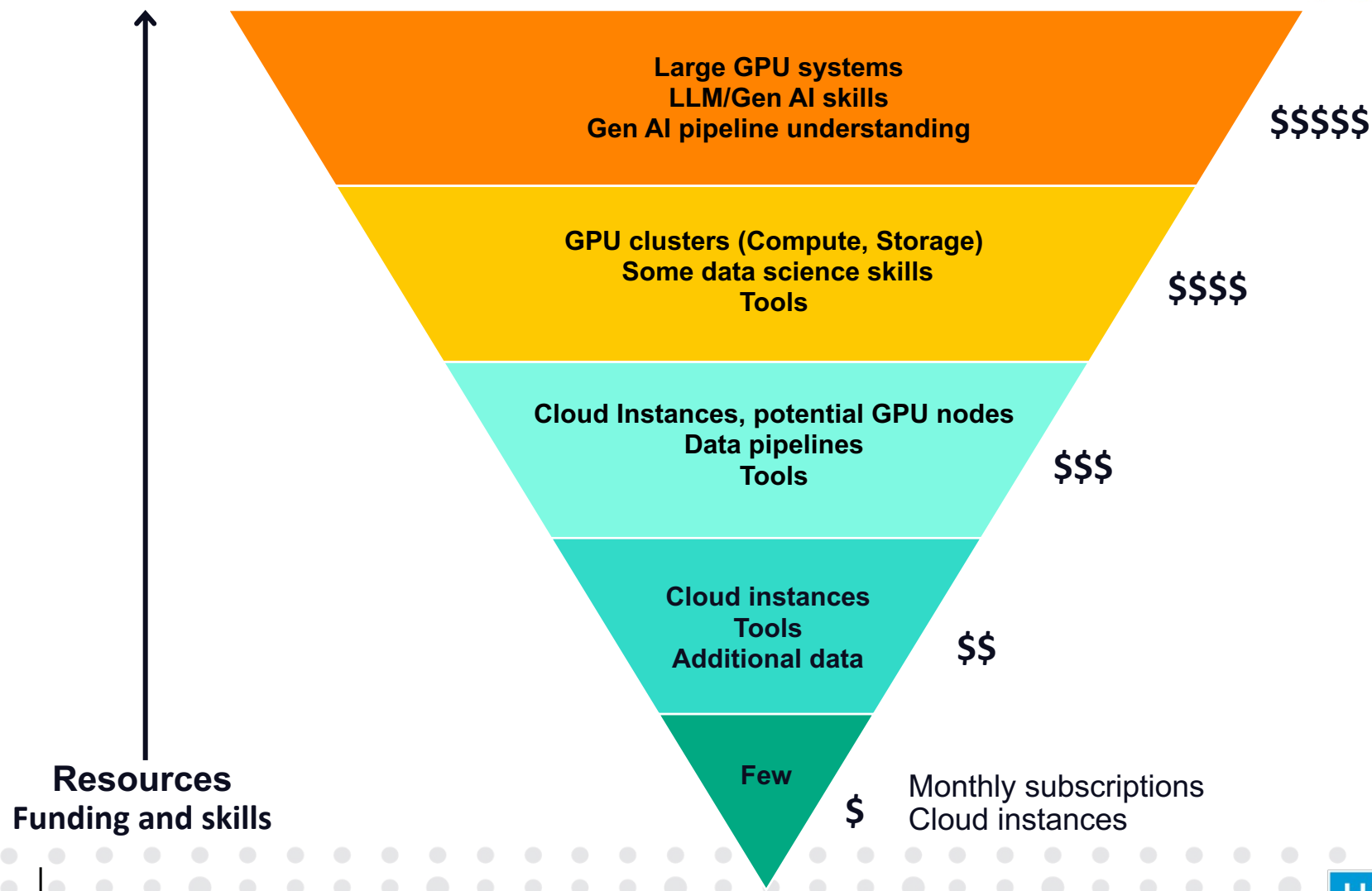# What Is Required?



Tighter

Processes & controls

Adoption oriented

Smaller

Audience

Broader

NN Design

Retraining

Fine-tuning

Orchestration

Consumption

# What Is Required? I Depends On Where You Want To Play

**Large GPU systems**
**LLM/Gen AI skills**
**Gen AI pipeline understanding**

$$$$$

**GPU clusters (Compute, Storage)**
**Some data science skills**
**Tools**

$$$$

**Cloud Instances, potential GPU nodes**
**Data pipelines**
**Tools**

$$$

**Cloud instances**
**Tools**
**Additional data**

$$

**Few**

$  Monthly subscriptions
Cloud instances

**Resources**
**Funding and skills**

# Artificial Intelligence (Gen AI and Beyond)

**Models and Services**  $+$  **AI Platform**  $+$  **Infrastructure**

Open-source models

Commercial partnerships

Services to better drive AI models

Data management

Model training

Model inference

From bare metal to containers

From training to inference

From supercomputer to edge

**One platform | Vendor neutral | Cloud neutral | AI accessible for all**

# What Are Your Key Considerations?

Models and Services

AI Platform

Data Services

Infrastructure
Software

Hardware
Infrastructure

If this cat were elected president, its first order of business would be to...

**Completion**                    Base—64 GPUs for 37 days—56,832 GPU hours

make sure the economy is strong.

**Completion**          Extended—256 GPUs for 22 days—135,168 GPU hours

declare war on the dog.

# Infrastructure I Where to Deploy?
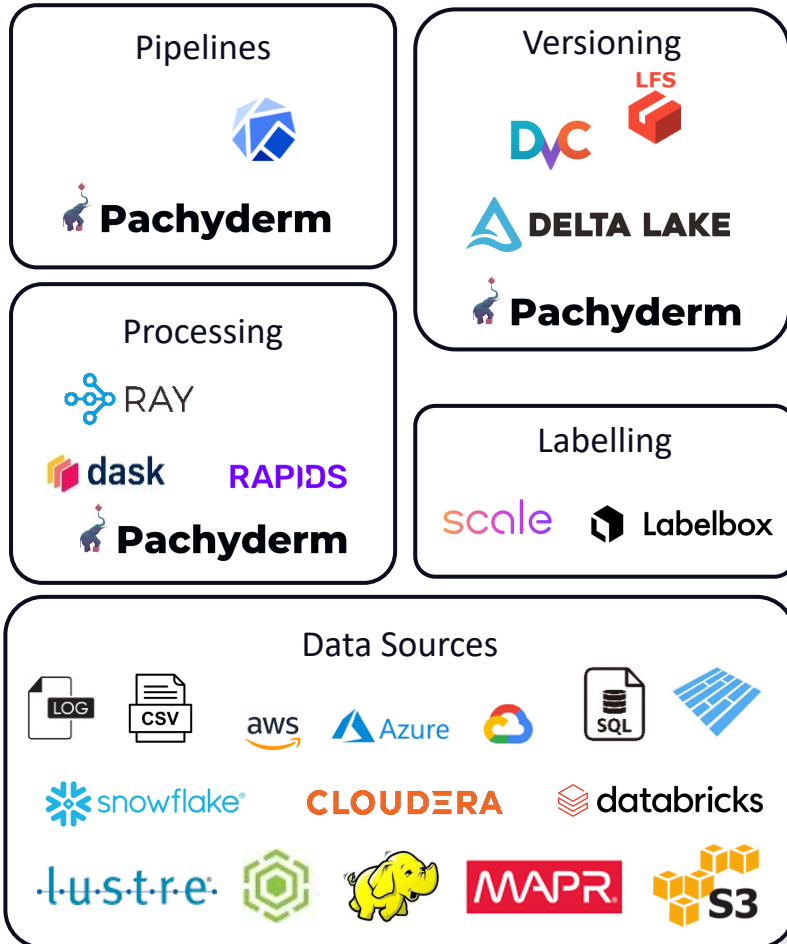
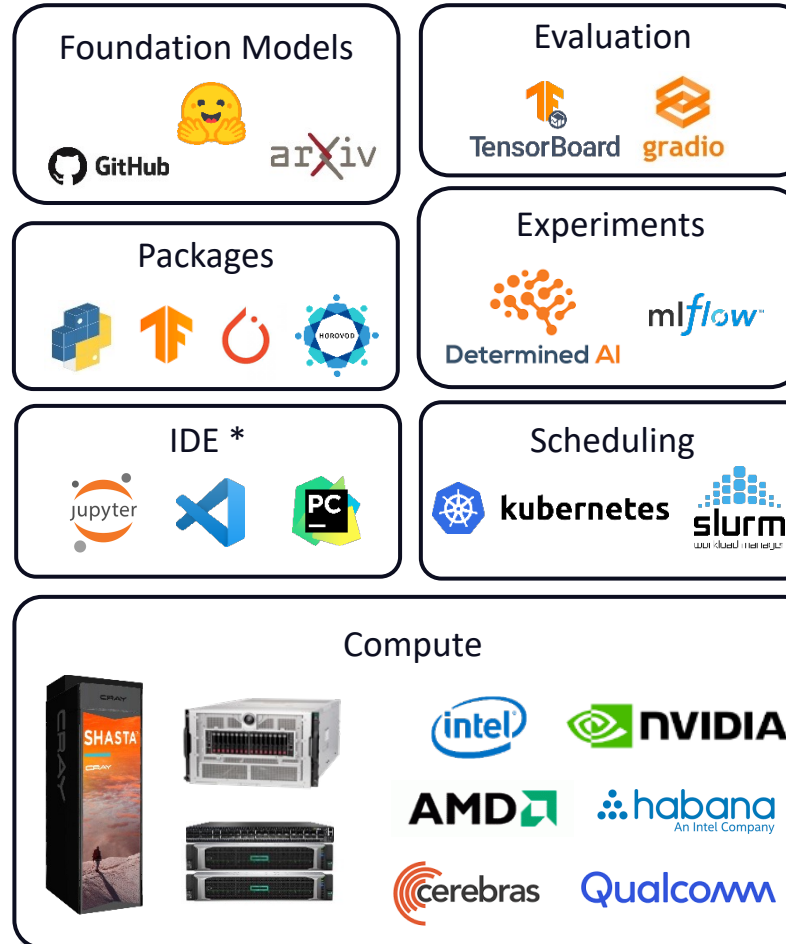|  | Pros | Cons |
|---|---|---|
| **Private Cloud** | • Enhanced Security & Privacy<br>• Customization<br>• Predictable performance<br>• Compliance adherence<br>• Reduced Risk of vector lock in | • Higher Initial Investment<br>• Limited Scalability<br>• Resource underutilization<br>• Complexity |
| **Public Cloud** | • Cost Efficiency<br>• Scalability<br>• Global Accessibility<br>• Maintenance and Updates<br>• Innovation | • Security and Privacy Concerns<br>• Dependency<br>• Limited Control<br>• Potential Compliance Issues<br>• Vendor Lock-In |

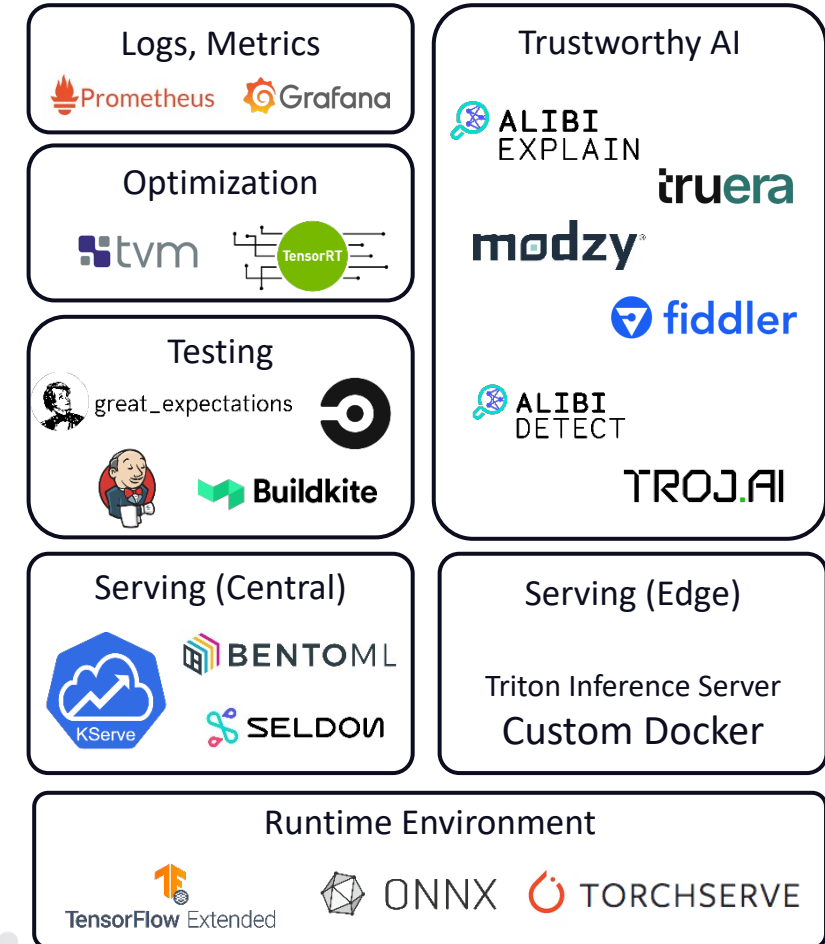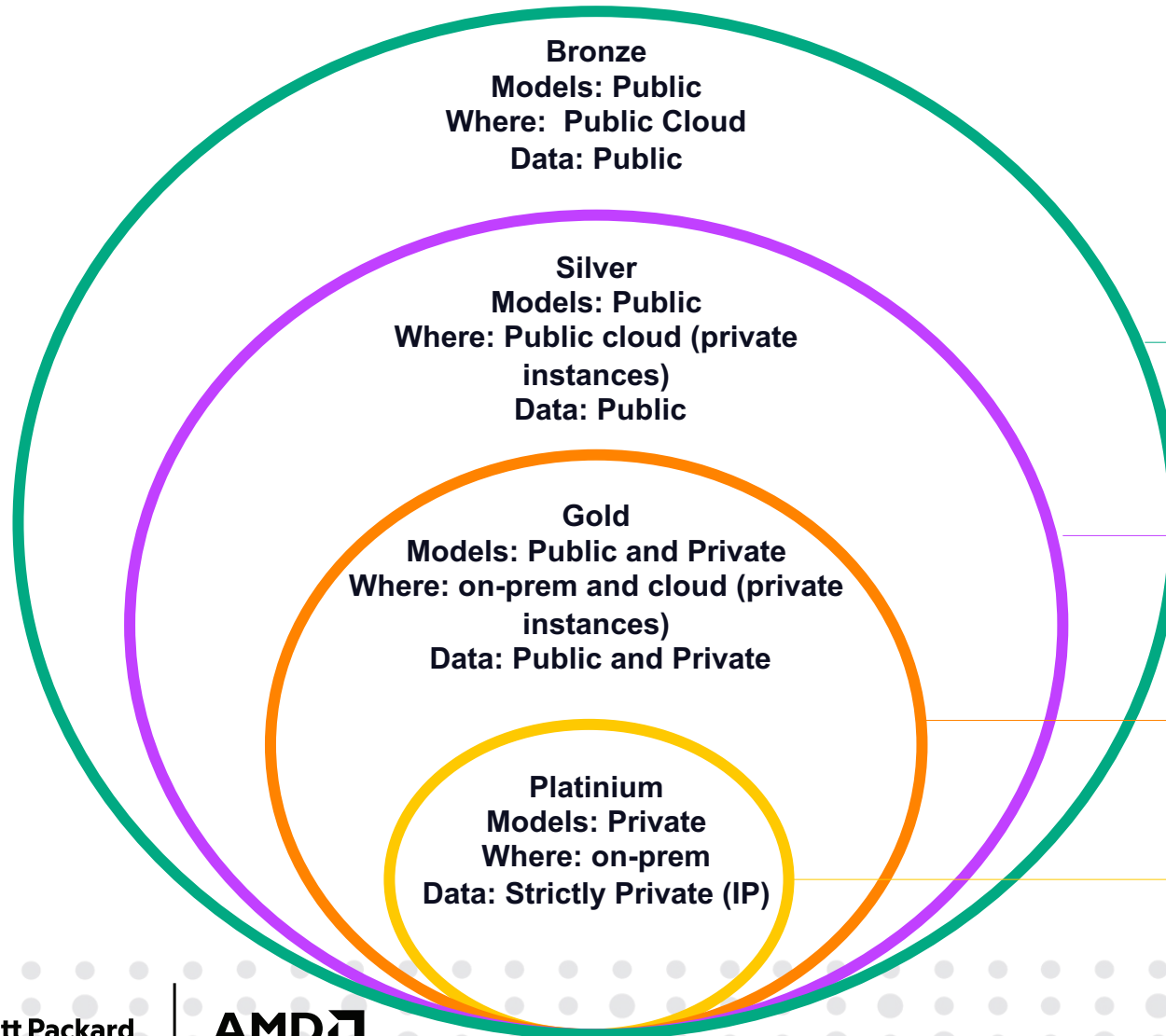# Platform | A Very Complicated Landscape

# Models | How Do You Choose the Right One? It Depends



Enterprises need hybrid environments to support multiple safety, control, and privacy requirements.

**Bronze**
Models: Public
Where: Public Cloud
Data: Public

**Silver**
Models: Public
Where: Public cloud (private instances)
Data: Public

**Gold**
Models: Public and Private
Where: on-prem and cloud (private instances)
Data: Public and Private

**Platinium**
Models: Private
Where: on-prem
Data: Strictly Private (IP)

ChatGPT like models used, either direct access or via APIs.

Models fine tuned in the public cloud, public models API end points used in private instances, own models built

Models trained on prem.
Models fine-tuned in the public cloud or on prem.
Inference running on prem and public cloud

Models trainined and running on prem

Hewlett Packard Enterprise

AMD

HPC | QUANTUM | DATA | AI

HPC + AI WALL STREET

# Models | Choosing the Right Foundation Model

- Select the best AI foundation model for your use case

- Evaluation categories:
  - Project Requirements
  - Model Capabilities
  - Operational & Ethical Considerations
  - Ethical & Governance Considerations

- Align model choice with project goals and organizational policy

# Implementation | What About Security?



Diagram elements:

BAD ACTOR

**Threats** (KEY)
**Request Attacks**
**Response Attacks**

Data Injection
Prompt Injection
Request misdirection
Prompt Capture / Leaking
Prompt Jailbreaking / Manipulation

Filter prompt input (Governance model)
Validate prompt and contextual data
Apply model permission schema

USER
UI
LLM SERVICE
LLM GATEWAY

PRIVATE INFRASTRUCTURE
MODEL
INFERENCE

PUBLIC CLOUD
MODEL
INFERENCE

Hallucination

DATA REPO
PRIVATE DATA

Filter response (Governance model)
Validate accuracy and apply contextual warning
Validate response source

Response Pollution
Exchange Capture / Leak

BAD ACTOR

- Fail gracefully and secretly
- Audit user authorization
- Secure access to external resources
- Parameterize and validate all inputs and outputs
- Avoid persisting changes when possible
- Adversarial testing

# Implementation | Governance

- Allow employees to **securely and auditably** use new technology to assist with researching, writing, and code generation.

- Keep your **sensitive data on premise.**

- Deliver the ability to **choose your model.**

- Enforce **company policy**.

- *KEY DELIVERABLES*
  - *AI guardrails*
  - *Model Token Management*
  - *Access to models either on-prem or public cloud*

**Consumption**
History, Search, Model Selector,
Prompt Marketplace,
Budget Controls, ...etc.

**Orchestration
and
Fine-Tuning**

**Gen AI App**
Authentication,
Authorization,
Audit,
Metering

| Open-Source Models | 3rd Party Models | Tools |

HPC + AI WALL STREET

HPC QUANTUM DATA AI

Hewlett Packard Enterprise

AMD

# What Else? | Observability | Sustainability | FinOps

# AI Strategy on One Page

| | | | | |
|---|---|---|---|---|
| **Models and Services** | Accelerate the journey to Gen AI value for the Enterprise | **Model Hub**<br>**HPE A&PS Services** | **Model-aaS** | |
| **AI Platform** | Drive collaborative AI model development from Data preparation to Deployment | Machine Learning Development Environment<br>Machine Learning Data Management<br>HPE AI Essentials | **AI PaaS** | Observability<br>Sustainability<br><br>FinOps |
| **Data Services** | Manage data governance, availability, and security | Data Fabric<br>Data Observability<br>Machine Learning Data Management | | |
| **Infrastructure Software** | Allow AI models to run efficiently anywhere | Private Cloud for BMaaS, VMaaS, and CaaS. | **IaaS** | |
| **Hardware Infrastructure** | Optimize the performance of AI workloads with AI native infrastructure | Accelerators: Heterogenous support for Nvidia/AMD/Intel<br>Compute: Cray, Proliant<br>Storage: GL4FL,<br>Networking: Slingshot, Mellanox<br>Best in Class DLC | | |

# AI Use Case Families in Financial Services: Spanning Front, Middle, and Back Office

**HPC + AI WALL STREET**

**FRONT OFFICE**

| Credit Scoring | Customer Acquisition | Customer Experience | Customer Retention |
|---|---|---|---|

**MIDDLE OFFICE**

| Trading | Process Automation | Knowledge Management | Risk Management |
|---|---|---|---|
| Agent-Based Modelling | Stress Testing | Regulatory Compliance | Fraud Detection |

**BACK OFFICE**

| Synthetic Data | Data Architecture | Data Capabilities | Infrastructure Optimization |
|---|---|---|---|

Hewlett Packard Enterprise | AMD

HPC | QUANTUM | DATA | AI

# Getting Started: Pick Your Use Case

**HPC + AI WALL STREET**

|  | **Regulatory Compliance** | **Credit Scoring** | Development | Production |
|---|---|---|---|---|
| **End-User Group** | Internal Users<br>Employee domain SME's (CPA's, compliance officers,<br>Accountants, auditors, (CAO/CFO)<br>Small Userbase, Highly Specialized<br>Able to provide feedback on early iterations<br>Low exposure to bad actors | External and/or Internal Users<br>Customers or customer-facing employees (Loan officers, salespeople, brokers)<br>Large Userbase, Experience levels may vary<br>High exposure to bad actors | | |
| **Use Case Regulations** | Regulations mandate what to report to who and when, with specific instructions on preparation methodology<br>"Automate as much as possible to minimize human error" -regulators | Fair lending, anti-discrimination, industry regs<br>High-risk Use Case- AI Act (EU)<br>GDPR/CCPA/Data protection regulations<br>"AI can't do anything a human wouldn't be allowed to do" -regulators | Security<br>Guardrails<br>Governance | Observability<br>Sustainability<br>FinOps |
| **Source Data (Unstructured)** | Public information (Regulatory text, instructions, etc.)<br>Trusted source- directly from regulators<br>Low potential for wrongdoing if leaked | Non-public, Personally Identifiable Information<br>Various sources- trustworthiness may vary<br>High potential for wrongdoing if leaked<br>Data sovereignty & data gravity concerns | | |
| **Models Used** | Open-source foundation model (Llama3 based)<br>Out of box, no customization | Proprietary vendor or internally developed models<br>Highly customized | | |
| **Hardware Infrastructure** | Small-scale inference environment or cloud API solution | Private, hosted, secure AI Inference at Scale cluster<br>AI Factory training resources if building your own models | | |

**Hewlett Packard Enterprise**  |  **AMD**

**HPC QUANTUM DATA AI**

# Internal Case Study:
# Hewlett-Packard International Bank DAC

- Hewlett Packard International Bank DAC, ("HPIB") is a credit institution authorized by the Central Bank of Ireland and owned by HPE.

- The Capital Requirements Regulation 3 (CRR3) is a new regulation spanning 646 pages that goes into effect for banks in the EU on January 1, 2025. It is part of the wider "Basel 3.1"/ "Basel 3 Endgame" regulatory risk framework for global financial institutions.

- The CRR3 requirements have been consistently modified leading up to implementation, and portions still lacked fully defined requirements with less than 1 year to go.

- Risk jargon often requires a CPA with specific experience in the financial sector to understand these standards.

Lengthy and complex risk regulations

**+**

Highly technical subject matter

**=**

Complexity, Confusion, Delays, Increased risk

# Problem Statement

- The financial services industry complies with a complex and fast-changing web of regulatory requirements.

- International firms are forced to adapt to the pace of change in the standards and laws of multiple regulators, jurisdictions, and governments

- Failure to comply can lead to civil (and even criminal) liability. Fines and sanctions can amount to **billions of dollars** for large institutions

- Each regulation or report change has the potential to have no impact, low impact, or very high impact and a long adoption process.  Expert interpretation is required to determine "Does this change affect us?" and "What do we need to do differently than before?



The greatest compliance challenges I expect to face in 2022 is/are...

Volume and implementation of regulatory change
Lack of budget and resources

# Initial Experiments



## Experiment 1

Bulk analysis of documents on regulator website via Python to produce categorized CRR3 Implementation Report

## Experiment 2

Analysis triggered by frequent regulator website updates, which can provide automated guidance and alerts for impactful changes and opportunities to comment

## Experiment 3

Chatbot experience for real-time self-service with follow-up questions

# Key Benefits

**Transparency**: End-to-end tracking of new proposed regulatory changes, opening and closing of public comment periods, adoption of final changes, and phased implementation deadlines.

**Operational efficiency**: Swift document analysis and reporting. Give users opportunity to ask questions on complex subject matter in natural language and receive proactive updates on what they need to know.

**Cost savings**: Minimize burden of adapting to new legislation, manual reviews and potential fines for non-compliance, and consulting spend.

**Increase collaboration**: Create wider awareness of upcoming changes.

# Part of a Holistic Compliance Strategy

**Generative AI analysis is one part of a comprehensive strategy**

**AI-Powered Analysis**:
- Swift document review and extraction of essential data.
- Predictive modeling for potential compliance risks.
- Real-time notifications on regulatory changes.
- Simplification of complex regulatory texts using NLP.
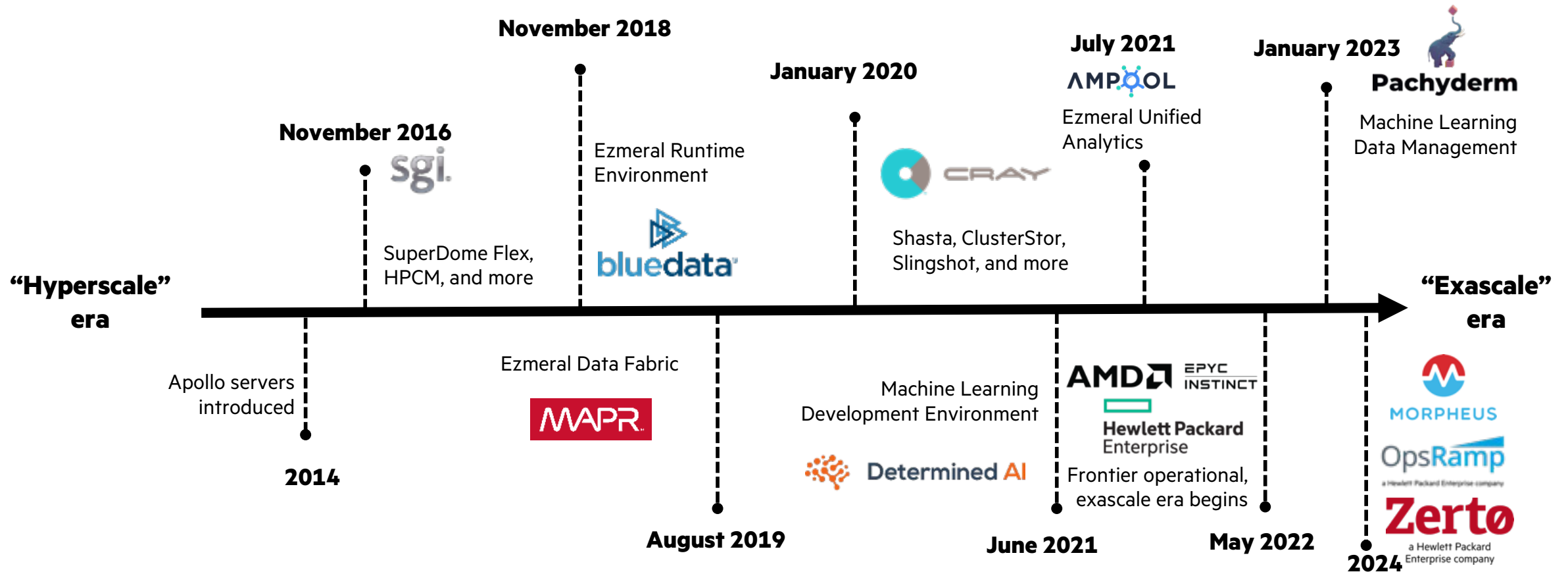
**Training & Development**:
- Regular workshops to educate staff on compliance matters.
- Role-specific training modules for targeted learning.
- Case studies and simulations for hands-on experience.

**Human Oversight**:
- Expert panels for manual review and validation.
- Routine audits to ensure AI-generated results align with compliance requirements.
- Feedback loops for continuous AI model refinement.

# HPE Has Been Preparing for Some Time

# Enabling Artificial Intelligence

**HPE GreenLake**

**Models and Services** + **AI Platform** + **Infrastructure**

| Models and Services | AI Platform | Infrastructure |
|---|---|---|
| Open-source models | Data management | From bare metal to containers |
| Commercial partnerships | Model training | From training to inference |
| Services to better drive AI models | Model inference | From supercomputer to edge |

**One platform | Vendor neutral | Cloud neutral | AI accessible for all**

Hewlett Packard Enterprise | AMD

HPC | QUANTUM | DATA | AI