# How AI
## *has - is - will*
# impact IT Service
# Management programs

**JOIN US FOR THIS WEBINAR**
**THURSDAY, MAY 16TH · 11AM ET**

Mark R. Hinkle
CEO
Peripety Labs

# Mark Hinkle, I Help Enterprises use Artificial Intelligence

25 Years of Executive Leadership and Nerding out over Emerging Technologies.

**MindSpring** — Director of Tech Support  [Nation ISP]

**zenoss** — VP of Developer Relations [Open Source Systems Management]

**cloud.com** — VP of Developer Relations [Open Source Cloud Computing]

**citrix** — Head of Open Source Business Office [Enterprise Software]

**THE LINUX FOUNDATION** — VP of Marketing [Open Source Enterprise Software]

**node JS Foundation** — Executive Director [Open Source Application Development]

**TRIGGERMESH** — CEO and Co-founder [Cloud and SaaS]

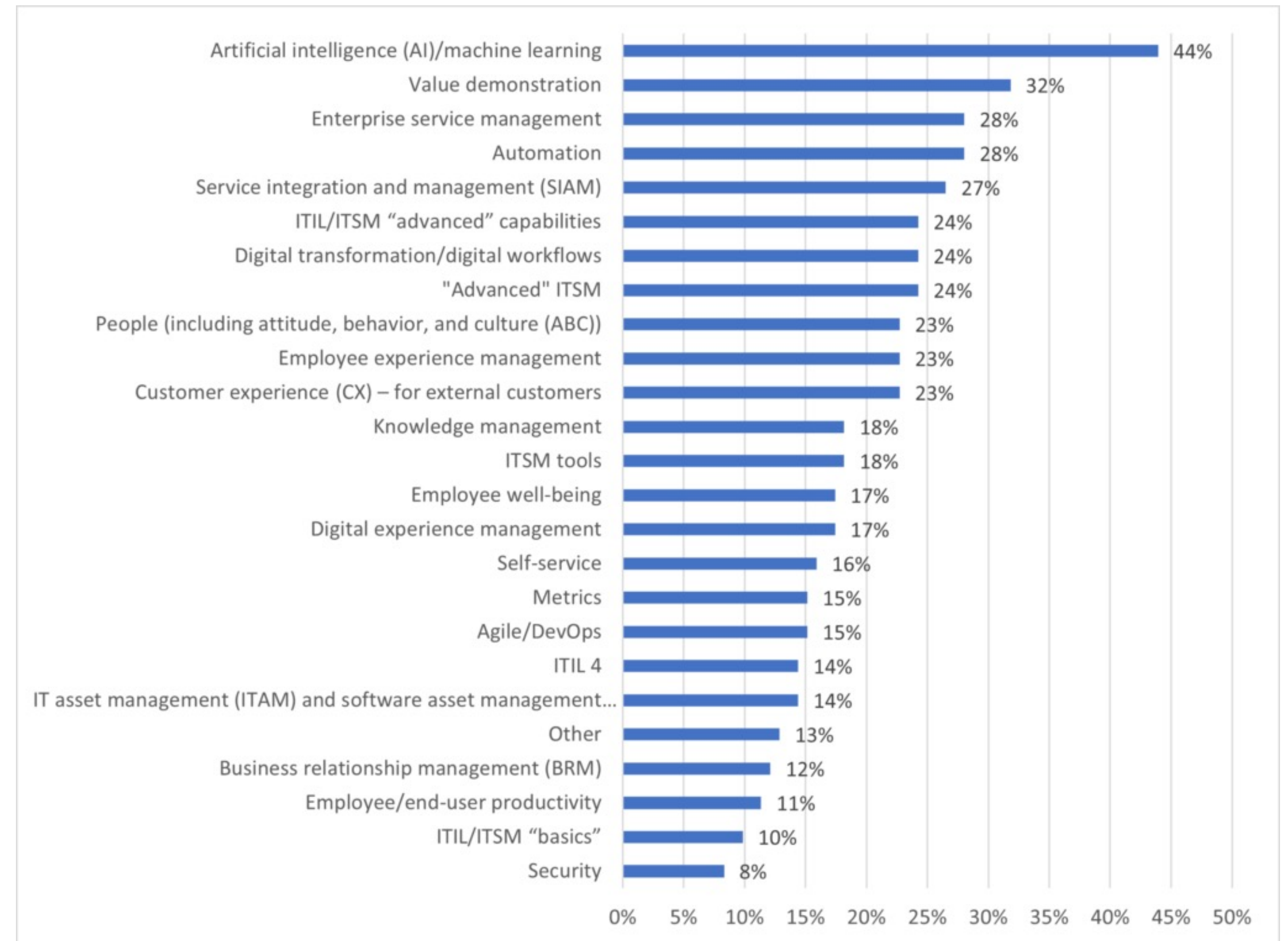**peripety labs** — Founder and CEO [Artificial Intelligence]

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# Five Hottest ITSM Trends for 2024

AI Moves from #5 in 2023 to #1 in 2024

The top five ITSM trends for 2024 are:

1. Artificial intelligence/machine learning (44%)

2. Value demonstration (32%)

3. Automation (28%)

4. Enterprise service management (28%)

5. Service integration and management (27%)

**ITSM TOOLS** — ITSM Trends for 2024



| Trend | Value |
|---|---|
| Artificial intelligence (AI)/machine learning | 44% |
| Value demonstration | 32% |
| Enterprise service management | 28% |
| Automation | 28% |
| Service integration and management (SIAM) | 27% |
| ITIL/ITSM "advanced" capabilities | 24% |
| Digital transformation/digital workflows | 24% |
| "Advanced" ITSM | 24% |
| People (including attitude, behavior, and culture (ABC)) | 23% |
| Employee experience management | 23% |
| Customer experience (CX) – for external customers | 23% |
| Knowledge management | 18% |
| ITSM tools | 18% |
| Employee well-being | 17% |
| Digital experience management | 17% |
| Self-service | 16% |
| Metrics | 15% |
| Agile/DevOps | 15% |
| ITIL 4 | 14% |
| IT asset management (ITAM) and software asset management… | 14% |
| Other | 13% |
| Business relationship management (BRM) | 12% |
| Employee/end-user productivity | 11% |
| ITIL/ITSM "basics" | 10% |
| Security | 8% |

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# How Generative AI is Going to Affect ITSM

# Enhanced Efficiency and Automation

Boosting Efficiency Through Automation

Generative AI can automate numerous ITSM tasks, such as ticketing and incident management, which traditionally require significant human intervention.

By automating these processes, generative AI reduces the workload on IT staff, allowing them to focus on more complex and strategic tasks.

**Examples**

- AI categorizes and prioritizes incoming service requests.
- Automatically responds to common queries.
- Predicts and addresses issues before they escalate.

# Improved User Experience

Elevating User Satisfaction

The integration of generative AI into ITSM tools significantly enhances the user experience. AI-powered systems can offer personalized service experiences by understanding user preferences and past interactions. This capability not only speeds up the resolution of issues but also enhances user satisfaction.

**Examples**

- Personalized service experiences.
- Faster issue resolution.
- Proactive service tailored to user preferences.

# Proactive Service Management

## Predictive and Preventive ITSM

Generative AI excels in predicting potential IT issues before they become critical, which is a significant shift from the reactive nature of traditional ITSM. By analyzing patterns and historical data, AI can forecast incidents and automate preventive measures, thereby minimizing downtime and improving system reliability.
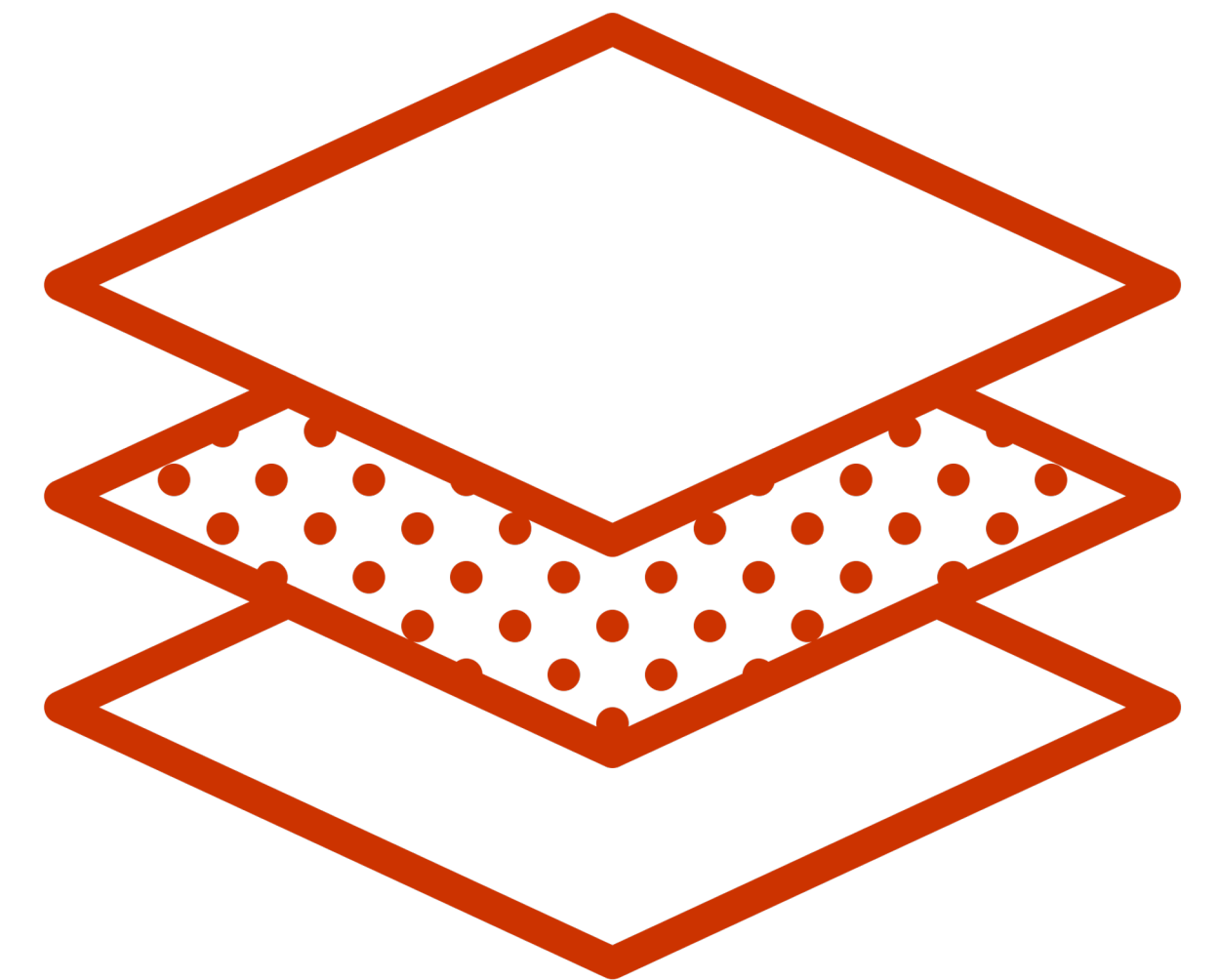
**Examples**

- Predicting incidents.
- Automating preventive measures.
- Minimizing downtime.

# Infrastructure Stacks Definition

The way we talk about groups of infrastructure

A **tech infrastructure stack**, often simply referred to as a "tech stack," is a combination of software tools, frameworks, and technologies used to build and run a digital application or service.

The term "stack" is used because these components often **build on top of one another**, much like how items are stacked on top of each other
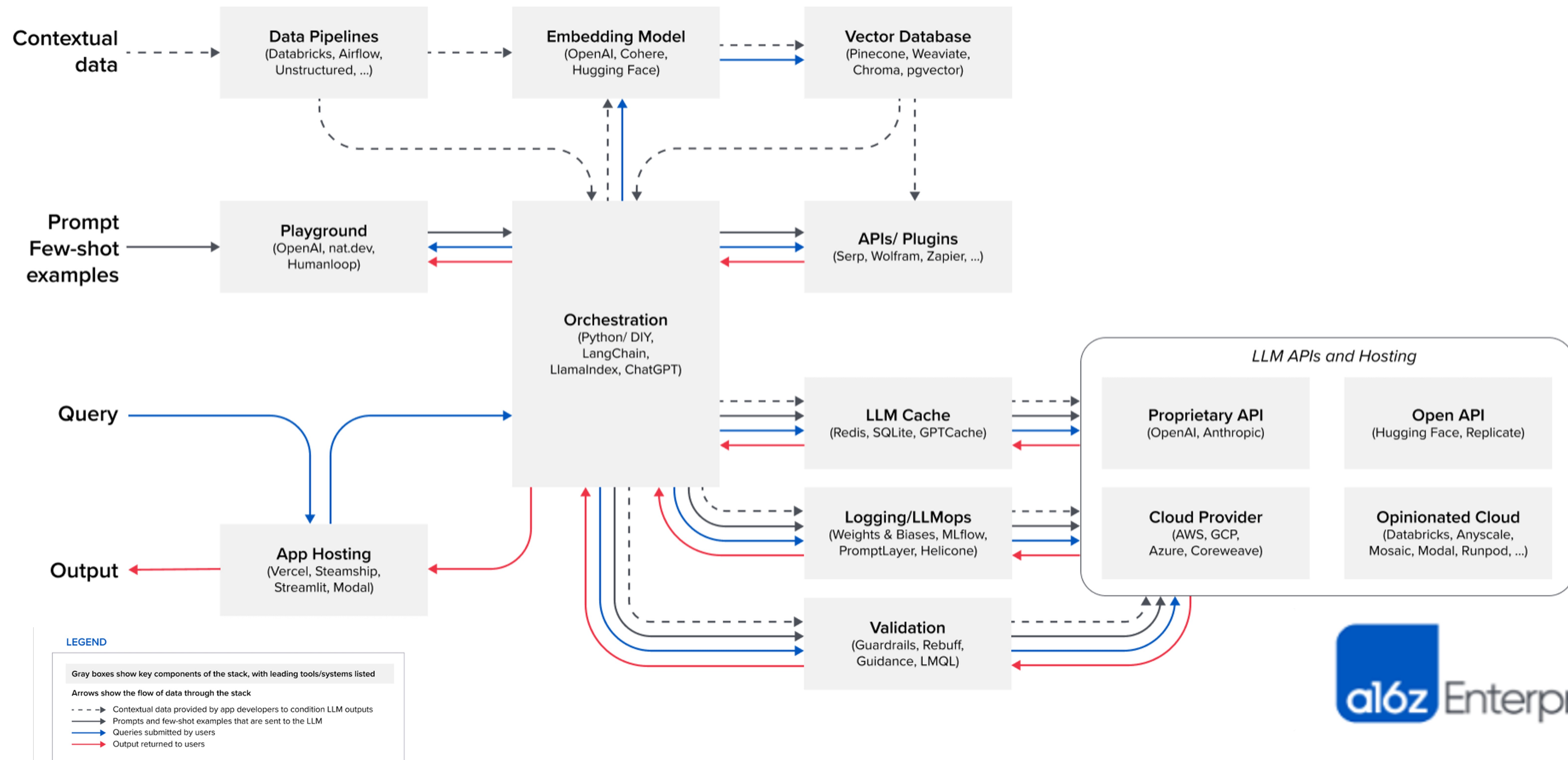
# An Ontology but not a Stack

A way to classify AI infrastructure

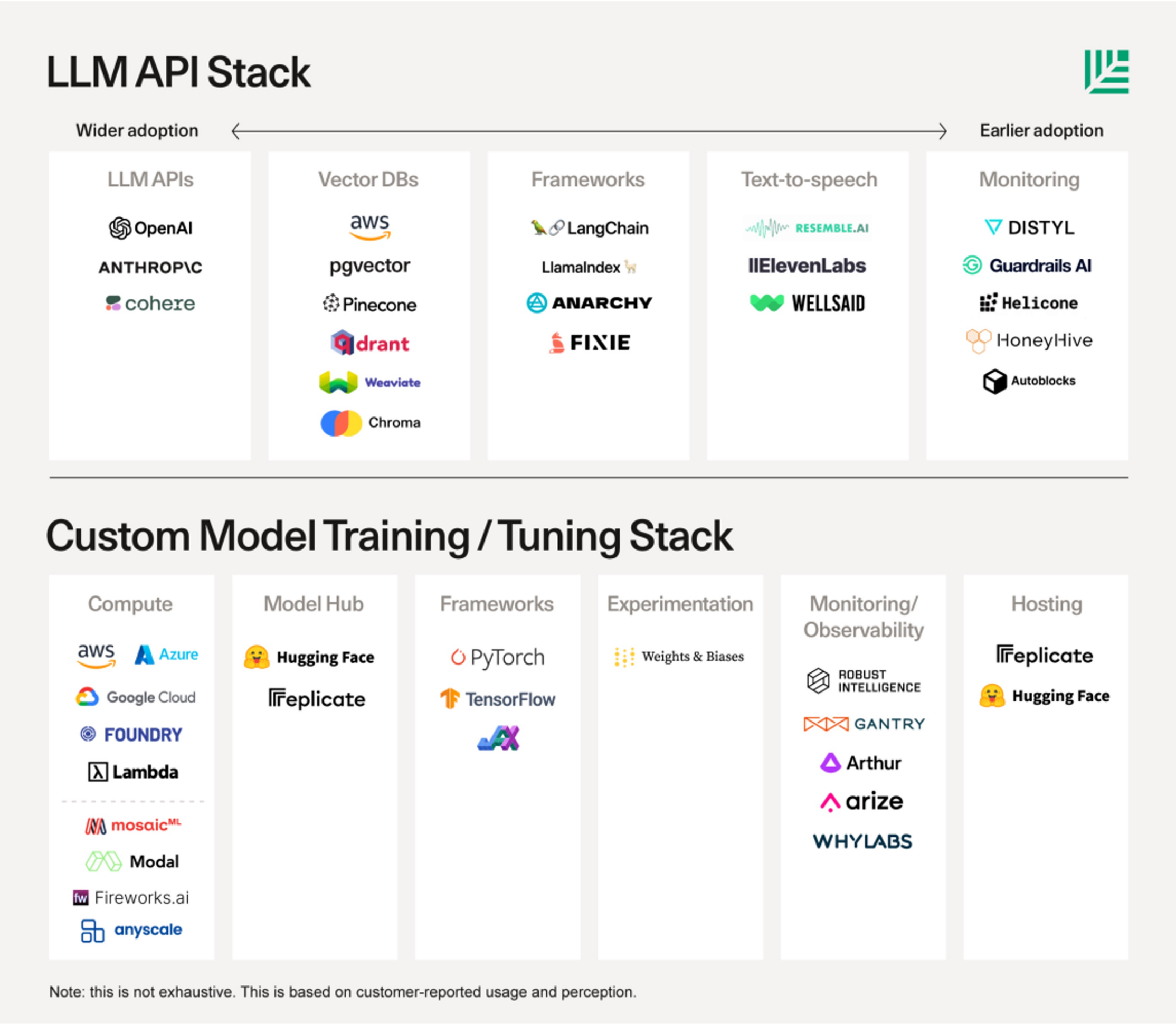| | |
|---|---|
| **Infrastructure Layer** | • **Hardware:** Specialized hardware like GPUs (Graphics Processing Units) and TPUs (Tensor Processing Units) that accelerate machine learning computations.<br>• **Cloud Services**: Platforms like AWS, Google Cloud, and Azure offer cloud-based infrastructure optimized for AI workloads. |
| **Data Layer** | • **Data Storage**: Databases and data lakes where raw data is stored.<br>• **Data Processing**: Tools like Apache Spark and Hadoop for processing large datasets.<br>• **Data Labeling**: Platforms like Amazon SageMaker Ground Truth or tools like Labelbox for annotating data. |
| **Machine Learning Frameworks** | • Deep Learning: Frameworks like TensorFlow, PyTorch, and Keras.<br>• **Traditional ML**: Libraries like Scikit-learn, XGBoost, and LightGBM. |
| **Model Development & Training** | • Development Environments: Jupyter notebooks, Google Colab, etc.<br>• **Training Platforms**: Platforms like Google AI Platform, AWS SageMaker for training models at scale. |
| **Deployment & Serving** | • Model Serving: Tools like TensorFlow Serving, NVIDIA Triton Inference Server.<br>• **Deployment Platforms**: AWS Lambda, Google Cloud Functions, Azure Machine Learning for deploying models as APIs. |
| **Management & Monitoring** | • **Model Management:** Tools like MLflow, TFX (TensorFlow Extended) for managing the lifecycle of ML models.<br>• **Monitoring**: Platforms to monitor the performance of deployed models and ensure they're working as expected. |
| **Application Layer** | • **Integrations:** How the AI models integrate with web services, apps, or other systems.<br>• **User Interfaces**: Dashboards, chatbots, or any other interface that interacts with the end-users. |
| **Ethics & Fairness** | • **Bias Detection:** Tools to detect and mitigate biases in AI models.<br>• **Explainability**: Tools like SHAP, LIME to explain model decisions. |

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# Example: Andressen-Horowitz Emerging LLM Stacks

Process-oriented stack



**Contextual data**

**Data Pipelines**
(Databricks, Airflow, Unstructured, ...)

**Embedding Model**
(OpenAI, Cohere, Hugging Face)

**Vector Database**
(Pinecone, Weaviate, Chroma, pgvector)

**Prompt Few-shot examples**

**Playground**
(OpenAI, nat.dev, Humanloop)

**APIs/ Plugins**
(Serp, Wolfram, Zapier, ...)

**Orchestration**
(Python/ DIY, LangChain, LlamaIndex, ChatGPT)

**Query**

**Output**

**App Hosting**
(Vercel, Steamship, Streamlit, Modal)

**LLM Cache**
(Redis, SQLite, GPTCache)

**Logging/LLMops**
(Weights & Biases, MLflow, PromptLayer, Helicone)

**Validation**
(Guardrails, Rebuff, Guidance, LMQL)

*LLM APIs and Hosting*

**Proprietary API**
(OpenAI, Anthropic)

**Open API**
(Hugging Face, Replicate)

**Cloud Provider**
(AWS, GCP, Azure, Coreweave)

**Opinionated Cloud**
(Databricks, Anyscale, Mosaic, Modal, Runpod, ...)

**LEGEND**

Gray boxes show key components of the stack, with leading tools/systems listed

Arrows show the flow of data through the stack

- - - → Contextual data provided by app developers to condition LLM outputs
——→ Prompts and few-shot examples that are sent to the LLM
——→ Queries submitted by users
——→ Output returned to users

a16z Enterprise

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# Example: Sequoia New Language Model Stack

Function-Oriented Stack



**LLM API Stack**

Wider adoption ← → Earlier adoption

| LLM APIs | Vector DBs | Frameworks | Text-to-speech | Monitoring |
|---|---|---|---|---|
| OpenAI | aws | LangChain | RESEMBLE.AI | DISTYL |
| ANTHROP\C | pgvector | LlamaIndex | ElevenLabs | Guardrails AI |
| cohere | Pinecone | ANARCHY | WELLSAID | Helicone |
| | drant | FIXIE | | HoneyHive |
| | Weaviate | | | Autoblocks |
| | Chroma | | | |

**Custom Model Training / Tuning Stack**

| Compute | Model Hub | Frameworks | Experimentation | Monitoring/ Observability | Hosting |
|---|---|---|---|---|---|
| aws  Azure | Hugging Face | PyTorch | Weights & Biases | ROBUST INTELLIGENCE | replicate |
| Google Cloud | replicate | TensorFlow | | GANTRY | Hugging Face |
| FOUNDRY | | | | Arthur | |
| Lambda | | | | arize | |
| mosaicML | | | | WHYLABS | |
| Modal | | | | | |
| Fireworks.ai | | | | | |
| anyscale | | | | | |

Note: this is not exhaustive. This is based on customer-reported usage and perception.

Sequoia's takeaway is that **AI is moving too quickly to have high confidence in the end-state stack**, but there was consensus that LLM APIs will remain a key pillar, followed in popularity by retrieval mechanisms and development frameworks like LangChain.

Open source and custom model training and tuning also seem to be on the rise. Other areas of the stack are important, but earlier in maturity.

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# Webs Not Stacks

Microservices and cloud architecture change the game

The traditional notion of a **tech infrastructure "stack" implies a linear, layered approach where one component is built on top of another** in a sequential manner.

As technology has evolved, especially with the rise of cloud services, microservices, and distributed systems, the architecture of modern infrastructure is becoming **more interconnected and less linear**.

This interconnectedness **makes the new infrastructure resemble more of a "web" than a traditional "stack**.



THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# A Simplified Framework For AI Operations

# The New AI Stack

Microservices and cloud architecture change the game

**The new AI Stack is a way to think about how the components of AI Infrastructure** that maps to the AI Operations tools to manage those component.

It looks **more like a web than a stack because it is more of a loose—coupling** of services that have have common traits rather than a linear stack.

## AI Infrastructure

**Data**
All the data that is inputed, analyzed, and outputted by AI models.

**Fine-Tuning and Training**
Process for incorporating enterprise data into Large Language Models.

**Narrow AI/ Autonomous Agents**
AI models trained for specific uses and agents that are goal oriented.

**Model Integration Frameworks**
Facilitate the development, integration, and deployment of applications powered by Large Language Models.

## AI Operations

**CI/CD**
How we dploy and update LLMS and other AI infrastructure.

**Monitoring and Observability**
Performance and insight into how your LLMs and other AI infrastructure is working.

**Configuration Management**
Systematic process of establishing and maintaining consistency in the model's performance

**Security**
Security for inputs and outputs of the LLM including training and fine-tuning as well as integration with AI - applications

APIs

APIs

Enterprise LLM

# The New AI Stack

The AI stack is a **loosely-coupled group of services that not only provide AI services** but the management of the services. **Relationships are not linear but rather one to many**.

# Artificial Intelligence Models

Models are algorithms that power services like ChatGPT, Google Bard,

An **AI model** is a program or algorithm that is trained on a set of data to make **predictions** or decisions. AI models are used in a wide variety of applications, including **natural language processing**, **image recognition**, and **fraud detection**.

# Size of AI Models

Parameters in AI models are the internal variables that enable machines to learn and make predictions

**Parameters** encapsulate the knowledge and behavior of the model, representing the weights and biases that govern its decision-making process.

During the training phase, these parameters are iteratively adjusted to minimize the disparity between the model's predictions and the desired outputs.

The number and complexity of parameters play a significant role, as they influence the model's capacity and computational requirements. **Ultimately, parameters form the backbone of AI models**, allowing them to generalize from training data and make accurate predictions on new, unseen inputs.

## NUMBER OF PARAMETERS IN AI MODELS

Open AI's ChatGPT-4 has nearly twice as many parameters with 100 trillion in their model versus second place contender Google PaLM2. Though once Microsoft and Google start rolling out their AI to their productivity suites, it's possible they will eclipse OpenAI. Though it's not necessarily about bigger but the quality of data also plays a roll in the effectiveness of the model.

| Model | |
|---|---|
| NVIDIA Megatron-LM | |
| Microsoft Turing NLG | |
| Meta AI LLAMA (largest) | |
| OpenAI GPT-3 | |
| Microsoft and NVIDIA Megatron NLG | |
| Google PaLM 2 | |

0   200,000,000,000   400,000,000,000   600,000,000,000

# AI Infrastructure – Enterprise LLMS

The Center of the Enteprise AI infrastructure

**Enterprise LLM**

An enterprise LLM can serve as a centralized knowledge hub, answering queries, offering insights, and assisting in decision-making processes.

- **Public Model:** Leveraging a pre-trained model like OpenAI's GPT models offers the advantage of a vast knowledge base. However, it might not be tailored to the specific nuances and terminologies of an enterprise.

- **Custom Foundation Model:** Training a foundation model on enterprise data ensures that the LLM understands the organization's specific context, jargon, and nuances. This approach is more resource-intensive but offers a higher degree of customization.

# AI Infrastructure

Classifying Infrastructure

The traditional notion of a **tech infrastructure "stack" implies a linear, layered approach** where one component is built on top of another in a sequential manner.
However, as technology has evolved, especially with the rise of cloud services, microservices, and distributed systems, the architecture of modern infrastructure is becoming more interconnected and less linear. This interconnectedness makes **the new infrastructure resemble more of a "web" than a traditional "stack**.

## AI Infrastructure

**Data**
All the data that is inputed, analyzed, and outputted by AI models.

**Fine-Tuning and Training**
Process for incorporating enterprise data into Large Language Models.

**Narrow AI/ Autonomous Agents**
AI models trained for specific uses and agents that are goal oriented.

**Model Integration Frameworks**
Facilitate the development, integration, and deployment of applications powered by Large Language Models.

# AI Infrastructure - Data

Data the new oil of AI

Data is undeniably the cornerstone of Large Language Models (LLMs) and the broader realm of artificial intelligence. The depth, breadth, and quality of this data play a pivotal role in determining the performance, accuracy, and applicability of these models.

- **Data Warehouses:** Centralized repositories where data from various sources is stored, consolidated, and made available for analysis and querying.

- **Data Pipelines:** Systems that automate the flow of data from diverse sources to storage or processing. They handle tasks like data extraction, transformation, and loading (ETL).

- **Data Provenance and Tracking:** Systems that track the source of data and its journey through the pipeline, ensuring transparency and traceability.

# Embeddings and Vector Databases

Long-term memory for artificial intelligence models

For machine learning to "learn" from natural language, the language is converted into **numerical representations called embeddings.** These embeddings can be stored in Vector Databases. Then the machine learning algorithm starts to work through creating probabilities.

## Vector Embeddings

Embeddings, in the context of machine learning and natural language processing (NLP), refer to the representation of words, sentences, or documents as dense vectors in a continuous vector space. These vectors capture semantic and syntactic relationships between different elements of the language.

## Vector Databases

A vector database, in the context of AI, refers to a specialized type of database that is designed to store and efficiently retrieve high-dimensional vectors. Vectors, in this context, represent numerical representations of data points or entities in a multi-dimensional space. These vectors can be generated from various types of data, such as text documents, images, audio files, or other forms of structured or unstructured data.

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# Tokenization Training Artificial Intelligence Models

Tokens how LLMs are trained using NLP

In machine learning **tokens** refer to the individual units of text or language that are processed by artificial intelligence models. In natural language processing (NLP) tasks, such as text generation or language translation, **AI models break down the input text into smaller units called tokens**. These tokens can be words, subwords, or even characters, depending on how the model is trained and the specific tokenization method used.

•Word-level tokens: ["I", "love", "cats", "and", "dogs"]

•Character-level tokens: ["I", " ", "l", "o", "v", "e", " ", "c", "a", "t", "s", " ", "a", "n", "d", " ", "d", "o", "g", "s"]

## Example: Open API's Billing is based partially on tokens

**What are the rate limits for our API?**

You can view the rate limits for your organization under the rate limits section of the account management page.

We enforce rate limits at the organization level, not user level, based on the specific endpoint used as well as the type of account you have. Rate limits are measured in two ways: **RPM** (requests per minute) and **TPM** (tokens per minute). The table below highlights the default rate limits for our API but these limits can be increased depending on your use case after filling out the Rate Limit increase request form.

The **TPM** (tokens per minute) unit is different depending on the model:

| TYPE | 1 TPM EQUALS |
|------|--------------|
| davinci | 1 token per minute |
| curie | 25 tokens per minute |
| babbage | 100 tokens per minute |
| ada | 200 tokens per minute |

In practical terms, this means you can send approximately 200x more tokens per minute to an `ada` model versus a `davinci` model.

# AI Infrastructure – Fine Tuning and Training
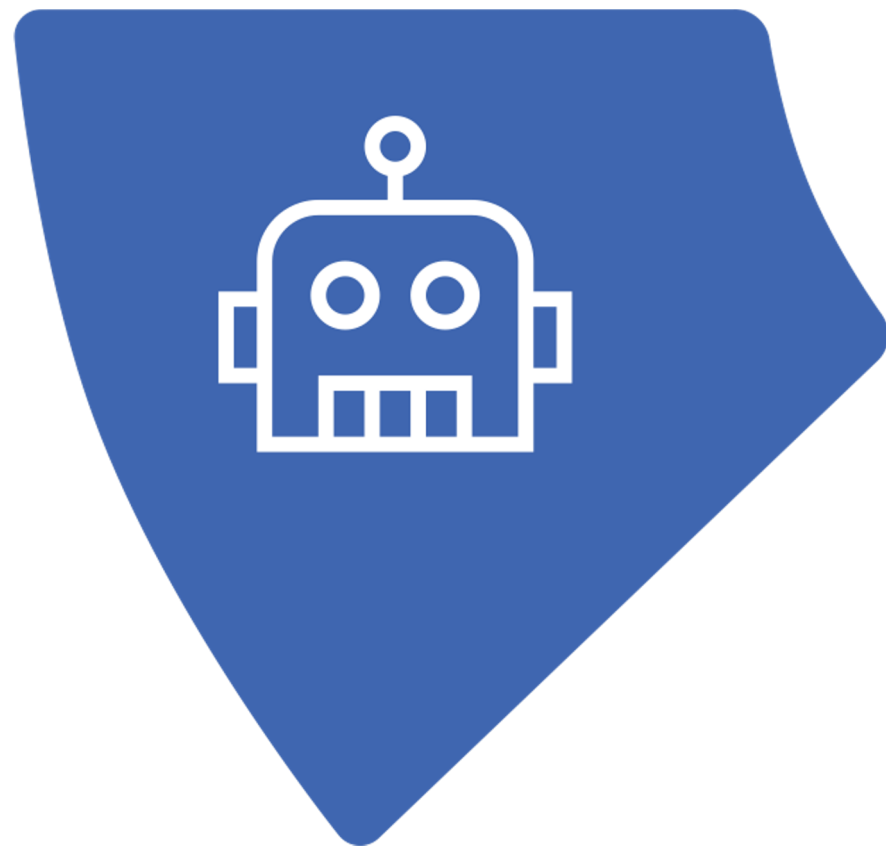
Adjusting models to work as you want

Fine-Tuning and Training form the crux of developing proficient Large Language Models (LLMs). At the heart of this process lies the foundational training, where models are exposed to vast datasets, enabling them to grasp the nuances of language and context.

- **Data Labeling Platforms:** Dedicated platforms that facilitate the annotation or marking of data. They often incorporate tools for automated labeling, manual review interfaces, quality checks, and collaboration among human annotators.
- **Machine Learning Backends:** These are the computational engines that power the training and fine-tuning processes. They handle tasks like model initialization, forward and backward propagation, optimization, and more. Often, they leverage specialized hardware like GPUs or TPUs for accelerated processing.

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# AI Infrastructure – Narrow AI and AAIA

Microservices and cloud architecture change the game

**Narrow AI**, often contrasted with general AI, is designed to perform a specific task without possessing the broad cognitive abilities that humans have. From chatbots answering customer queries to autonomous vehicles navigating the streets.

**Autonomous AI Agent (AAIA),** Autonomous artificial intelligence agents are an agent-driven type of AI that operates without consistent human intervention. They can solve various problems, make logical decisions, and handle numerous tasks without constant human input. As generative AI and foundation models continue to evolve, autonomous AI agents are growing in popularity and potential use cases.

# Autonomous Artificial Intelligence Agents

AAI Agents make Alexa and Siri look like primitive.

**Autonomous artificial intelligence** agents (AAIA) are software programs that can **operate independently** without human intervention. They are designed to learn from data, make decisions, and perform tasks without any human input. These systems are designed to be self-sufficient and **operate without the need for human oversight**.

# AI Infrastructure – Model Integration Frameworks

## Middleware for AI Infrastructure

**Langchain**, a data framework designed to aid in building applications with Large Language Models (LLMs). It provides essential tools that facilitate data ingestion, structuring, retrieval, and integration with various application frameworks.

**LlamaIndex,** a data framework designed to aid in building applications with Large Language Models (LLMs). It provides essential tools that facilitate data ingestion, structuring, retrieval, and integration with various application frameworks.

# AI For Ops

DevSecOps have typically spent their time with a toolbox that helps manage and automate the lifecycle of IT infrastructure with a certain set of tools.

By mapping those tools to new infrastructure it may make it easier for AI operations to existing systems.

This is what I call, "Ai For Ops"

# AIFOROPS

## AI Operations

**CI/CD**
How we dploy and update LLMS and other AI infrastructure.

**Monitoring and Observability**
Performance and insight into how your LLMs and other AI infrastructure is working.

**Configuration Management**
Systematic process of establishing and maintaining consistency in the model's performance

**Security**
Security for inputs and outputs of the LLM including training and fine-tuning as well as integration with AI - applications

# AI Infrastructure – CI/CD

Continuous Integration and Continuous Deployment of AI infrastructure

CI/CD stands for Continuous Integration and Continuous Deployment. Here's a breakdown of what CI/CD means for LLMs:

## Continuous Integration:

- **Model Training Integration:** As new data becomes available or as models are refined, they can be retrained. CI ensures that this retraining integrates smoothly with existing systems.

- **Automated Testing:** After integrating new changes, automated tests are run to ensure that the model still performs as expected. This might include tests for accuracy, bias, or other performance metrics.

- **Version Control:** As models are updated, it's essential to keep track of different versions. CI tools can help manage these versions, ensuring that developers can roll back to previous models if needed.

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# AI Infrastructure – CI/CD (cont'd)

Continuous Integration and Continuous Deployment of AI infrastructure

## Continuous Deployment (CD):

- **Model Deployment:** Once a model has been trained and tested, it needs to be deployed to a production environment. CD automates this deployment process, ensuring that new models can be used as soon as they're ready.

- **Scaling:** LLMs can be resource-intensive. CD tools can help deploy models across multiple servers or cloud instances, ensuring that they can handle real-world usage

.

- **Monitoring and Feedback Loop:** After deployment, CD tools react to the monitoring platform the model's performance in real-time. If the model starts to underperform or if issues are detected, the CD process can automatically roll back to a previous version or alert developers to the problem.

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# AI Infrastructure – Monitoring and Observability

Keeping an eye on performance and testing bias and outputs

Observability and monitoring, in the context of Large Language Models (LLMs), refer to the tools, practices, and methodologies used to gain insight into the performance, behavior, and health of these models during training, validation, and deployment.

## Observability

- **Performance Metrics:** Monitoring tools continuously track key performance metrics like accuracy, loss, and latency. This helps ensure that the LLM is performing as expected and allows for quick interventions if something goes awry.

- **Health Checks:** These are regular checks to ensure that the LLM is operational and responding. For deployed models, this might mean checking that the model's API endpoint is up and running.

- **Resource Utilization:** LLMs can be computationally intensive. Monitoring tools can track resource utilization, ensuring that the model isn't using more CPU, memory, or GPU resources than expected.

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# AI Infrastructure – Monitoring and Observability (cont'd)

Keeping an eye on performance and testing bias and outputs

Observability provides a deeper understanding of what's happening inside the model. For LLMs, this could mean understanding which parts of the model are activated for certain inputs, or how weights and biases change over time.

- **Internal Model Insights:** Observability provides a deeper understanding of what's happening inside the model.

- **Model Interpretability:** It's crucial to understand why an LLM produces a particular output.

- **Debugging and Troubleshooting:** If an LLM starts producing unexpected outputs or behaves erratically, observability tools can help pinpoint the cause, whether it's an issue with the model itself or with the data it's been trained on.

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# AI Infrastructure – Configuration Management

Settnings and consistency around AI LLMs

In the context of Large Language Models (LLMs), configuration management refers to the systematic process of establishing and maintaining consistency in the model's performance, functionalities, and attributes throughout its lifecycle.

| Model Versioning | Training Configuration | Deployment Configuration | Integration Configuration | Monitoring and Feedback | Backup and Recovery |
|---|---|---|---|---|---|
| • Tracking Changes<br>• Hyperparameter Management<br>• Storing and Retrieving | • Reproducibbility | • Environment Settings<br>• Scaling and Resource Allocation | • Data Connectors<br>• API Endpoints | • Metrics Tracking<br>• Alert Settings | • Snapshotting<br>• Migration |

# AI Infrastructure – Security

Security and Privacy for LLMs and AI Infrastructure

In the context of Large Language Models (LLMs), security pertains to the measures, practices, and protocols implemented to protect the model, its data, and its outputs from threats and vulnerabilities.

| Data Security | Model Security | Deployment Security | Output Security | Audit and Compliance |
|---|---|---|---|---|
| • Training Data Protection<br>• Data Privacy | • Model Inversion Attacks<br>• Model Poisoning<br>• Adversarial Attacks<br>• Model Copying | • Access Controls<br>• Rate Limiting<br>• Input Validation | • Information Leakage<br>• Output Filtering | • Logging<br>• Regulatory Compliance |

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# Measuring Performance of AI − Metrics

AI specific measurements



- **Accuracy:** The ratio of correctly predicted instances to the total instances in the dataset.

- **Data Drift:** The change in data distribution over time.

- **Model Explainability**: The degree to which a model's predictions can be understood and interpreted.

- **Model Inference Time:** The time it takes for the model to make a prediction once it's trained.

- **Model Inference Time:** The time it takes for the model to make a prediction once it's trained.

# Measuring Performance of AI – Metrics (cont'd)

AI specfic measurements

- **Model Training Time:** The amount of time it takes to train a model.

- **Precision:** The ratio of correctly predicted positive observations to the total predicted positives.

- **Recall (Sensitivity):** The ratio of correctly predicted positive observations to all the actual positives.

- **AL/ML Ops Metrics:** F1-Score, Area Under the Receiver Operating Characteristic Curve (AUC-ROC), Mean Absolute Error (MAE), Root Mean Square Error (RMSE), Log Loss, Confusion Matrix, Model Size, Concept Drift

# Measuring Performance of AI – Metrics (cont'd)

## AI Observability

Once we understand how to measure AI operations, the next step is the tooling to do so. Observability is the ability to understand the internal state of a system from its external outputs (the measurement of AI). In the context of AI-powered software, observability becomes crucial due to the inherent complexity and unpredictability of AI models.
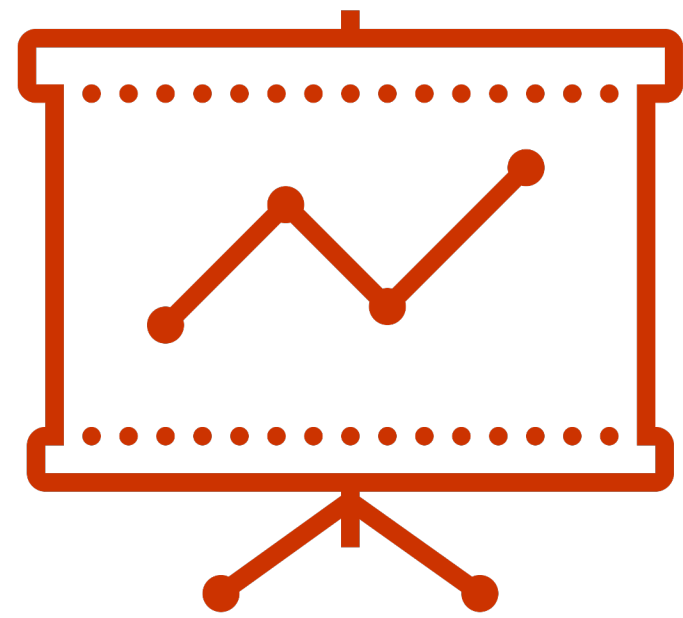
- **Model Monitoring:** AI models can drift over time as they encounter new and unseen data in production

- **System Health:** The entire software ecosystem must be monitored beyond just the AI model.

THE ARTIFICIALLY INTELLIGENT ENTERPRISE

# Measuring Performance of AI – Metrics (cont'd)

AI specfic measurements

## AI Observability

Once we understand how to measure AI operations, the next step is the tooling to do so. Observability is the ability to understand the internal state of a system from its external outputs (the measurement of AI). In the context of AI-powered software, observability becomes crucial due to the inherent complexity and unpredictability of AI models.

- **Model Monitoring:** AI models can drift over time as they encounter new and unseen data in production

- **System Health:** The entire software ecosystem must be monitored beyond just the AI model.

# The Takeaways

A framework for ops to talk about AI

- AI Brings Lots of New Infrastructure

- Stacks are not a good descriptor, webs are better for the cloud era

- Having a framework allows us to have a conversation for existing DevSecOps

- Tool categories are the same, way we apply them are slightly different

- Metrics are changing and will require new or updated tools

# Challenges and Strategic Implementation for ITSM

New technology same old change management challenge

Organizations may face issues such as:

- Integration complexities with existing IT infrastructures,
- The need for significant upfront investment,
- Potential resistance from IT staff due to fears of job displacement.

However, these challenges can be mitigated through strategic planning, continuous training, and phased implementation of AI technologies
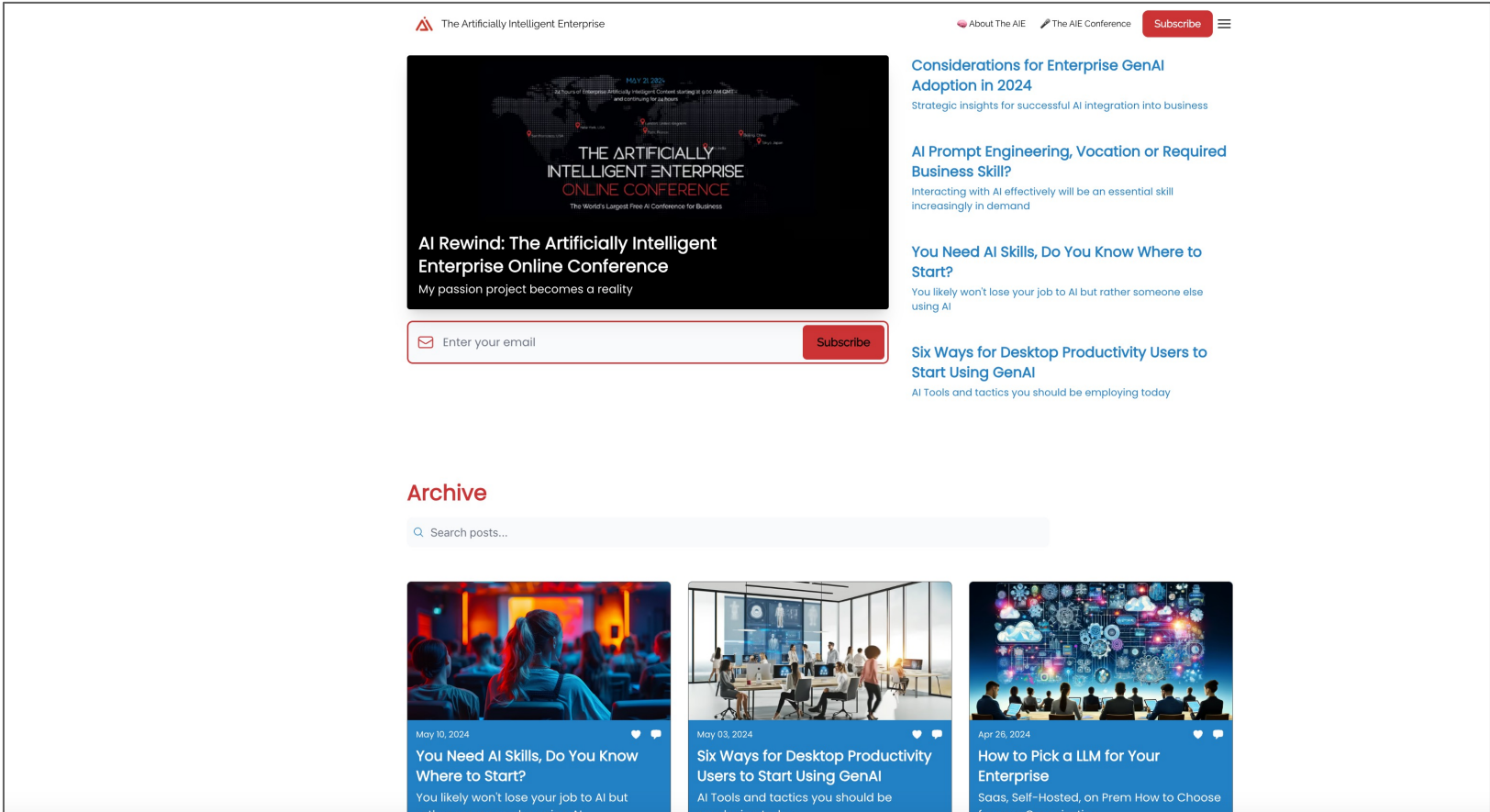
# The GenAI Future

Small Challenges, Big Opportunities

Looking ahead, the role of generative AI in ITSM is expected to grow, with advancements likely to introduce more sophisticated capabilities such as multimodal interactions and deeper integration with enterprise systems.

This evolution will further enhance the agility and efficiency of IT services, making businesses more resilient and adaptive to technological changes

# Get Free Access to GenAI Training and Knowledge

**GenAI Strategy, Tips, and Tricks from Experts**



## The AIE Newsletter
https://theaienterpise.io

Weekly email with GenAI strategy, information, and tips.



## The AIE Conference
https://theaienterprise.online

May 21st, 2024 will be a free conference

# Thank You!