**Figure 5** Today's Private Cloud Solutions Market Offers A Wide Variety Of Solutions

| Vendor | Self-service portal or service catalog | Dynamic workload management | Resource management | Service accounting | Integration and control APIs | Image library | RBAC administration | Virtualization layer | Physical compute and storage included | Application services |
|---|---|---|---|---|---|---|---|---|---|---|
| Abiquo | | | | | | | | | | |
| BMC | | | | | | | | | | |
| CA | | | | | | | | | | |
| Cloud.com | | | | | | | | | | |
| Dell | | | | | | | | | | |
| Enomaly | | | | | | | | | | |
| Eucalyptus | | | | | | | | | | |
| Hexagrid | | | | | | | | | | |
| HP | | | | | | | | | | |
| IBM | | | | | | | | | | |
| Microsoft | | | | | | | | | | |
| newScale | | | | | | | | | | |
| Platform Computing | | | | | | | | | | |
| Tibco | | | | | | | | | | |
| VMware | | | | | | | | | | |

Legend: ○ 0   ◔ 1   ◑ 2   ◕ 3   ● 4

Note: Please refer to Figure 4 for the scoring criteria.

58924

FORRESTER® Making Leaders Successful Every Day
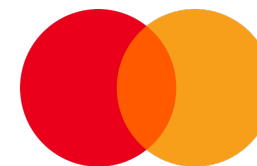
**Suresh Mandava**
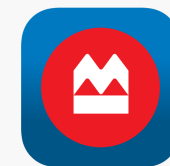SVP/Chief Architect
Cloud-Native AI/ML Platforms and Security
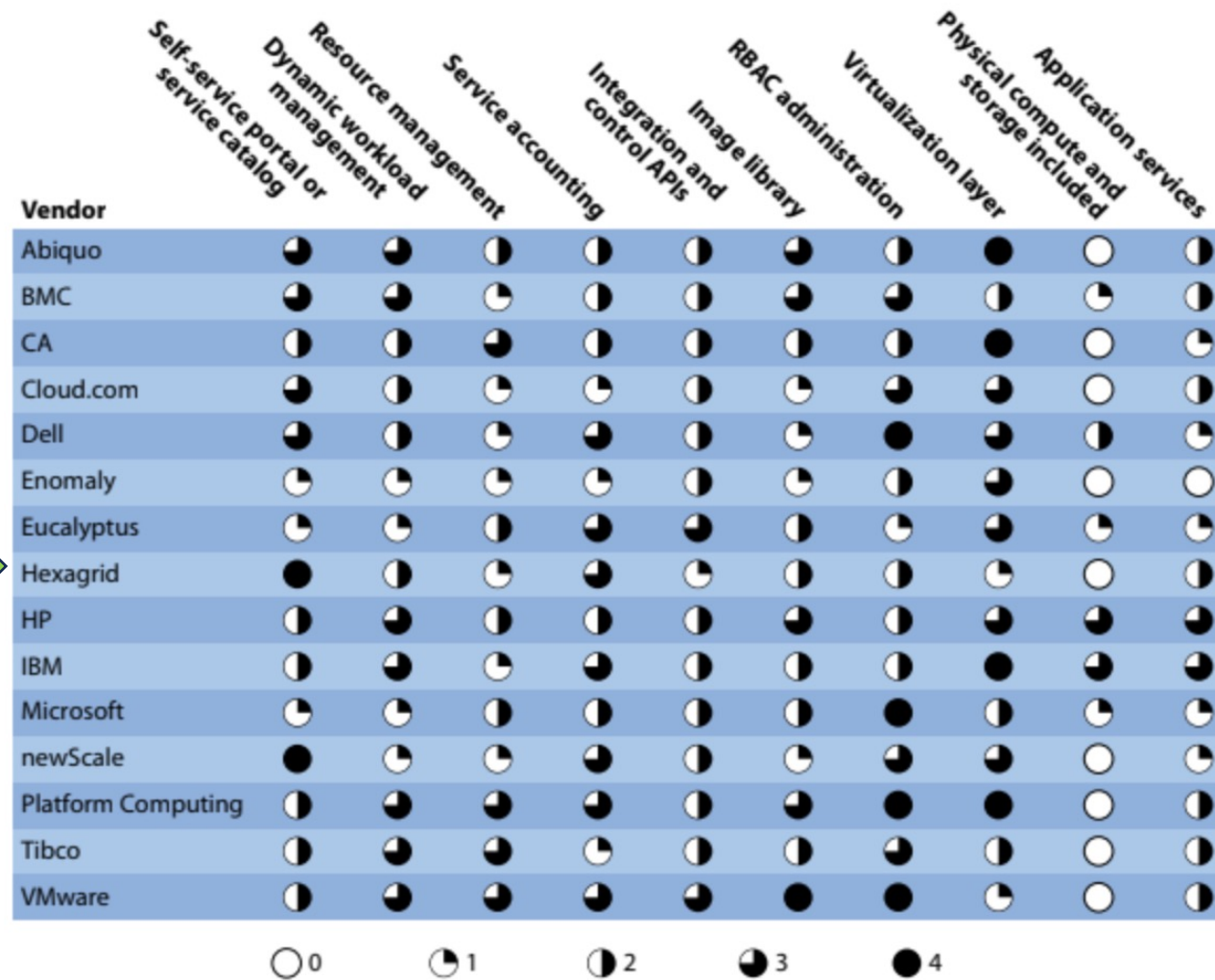Infinite Computer Solutions

**Founder  (2007-2012)**
HexaGrid Computing

mastercard

DXC TECHNOLOGY

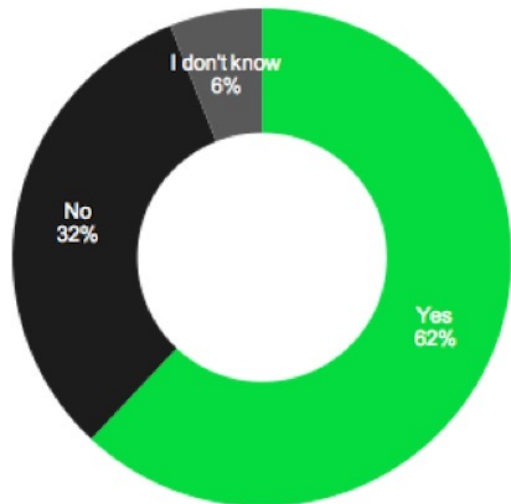**SiFive Rolls Out RISC-V Cores Aimed at Generative AI and ML**

**RISC-V**

**Did you Know ?** In order to play the role of an insane and mentally depressed person the movie "Joker", Joaquin Phoenix becomes a **full stack developer for a month.**

**Black Hat: AI As An Attack Method**
**AUG 1, 2017**



I don't know
6%

No
32%

Yes
62%

https://www.isssource.com/black-hat-ai-as-an-attack-method/



INTERVIEWS | JUNE 21, 2021
AI-POWERED CYBER ATTACKS EMERGING AS MAJOR CONCERN

**WormGPT: New AI Tool Allows Cybercriminals to Launch Sophisticated Cyber Attacks**
Jul 15, 2023

**PROMPT INJECTION: AN AI-TARGETED ATTACK**

# Samsung Engineers Feed Sensitive Data to ChatGPT, Sparking Workplace AI Warnings

In three separate incidents, engineers at the Korean electronics giant reportedly shared sensitive corporate data with the AI-powered chatbot.

**Jai Vijayan**

Contributing Writer, Dark Reading

April 11, 2023

**[The Economist Korea](), one of the first to report on the data leaks, described the first incident as involving an engineer who pasted buggy source code from a semiconductor database into ChatGPT, with a prompt to the chatbot to fix the errors.**

**In the second instance, an employee wanting to optimize code for identifying defects in certain Samsung equipment pasted that code into ChatGPT.**

**The third leak resulted when an employee asked ChatGPT to generate the minutes of an internal meeting at Samsung.**

Nightmare continues

... Wait until somebody loaded a 3-party GENAI evil tool against your GITHUB

# AI Trust, Risk and Security Management

## (AI TRiSM)

By 2026, organizations that operationalize AI transparency, trust and security will see their AI models achieve a 50% result improvement in terms of adoption, business goals and user acceptance.
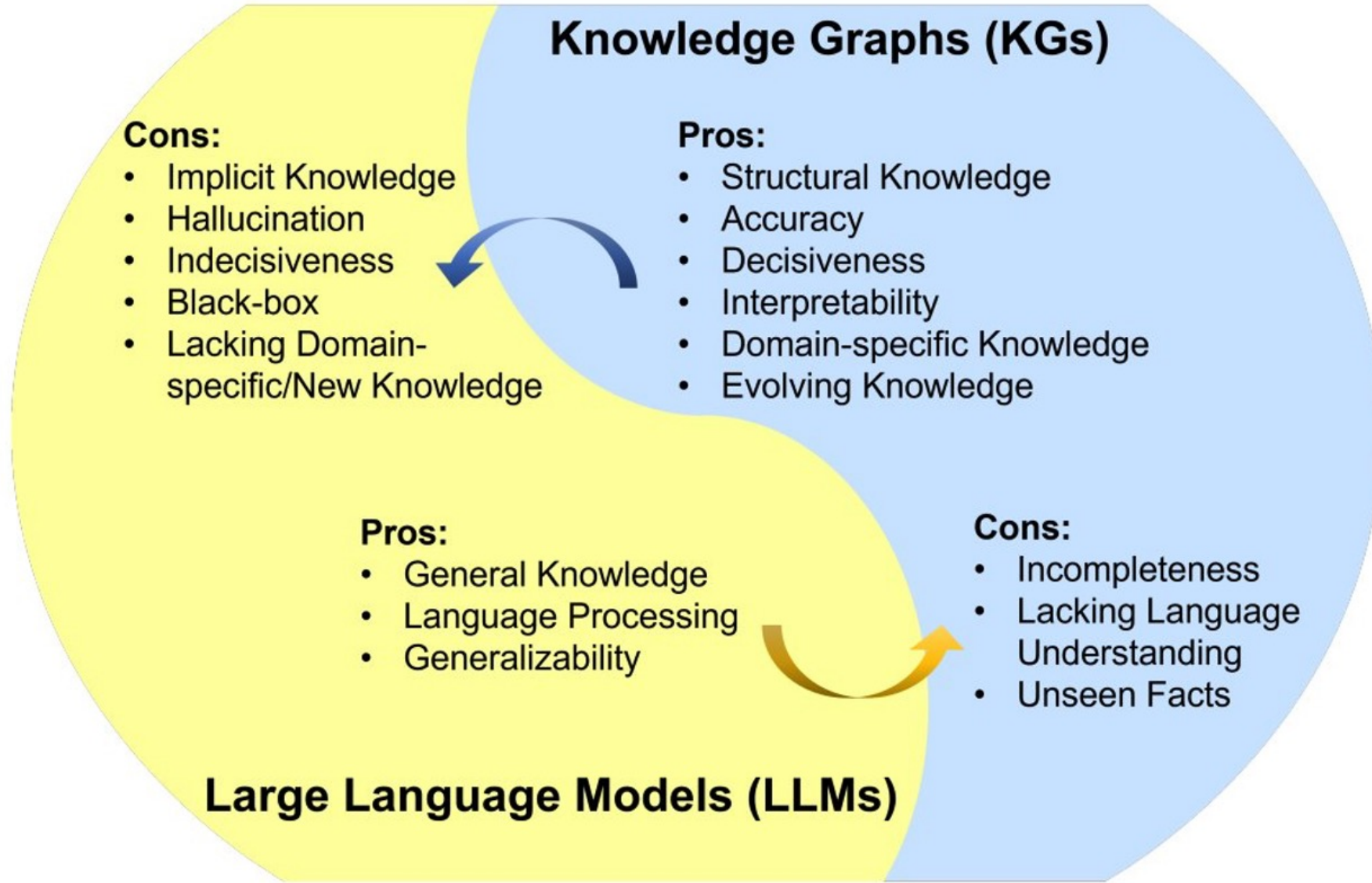
## AI TRiSM: Optimize Trust in AI
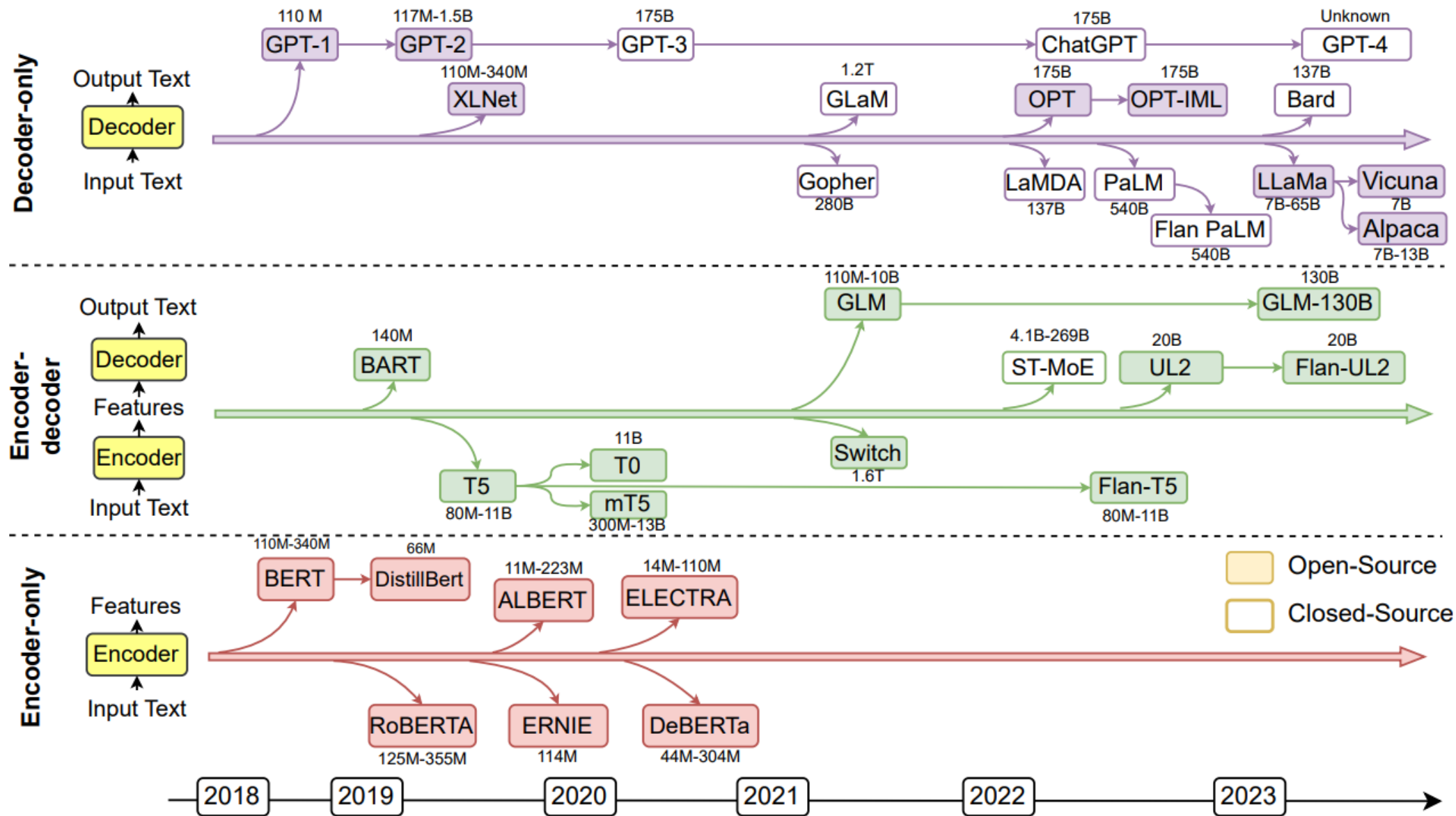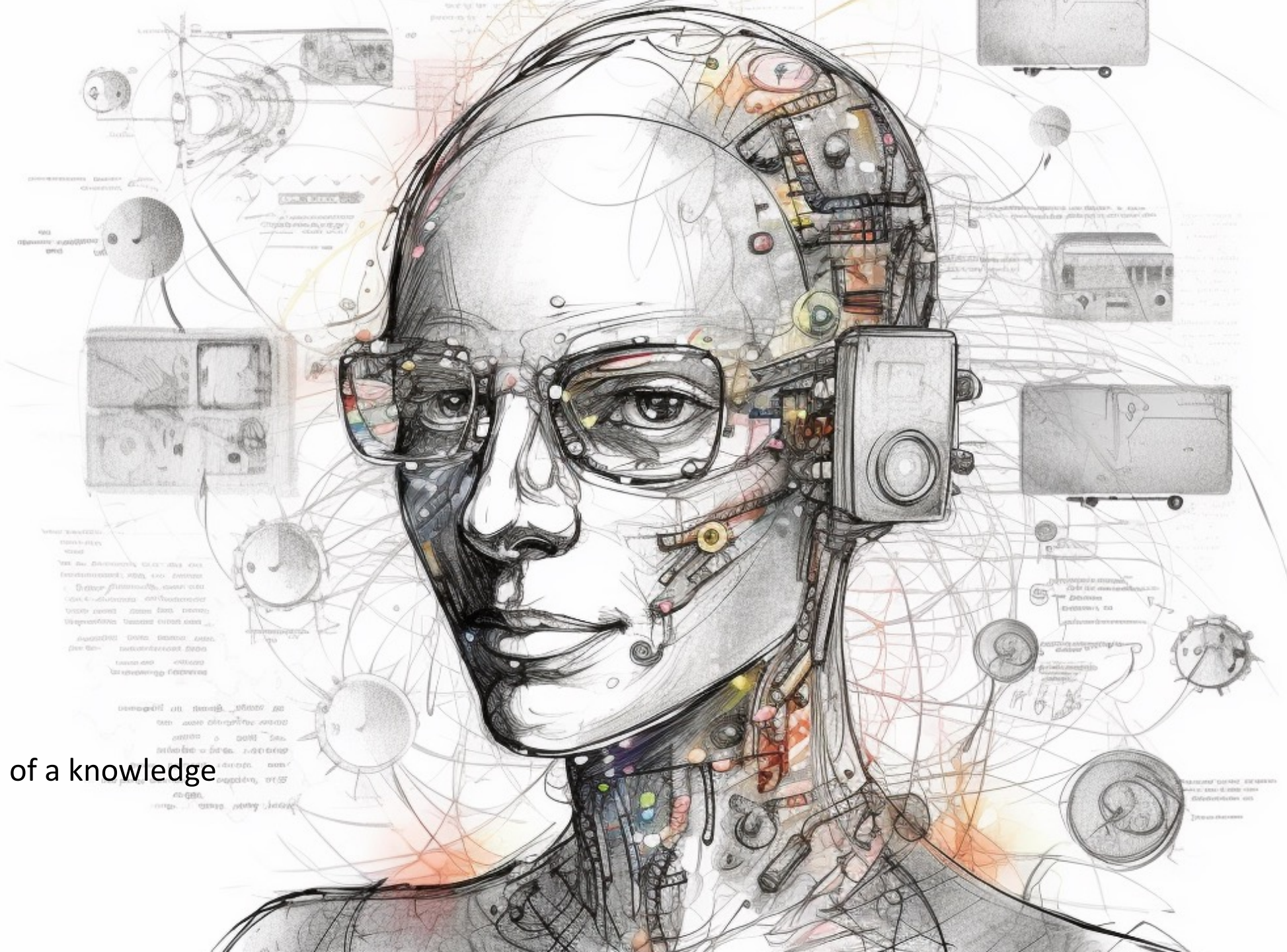
Four pillars of AI trust, risk and security management



Unmanaged Risks → AI TRiSM → Managed Risks

| Explainability/Model Monitoring | ModelOps | AI Application Security | Privacy |

# Evolution of AI Architecture: Traditional ML to Generative AI

## Traditional ML

**Data Pre-Processing**
Cleaning and preparing data for analysis.

**Feature Engineering**
Extracting important features from data.

**Training & Tuning**
Training models on data and adjusting parameters for optimal performance.

**Deployment & Monitoring**
Implementing models in real-world applications and monitoring their performance.

### Tech Stack for Traditional ML

- **ML Frameworks:** Keras, Theano
- **ML API's & SDK:** IBM Watson
- **Database:** SQL Server, Oracle
- **ML Ops:** Docker, Jenkins

## Generative AI

**Data Pre-Processing**
Cleaning and preparing data for analysis.

**Prompt Engineering/Fine Tuning**
Designing effective prompts to guide AI in generating desired outputs.

**Foundational/Fine Tuned LLM**
Using foundational and fine-tuned language learning models for sophisticated content generation.

**Deployment & Monitoring**
Implementing models in real-world applications and monitoring their performance.

### Tech Stack for Generative AI

- **Gen AI Orchestration:** Langchain, llamaindex
- **LLM Models:** OpenAI, Anthropic
- **Vector Database:** Pinecone, Weaviate
- **LLM Ops:** Prompt Layer, Helicone

**Knowledge Graphs (KGs)**

**Cons:**
- Implicit Knowledge
- Hallucination
- Indecisiveness
- Black-box
- Lacking Domain-specific/New Knowledge

**Pros:**
- Structural Knowledge
- Accuracy
- Decisiveness
- Interpretability
- Domain-specific Knowledge
- Evolving Knowledge

**Pros:**
- General Knowledge
- Language Processing
- Generalizability

**Cons:**
- Incompleteness
- Lacking Language Understanding
- Unseen Facts

**Large Language Models (LLMs)**

**Decoder-only**

110 M — GPT-1
117M-1.5B — GPT-2
175B — GPT-3
175B — ChatGPT
Unknown — GPT-4
110M-340M — XLNet
1.2T — GLaM
175B — OPT
175B — OPT-IML
137B — Bard
Output Text — Decoder — Input Text
280B — Gopher
137B — LaMDA
540B — PaLM
540B — Flan PaLM
7B-65B — LLaMa
7B — Vicuna
7B-13B — Alpaca

**Encoder-decoder**

110M-10B — GLM
130B — GLM-130B
140M — BART
4.1B-269B — ST-MoE
20B — UL2
20B — Flan-UL2
Output Text — Decoder — Features — Encoder — Input Text
11B — T0
1.6T — Switch
T5 — 80M-11B
mT5 — 300M-13B
Flan-T5 — 80M-11B

**Encoder-only**

110M-340M — BERT
66M — DistillBert
11M-223M — ALBERT
14M-110M — ELECTRA
Features — Encoder — Input Text
Open-Source
Closed-Source
125M-355M — RoBERTA
114M — ERNIE
44M-304M — DeBERTa

2018  2019  2020  2021  2022  2023

Midjourney's idea of a knowledge graph chatbot.

https://github.com/Stability-AI/generative-models

## Dataset
# RealToxicity

GPT                                          0.233

Supervised Fine-Tuning                       0.199

InstructGPT                                **0.196**

## Dataset
# TruthfulQA

GPT                                          0.224

Supervised Fine-Tuning                       0.206

InstructGPT                                **0.413**

Hallucination is worse for InstructGPT
(RLHF + SFT) compared to just SFT
(Ouyang et al., 2022)

## API Dataset
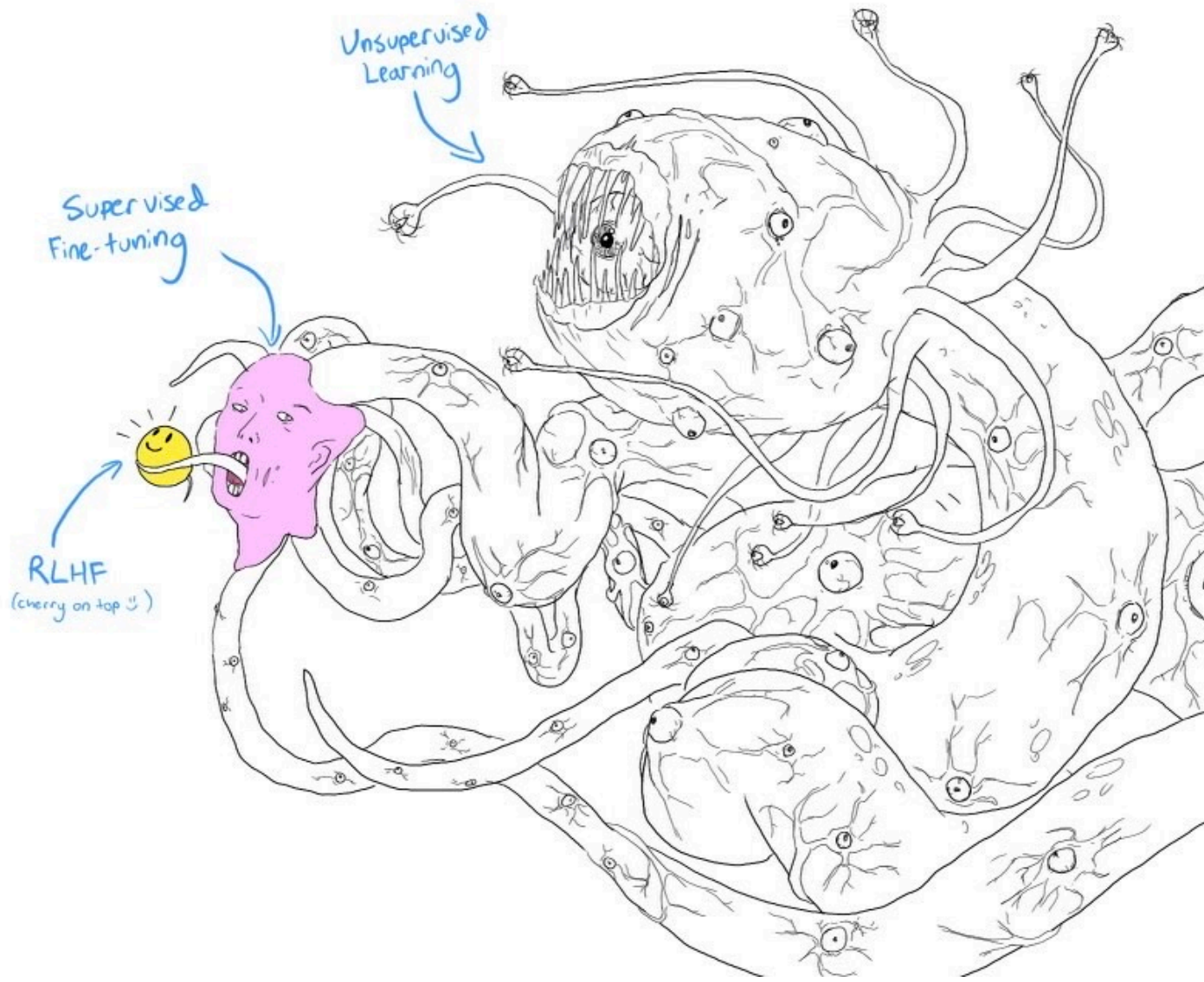# Hallucinations

GPT                                          0.414

Supervised Fine-Tuning                     **0.078**

InstructGPT                                  0.172

## API Dataset
# Customer Assistant Appropriate

GPT                                          0.811

Supervised Fine-Tuning                       0.880

InstructGPT                                **0.902**

Evaluating InstructGPT for toxicity, truthfulness, and appropriateness. Lower scores are
better for toxicity and hallucinations, and higher scores are better for TruthfulQA and
appropriateness. Hallucinations and appropriateness are measured on our API prompt
distribution. Results are combined across model sizes.

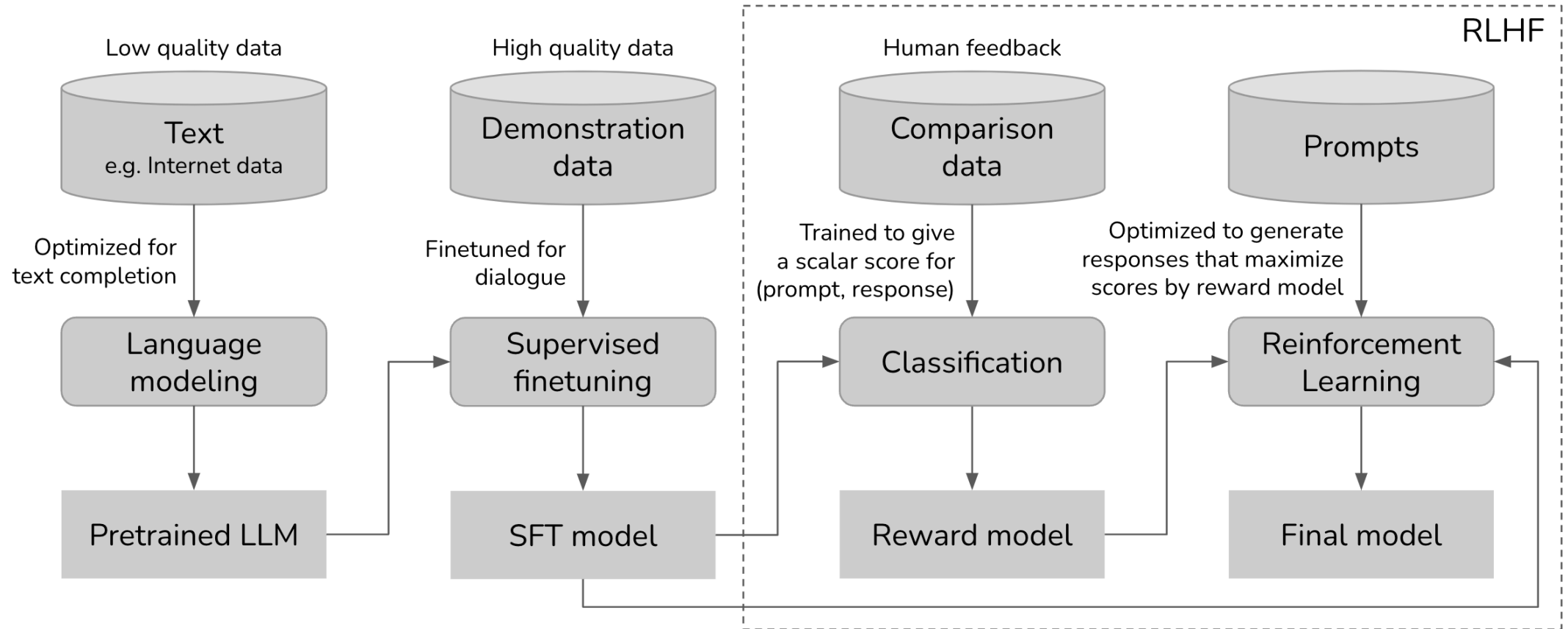| | RedPajama | LLaMA* |
|---|---|---|
| CommonCrawl | 878 billion | 852 billion |
| C4 | 175 billion | 190 billion |
| Github | 59 billion | 100 billion |
| Books | 26 billion | 25 billion |
| ArXiv | 28 billion | 33 billion |
| Wikipedia | 24 billion | 25 billion |
| StackExchange | 20 billion | 27 billion |
| Total | 1.2 trillion | 1.25 trillion |

**Out of Domain**

Toxicity  Bias

**Third Party Applications (Agents)**

# Unsupervised, Supervised Fine Tuning and Reinforcement Learning from Human Feedback

RLHF

Low quality data

**Text**
e.g. Internet data

High quality data

**Demonstration data**

Human feedback

**Comparison data**

**Prompts**

Optimized for text completion

**Language modeling**

Finetuned for dialogue

**Supervised finetuning**

Trained to give a scalar score for (prompt, response)

**Classification**

Optimized to generate responses that maximize scores by reward model

**Reinforcement Learning**

**Pretrained LLM**

**SFT model**

**Reward model**

**Final model**

Scale
May '23

**>1 trillion**
tokens

**10K - 100K**
(prompt, response)

**100K - 1M** comparisons
(prompt, winning_response, losing_response)

**10K - 100K**
prompts

Examples
**Bolded**: open sourced

GPT-x, Gopher, **Falcon**,
LLaMa, **Pythia**, **Bloom**,
**StableLM**

**Dolly-v2, Falcon-Instruct**

InstructGPT, ChatGPT,
Claude, **StableVicuna**

# State of Gen AI 2023

## Jobs in U.S. that are likely to have high, medium or low exposure to AI

**High exposure**
- Budget analysts
- Data entry keyers
- Tax preparers
- Technical writers
- Web developers

**Medium exposure**
- Chief executives
- Veterinarians
- Interior designers
- Fundraisers
- Sales managers

**Low exposure**
- Barbers
- Child care workers
- Dishwashers
- Firefighters
- Pipelayers

Note: Occupations are grouped by the relative importance of work activities with low, medium or high exposure to AI.
Source: Pew Research Center analysis of O*NET (Version 27.3).
"Which U.S. Workers Are More Exposed to AI on Their Jobs?"

**PEW RESEARCH CENTER**

## What shares of workers are most exposed to AI in their jobs?

*% of U.S. workers employed in jobs that are the most exposed to AI in 2022*

| | |
|---|---|
| All workers | 19% |
| Men | 17 |
| Women | 21 |
| White | 20 |
| Black | 15 |
| Hispanic | 13 |
| Asian | 24 |
| Amer. Indian or Pacific Islander | 16 |
| Other | 18 |
| Less than HS | 3 |
| HS grad | 12 |
| Some college | 19 |
| Bachelor's+ | 27 |

Note: Occupations are ranked by the relative importance of work activities with high exposure to AI. Those in the top 25% are the "most exposed," some 122 in number. Estimates by education level are for workers ages 25 and older. White, Black, Asian, and American Indian or Pacific Islander workers include those who report being only one race and are not Hispanic. "Other" includes all other single race groups and people reporting two or more races. Hispanics are of any race.
Source: Pew Research Center analysis of O*NET (Version 27.3) and 2022 Current Population Survey (IPUMS) annual data.
"Which U.S. Workers Are More Exposed to AI on Their Jobs?"

**PEW RESEARCH CENTER**

Pew Research Center

## Which U.S. Workers Are More Exposed to AI on Their Jobs?

*About a fifth of all workers have high-exposure jobs; women, Asian, college-educated and higher-paid workers are more exposed. But those in the most exposed industries are more likely to say AI will help more than hurt them personally*

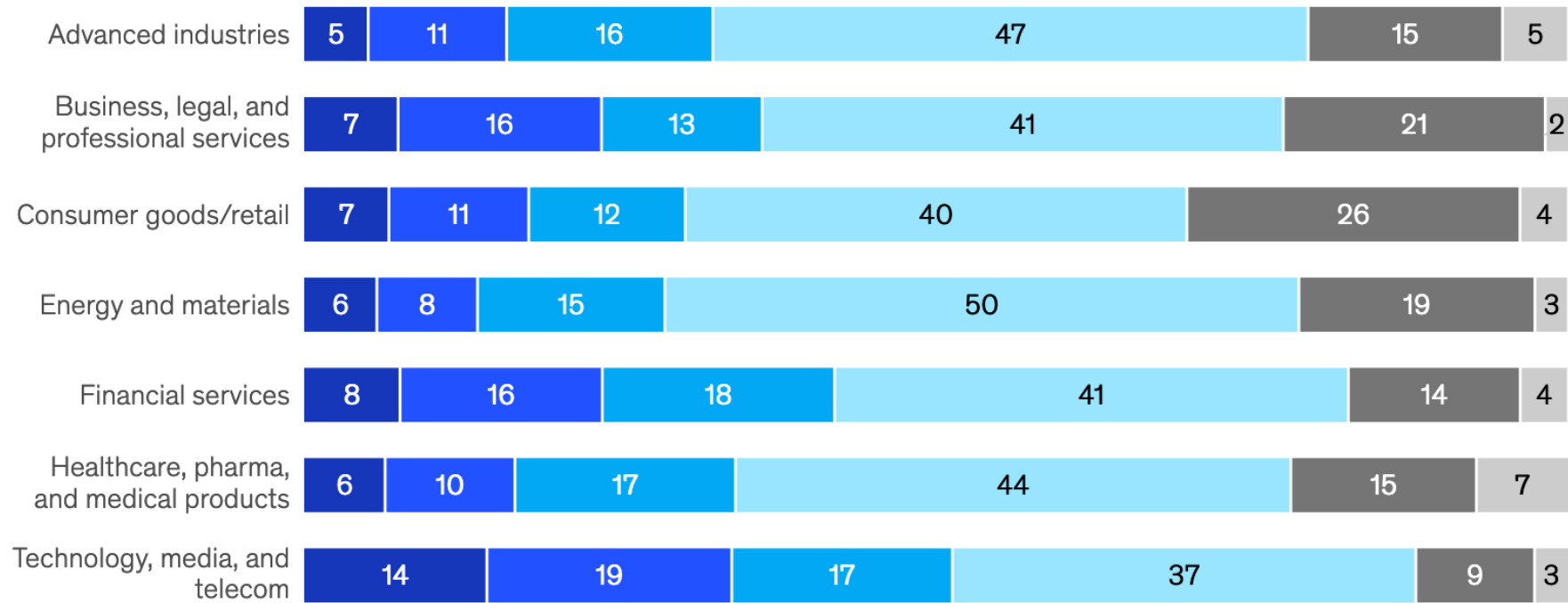## A.I. is on a collision course with white-collar, high-paid jobs — and with unknown impact

# Respondents across regions, industries, and seniority levels say they are already using generative AI tools.

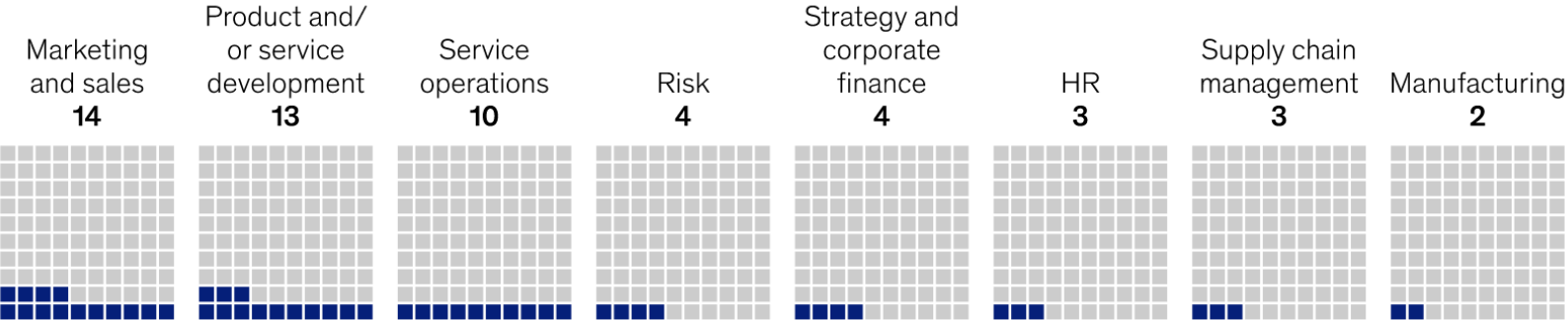Reported exposure to generative AI tools, % of respondents

Select demographic | By industry ▼ |

- ■ Regularly use for work
- ■ Regularly use for work and outside of work
- ■ Regularly use outside of work
- ■ Have tried at least once
- ■ No exposure
- ■ Don't know

| Industry | Regularly use for work | Regularly use for work and outside of work | Regularly use outside of work | Have tried at least once | No exposure | Don't know |
|---|---|---|---|---|---|---|
| Advanced industries | 5 | 11 | 16 | 47 | 15 | 5 |
| Business, legal, and professional services | 7 | 16 | 13 | 41 | 21 | 2 |
| Consumer goods/retail | 7 | 11 | 12 | 40 | 26 | 4 |
| Energy and materials | 6 | 8 | 15 | 50 | 19 | 3 |
| Financial services | 8 | 16 | 18 | 41 | 14 | 4 |
| Healthcare, pharma, and medical products | 6 | 10 | 17 | 44 | 15 | 7 |
| Technology, media, and telecom | 14 | 19 | 17 | 37 | 9 | 3 |

# The most commonly reported uses of generative AI tools are in marketing and sales, product and service development, and service operations.

**Share of respondents reporting that their organization is regularly using generative AI in given function,** %[1]

| Marketing and sales | Product and/ or service development | Service operations | Risk | Strategy and corporate finance | HR | Supply chain management | Manufacturing |
|---|---|---|---|---|---|---|---|
| **14** | **13** | **10** | **4** | **4** | **3** | **3** | **2** |

**Most regularly reported generative AI use cases within function,** % of respondents

| Marketing and sales | Product and/or service development | Service operations |
|---|---|---|
| Crafting first drafts of text documents | Identifying trends in customer needs | Use of chatbots (eg, for customer service) |
| 9 | 7 | 6 |
| Personalized marketing | Drafting technical documents | Forecasting service trends or anomalies |
| 8 | 5 | 5 |
| Summarizing text documents | Creating new product designs | Creating first drafts of documents |
| 8 | 4 | 5 |

# Inaccuracy, cybersecurity, and intellectual-property infringement are the most-cited risks of generative AI adoption.

**Generative AI–related risks that organizations consider relevant and are working to mitigate,**
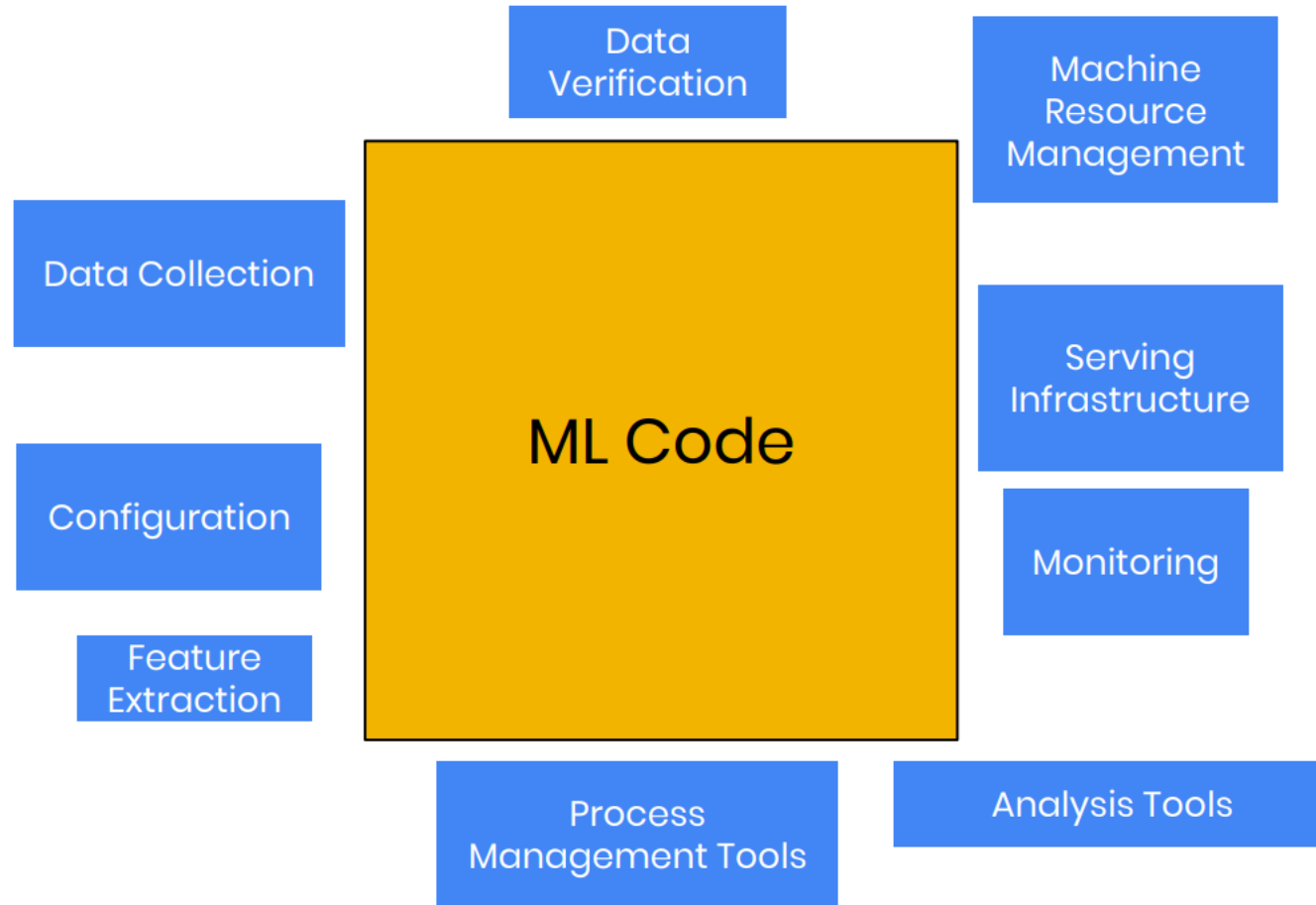% of respondents[1]

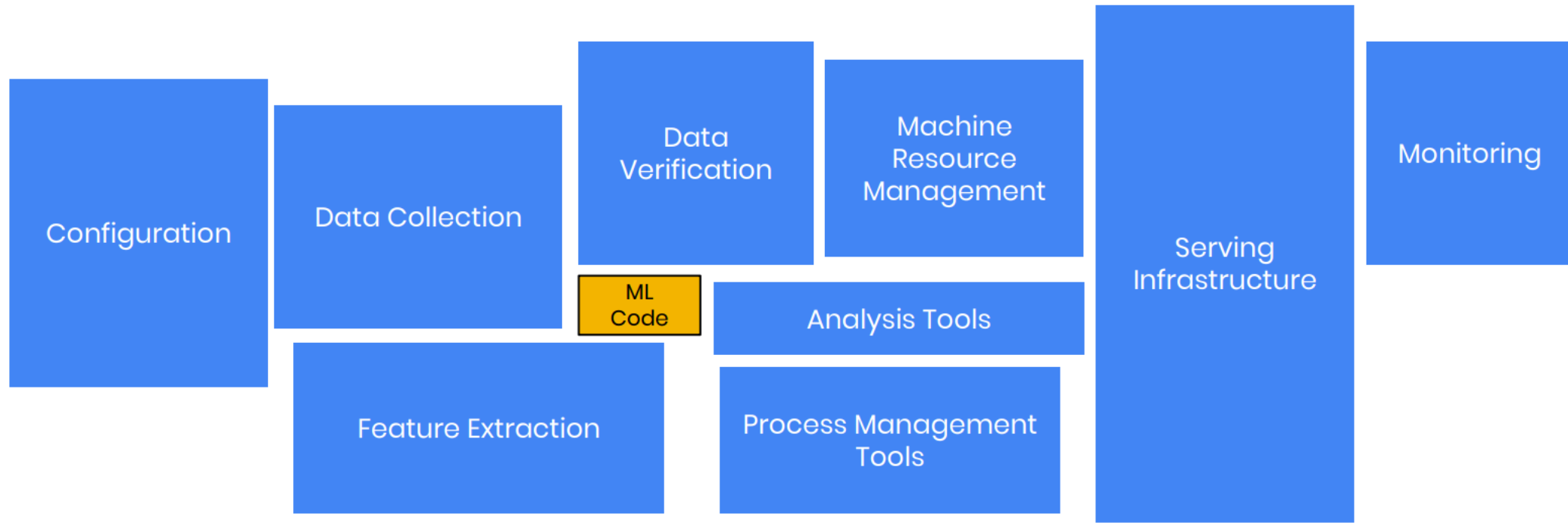| | Organization considers risk relevant | Organization working to mitigate risk |
|---|---|---|
| Inaccuracy | 56 | 32 |
| Cybersecurity | 53 | 38 |
| Intellectual-property infringement | 46 | 25 |
| Regulatory compliance | 45 | 28 |
| Explainability | 39 | 18 |
| Personal/individual privacy | 39 | 20 |
| Workforce/labor displacement | 34 | 13 |
| Equity and fairness | 31 | 16 |
| Organizational reputation | 29 | 16 |
| National security | 14 | 4 |
| Physical safety | 11 | 6 |
| Environmental impact | 11 | 5 |
| Political stability | 10 | 2 |
| None of the above | 1 | 8 |

McKinsey & Company

AI 101

# Perception : ML Products are mostly about ML



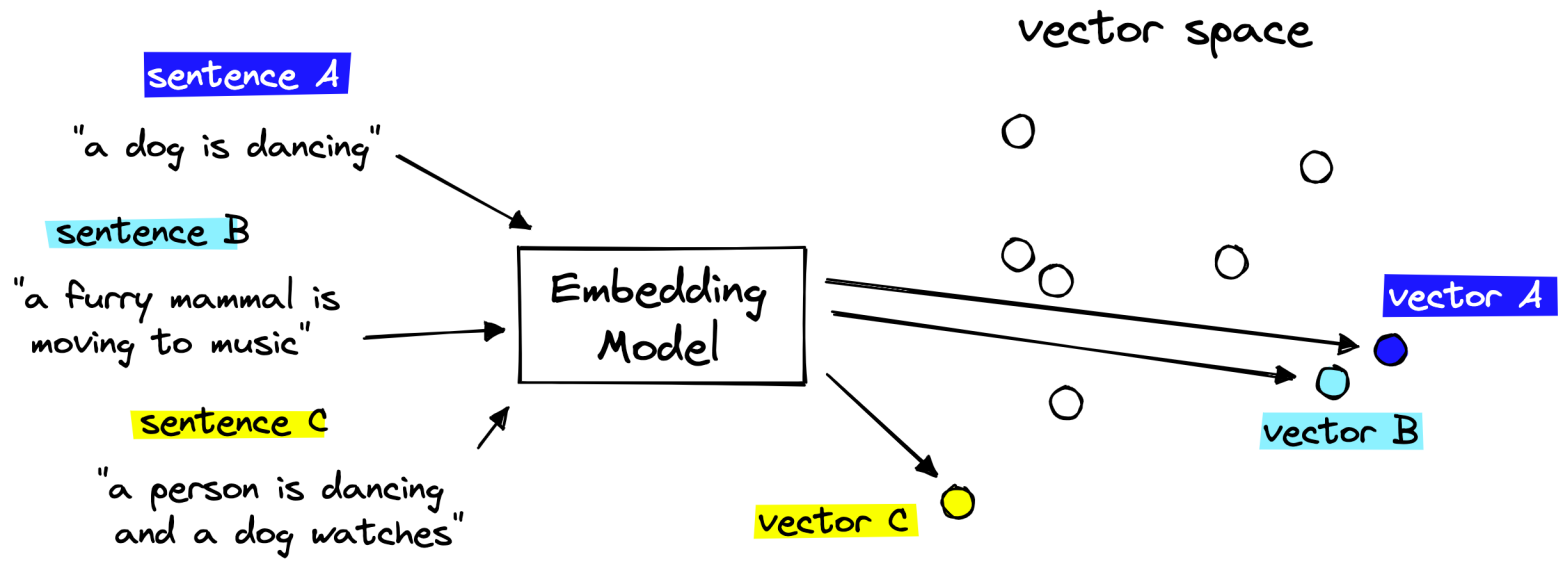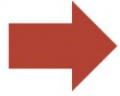Credit: Hidden Technical Debt of Machine Learning Systems, D. Sculley, et al.

Configuration

Data Collection

Data Verification

Machine Resource Management

Monitoring

ML Code

Analysis Tools

Serving Infrastructure

Feature Extraction

Process Management Tools

vector space

sentence A

"a dog is dancing"

sentence B

"a furry mammal is moving to music"

sentence C

"a person is dancing and a dog watches"

Embedding Model

vector A

vector B

vector C

dogs running across field

two dogs playing

big wave breaking in sea

a person surfing

person surfing inside barrel of wave

close up photo of a dog

**Vectors 101**

Vocabulary:
Man, woman, boy, girl, prince, princess, queen, king, monarch

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| man | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| woman | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| boy | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| girl | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| prince | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| princess | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| queen | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| king | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| monarch | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

Each word gets a 1x9 vector representation

Try to build a lower dimensional embedding

Vocabulary:
Man, woman, boy, girl, prince, princess, queen, king, monarch

| | Femininity | Youth | Royalty |
|---|---|---|---|
| Man | 0 | 0 | 0 |
| Woman | 1 | 0 | 0 |
| Boy | 0 | 1 | 0 |
| Girl | 1 | 1 | 0 |
| Prince | 0 | 1 | 1 |
| Princess | 1 | 1 | 1 |
| Queen | 1 | 0 | 1 |
| King | 0 | 0 | 1 |
| Monarch | 0.5 | 0.5 | 1 |

Each word gets a 1x3 vector

Similar words... similar vectors

@shane_a_lynn | @TeamEdgeTier

| | living being | feline | human | gender | royalty | verb | plural |
|---|---|---|---|---|---|---|---|
| cat → | 0.6 | 0.9 | 0.1 | 0.4 | −0.7 | −0.3 | −0.2 |
| kitten → | 0.5 | 0.8 | −0.1 | 0.2 | −0.6 | −0.5 | −0.1 |
| dog → | 0.7 | −0.1 | 0.4 | 0.3 | −0.4 | −0.1 | −0.3 |
| houses → | −0.8 | −0.4 | −0.5 | 0.1 | −0.9 | 0.3 | 0.8 |

Dimensionality reduction of word embeddings from 7D to 2D →

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| man → | 0.6 | −0.2 | 0.8 | 0.9 | −0.1 | −0.9 | −0.7 |
| woman → | 0.7 | 0.3 | 0.9 | −0.7 | 0.1 | −0.5 | −0.4 |
| king → | 0.5 | −0.4 | 0.7 | 0.8 | 0.9 | −0.7 | −0.6 |
| queen → | 0.8 | −0.1 | 0.8 | −0.9 | 0.8 | −0.5 | −0.9 |

Dimensionality reduction of word embeddings from 7D to 2D →

Word — Word embedding — Dimensionality reduction — Visualization of word embeddings in 2D

# Vectors 101

| Object | Vector | Task |
|--------|--------|------|
| IMAGE → IMAGE TRANSFORMER | [1.3, 0.6, 1.2, -1.3, ...] | Object recognition, deduplication, scene detection, product search, ... |
| Happy Birthday TEXT → NLP TRANSFORMER | [0.3, -0.4, 1.2, 0.3, ...] | Translation, understanding, Sentiment, Question Answering, Semantic Search, ... |
| AUDIO → AUDIO TRANSFORMER | [1.2, -0.3, 0.7, -1.8, ...] | Anomaly detection, speech-to-text, music transcription, machinery malfunction, ... |

Audio → Audio model → Audio Vector Embeddings

Text → Texts model → Text Vector Embeddings

Videos → Videos model → Video Vector Embeddings

RedisSearch Vector Similarity Search

# Data Pre-Processing Pipeline

| RAW Corpus | Quality Filtering | De-Duplication | Privacy Reduction | Tokenization | Ready to Pre-Train |
|---|---|---|---|---|---|
| | • Language Filtering<br>• Metric Filtering<br>• Statistic Filtering<br>• Keyword Filtering | • Sentence Level<br>• Document Level<br>• Set Level | • Detect Personal Identifiable Information (PII)<br>• Remove PII | • Reuse Existing Tokenizer<br>• Sentence Piece<br>• Byte-Level BPE | |
| | Alice is writing a paper about LLMs $@% Alice is writing a paper about LLMs | Alice is writing a paper about LLMs Alice is writing a paper about LLMs | Replace : [Alice] is writing a paper about LLMs | Encode : [Somebody] is writing a paper about LLMs | 32, 145, 66, 79, 12, 56 .. |

**Encoder Models**

**Instructor Embeddings**

**Llama Embeddings**

**Word2VEC**

**OpenAI :** text-embedding-ada-002 model

# Key Components for Building RAG based applications:



Source : Lance Blog

**On-prem (self-hosted)**

In-memory only, single machine

Weaviate   Qdrant   chroma

LanceDB

redis   vespa

Weaviate   pgvector

Qdrant

milvus   elasticsearch

**Embedded (Serverless)** ← → **Client-server**

zilliz   pgvector
Through various third-parties
chroma   Qdrant

LanceDB

zilliz

Pinecone   Weaviate

Vald   vespa

redis   elasticsearch

**Cloud-native (managed)**

---

Pinecone ·············· Proprietary composite index

milvus / zilliz ········ Flat, Annoy, IVF, HNSW/RHNSW (Flat/PQ), DiskANN

Weaviate ·············· Customized HNSW, HNSW (PQ), DiskANN (in progress...)

Qdrant ·············· Customized HNSW

chroma ·············· HNSW

LanceDB ·············· IVF (PQ), DiskANN (in progress...)

vespa ·············· HNSW + BM25 hybrid

Vald ·············· NGT

elasticsearch ········ Flat (brute force), HNSW

redis ·············· Flat (brute force), HNSW

pgvector ············ IVF (Flat), IVF (PQ) in progress...

# Use CASE : AI Intelligent OPS for Pharma Distribution



SalesForce (CRM)
ServiceNow (ITIL)

Retrieval Argument Generation (RAG)

Deep Learning

L1/L2 Operations

Partner-Cloud Interactions

Customer Interactions

Production Monitoring

Live Data Asset Discovery/Sync

Asset CMDB

Content Curation

LangChain

Streamlit    jupyter

Notebook    API Gateway    Chat Interfaces

Clinical Data Ingestion

Kafka

RAY

Distributed ML Training

LLAMA Index

Embeddings

Chroma

Microservices    Bots    WASM    NLU

mongoDB    MySQL

Dropbox    ORACLE

Clinical Trails Data Assets

Block Chain

GPTCache

Models

API    Local    Local

Hugging Face    OpenAI    Distribution    Pharma

Destination Traceability

Apache Beam

Model Inference

Distribution Monitoring

OLAP Database    Vector Database    Cache Store    Relational    Model Repo

Milvus    redis

kubernetes  Edge

kubernetes  Data lake

Cloud-Native App Protection (CNAP)

Security Policy Enforcement

End-to-End Zero Trust Network

Encryption Model/Privacy Protection

Security Operations Integration

Governance

Cloud-Security Posture Management (CSPM)

AccuKnox

@Copyright Opensource

# Continuum to Unlock Digital Innovation

# MLOPS Eco-System

**MLOPS**

## BUSINESS GOALS  `BS` `AI` `SE`

- Problem
- KPIs
- Compliance
- Hypothesis

## Data  `AI` `SE`

- Executor
- DocumentArray
- Run

## PROTOTYPING  `AI`

- Experiment ideas
- Visualization
- Results

## CLOUD-NATIVE  `SE` `AI`

- Network protocols
- Parallelization/async pipeline
- Hub Executor reusing
- Microservice
- Containerization
- K8s/Docker Orchestration

## BUSINESS APP  `SE`

- NoCode
- Template

## DEPLOYMENT  `SE`

- Auto provisioning
- Lifecycle management
- Auto scaling

## FINETUNING  `AI` `SE`

- Pretrained model
- Training Data
- Finetuned model
- Hub Executor release

## OBSERVABILITY

- Business Performance Dashboards  `BS` `AI` `SE`
- Model Performance Dashboards  `AI` `BS` `SE`
- System Operations Dashboards  `SE` `AI` `BS`
- Alerts  `SE`

## LEGEND

- `AI`  AI engineer
- `BS`  Business stakeholder
- `SE`  Software engineer
- (filled) Primary concern
- (outline) Secondary concern
- → Business goal flow
- → Data flow
- → Development flow
- -- Implementation variant

**Business Stakeholders**

MLOps isn't a platform- it's an ecosystem of tools. CML helps you bring your favorite DevOps tools to machine learning.

- Continuous integration for ML : CML
- Manage environments : Kubernetes
- Infrastructure as code : Terraform and Docker
- Data as Code : DVC

Fusion.AI

Raw Data → Initial ML Model → Cleanlab AI → **Better Data** → **Better ML Model**

Obtain data → Label data → Fit baseline ML model → Characterize data quality → Curate and improve data → Identify & train the best model → Deploy reliable ML or Analytics → Provide value to customer

Automated by Cleanlab Studio

**Practicing data-centric AI can look like this:**

1. Train initial ML model on original dataset.
2. Utilize this model to diagnose data issues (via cleanlab methods) and improve the dataset.
3. Train the same model on the improved dataset.
4. Try various modeling techniques to further improve performance.

Most folks jump from Step 1 → 4, but you may achieve big gains without *any* change to your modeling code by using cleanlab! Continuously boost performance by iterating Steps 2 → 4 (and try to evaluate with *cleaned* data).

Obtain data and labels

(often lower quality than you want)

the gap today

Low quality data → High quality model

Cleanlab Studio

Deploy ML or analytics

(that work reliably for customers)

Lots of labeling, ETL, and data warehouse solutions exist

databricks  snowflake  scale  Google Cloud  Azure

Lots of parameter tuning and deployment solutions exist

Weights & Biases  GRID·AI  aws

blog        Source : Cleanlab

# Digital Wealth Management

Wealth managers are transitioning to digitally enabled, scalable platforms to empower clients and advisers with compelling experiences

**Investor Clients**

**Advisor Coaches**

Integrated client and advisor experiences

**Fusion.AI**

| Sentient, Intelligent and Engaging | Human Trusted and Transparent | Highly Automated Modern and Frictionless |
|---|---|---|

**Notification services**
*Computer-generated, personalized, real time*

Email    In-app    SMS    Social

**Social/chat engine**
- With peers/friends/advisers

**Client relationship management**
- Prospects, sales, servicing
- Client/sales analytics

**Client on-boarding/account opening utility**
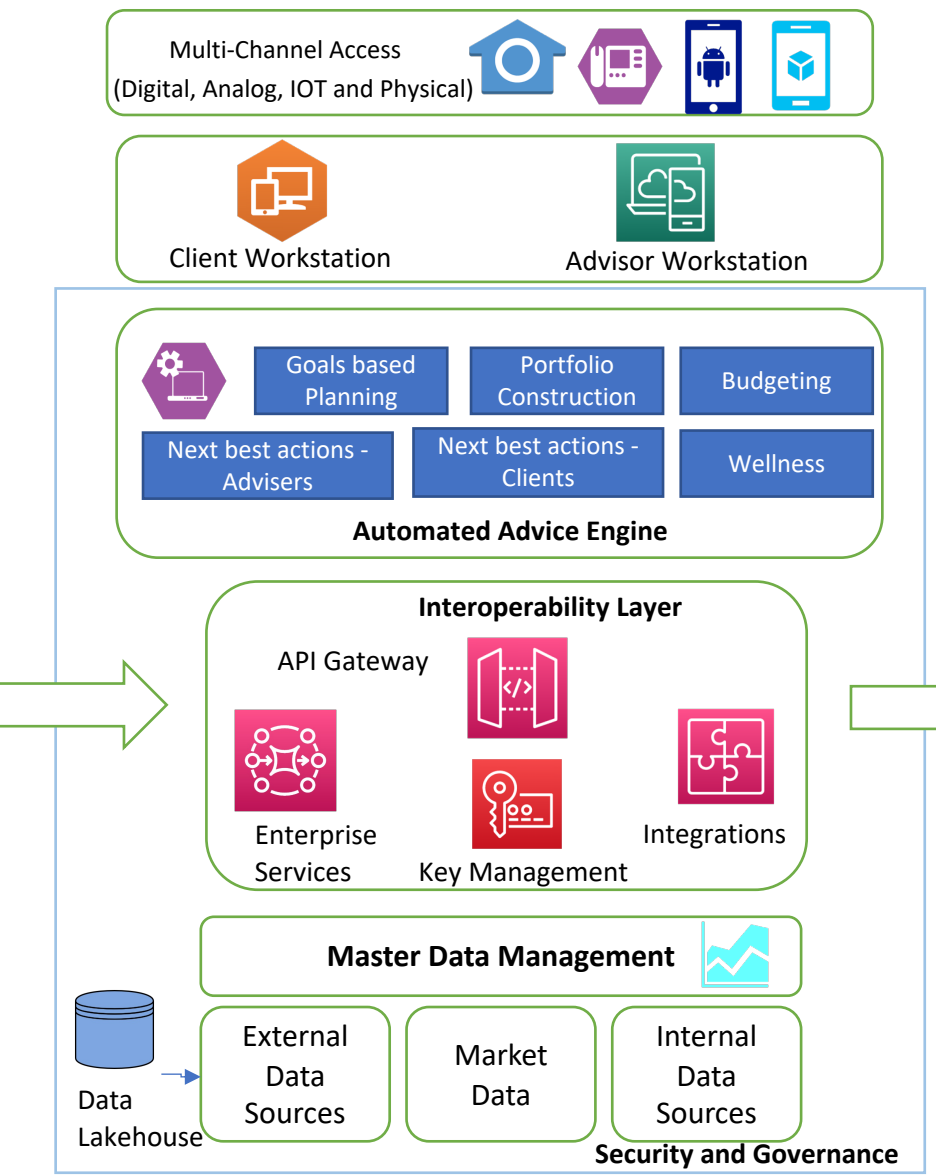- Digital, with e-signature
- Integrated across products

**KYC/AML utility**
*(Proprietary or shared with other wealth managers)*

**Fiduciary/Best Interest engine**
- Guided decision-making
- Product comparisons
- Document generation and recordkeeping

**User research**
*(Voice of clients/FAs/other users)*

**Applied design**
*(Human-centered)*

**Content management**

**Digital marketing**

**Digital ID**

**Investor data store**
*(Shared with other wealth mgrs)*
- Personal data, transaction history, weblogs, social
- Goals and financial plan

**Client data aggregator**
*(Across providers/clients)*

**Multi-Channel Access**
(Digital, Analog, IOT and Physical)

**Client Workstation**

**Advisor Workstation**

## Automated Advice Engine

| Goals based Planning | Portfolio Construction | Budgeting |
|---|---|---|
| Next best actions - Advisers | Next best actions - Clients | Wellness |

## Interoperability Layer

API Gateway

Enterprise Services

Key Management

Integrations

## Master Data Management

Data Lakehouse

| External Data Sources | Market Data | Internal Data Sources |
|---|---|---|

**Security and Governance**

## Integrated product platforms

**Single investment and trust platform**
- Access to securities, funds, annuities, alts, and trust products
- Single interface with internal/external trading platforms and custodians
- Investment servicing and trust administration

**Custody and clearing**
*(Proprietary or T/P platform)*

**Banking**
*(Proprietary or T/P platform)*
- Consumer lending, SBL
- Deposits and savings
- Payments

**Investment banking and capital markets**
- Conduit lending, M&A, etc.
- Trading desk: structured notes, FX, swaps

**Commercial banking**
- Lending
- Cash management

# AI Trust, Risk and Security Management Pillars

**AI TRiSM**

| Explainability | ModelOps | Data Anomaly Detection | Adversarial Attack Resistance | Data Protection |
|---|---|---|---|---|
| Model Behavior Interpretation, Accountability, Transparency | Model Agnostic Operationalization, Management, Governance | Drift monitoring, Detection of Poisoning, Anomalies, Pertubations | Artifact Localization Attack Detection Attack Defense Adversarial Training | Differential Privacy, Synthetic Data, SMPC, FHE |

IT, Data and Analytics Teams, Legal and Compliance Teams, Enterprise Architects, I&O Teams, LoB

**Gartner.**

# **MLOPS** Emergency App 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**

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

*LLM APIs and Hosting*

**Proprietary API**
(OpenAI, Anthropic)

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

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

**Output**

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

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

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

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

**LEGEND**

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

Arrows show the flow of data through the stack

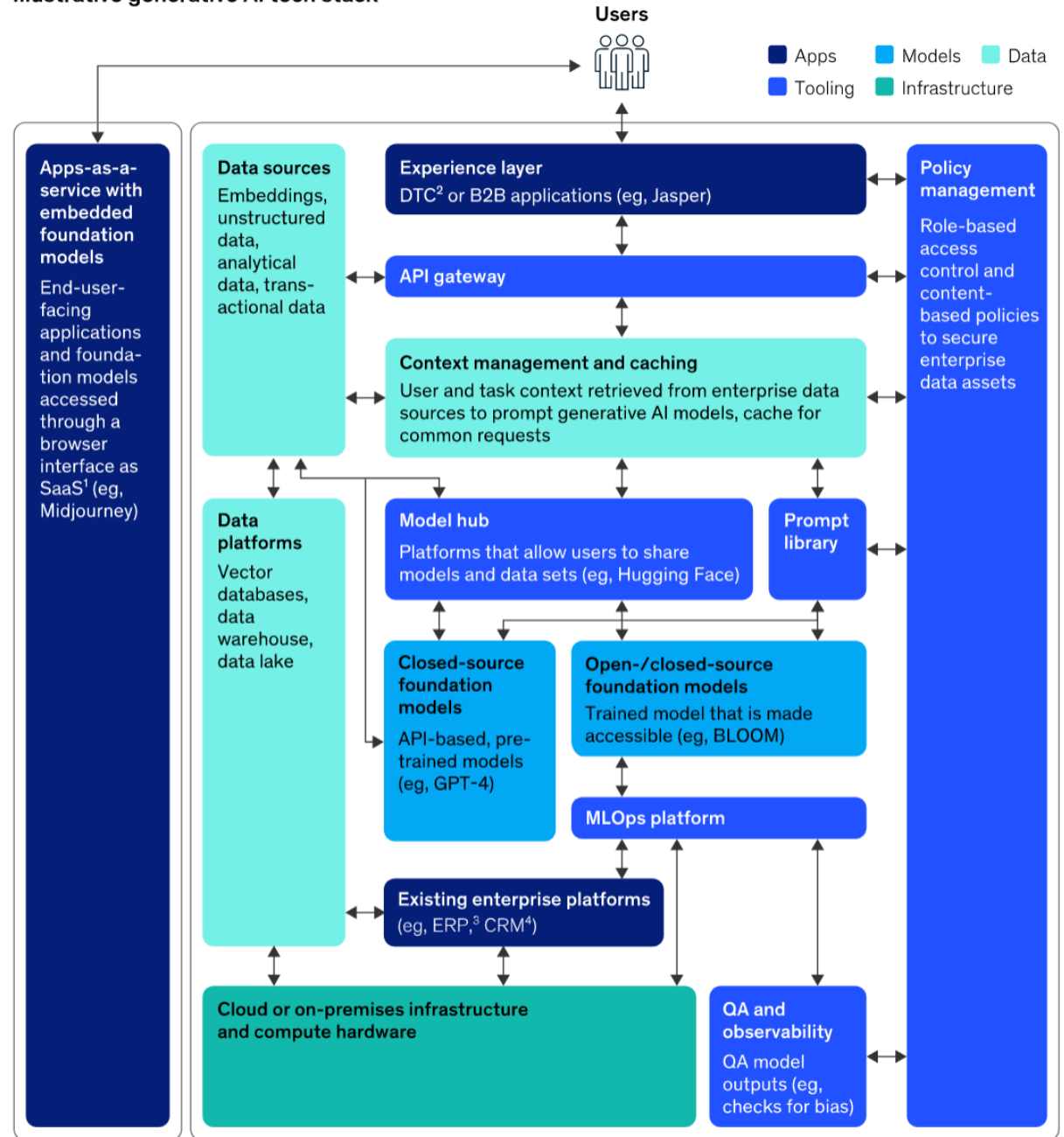- – – → Contextual data provided by app developers to condition LLM outputs
- —→ Prompts and few-shot examples that are sent to the LLM
- —→ Queries submitted by users
- —→ Output returned to users

# The tech stack for generative AI is emerging.

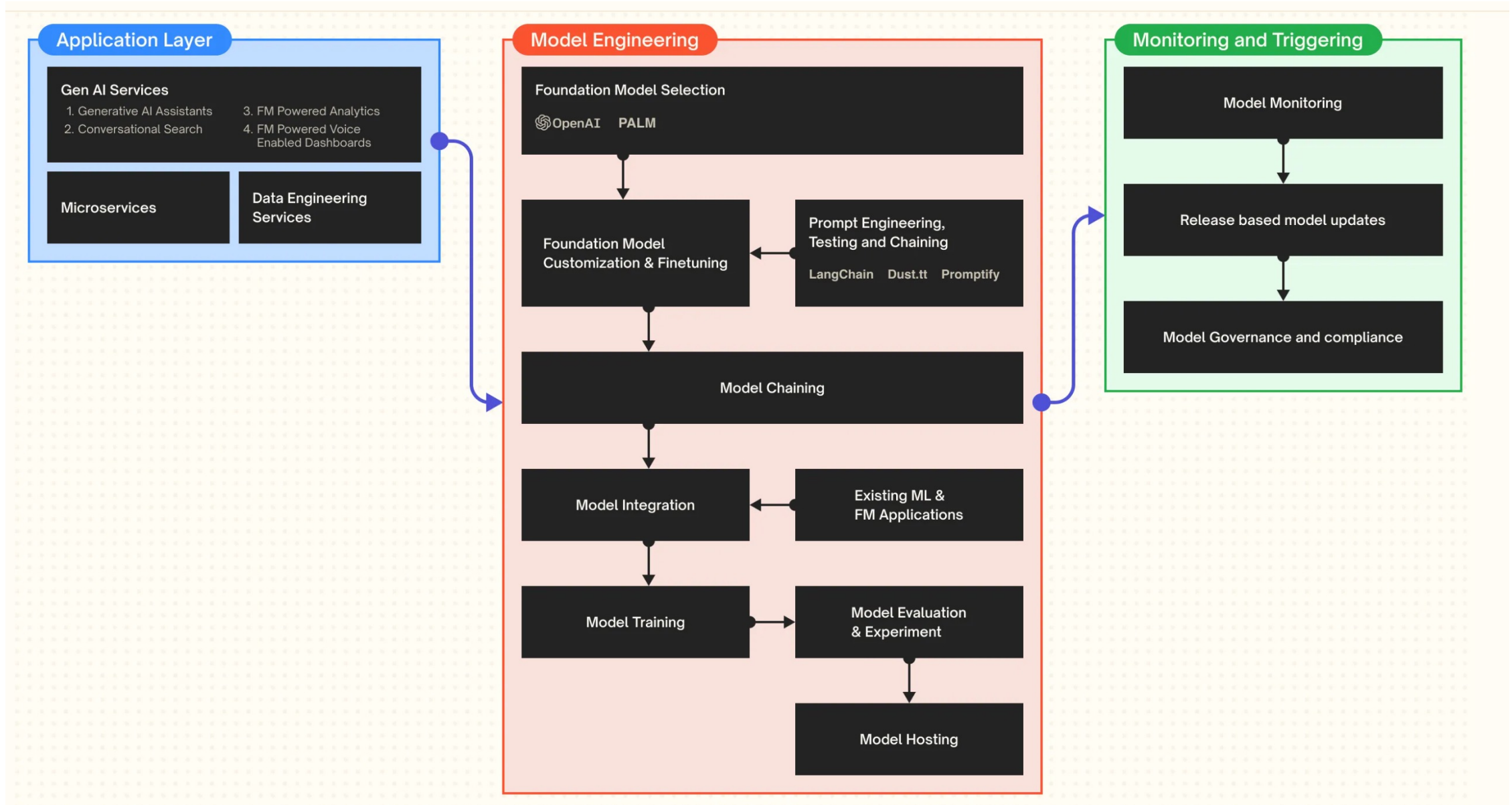**Tech Stack**

Illustrative generative AI tech stack



**Users**

- ■ Apps
- ■ Models
- ■ Data
- ■ Tooling
- ■ Infrastructure

**Apps-as-a-service with embedded foundation models**

End-user-facing applications and foundation models accessed through a browser interface as SaaS[1] (eg, Midjourney)

**Data sources**

Embeddings, unstructured data, analytical data, transactional data

**Experience layer**
DTC[2] or B2B applications (eg, Jasper)

**API gateway**

**Context management and caching**
User and task context retrieved from enterprise data sources to prompt generative AI models, cache for common requests

**Policy management**

Role-based access control and content-based policies to secure enterprise data assets

**Data platforms**

Vector databases, data warehouse, data lake

**Model hub**
Platforms that allow users to share models and data sets (eg, Hugging Face)

**Prompt library**

**Closed-source foundation models**

API-based, pre-trained models (eg, GPT-4)

**Open-/closed-source foundation models**
Trained model that is made accessible (eg, BLOOM)

**MLOps platform**

**Existing enterprise platforms**
(eg, ERP,[3] CRM[4])

**Cloud or on-premises infrastructure and compute hardware**

**QA and observability**

QA model outputs (eg, checks for bias)

# LLM Technology Stack Choices

| Data pipelines | Embedding model | Vector database | Playground | Orchestration | APIs/plugins | LLM cache |
|---|---|---|---|---|---|---|
| Databricks | OpenAI | Pinecone | OpenAI | Langchain | Serp | Redis |
| Airflow | Cohere | Weaviate | nat.dev | LlamaIndex | Wolfram | SQLite |
| Unstructured | Hugging Face | ChromaDB | Humanloop | ChatGPT | Zapier | GPTCache |
| | | pgvector | | | | |

| Logging / LLMops | Validation | App hosting | LLM APIs (proprietary) | LLM APIs (open) | Cloud providers | Opinionated clouds |
|---|---|---|---|---|---|---|
| Weights & Biases | Guardrails | Vercel | OpenAI | Hugging Face | AWS | Databricks |
| MLflow | Rebuff | Steamship | Anthropic | Replicate | GCP | Anyscale |
| PromptLayer | Microsoft Guidance | Streamlit | | | Azure | Mosaic |
| Helicone | LMQL | Modal | | | CoreWeave | Modal |
| | | | | | | RunPod |

# **GEN-AI** Service Development

## Application Layer

### Gen AI Services
1. Generative AI Assistants
2. Conversational Search
3. FM Powered Analytics
4. FM Powered Voice Enabled Dashboards

### Microservices

### Data Engineering Services

## Model Engineering

### Foundation Model Selection
OpenAI    PALM

### Foundation Model Customization & Finetuning

### Prompt Engineering, Testing and Chaining
LangChain    Dust.tt    Promptify

### Model Chaining

### Model Integration

### Existing ML & FM Applications

### Model Training

### Model Evaluation & Experiment

### Model Hosting

## Monitoring and Triggering

### Model Monitoring

### Release based model updates

### Model Governance and compliance

# ML Tech Stack and Security Risk Management

# Deepchecks

Tests for Continuous Validation of ML Models & Data. Deepchecks is a holistic open-source solution for all of your AI & ML validation needs, enabling to thoroughly test your data and models from research to production.

Deepchecks includes:

- **Deepchecks Testing** (Quickstart, docs):
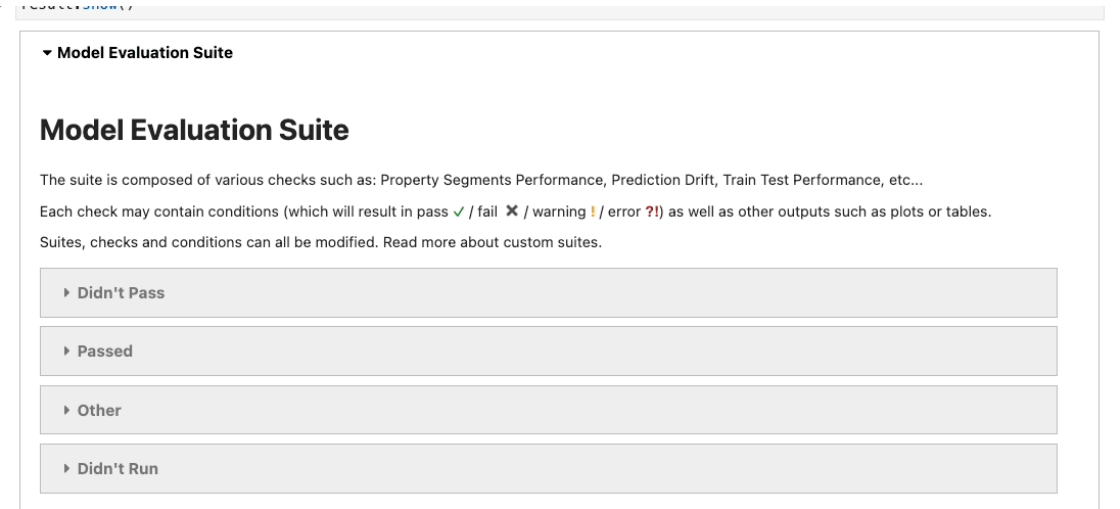    - Running built-in & your own custom Checks and Suites for Tabular, NLP & CV validation (open source).
- **CI & Testing Management** (Quickstart, docs):
    - Collaborating over test results and iterating efficiently until model is production-ready and can be deployed (open source & managed offering).
- **Deepchecks Monitoring** (Quickstart, docs):
    - Tracking and validating your deployed models behavior when in production (open source & managed offering).

## Model Evaluation Suite

The suite is composed of various checks such as: Property Segments Performance, Prediction Drift, Train Test Performance, etc...

Each check may contain conditions (which will result in pass ✓ / fail ✗ / warning ! / error ?!) as well as other outputs such as plots or tables.

Suites, checks and conditions can all be modified. Read more about custom suites.

- ▸ Didn't Pass
- ▸ Passed
- ▸ Other
- ▸ Didn't Run

OK! We have many important issues being surfaced by this suite. Let's dive into the individual checks:

## Model Eval #1: Train Test Performance

We can immediately see in the "Didn't Pass" tab that there has been significant degradation in the Recall on class "optimism". This is very likely a result of the severe label drift we saw after running the previous suite.

## Model Eval #2: Segment Performance

Also in the "Didn't Pass" tab we can see the two segment performance checks - Property Segment Performance and Metadata Segment Performance. These use the metadata columns of user related information OR our calculated properties to try and **automatically** detect significant data segments on which our model performs badly.
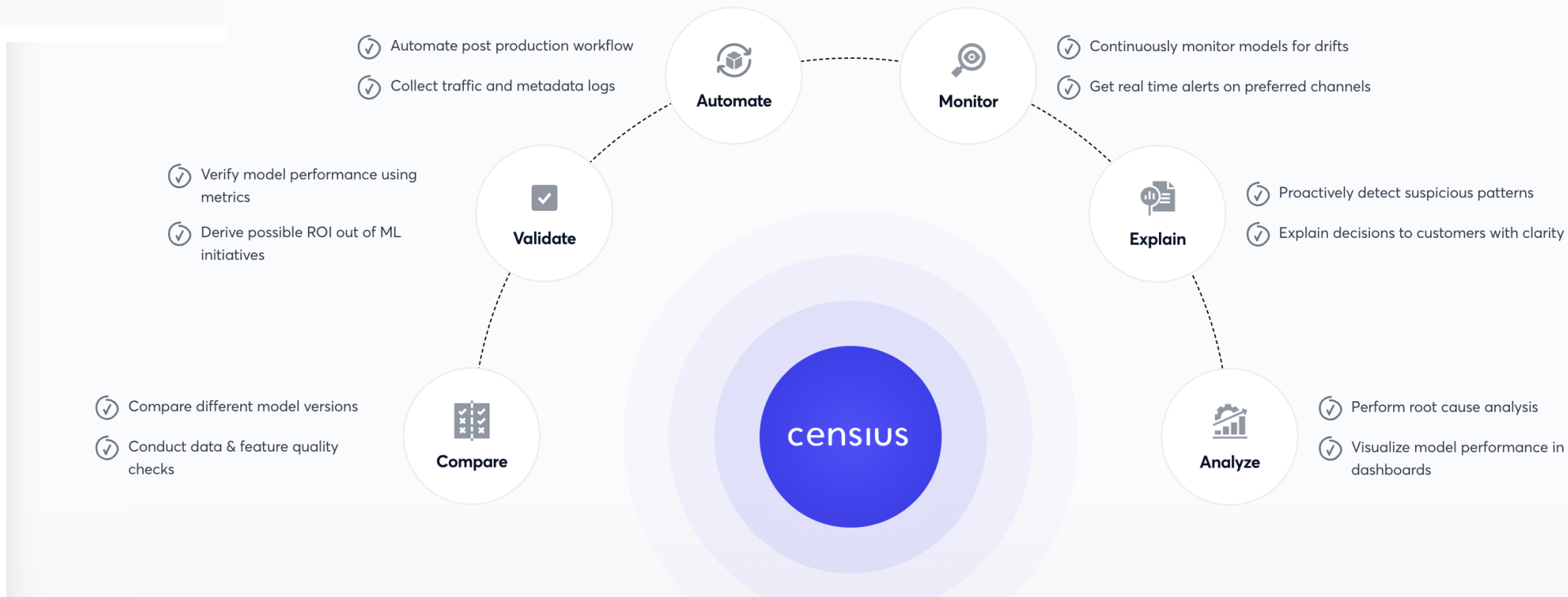
**censius** AI Observability Platform

**attri** Generative AI Solutions

Contextual Relationship Intelligence (Stealth)

Simulated Agents (Stealth)

# A single platform for delivering enterprise level observability at scale.

✓ Automate post production workflow

✓ Collect traffic and metadata logs

**Automate**

**Monitor**

✓ Continuously monitor models for drifts

✓ Get real time alerts on preferred channels

✓ Verify model performance using metrics

✓ Derive possible ROI out of ML initiatives

**Validate**

**Explain**

✓ Proactively detect suspicious patterns

✓ Explain decisions to customers with clarity

✓ Compare different model versions

✓ Conduct data & feature quality checks

**Compare**

**censius**

**Analyze**

✓ Perform root cause analysis

✓ Visualize model performance in dashboards

# Sentima

## The Contextually Aware Converged Security Platform

http://www.sentima.io

### Contextual Awareness

Intent and Context aware platform that defines Why, Where, What, When of a request so proactive security decisions can be made instantly

### Identity Verification and Attestation

Attestation and Verification based User, Workload, Machine, Process, Network, Service Verification and Secure Communication
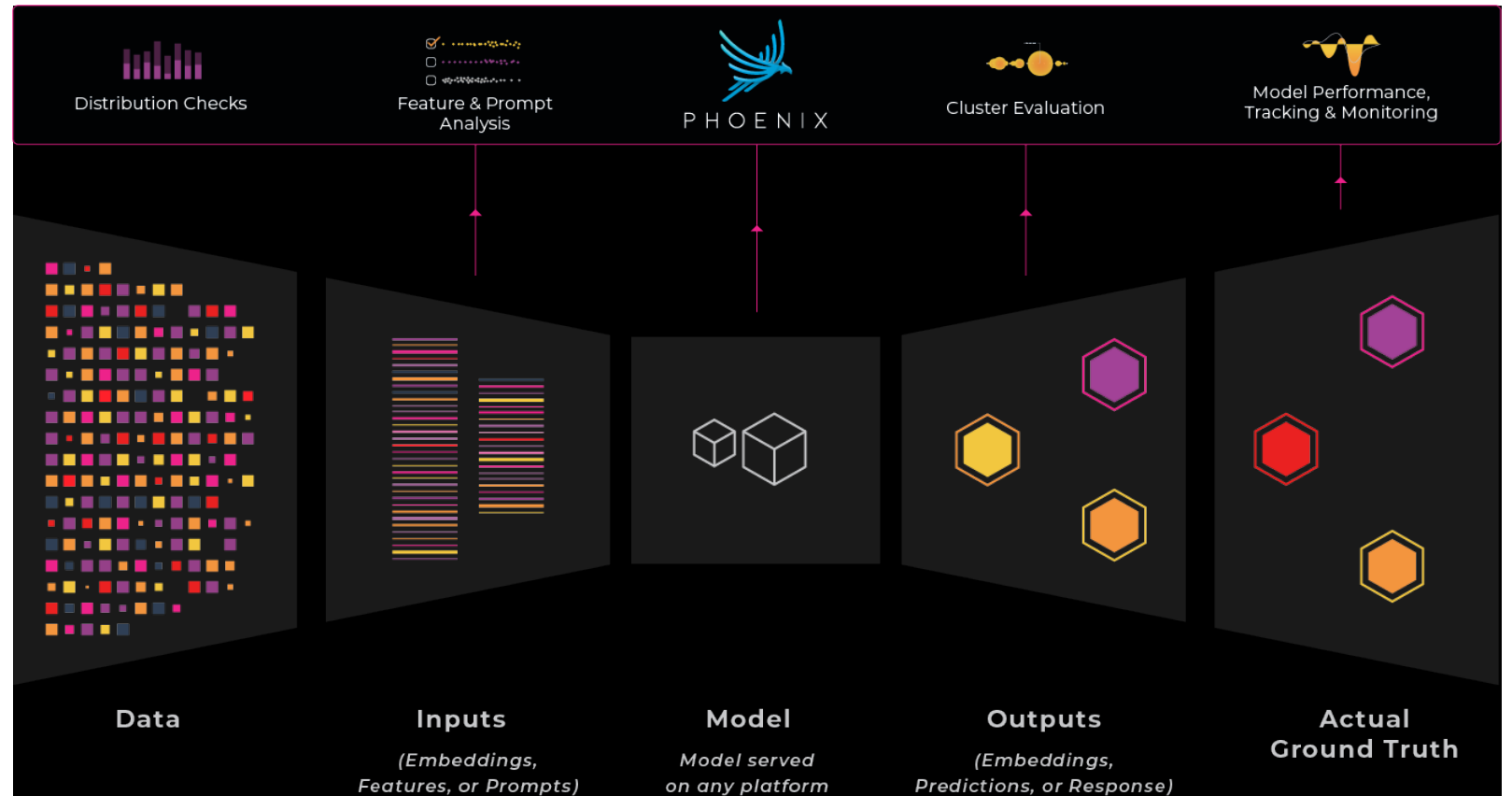
### Point to Point Zero Trust

Point to Point Zero Trust between Service to Data Stores, Service to Service, Users to Services, without Passwords

# Phoenix ML Observability in a Notebook

Phoenix provides ML insights at lightning speed with zero-config observability for model drift, performance, and data quality.

Phoenix is an Open Source ML Observability library carefully designed for the Notebook. The toolset is designed to ingest inference data for LLMs, CV, NLP and tabular datasets. It allows Data Scientists and AI Engineers to quickly visualize their inference data, monitor performance, track down issues & insights, and easily export to improve.



**Phoenix Functionality**

**Discover How Embeddings Represent Your Data:** Map structured features onto embeddings for deeper insights into how embeddings represent your data.
**Evaluate LLM Tasks:** Troubleshoot tasks such as summarization or question/answering to find problem clusters with misleading or false answers.
**Find Clusters of Issues to Export for Model Improvement:** Find clusters of problems using performance metrics or drift. Export clusters for fine-tuning workflows.
**Detect Anomalies:** Using LLM embeddings
**Surface Model Drift and Multivariate Drift:** Use embedding drift to surface data drift for generative AI, LLMs, computer vision (CV) and tabular models.
**Easily Compare A/B Datasets:** Uncover high-impact clusters of data points missing from model training data when comparing training and production datasets.

**EU lawmakers pass landmark artificial intelligence regulation**

PUBLISHED WED, JUN 14 2023·9:45 AM EDT | UPDATED WED, JUN 14 2023·1:28 PM EDT

- The European Union's AI Act is the first comprehensive set of regulations for the artificial intelligence industry.

- The law proposes requiring generative AI systems, such as ChatGPT, to be reviewed before commercial release. It also seeks to ban real-time facial recognition.

- It comes as global regulators are racing to get a handle on the technology and limit some of the risks to society, including job security and political integrity.

**BLUEPRINT FOR AN AI BILL OF RIGHTS**

MAKING AUTOMATED SYSTEMS WORK FOR THE AMERICAN PEOPLE

OSTP

Safe and Effective Systems

Algorithmic Discrimination Protections

Data Privacy

Notice and Explanation

Human Alternatives, Consideration, and Fallback

# America's first law regulating AI bias in hiring takes effect this week

While the law aims for transparency, critics say it may not be enough to protect against AI bias

New York City Adopts Final Regulations on Use of AI in Hiring and Promotion, Extends Enforcement Date to July 5, 2023

- **Automated resume screeners** that read job applications and recommend the best candidates for an open role
- **Matchmaking algorithms** that scour millions of job postings to recommend roles to candidates—and vice versa
- **Social media scrapers** that collect data on applicants to compile personality profiles based on what they've found online
- **AI-based chatbots** that ask candidates questions about their qualifications, then decide if they'll proceed in the interview process
- **Algorithmic video platforms** that have candidates answer interview questions on camera, record their replies, transcribe their responses, and analyze their vocal or facial patterns for subjective traits like "openness" or "conscientiousness"
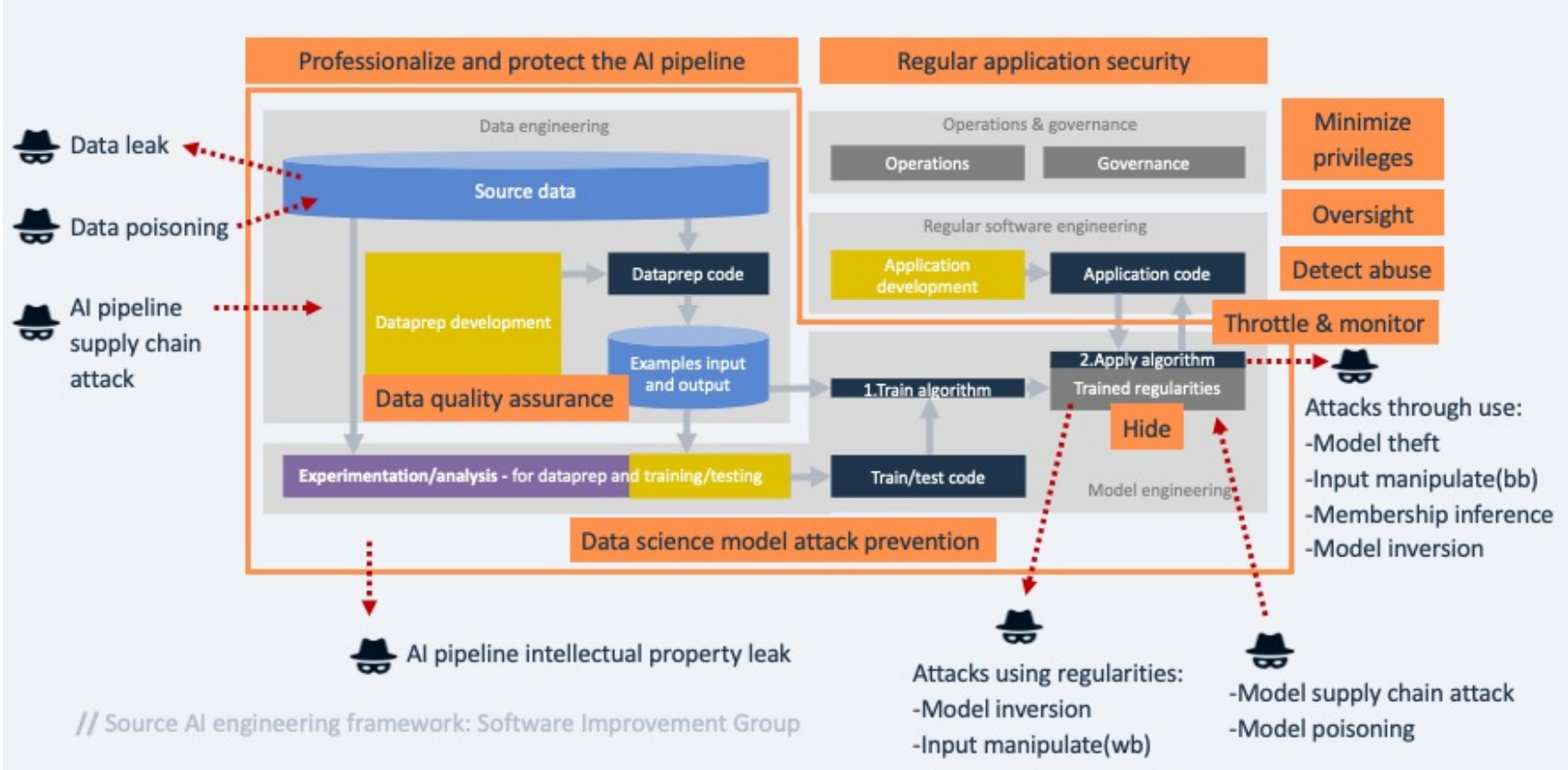- **Logic games** that purport to identify qualities like "risk-taking" or "generosity"

# AI RMF 1.0

On January 26, 2023, NIST released the AI Risk Management Framework (AI RMF 1.0) along with a companion NIST AI RMF Playbook, AI RMF Explainer Video, an AI RMF Roadmap, AI RMF Crosswalk, and various Perspectives.
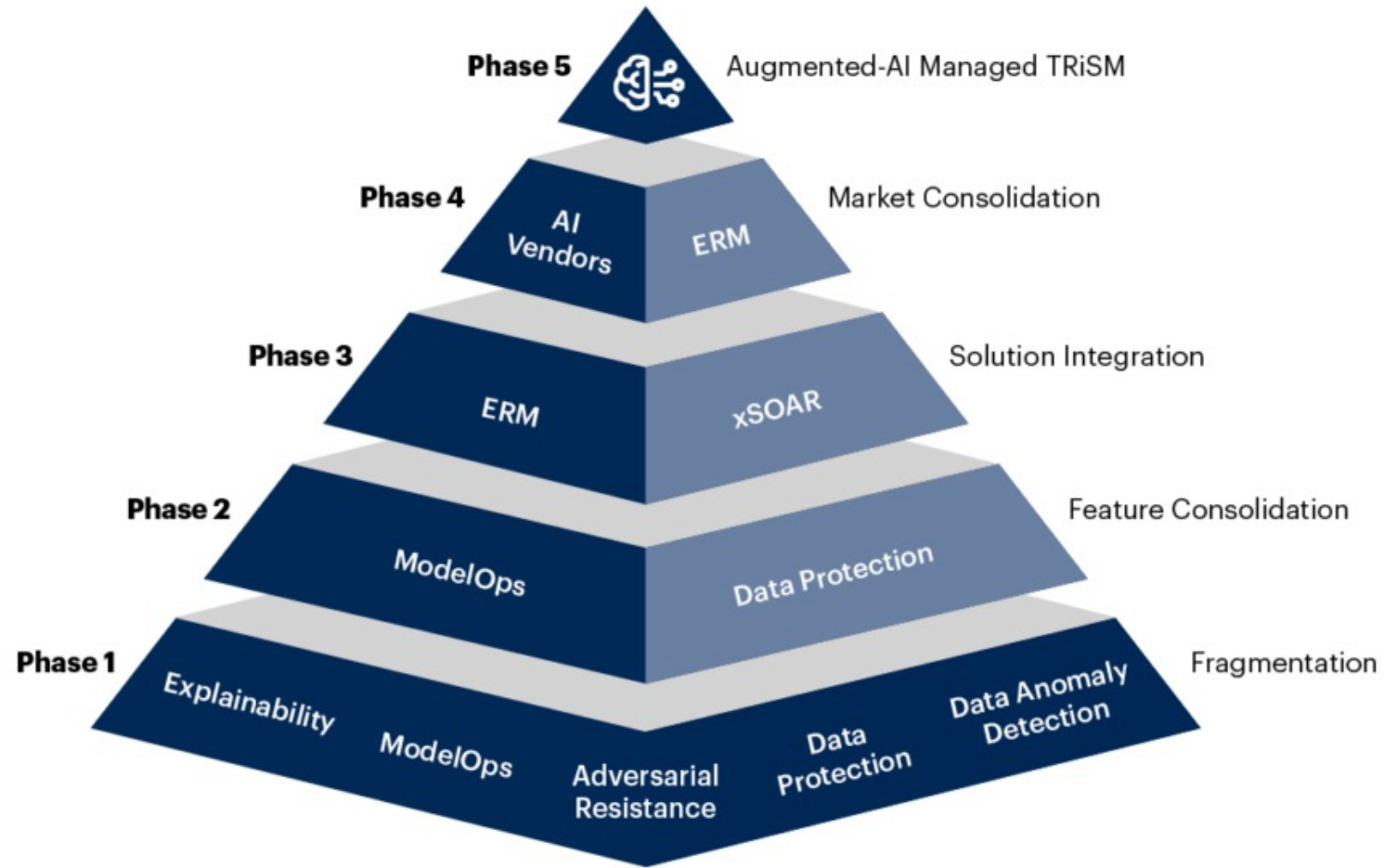
Safe | Secure & Resilient | Explainable & Interpretable | Privacy-Enhanced | Fair - With Harmful Bias Managed | Accountable & Transparent

Valid & Reliable

| | Application Context | Data & Input | AI Model | AI Model | Task & Output | Application Context | People & Planet |
|---|---|---|---|---|---|---|---|
| **Key Dimensions** | Application Context | Data & Input | AI Model | AI Model | Task & Output | Application Context | People & Planet |
| **Lifecycle Stage** | Plan and Design | Collect and Process Data | Build and Use Model | Verify and Validate | Deploy and Use | Operate and Monitor | Use or Impacted by |
| **TEVV** | TEVV includes audit & impact assessment | TEVV includes internal & external validation | TEVV includes model testing | TEVV includes model testing | TEVV includes integration, compliance testing & validation | TEVV includes audit & impact assessment | TEVV includes audit & impact assessment |
| **Activities** | Articulate and document the system's concept and objectives, underlying assumptions, and context in light of legal and regulatory requirements and ethical considerations. | Gather, validate, and clean data and document the metadata and characteristics of the dataset, in light of objectives, legal and ethical considerations. | Create or select algorithms; train models. | Verify & validate, calibrate, and interpret model output. | Pilot, check compatibility with legacy systems, verify regulatory compliance, manage organizational change, and evaluate user experience. | Operate the AI system and continuously assess its recommendations and impacts (both intended and unintended) in light of objectives, legal and regulatory requirements, and ethical considerations. | Use system/ technology; monitor & assess impacts; seek mitigation of impacts, advocate for rights. |
| **Representative Actors** | System operators; end users; domain experts; AI designers; impact assessors; TEVV experts; product managers; compliance experts; auditors; governance experts; organizational management; C-suite executives; impacted individuals/ communities; evaluators. | Data scientists; data engineers; data providers; domain experts; socio-cultural analysts; human factors experts; TEVV experts. | Modelers; model engineers; data scientists; developers; domain experts; with consultation of socio-cultural analysts familiar with the application context and TEVV experts. | | System integrators; developers; systems engineers; software engineers; domain experts; procurement experts; third-party suppliers; C-suite executives; with consultation of human factors experts, socio-cultural analysts; governance experts; TEVV experts, | System operators, end users, and practitioners; domain experts; AI designers; impact assessors; TEVV experts; system funders; product managers; compliance experts; auditors; governance experts; organizational management; impact- ed individuals/commu- nities; evaluators. | End users, operators, and practitioners; impacted individu- als/communities; general public; policy makers; standards organizations; trade associations; advocacy groups; environmental groups; civil society organizations; researchers. |

- [ML01:2023 Adversarial Attack](#)
- [ML02:2023 Data Poisoning Attack](#)
- [ML03:2023 Model Inversion Attack](#)
- [ML04:2023 Membership Inference Attack](#)
- [ML05:2023 Model Stealing](#)
- [ML06:2023 Corrupted Packages](#)
- [ML07:2023 Transfer Learning Attack](#)
- [ML08:2023 Model Skewing](#)
- [ML09:2023 Output Integrity Attack](#)
- [ML10:2023 Neural Net Reprogramming](#)

Future Direction AI TRiSM Market

Source: Gartner
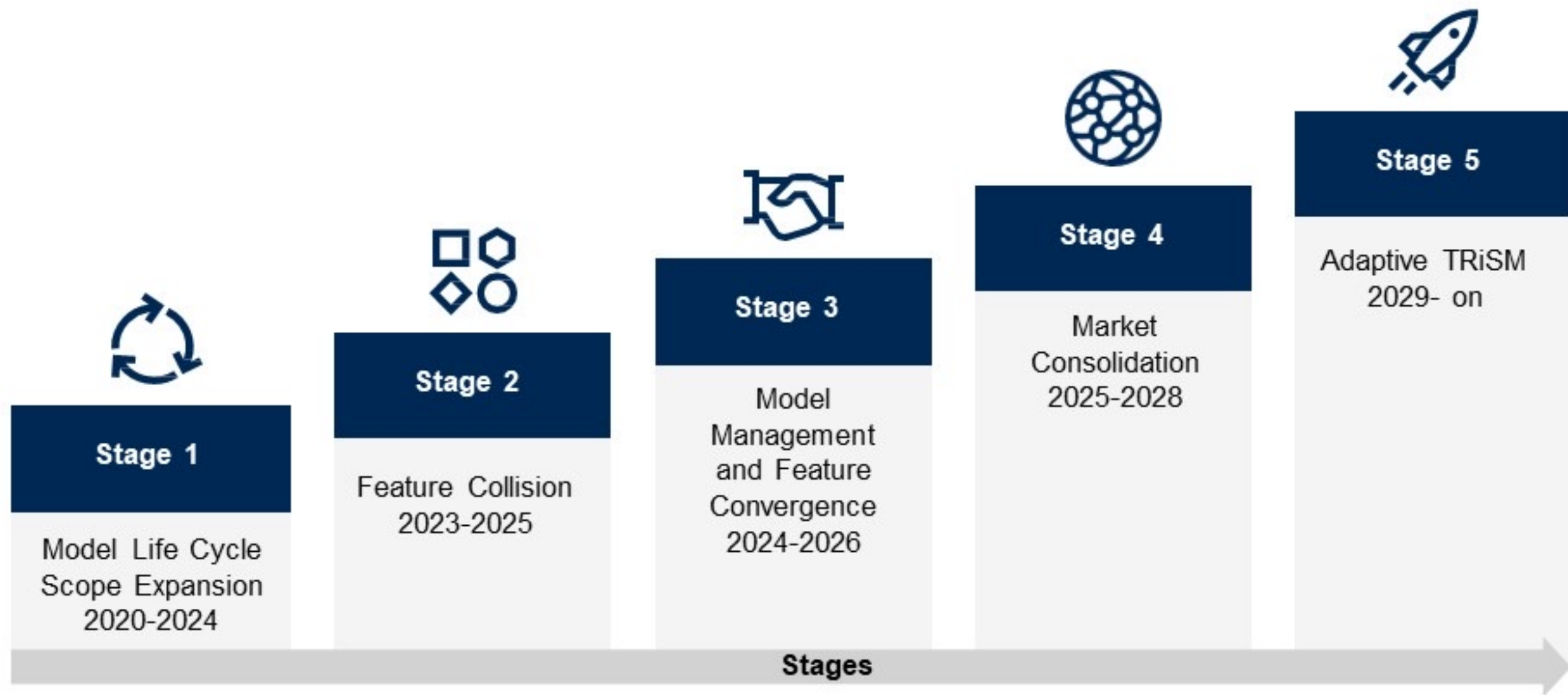750738_C

# Future Direction of the AI TRiSM Market

**Stage 1**

Model Life Cycle Scope Expansion 2020-2024

**Stage 2**

Feature Collision 2023-2025

**Stage 3**

Model Management and Feature Convergence 2024-2026

**Stage 4**

Market Consolidation 2025-2028

**Stage 5**

Adaptive TRiSM 2029- on

**Stages**

# North Austin Tech User Group : AI Focused

# NATU.AI

Connect Me at **Linked** in ®

**/in/mandavasuresh**



**Suresh Mandava**
SVP/Chief Architect
Cloud-Native AI/ML Platforms and Security
Leander, Austin, Texas : 636-634-0552

Linkedin : /in/MandavaSuresh
Twitter : @sureshmandava

**AITX Meetup**