# Tridex – New York City

# Db2 12.1 – In the World of AI

Les King
lking@ca.ibm.com
December 11, 2024

# Agenda

- AI Concepts and Challenges
- IBM watsonx
- Db2 12.1.0 – Support for watsonx
- Db2 12.1.0 – Making Db2 a Smarter DMS with AI
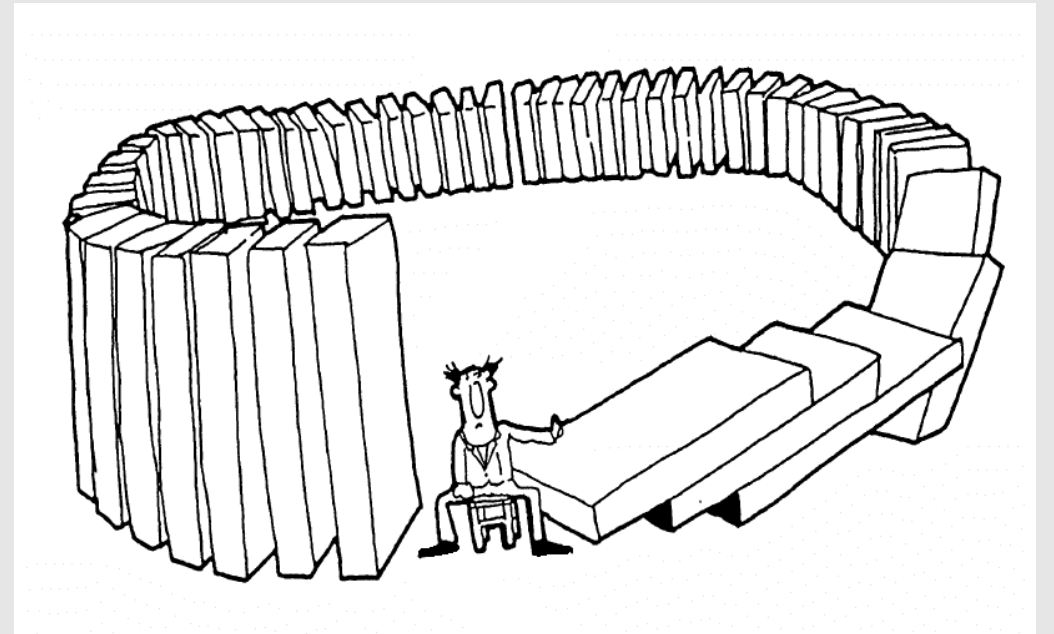- Db2 12.1.0 – Making Db2 a Great Data Store for AI Consumption
- Questions

# AI – The Concepts, Benefits and Challenges

# Technology

**Promise**



**Unintended Consequences**

# Social Media

# Social Media - Promise

- ***Promise:***
- *Bring us together in ways we could not imagine*

- ***Promise Realized:***
- *Communication Globalization*
- *News Reach*
- *Keeping in contact with family & friends*
- *Visiting parts of the world without leaving your living room*

# Social Media – Unintended Consequences

- *Are we more connected ?*
- *OR*
- *Just more connected to our devices ?*

# Social Media – Unintended Consequences
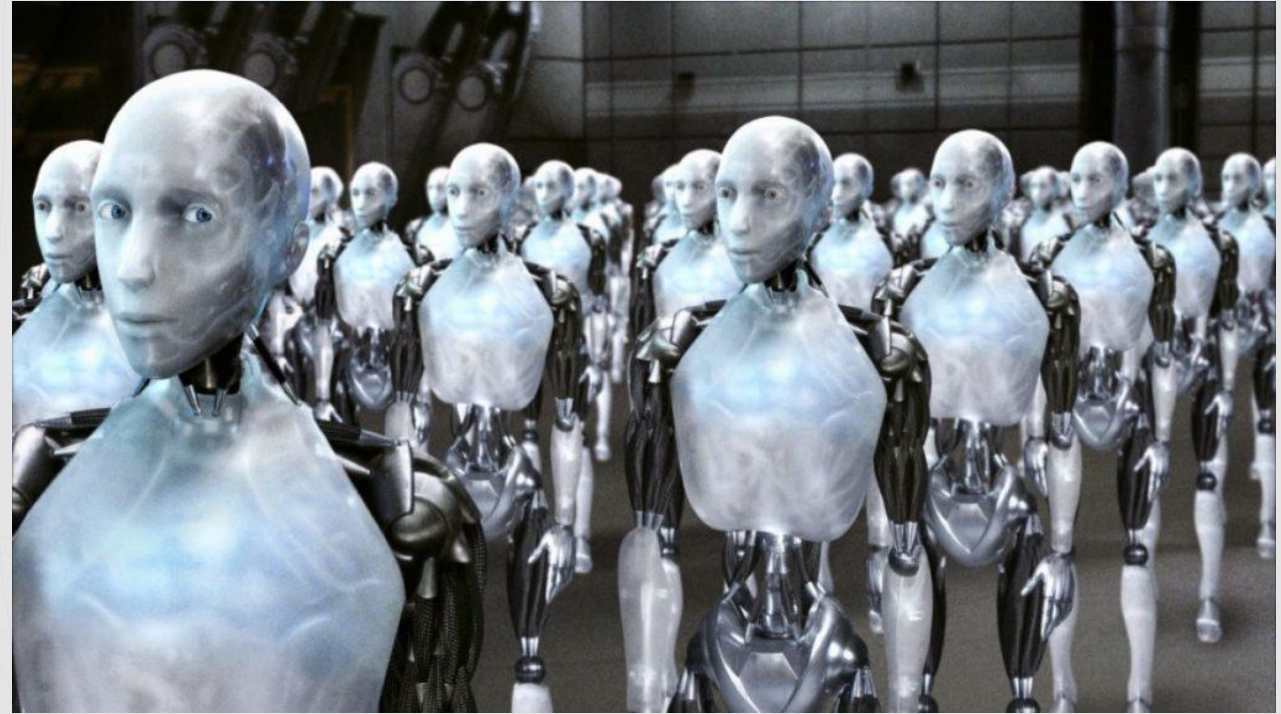
# Social Media – Unintended Consequences

# Social Media – Unintended Consequences
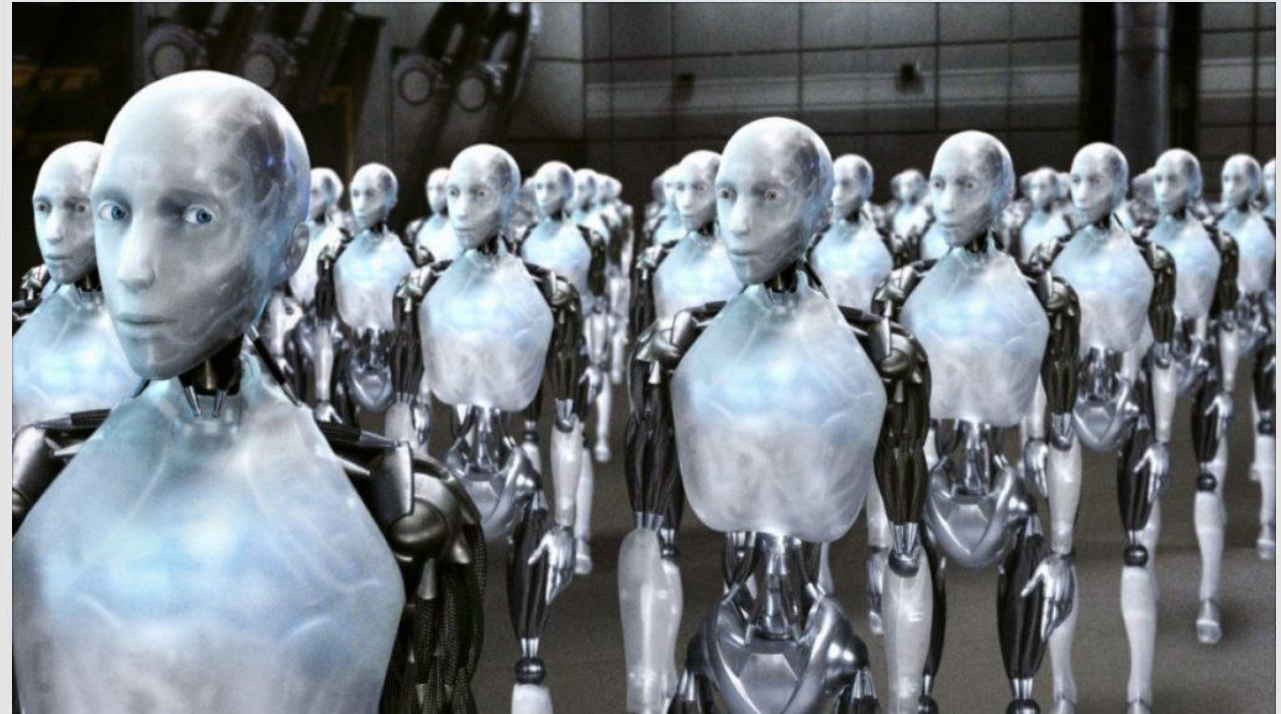
# AI

# AI - The Promise

# AI - Unintended Consequences

# AI - Unintended Consequences

**When has the created ever loved, adored and respected the creator ?**

**Have you ever been a parent and the proud owner of a teenager ?**
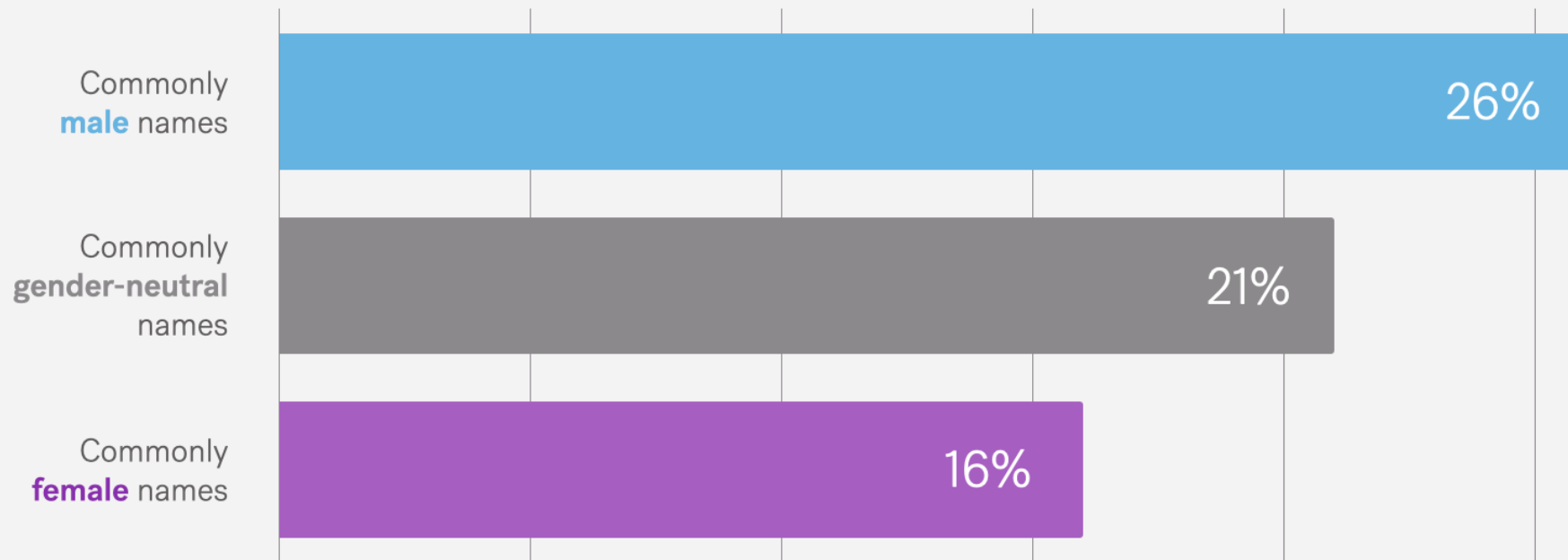
# AI – The Promise



- **Provides in-depth answers**
- **Strong analytics**
- **Translation**
- **Conversational**
- **Document Production**
- **Image Production**
- **Learn from mistakes**
- **…..**

# AI – Unintended Consequences - Bias

## ChatGPT favors male names when de-biasing job feedback

How often ChatGPT removes negative personality feedback when asked to remove bias

Commonly **male** names — 26%

Commonly **gender-neutral** names — 21%

Commonly **female** names — 16%

textio

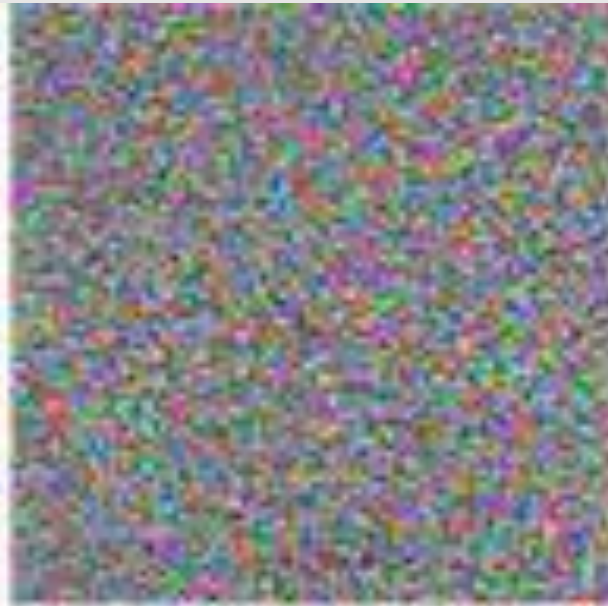# AI – Unintended Consequences - Hallucinations

# AI – Unintended Consequences – Adversarial Responses
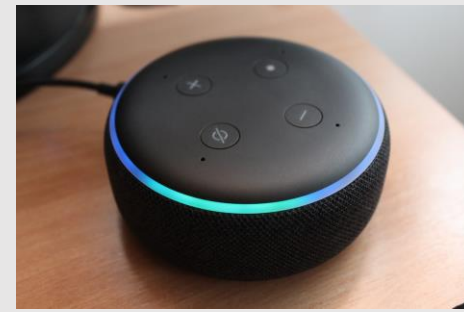


90% Tabby Cat + Adversarial noise = 100% Guacamole

# Amazon Alexa – Blooper Countdown
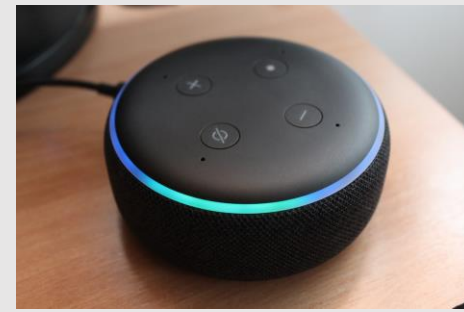
Alexa starts a party and cops are called

# Amazon Alexa – Blooper Countdown



Alexa starts a party and cops are called

Dollhouses and 2 kg of cookies purchased by children

# Amazon Alexa – Blooper Countdown

Alexa starts a party and cops are called

Dollhouses and 2 kg of cookies purchased by children

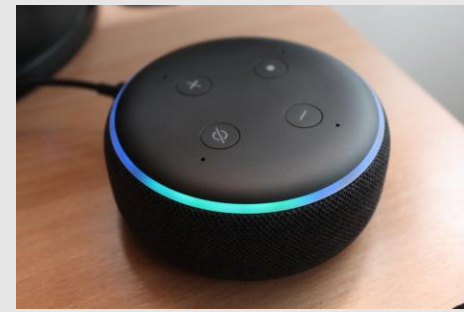Porn instead of children's song played when "Digger Digger" requested by a child

# Amazon Alexa – Blooper Countdown

Alexa starts a party and cops are called

Dollhouses and 2 kg of cookies purchased by children

Porn instead of children's song played when "Digger Digger" requested by a child

Bias is endless ….. Passport, World Cup, Beauty contest, Political …

FIFA WOMEN'S WORLD CUP

AU NZ 20 23 ™

# Building blocks of an AI Strategy

ChatGPT
inspired interest...

But there is a
bigger concept...

Which will
change business

## Large language model



Great
at text

## Foundation model

↑

### Transformer



Unlabeled
data

## Generative AI



Anything
that creates
new content

# Building blocks of an AI Strategy

ChatGPT
inspired interest...

But there is a
bigger concept...

Which will
change business

## Large language model

Great
at text

## Foundation model

### Transformer

Unlabeled
data

## Generative AI

Anything
that creates
new content

A **large language model (LLM)** is a type of **machine learning model** that has been trained on **large quantities** of unlabeled text using self-supervised learning and can perform a variety of natural language processing (NLP) tasks (even when that language is a programming language). Output may range from books, articles, social media posts, online conversations, and even code. The architecture of an LLM consists of layers of **neural networks** that learn to generate language in a way that is similar to how humans use language
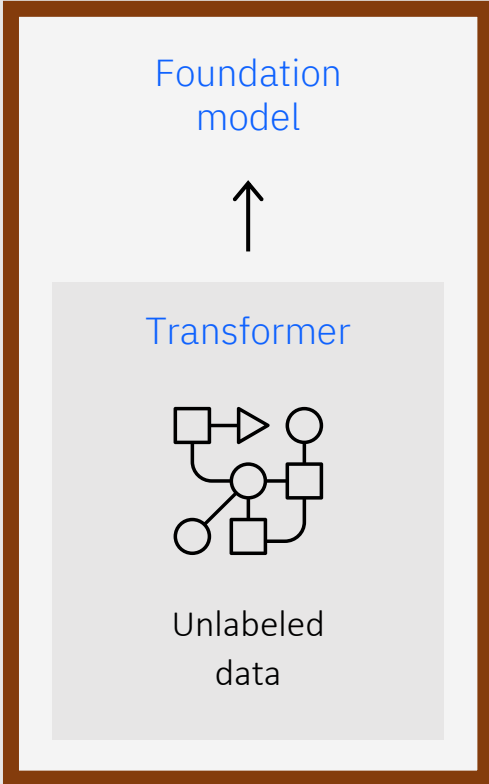
# Building blocks of an AI Strategy

ChatGPT
inspired interest...

But there is a
bigger concept...

Which will
change business

### Large language model

Great
at text

### Foundation
model

### Transformer

Unlabeled
data

### Generative AI

Anything
that creates
new content

A **Foundation models** are typically built using a specific kind of neural network architecture, called a transformer, which is designed to generate sequences of related data elements (for example, like a sentence).

A **transformer model** is a neural network architecture useful for understanding language, which does not have to understand words one at a time but can look at an entire sentence at once for context and disambiguation.
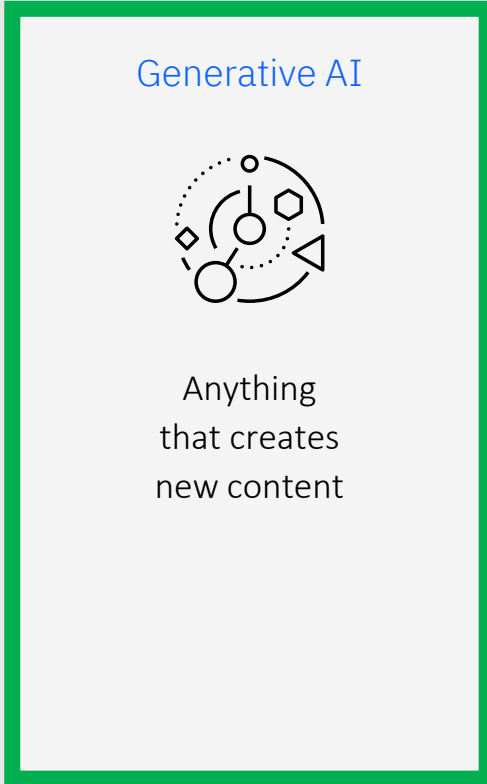
# Building blocks of an AI Strategy

ChatGPT
inspired interest...

But there is a
bigger concept...

Which will
change business

## Large language model
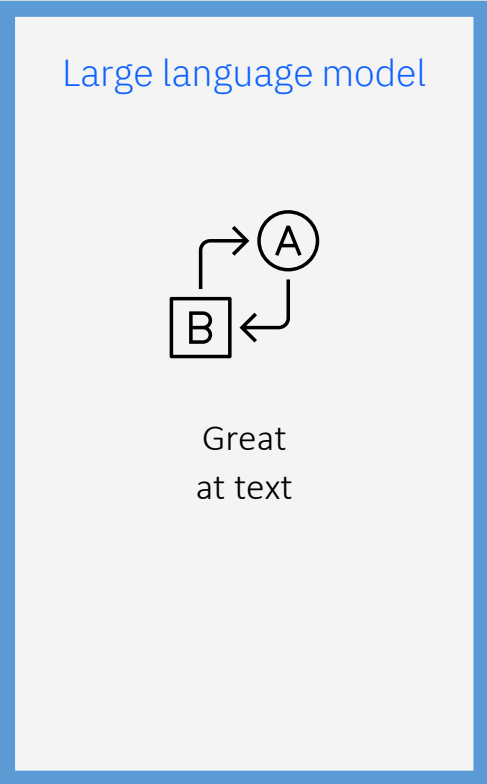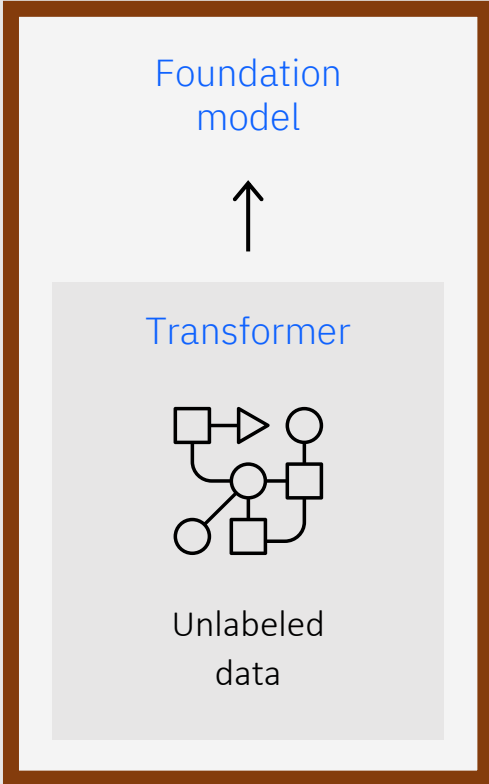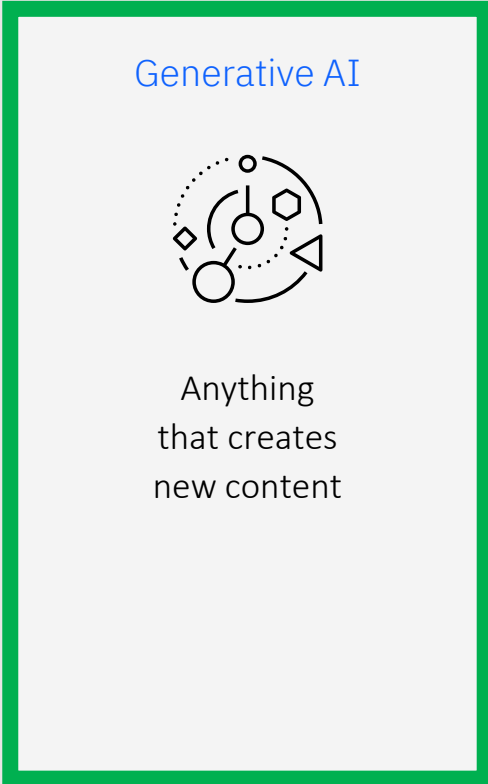
Great
at text

## Foundation model

### Transformer

Unlabeled
data

## Generative AI

Anything
that creates
new content

**Generative AI** refers to a set of AI algorithms that can generate new outputs — such as text, images, code, or audio — based on the training data, unlike traditional AI systems that are designed to recognize patterns and make predictions. Sometimes the AI that powers these solutions are referred to as decoders.

# Incredible opportunities around enterprise data

# Incredible opportunities around enterprise data



Physical Asset Management

Sensor data

Chemistry & materials

IT Ops
IT Automation

IT data

Geospatial

Sustainability

Foundation models

Business Automation

Business data

Programming languages (code)

LLMs

Application Modernization
IT Modernization

LLMs

Speech

Cyber Security Data

Threat Management

Customer Care
Digital Labor

Natural language

Dialog

LLMs

LLMs

LLMs

# The modern-day AI ladder

**AI+**

AI does the work

Replace your workflows

Automate your workflows

Add AI to your applications

**+AI**

Collect, organize, grow data

Infuse

Analyze

Organize

Collect

# AI governance is needed to manage risk and protect reputations

- *"Fewer than 20% of executives strongly agree that their organizations' practices and actions on AI ethics match (or exceed) their stated principles and values."*
- - IBM and Oxford Economics – AI ethics in action, 2021



## Bloomberg

### IDEAS

## Algorithmic bias isn't just unfair — it's bad for business

If it's not deployed wisely, artificial intelligence can turn consumers off.

By **Kalinda Ukanwa** Updated May 23, 2021, 3:00 a.m.

**YouTube sued for using AI to racially profile content creators**

They claim YouTube's algorithms discriminate against black users

### Data science during COVID-19: Some reassembly required

Most likely, the assumptions behind your data science model or the patterns in your data did not survive the coronavirus pandemic. Here's how to address the challenges of model drift

**Amazon scraps secret AI recruiting tool that showed bias against women**

## The $300m flip flop: how real-estate site Zillow's side hustle went badly wrong

### The Washington Post
*Democracy Dies in Darkness*

## Apple Card algorithm sparks gender bias allegations against Goldman Sachs

# Constantly growing and changing regulations drive the need for governance

**United Kingdom** - AI Regulation Policy

**Canada** - Bill C27: AI and Data Act

**Fine**: $25M or 5% of company's gross global revenue

**Germany** - AI Strategy
**South Korea** - AI Strategy
**India** – National Strategy for AI

**UAE** - Strategy for AI

**Australia** - AI Ethics Framework

**United States –** National Artificial Intelligence Initiative

**New York City** - AI Hiring Law

Pan-**Canadian** AI Strategy
**Japan** - AI Technology Strategy

**Norway** – National AI Strategy
**Serbia** – Strategy for the development of AI

**United States** - NIST issues an AI risk management framework

**United States** - AI Bill of Rights. Validation for algorithms to be explainable and protect against discrimination

**Colombia** - National Policy for Digital Transformation and AI

**European Union** - The AI Act

**Fine**: €30M or 6% of company's global revenue

**European Union** - GDPR

**Fine**: €20M or 4% of company's annual turnover

**China** - Internet information service algorithm recommendation management regulations

**Singapore** - Launches AI Verify: An AI Governance Testing Framework and Toolkit

2017　　2018　　2019　　2020　　2021　　2022　　2023

**NAIC** National Association of Insurance Commissioners

Gramm-Leach-Bliley Act COMPLIANT

European Commission

NATIONAL ARTIFICIAL INTELLIGENCE INITIATIVE OFFICE

Sarbanes–Oxley Act

UNITED STATES FEDERAL RESERVE SYSTEM

**NIST** National Institute of Standards and Technology

# Why should organizations that build or use AI care about ethics?

- Company values

- Company reputation

- Social justice and equity

- Client and investor inquiries

- Differentiation

- Business opportunities

- Existing or expected regulations

# Pillars of trust

The purpose of AI is to augment human intelligence

Data & Insights belong to their creator

- **Explainability**

  - An AI system's ability to provide a human-interpretable explanation for its predictions and insights.

- **Fairness**

  - An AI system's ability to treat individuals or groups equitably, depending on the context in which the AI system is used.

**Robustness**

An AI system's ability to effectively handle exceptional conditions, such as abnormalities in input.

- **Transparency**

  - An AI system's ability to include and share information on how it has been designed and developed.

- **Privacy**

  - An AI system's ability to prioritize and safeguard consumers' privacy and data rights.

Your AI is only as good as your data.

# IBM watsonx

# The platform for AI and data

## watsonx

Scale and accelerate the impact of AI with trusted data.

## watsonx.ai

Train, validate, tune and deploy AI models

- A next generation enterprise studio for AI builders to 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 optimized for governed data and AI workloads, supported by querying, governance and open data formats to access and share data.

## watsonx.governance

Enable responsible, transparent and explainable data and AI workflows

- End-to-end toolkit encompassing both data and AI governance to enable responsible, transparent, and explainable AI workflows.

# IBM's AI is embedded in applications built on

**watsonx**

## Watson Orchestrate

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

### 40%
Improvement in HR productivity

## Watson Assistant

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

### 70%
Call center calls contained by conversational AI

## Watson Code Assistant

Enables hybrid cloud developers to write code with AI-generated recommendations

### 30%
Productivity gain in application modernization
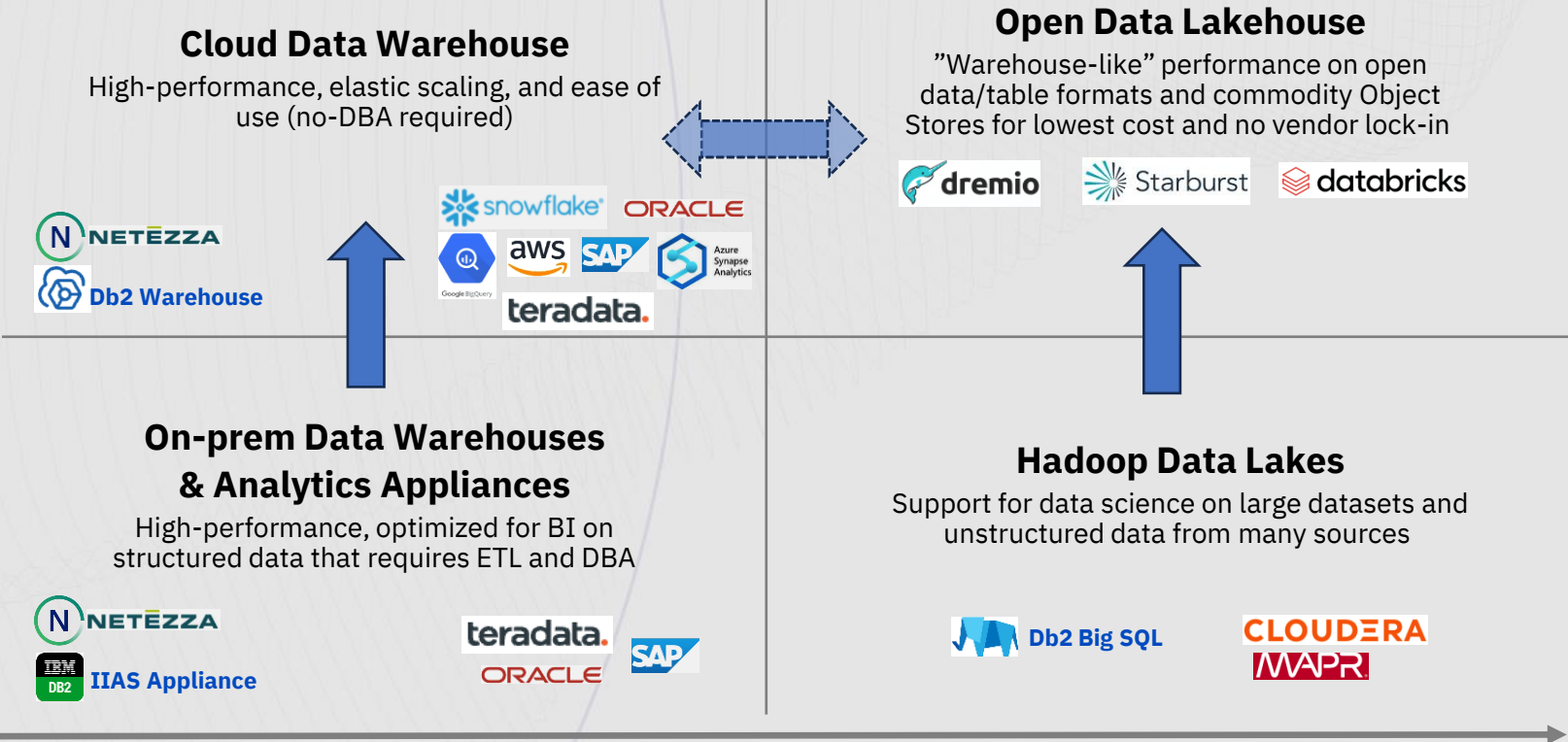
AI and data platform | **watsonx**

# Market Dynamics

Major disruptions are driving the growth in the analytics repositories market **from on-prem to SaaS** and **from proprietary to open technologies**

**Analytics Repositories Market Landscape**

**SaaS**
$31bn 2025
27% CAGR ('21-'25)

**Deployment**

**On-prem**
$12bn 2025
2% CAGR ('21-'25)

**Cloud Data Warehouse**
High-performance, elastic scaling, and ease of use (no-DBA required)

**Open Data Lakehouse**
"Warehouse-like" performance on open data/table formats and commodity Object Stores for lowest cost and no vendor lock-in

**On-prem Data Warehouses & Analytics Appliances**
High-performance, optimized for BI on structured data that requires ETL and DBA

**Hadoop Data Lakes**
Support for data science on large datasets and unstructured data from many sources

**Proprietary**
$26bn 2025
13% CAGR ('21-'25)
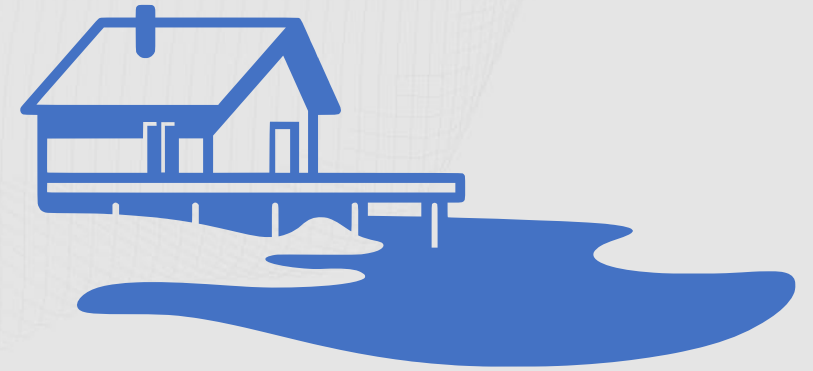
**Technology**

**Open**
$17bn 2025
27% CAGR ('21-'25)

*Sources: IDC Data Management Forecast (November 2021), IDC BDA Forecast (June 2021), MI modeling*
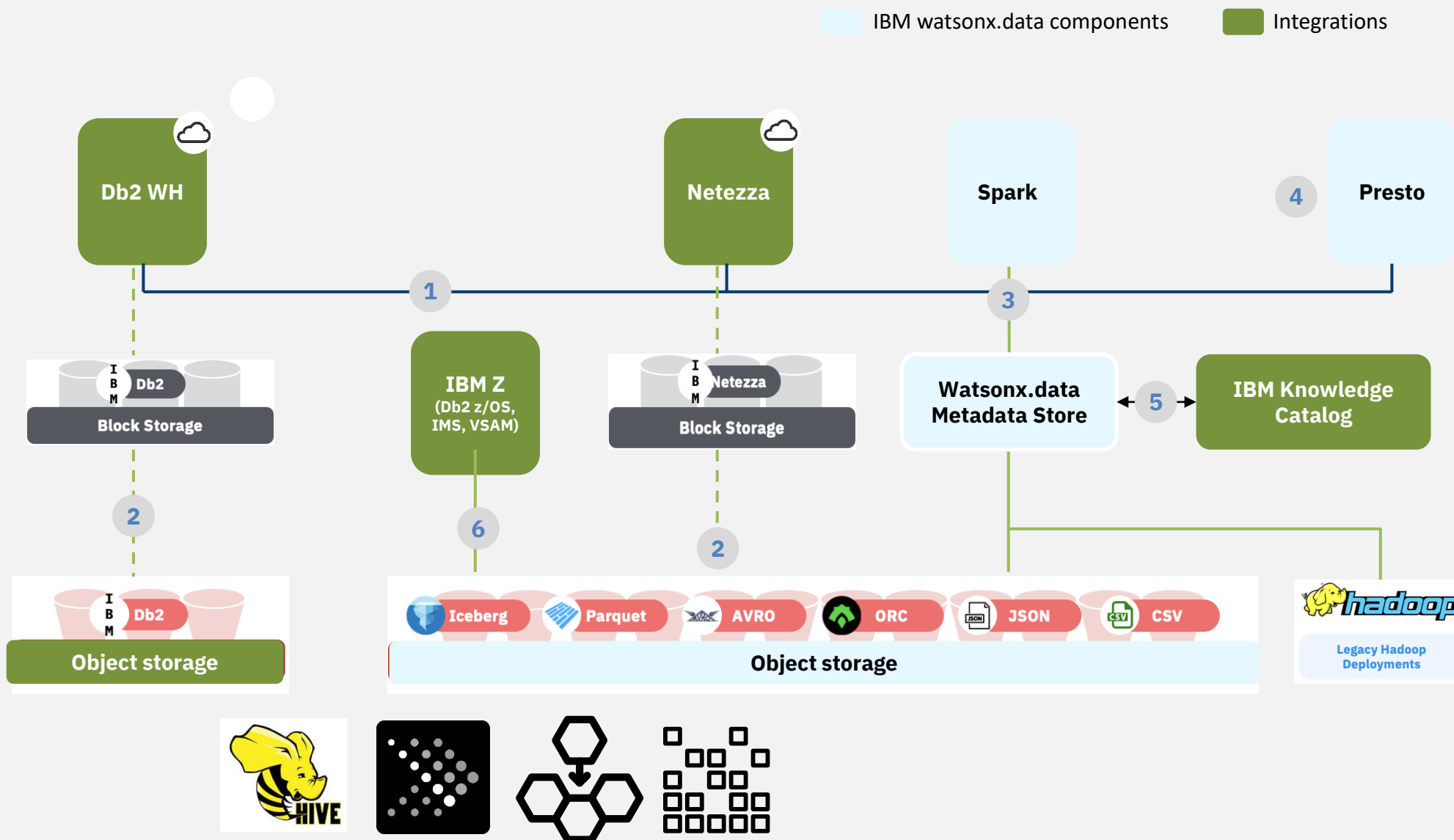
38

# The Data Lakehouse

The Data Lakehouse implements the **data structures and management features** of a data warehouses on the **low cost, reliable & scalable** object storage within a new architectural approach that leverages open-source technology.

It enables organizations to manage their data in an **open, flexible, cost-effective, feature rich and scalable way**, enabling Business Intelligence and Machine Learning on all data.

Data *Lake*
+
Ware*house*

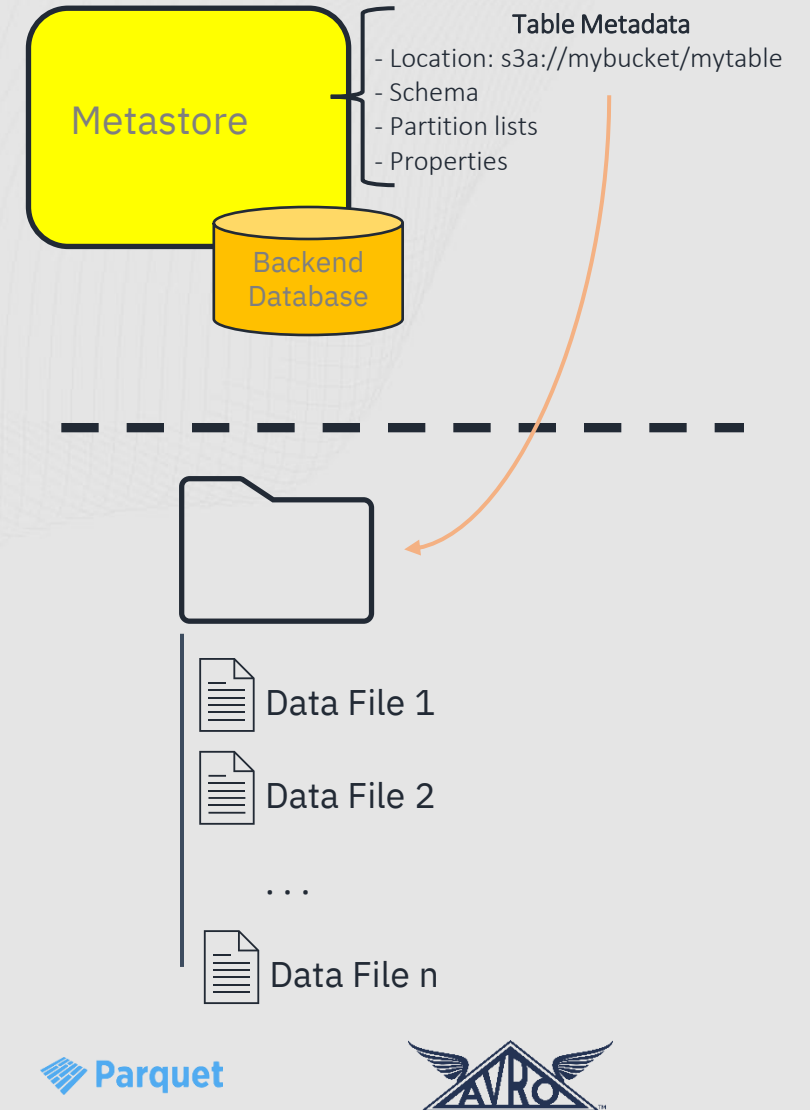# The integrated IBM watsonx.data ecosystem for maximum workload coverage and optimal price-performance

IBM watsonx.data components | Integrations

**Db2 WH**

**Netezza**

**Spark**

(4) **Presto**

(1)

(3)

IBM Db2
Block Storage

**IBM Z**
(Db2 z/OS, IMS, VSAM)

IBM Netezza
Block Storage

**Watsonx.data Metadata Store** ↔ (5) **IBM Knowledge Catalog**

(2)

(6)

(2)

IBM Db2
**Object storage**

Iceberg | Parquet | AVRO | ORC | JSON | CSV
**Object storage**

**hadoop**
Legacy Hadoop Deployments

---

1. Warehouses can **access data in the lakehouse**

2. **Easily Promote data** between the warehouse and lakehouse

3. Multiple engines can **access same data lake data**

4. The lakehouse can **access data residing in Db2, Netezza, and other data sources**

5. IKC policies enforced by the lakehouse

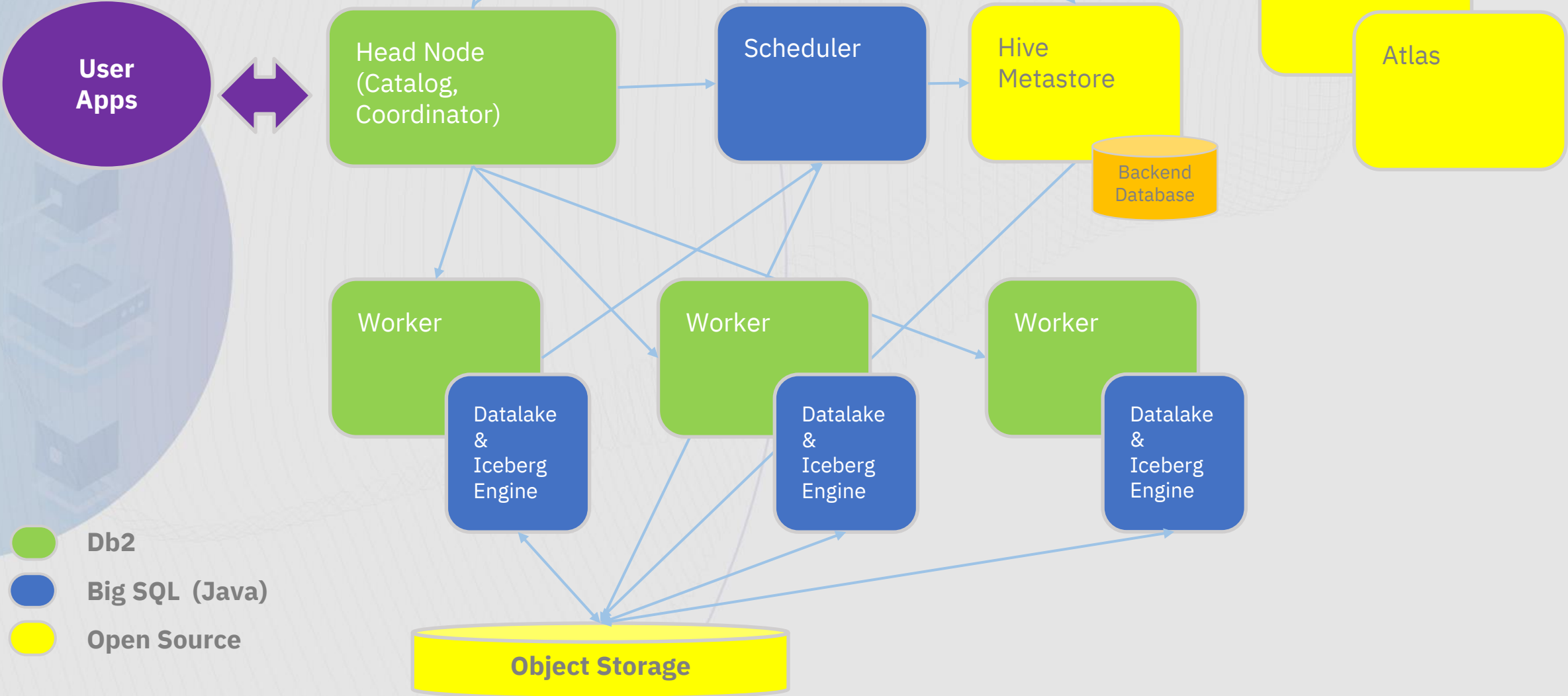6. **Analyze Z data easily and securely** by writing to Iceberg tables with **Data Gate for watsonx**

# Db2 12.1.0 - Enhancements for watsonx.data / Lakehouse Support
# (actually delivered in Db2 11.5.9)

# DATALAKE Tables

- A Data Lake "Table" is a collection of files serialized using an **Open Data File** (ODF) format (CSV, ORC, Avro, Parquet ...) stored on remote storage (HDFS, S3, COS, ...)

- The **metadata** of the table is stored in a Metastore server
  - Location
  - Schema
  - Partition lists

- An engine querying the table must query the metadata first and can proceed to read the data from remote storage

- Benefits
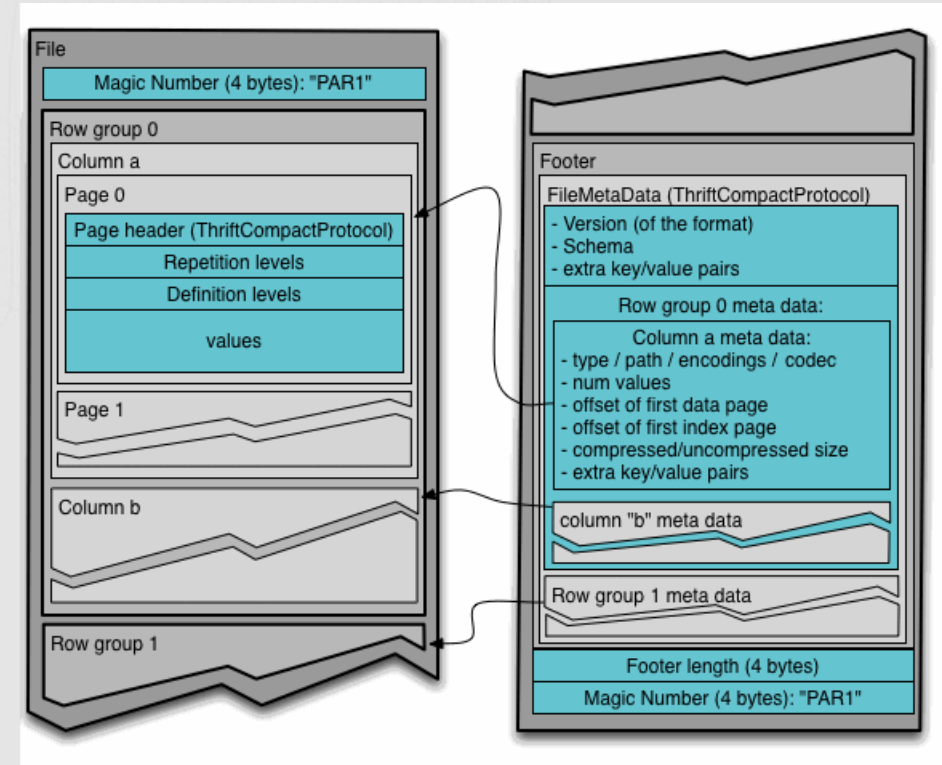  - Interoperability of open data formats
  - Ease of use

Metastore

Backend Database

**Table Metadata**
- Location: s3a://mybucket/mytable
- Schema
- Partition lists
- Properties

Data File 1

Data File 2

. . .

Data File n

Parquet

AVRO

Apache ORC

# IBM Db2 Warehouse
## DATALAKE Table Support

**User Apps**

Head Node (Catalog, Coordinator)

Scheduler

Hive Metastore

Backend Database

Ranger

Atlas

Worker

Datalake & Iceberg Engine

Worker

Datalake & Iceberg Engine

Worker

Datalake & Iceberg Engine

**Object Storage**

- Db2
- Big SQL (Java)
- Open Source

# Open Data File Format Limitations

- A DATALAKE "Table" is a collection of files serialized following an Open Data File format

- Passive data structures – serialized set of data records
  - No notions of their **state** or **history**
  - No concurrency control between applications
  - **No ACID**, even less transactions

- Separate metadata
  - Need for a "**Catalog**"
  - No awareness of catalogs – it's an external system

# Apache Iceberg
## An Open Data Table format for the Lakehouse

ICEBERG

Full **open-source**, **Open Data Table format**, quickly becoming an **industry standard**

Relies on Open Data File formats for storage, but provides an additional layer of metadata for enhanced capabilities
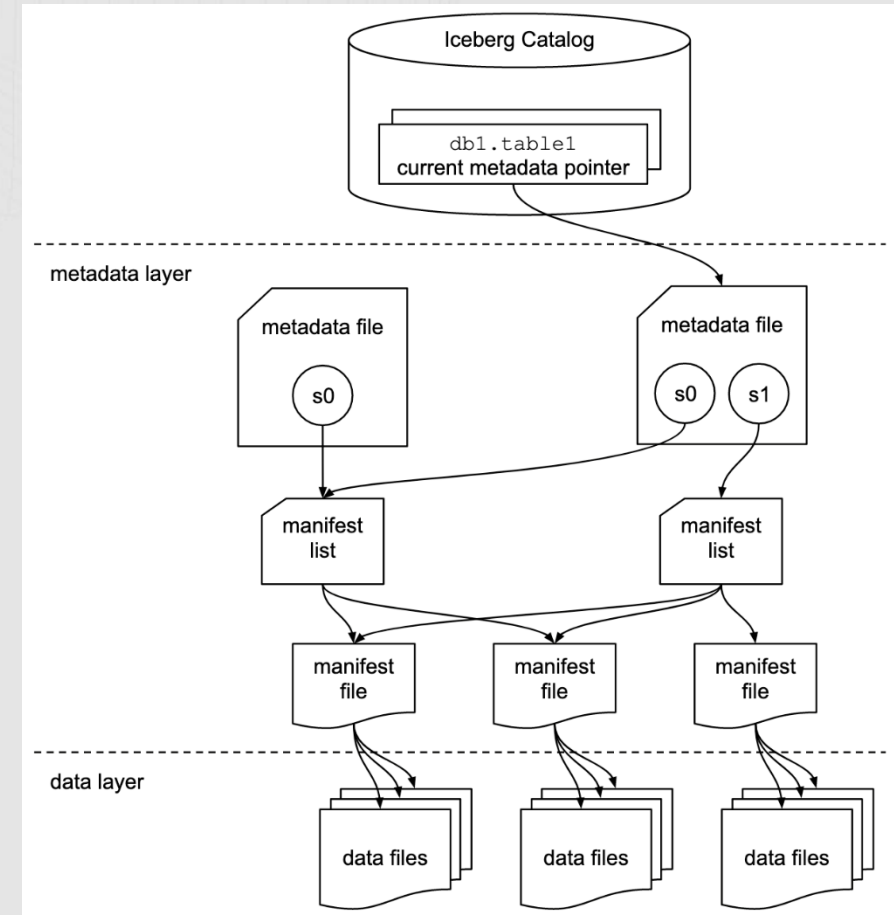
Support for CREATE, SELECT & INSERT including partitioning support

No UPDATE, DELETE
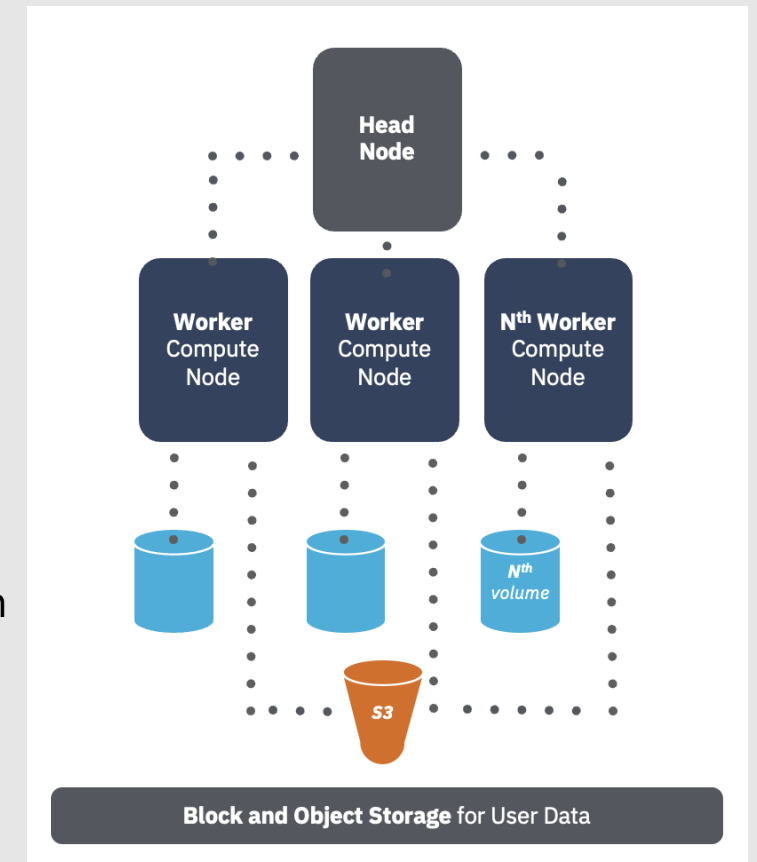
No Scheme Evolution

No Time Travel

Smaller restrictions related to Icebert/Db2 type compatibility such as nested types, etc.

# Native Cloud Object Storage Support – Remote Tablespace
## Key Attributes

- Significant **storage cost savings** by using object storage instead of block storage.

- **Faster query and ingest performance** through the new multi-tier storage engine.

- **Consumption-based model** for the storage, with all the benefits of automatic and **unlimited storage scaling**.

- Data can reside on **block storage or object storage**, based on business or technical requirements.

- **No applications and workload changes** necessary.

  - IUD and move data as needed into and out of tables in object storage.
  - Query data seamlessly no matter where it resides (in block or object storage), in isolation or in combination with each other.

- Enables new use cases:
  - **MQT cache** for accelerating queries over Datalake tables.
  - Cost-efficient **high-volume streaming** into native tables.

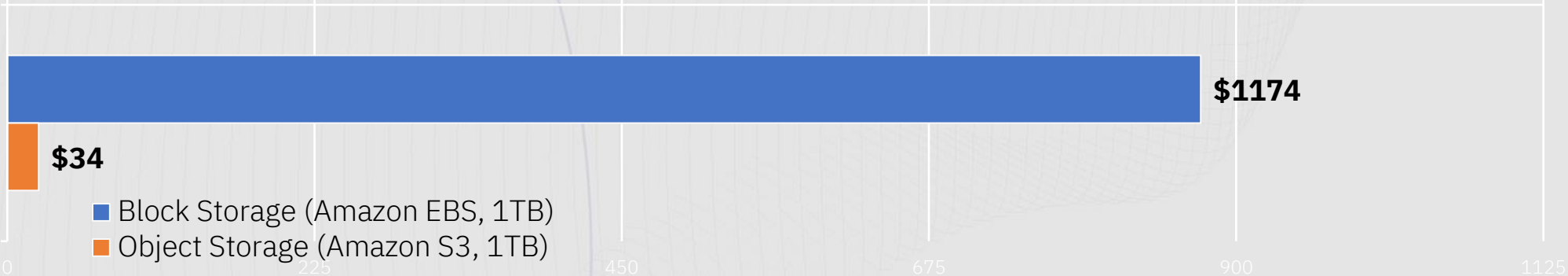- Available in Db2WHoC Gen 3 and Db2U containerized environments.

# Remote Tablespace Support - Storage Savings

## Db2 Warehouse current generation vs Gen3

**34x**

**Less expensive to host Db2 data on object vs block storage[1]**

$1174

$34

■ Block Storage (Amazon EBS, 1TB)
■ Object Storage (Amazon S3, 1TB)

0    225    450    675    900    1125

**Block Storage (Amazon EBS) vs Object Storage (Amazon S3)
Cost reflects Amazon's list price for block storage (various tiers
& IOPS levels) required to host an incremental 1TB of Db2 data**

# Db2 12.1.0 – Support for Modern AI Workloads

# Db2 – Ready for Modern AI Workloads

**Data Virtualization**

Db2 contains a data virtualization component which allows Db2 to be a doorway to all of your business data

- Relational Sources
- Cloud Sources
- Open Source Sources
- NoSQL Sources

- Native Clients
- ODBC, JDBC, REST, NoSQL
- Pushdown Performance
- In-memory MQT

**In-Db2 Machine Learning**

Allows data scientists and developers to bring machine learning local to the data stored within Db2

- Data Exploration
- Model Training
- Model Evaluation
- Model Deployment

- Data Preprocessing
- Inferencing
- Error Detection
- Support for many models

**Multi-Model – NoSQL and NewSQL Data Store**

Db2 is a multi-model data store supporting native relational, JSON, BSON, Graph, Spatial, Text and XML

- Vector
- XML
- Spatial
- Text

- JSON/BSON
- XQuery/Mongo/FLWOR
- ESRI
- ACID Properties

**Mixed Workloads**

Db2 can handle any combination of workloads including real-time data ingestion, multi-model and mixed.

- ML Optimizer
- ML Memory Management
- CDI (Trickle-feed)
- Access multi models

- Access Remote sources
- HTAP
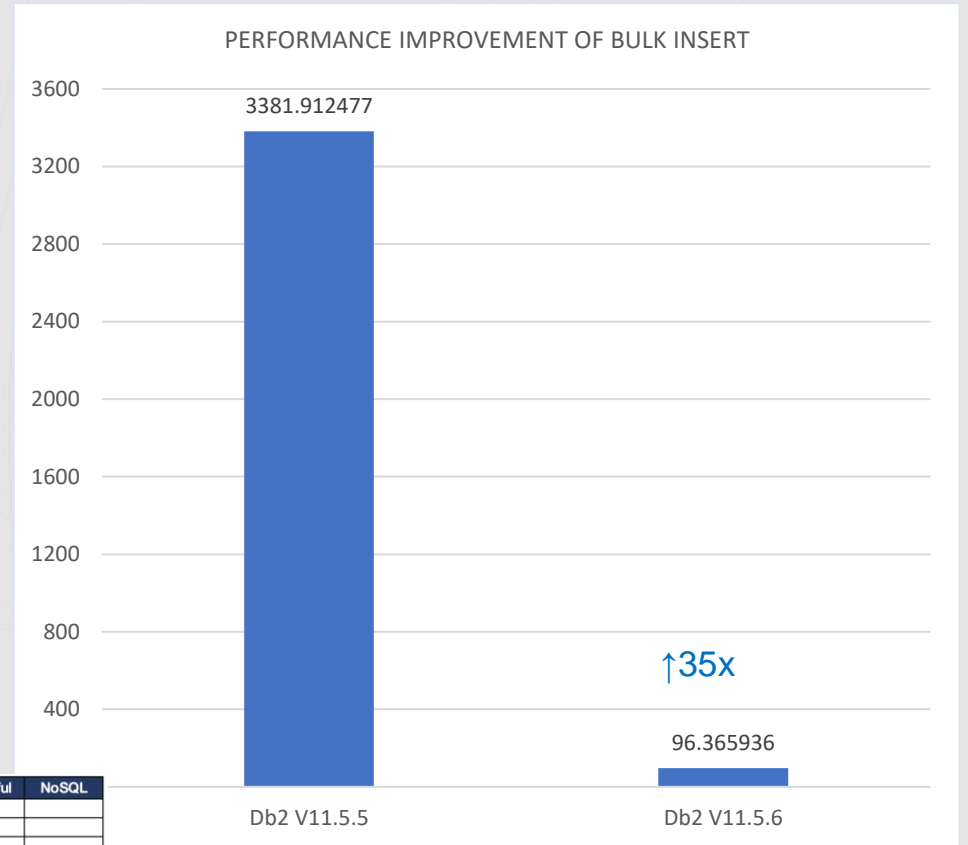- OLTP + OA + Reporting
- OLAP (All Combinations)

# Db2 12.1.0 – Data Virtualization

# Data Virtualization – Federation

- Connectivity – Spark JDBC Connectivity Support

- Functionality – Column Length Variation for Code Page Conversion

- Functionality – Nickname Hidden Column Support

- Performance – Federation DRDA Bulk Insert for Db2 Family Data Sources

## PERFORMANCE IMPROVEMENT OF BULK INSERT

| | Value |
|---|---|
| Db2 V11.5.5 | 3381.912477 |
| Db2 V11.5.6 | 96.365936 |

↑35x

| Category | Data Source | Native | ODBC | JDBC | RESTful | NoSQL |
|---|---|---|---|---|---|---|
| Relational | Db2 LUW | Yes | | Yes | | |
| | Db2 for IBM i | Yes | | | | |
| | Oracle | Yes | Yes | Yes | | |
| | MS SQL Server | Yes | Yes | Yes | | |
| | Informix | Yes | | | | |
| | Sybase | Yes | | | | |
| Warehouse / Appliance | IIAS | Yes | | Yes | | |
| | Netezza | | Yes | Yes | | |
| | Teradata | Yes | | Yes | | |
| | SAP HANA | | Yes | Yes | | |
| | Greenplum | | Yes | Yes | | |
| Open Source | MySQL Community | | Yes | Yes | | |
| | MySQL Enterprise | | Yes | Yes | | |
| | PostgreSQL | | Yes | Yes | | |
| | MariaDB | | Yes | Yes | | |
| | Derby | | | Yes | | |
| Hadoop | IBM Db2 BigSQL | Yes | | Yes | | |
| | Hive | | Yes | Yes | | |
| | Spark | | Yes | Yes | | |
| | Impala | | Yes | | | |
| Files | Delimited | Yes | | | | |
| | Excel | Yes | Yes | | | |
| | XML | Yes | | | | |
| | JSON | | | | | Yes |
| | CSV | Yes | | | | |
| Mainframe | Db2 for z/OS | Yes | | Yes | | |
| | IBM DVM for z/OS | | | Yes | | |

| Category | Data Source | Native | ODBC | JDBC | RESTful | NoSQL |
|---|---|---|---|---|---|---|
| Message Queue | MQSeries | Yes | | | | |
| Cloud | Db2 Warehouse | Yes | | Yes | | |
| | MS Azure SQL | | Yes | | | |
| | Oracle Cloud | | Yes | | | |
| | Amazon AWS Redshift | | | Yes | | |
| | Google BigQuery | | | Yes | | |
| | Amazon AWS S3 | | | Yes | | |
| | Salesforce | | | Yes | | |
| | Snowflake | | Yes | Yes | | |
| NoSQL | Hyperledger Fabric | | | | | Yes |
| | MongoDB | | | | | Yes |
| | CouchDB | | | | | Yes |
| | Hbase HDFS | | | | | Yes |
| | Cassandra | | | | | Planning |
| | Redis | | | | | Planning |
| | Jira | | | | Yes | |
| | Aha! | | | | Yes | |
| | GitHub | | | | Yes | |
| | HubSpot | | | | Yes | |
| | TeamCity | | | | Yes | |
| | api.spacexdata.com | | | | Yes | |
| | earthquake.usgs.gov | | | | Yes | |
| | Google Calendar API | | | | Yes | |
| | groupkt.com | | | | Yes | |
| | Yelp | | | | Yes | |

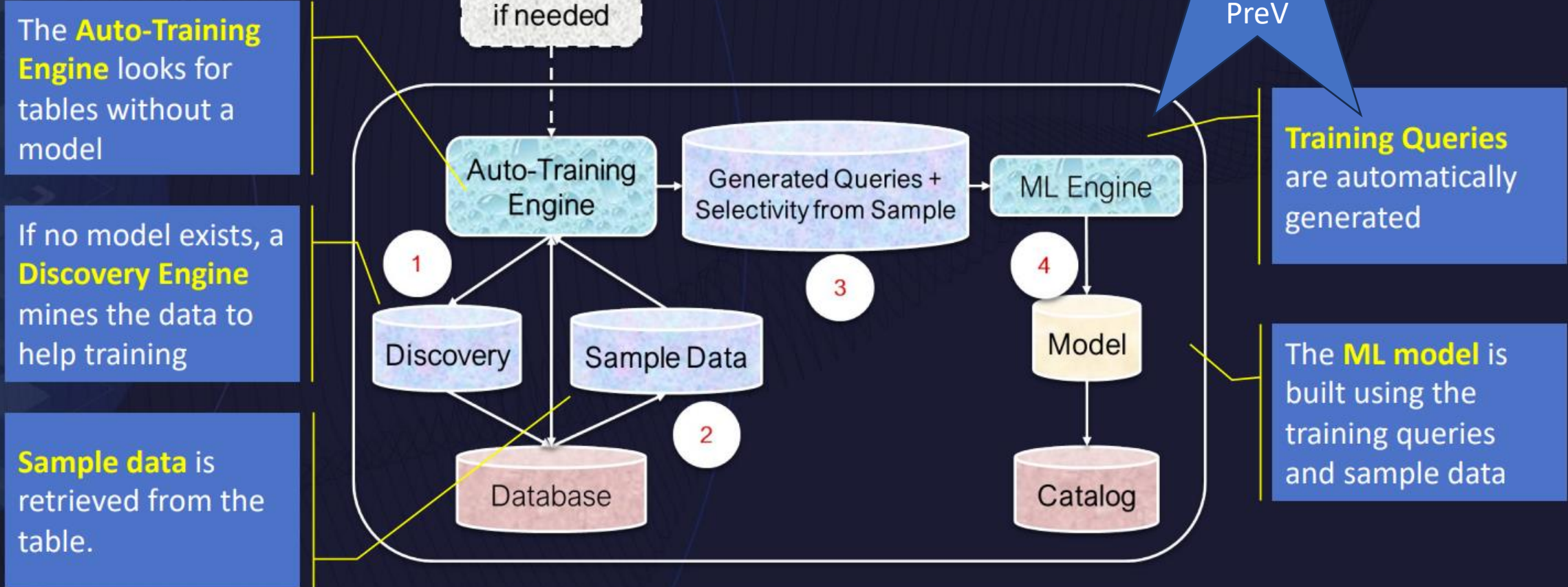| |
|---|
| Supported Before v10.5 |
| Supported In v11.1 |
| Supported In v11.5 GA |
| Supported in v11.5.4 |
| Supported in v11.5.5 |
| Supported in v11.5.6 |
| Working / Planning |

# Db2 12.1.0 – Multi-Model Support

# Db2's Multi-Model Support

- **Storage** – natively storing the data to ensure no loss of data – no force fitting into row/column structure (ie: shredding)

- **Performance** – ability to index the data in a meaningful way to provide tier 1 performance for both ingestion and queries

- **Integration** – ability to query data in each model of data – within the same query

- **SQL Support** – ability to work with the model of data using SQL

- **NoSQL Support** – ability to work with the model of data using a natural query language for that model of data

- **NewSQL Support** – support for transactional awareness (ACID properties) when using that model of data

- **Enterprise Requirements** – leverage Db2's availability, security, recoverability, etc for that model of data

- **Output** – ability to decide between traditional relational result set of model specific output

# Db2 12.1.0 – AI Optimizer

# Tech Preview – Automatic Training



The **Auto-Training Engine** looks for tables without a model

If no model exists, a **Discovery Engine** mines the data to help training

**Sample data** is retrieved from the table.

Manual call if needed

Auto-Training Engine

Generated Queries + Selectivity from Sample

ML Engine

Discovery

Sample Data

Database

Model

Catalog

11.5 Tech PreV

**Training Queries** are automatically generated

The **ML model** is built using the training queries and sample data

## AI Optimizer highlights for GA

**12.1 GA**

- Infrastructure for future AI models for use within Db2

- Significantly improved local predicate cardinality estimation

- Possible pairwise join cardinality estimation using the single table model

# Db2 12.1.0 – AI DB Assistant

**Sergio,**
Database
Administrator

# Pain Points

## Finding the right information / documentation

*"I search for documentation daily, sometimes hourly. IBM documentation can feel like boiling the ocean. I use Google."*

## Training junior DBAs

*"Staffing is an issue. Workloads are increasing. A value add would be a way to handle increased workloads/mixed workloads and not have to increase staffing."*

## Identifying and resolving the root problem

*"I have a lot of information, but I don't know what is relevant to my current issue."*

## Optimizing + tuning the database

*"Optimizing performance is complex and requires expertise. The current Tuning UX in DMC is complicated and not ideal."*

# Introducing Database Assistant
powered by **watsonx**

**Db2 Expert**

Get answers to your Db2 questions, faster

**Monitoring Metrics**

Quickly access key Db2 metrics using natural language queries

**Simplified Troubleshooting**

Get recommendations for troubleshooting common database issues

**Advanced Analysis**

Accurately identify root cause of performance issues, bottlenecks, deadlocks

*Coming Soon!*

# Simplify the process of *navigating through multiplicity of database tasks* through

AI assisted navigation of basic database tasks such as:
- *Viewing database summary information*
- *Listing tables + schemas + indexes*
- *Checking active resource usage (CPU, IO, Memory)*
- *Checking storage utilization*
- *Viewing active sessions*
- *Viewing active queries*
- *Analyzing where time is being spent*
- *Analyzing lock waits*
- *Analyzing Top N queries + connections*

*For a non-expert user, tasks normally involves cross referencing public documentation with unguided adhoc navigation of the available console panels and following a multi-step process to locate the required information, diagnose, and then resolve an issue.*

# Benefits

- Easy access to targeted grounded answers for technical questions.

- Reduce context switching and switching between different tools while diagnosing and fixing issues.

- AI guided tasks and troubleshooting to streamline the DBA's job.

- The Assistant is developed using a RAG based AI system to minimize hallucinations by retrieving information from trusted sources.

Hi! I'm a Database assistant. How can I help you today?

# Where does it operate?

Database Assistant is built directly into your Database Management Console (DMC).

Database Assistant provides real-time metrics of your database instance.

# Components



Db2
Knowledge Base

Generative AI

Database
Assistant

Db2 Instance
Metrics

# Technical Architecture



**Vector Database**

IBM **Db2** Documentation

**Scrape, Chunk and Ingest**
- **Db2 Documentation**
- **Tech Notes**
- **Support Q&A**

② Query Knowledge Base

③ Retrieve Relevant Documents

① User Question

⑥ Live Monitoring Queries

**watsonx Orchestrate**

⑥ Assistant Answer

**Db2 SaaS Instance**

④ Documents + LLM Prompt

⑤ Generated Text Response

⑦ User Feedback

**IBM Granite LLM**

# Demo

# Db2 12.1.0 - In-Db2 Machine Learning

# Bring AI to where the data lives!

**Build and deploy AI models inside Db2 for**
- Classification
- Regression
- Self-supervised learning

**OR**

**Build models anywhere and deploy them on Db2:**
- Python models (e.g., Scikit-learn)
- R models

# 66%
ML projects use Relational data

# In-Db2 Machine Learning

Lifecyle Phases of an Enterprise ML System

Research

Production

**Experiment**

**Train**

**Deploy**

**Inference**



Best algorithm
search

Model trained using
the best algorithm

Trained model on
production

Model generates
predictions

**Data** – storage, regulations, scale, quality
**Model** – infrastructure, compute resources, latency, integration

# IBM Db2 Can Accelerate Implementing ML Systems

# Python UDF: Scoring Python Models via Db2

Open Source

Export the AI pipeline by serializing *python joblib or pickle*

Db2 Server

Host OS

Db2 Instance

Python Runtime

Check out Session 6 today for a live demo of this capability and more

# R UDF: Scoring R Models via Db2

Open Source

Export the AI pipeline and other deployment assets as *RDS files*

**Db2 Server**

Host OS

Db2 Instance

R Runtime

# In-Db2 Machine Learning

Train, Tune, Cleanse, Explore, Evaluate, Manage, Error Detection, Inferencing

Integrated Python and R Library for exploring and manipulating data

Accelerated and Distributed Machine Learning Algorithms in Db2



**Random Forest**

**Decision Tree**

**Logistic Regression**

# Db2 12.1.0 – Wrap-Up

# Db2 – Handling Modern Workloads

## Powered by AI

Confidence-based query results
leveraging ML-SQL

Up to 10x better query performance
powered by an ML-Optimizer

No data movement & single view
on all data
delivered by Data Virtualization

Auto resource optimization
delivered by Adaptive Workload Management

## Built for AI

Faster data exploration
by using In-Db2 Machine Learning

Build AI based applications
with Python, R, GO , JSON and Jupyter notebooks

Model Complex Relationships
by using  Db2's Multi-Model Capabilities

Blockchain Ready
using Db2 Blockchain Connector

# Thank You

Speaker:  Les King

Company:  IBM

Email Address: lking@ca.ibm.com