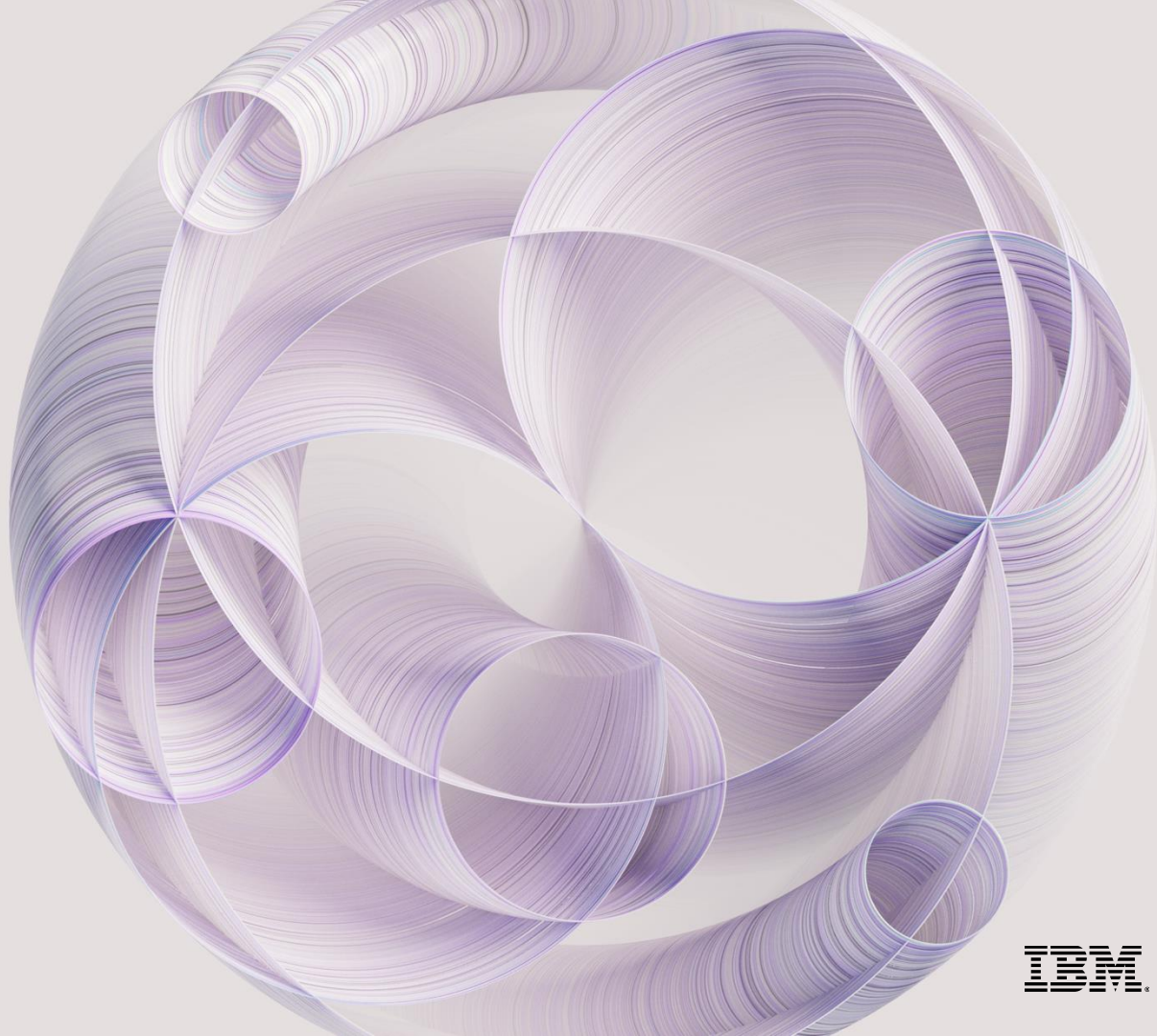# Onboarding Workshop

# watsonx

IBM

# Agenda

- 12:00 PM - 12:10 PM Introductions

- 12:10 PM - 12:30 PM Account Onboarding / Logistics

- 12:30 PM - 1:00 PM watsonx Overview

- 1:00 PM - 2:00 PM watsonx.ai Labs

**Lab Materials**

Yann LeCun ✔ ∞
@ylecun
...

Big Tech is not the problem.
Closed and proprietary is the problem.
Meta, IBM are Big Tech and open.
Google, Apple are Big Tech and closed.
OpenAI, Anthropic are Small Tech and closed.
Hugging Face, Mistral are Small Tech and open.

2:02 PM · 9/23/23 · **12.1K** Views

# The platform for AI and data

## watsonx

Scale and accelerate the impact of AI across your business

### watsonx.ai

Build, train, validate, tune and deploy AI models

A next generation enterprise studio for AI builders to build, train, validate, tune, and deploy both traditional machine learning and new generative AI capabilities powered by foundation models. It enables you to build AI applications in a fraction of the time with a fraction of the data.

### watsonx.data

Scale AI workloads, for all your data, anywhere

Fit-for-purpose data store, built on an open lakehouse architecture, supported by querying, governance and open data formats to access and share data.
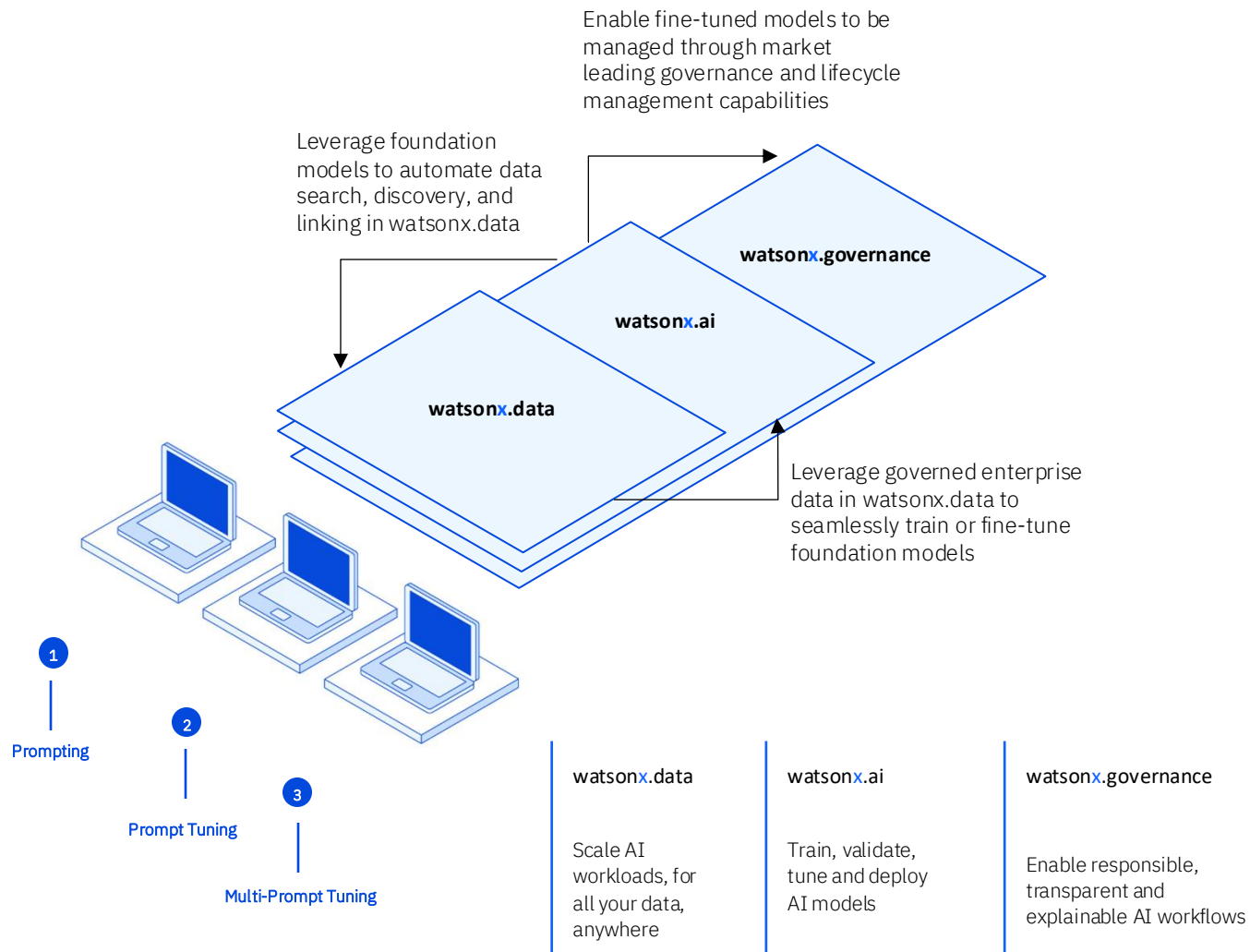
### watsonx.governance

Accelerate responsible, transparent and explainable AI workflows

End-to-end toolkit for AI governance across the entire model lifecycle to accelerate responsible, transparent, and explainable AI workflows

# watsonx

## Scale and accelerate the impact of AI with trusted data.

Enable fine-tuned models to be managed through market leading governance and lifecycle management capabilities

Leverage foundation models to automate data search, discovery, and linking in watsonx.data

**watsonx.governance**

**watsonx.ai**

**watsonx.data**

Leverage governed enterprise data in watsonx.data to seamlessly train or fine-tune foundation models

**1** Prompting

**2** Prompt Tuning

**3** Multi-Prompt Tuning

| watsonx.data | watsonx.ai | watsonx.governance |
|---|---|---|
| Scale AI workloads, for all your data, anywhere | Train, validate, tune and deploy AI models | Enable responsible, transparent and explainable AI workflows |

What IBM offers

# AI assistants

# watson**x**

## watson**x** Orchestrate

Harness the power
of AI and automation
to free up individuals
from tedious tasks

Enable employees to
quickly offload time-
consuming work to tackle
more of the work only they
can do. Business users can
delegate common and
complex tasks such as
creating a job description,
pulling a report in
Salesforce or SAP
SuccessFactors, sourcing
candidates, and more
using natural language.

## 40%
improvement in
HR productivity[1]

## watson**x** Assistant

Build better virtual agents,
to deliver consistent and
intelligent customer care

Understand customers in the
right context, and provide
fast, consistent, and accurate
answers, and self-service
support across any application,
device, or channel. The
intuitive build experience
empowers everyone in the
organization to build and
deploy AI-powered virtual
agents without writing a
line of code.

## >90%
customer inquiries
handled by AI assistant[2]

## watson**x** Code Assistant

Accelerate development,
application modernization,
and assist with IT Operations

Increase developer
productivity, reduce coding
complexity, and accelerate
developer onboarding.
Purpose-built for targeted
use cases, watsonx Code
Assistant uses AI to support
application modernization
and IT automation.

## 60%
software development
content automatically
generated by AI[3]

[1] IBM HR use case
[2] Vodafone Case Study in partnership with IBM
and Genesys
[3] IBM CIO case study based on limited internal test

# IBM's generative AI technology and expertise

| | | |
|---|---|---|
| **AI assistants** | Empower individuals to do work without expert knowledge across a variety of business processes and applications. | **watsonx** Code Assistant<br>**watsonx** Assistant<br>**watsonx** Orchestrate<br>**watsonx** Orders |
| **SDKs & APIs** | Embed watsonx platform in third party assistants and applications using programmatic interfaces. | **Ecosystem integrations** |
| **AI & data platform** | Leverage generative AI and machine learning — tuned with your data — with responsibility, transparency and explainability. | **watsonx**<br>watsonx.ai<br>watsonx.governance<br>watsonx.data  **Foundation models**<br>Granite │ *IBM*<br>Open Source │ *Hugging Face*<br>Llama 2 │ *Meta*<br>Geospatial │ *IBM + NASA*<br>... |
| **Data services** | Define, organize, manage, and deliver trusted data to train and tune AI models with data fabric services. | **Cloud Pak for Data**<br>**watsonx** Discovery |
| **Hybrid cloud AI tools** | Build on a consistent, scalable, foundation based on open-source technology. | **Red Hat** OpenShift AI<br>(*e.g.,* Ray, Pytorch) |

**Consulting**
Generative AI strategy, experience, technology, operations

**Ecosystem**
System Integrators, Software and SaaS partners, Public Cloud providers

# watsonx

Model strategy →

## Multi-model

**One model doesn't fit all use cases.** We offer IBM-developed, open-source, third party, and BYOM.

**Bigger is not always better.** Specialized models can outperform general-purpose models with lower infrastructure requirements.

## Hybrid, multi-cloud

**Hybrid deployments.** We provide the flexibility to deploy models on the platform of choice.



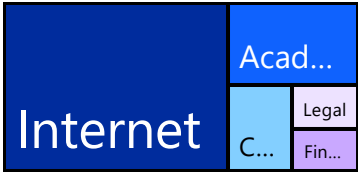*granite.20b.code is delivered through watsonx Code Assistant*

# What is IBM Granite ?

➢ Granite is IBM's flagship series of LLM foundation models based on decoder-only transformer architecture.

➢ Granite language models are trained on trusted enterprise data spanning internet, academic, code, legal and finance.

By 2027 more than 50% of the Gen AI models that enterprises use will be domain-specific — specific to either an industry or business function — up from approximately 1% in 2023.

Gartner Report, Predicts 2024: The Future of Generative AI Technologies.

## Trusted, Performant, Cost-effective AI foundation models purpose built for enterprises.

### granite-13b-v2 ( English LLM )
*-chat-v2.1, -instruct-v2*

**13B** parameters in size
**2.5T** tokens of data



(v1 breakdown)

➢ Chat derivative model is optimized for dialogue use cases and works well with virtual agent and chat applications.

➢ Instruct derivative model was designed to perform well on natural language tasks and can be customized for specific industries and domains via prompt-tuning.

### granite-20b-multilingual

**20B** parameters in size
**2.6 T** tokens of Data



### granite-7b-lab
*Open-source*

**7B** parameters in size

➢ Tuned using IBM's large-scale alignment of chatbots(LAB).

### granite-8b-japanese

**8B** parameters in size
**1.6T** tokens of Data



### granite-code
*Open-source*

**3B, 8B, 20B, 34B** parameters in size

➢ A family of models trained in 116 programming languages



Granite-code-20b

IBM open-source models : https://huggingface.co/ibm-granite

# watsonx

## Operating Costs →

### Inferencing Costs (aaS)

**Price Differentiation:** Model costs per token in managed infrastructure can vary significantly by model type and service provider.

### Infrastructure Scale

**Capacity and Performance.** Model performance can also differ greatly in infrastructure requirements and speed of inferencing.

## IBM Granite 13B models operating at up to 62 times lower cost than GPT4.

### Inference costs for a customer summarizing 80 million chat sessions.

| Cost $ | Price per 1K Input Tokens | Price per 1K Output Tokens | Avg. price per 1K Tokens* | Cost per 14K Token chat session** | Costs for 80M chat sessions per year |
|---|---|---|---|---|---|
| GPT 4 | 0.03 | 0.06 | 0.039 | 0.546 | $43.7M |
| Llama2 70B | 0.0018 | 0.0018 | 0.0018 | 0.0252 | $2.0M |
| GPT 3.5-Turbo | 0.0005 | 0.0015 | 0.0008 | 0.0112 | $0.9M |
| Granite 13B | 0.0006 | 0.0006 | 0.0006 | 0.0084 | $0.7M |

33

* Average of 70% input and 30% output tokens
**A typical session is 10,500 words or 14K tokens

Link to IBM Pricing
Link to OpenAI Pricing

62X
3X
1.34X

Granite 13b   GPT 3.5 Turbo   Llama2 70B   GPT 4

### Significant cost impact as you scale



Cost ($M) annually

GPT 4
$44M

Llama2 70B       $2M
GPT 3.5
$0.9M
Granite 13b      $0.7M

Volume of chatbot conversations

$43M Savings

$0.2-$1.3M Savings

# IBM watsonx.ai architecture

## AI Tooling

### Prepare Data

**Synthetic Data Generator**
Generate synthetic tabular data

**Data Refinery**
Prepare and visualize data

### Work with models

**Prompt Lab**
Experiment with foundation models and build prompts

**Tuning Studio**
Tune a foundation model with labeled data

**SPSS Modeler**
Build models as a visual flow

**RStudio IDE**
Work with data and models in R

**Jupyter notebook editor**
Work with data and models in Python or R notebooks

**AutoAI**
Build ML models automatically

**Decision Optimization**
Solve optimization problems

**Federated Learning**
Train models on distributed data

### Automate model lifecycles

**Orchestration Pipelines**
Automate model lifecycles

## Common core services

**Projects**
Collaborate with others to work with data and build AI assets

**Deployment Spaces**
Deploy and monitor your assets

## Hybrid Cloud Platform

**Red Hat OpenShift AI**

**IBM Cloud Infrastructure**

**Third party infrastructure**

## Common core services

- Collaborative projects
- Deployment spaces
- Jobs
- Notifications
- Common connectivity
- Access and Authentication
- Resource management
- Central asset management system

# watsonx.ai: Prompt Lab
## Experiment with foundation models and build prompts

Interactive prompt builder

Includes prompt examples for various use cases and tasks

Experiment with different prompts, save and reuse older prompts, use different models and vary different parameters

Experiment with zero-shot, one-shot, or few-shot prompting to get the best results

Experiment with prompt engineering

Choice of foundation models to use based on task requirements

Prevent the model from generating repeating phrases

Number of min and max new tokens in the response

Stop sequences – specifies sequences whose appearances should stop the model

# watson**x**.ai: Tuning Studio
## Tune your foundation models with labeled data

Summary:

- Tool for performing PEFT and fine-tuning training techniques to optimize FM task performance
- Tuned model can be deployed and inferenced via the API or Prompt Lab

Prompt-tuning:

- **How it works:** creates an optimized sequence of values (called a soft-prompt vector) to add as a prefix to FM prompt to improve task performance
- **Technical origins:** The Power of Scale for Parameter-Efficient Prompt Tuning
- Subset of PEFT, similar to P-Tuning, LoRA, etc.



Product documentation

# **watsonx.ai**: Synthetic Data Generator
## Generate synthetic tabular data to address your data gaps

### Create synthetic data at scale

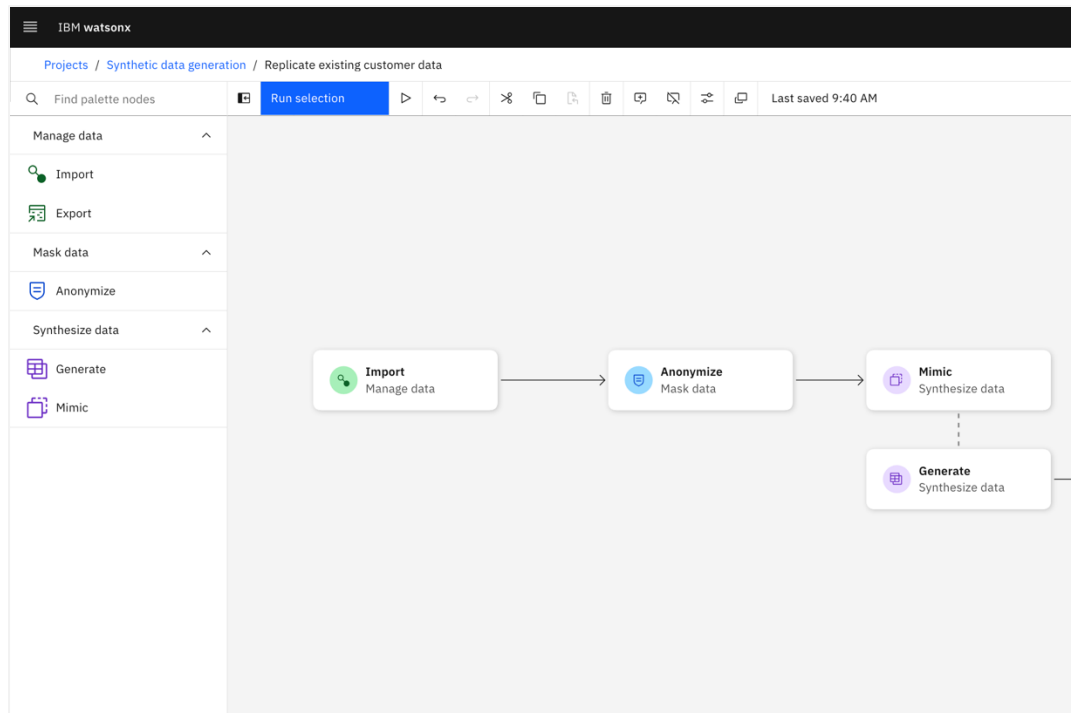Unlock your valuable insights by using synthetic data.

Create synthetic data using your existing data in a database or by uploading a file. If no data exists or can't be accessed, you can design your own data schema.

Address data gaps and create synthetic edge cases to expedite classical AI model training.

### Select your model & privacy needs

Depending on your cost, fidelity, application, or data needs, you can select from multiple IBM models* to create your synthetic tabular data.

When using existing data, IBM models apply differential privacy to minimize your privacy risk and give you control over the level of privacy protection required for your organization.



*Evaluation metrics available in Q3 2024*

# watsonx.ai: Data Science and MLOps
## Build machine learning models automatically in the studio

### Model training and development

Build experiments quickly and enhance training by optimizing pipelines and identifying the right combination of data
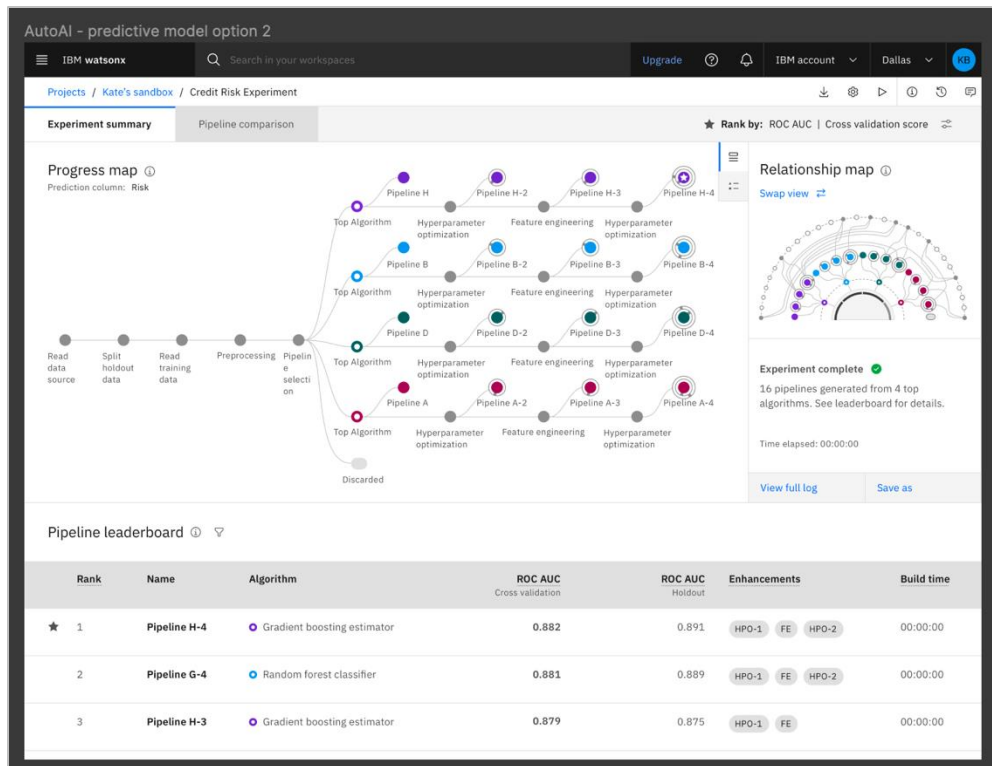
AutoAI, including preparing data for machine learning and generating and ranking candidate model pipelines

Use predictions to optimize decisions, create and edit models in Python, in OPL or with natural language

### Integrated visual modeling

Prepare data quickly and develop models visually to help visualize and analyze enterprise data to identify patterns and trends, explore opportunities, and make informed, insightful business decisions

- Uncover correlations
- Insight for hypotheses
- Find relationships and connections within the data

# watsonx.ai Embeddings API

**What does it do?**
- Converts input text into embeddings, which are dense vector representations of the input text
- Embeddings capture nuanced semantic and syntactic relationships between words and passages in vector space
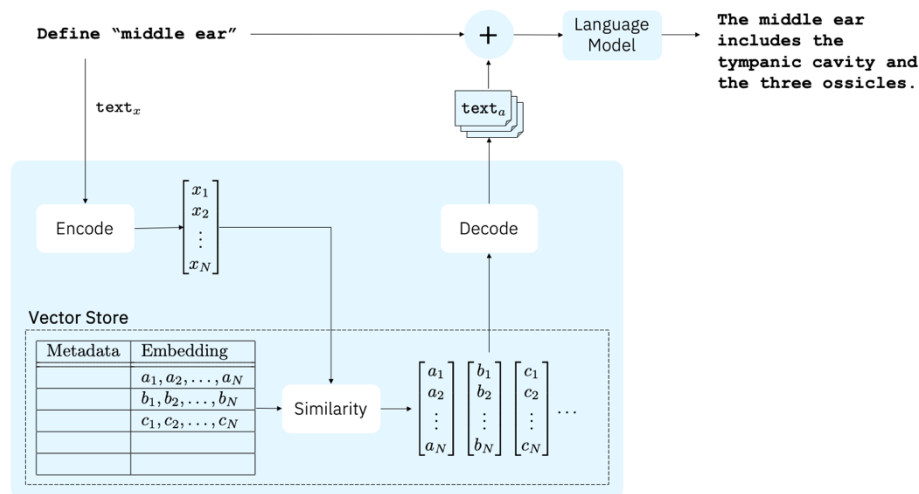
**Customer value**
- Embeddings provide a more semantically faithful representation of the supplied text, especially when compared to basic keyword-based alternatives in classic NLP modeling
- The efficient storage and compute profiles of embeddings make them easily infusible into generative AI application
- Retrieval Augmented Generation (RAG) patterns utilize embedding models for query and passage vectorization, enabling contextual grounding

**IBM differentiators**
- Performance matching or exceeding market leaders in retrieval benchmarks
- Exclusively trained on legally approved and commercially viable data to enable enterprise usage
- Scalable, fully-managed, and integrated with the broader watsonx portfolio

# watson**x**.governance

Accelerate responsible, transparent and explainable AI

*One unified, integrated AI Governance platform to govern generative AI and predictive ML*

## Lifecycle Governance

Govern across the AI lifecycle. Automate and consolidate tools, applications and platforms. Capture metadata at each stage and support models built and deployed in 3rd party tools.

## Risk Management

Manage risk & protect reputation by automating workflows to ensure quality and better detect bias and drift.

## Regulatory Compliance

Adhere to regulatory compliance by translating growing regulations into enforceable policies.

**Comprehensive**
Govern the end-to-end AI lifecycle with metadata capture at each stage

**Open**
Support governance of models built and deployed in 3rd party tools.

**Automatic metadata recording**
and data transformation/lineage capture though Python notebooks.

# watsonx.governance



**Model Risk Governance**

- Consolidated view of models from multiple platforms

- View development status, model performance and alerts or emerging issues

- Monitor and trigger workflows for model validation, retraining and performance issues



**Model Request**
Management approvals
Use case definition
Risk assessment

**Model Development**
Coding / Training
Experiments
Prompt-Tuning / Fine Tuning
Data manipulation and transformation
Fact collection

**Model Monitoring**
Performance tracking
Issue management
Change requests

**Model Validation**
QA Testing
Performance checks
Accuracy, Bias, Drift, Explainability

**Model Retirement**
Justification
Documentation

**Model Deployment**
MLOps
Provisioning
Endpoint creation

**Model Approval**
Review validation report
Regulatory check
Final sign-off

# watsonx.governance

## AI Governance object model
*Much more than just "a model"*



Model Inventory and lifecycle Tracking

# watsonx.governance



**Performance Monitoring**
- Ongoing health monitoring of AI Models during runtime
- Trace and explain AI predictions
- Document metrics and track metric values over time
- Bias detection and mitigation
- Notification of issues when quality thresholds or business KPIs are violated

**Change/Issue Management**
-Automatic deployment and execution of validation tests for AI Models
- Track issues and incidents related to models in OpenPages included Issue Management Solution
- Workflow to document and approve changes to models

Evaluation and Monitoring

# Thank you