# *Generative AI in Data Management and Analytics*
*– A New Era of Assistance, Productivity and Automation*

Mike Ferguson

CEO, Intelligent Business Strategies

Adept Events Data Warehousing & BI Summit

Utrecht, March 2024

𝕏 *@mikeferguson1*

---

## Who Is Mike Ferguson? – A Leading Analyst In Data Management & Analytics

**Codd & Date**

Edgar F Codd

Turing
Award
Winner

Relational Model

All relational DBMSs | SQL Language | Data Normalisation

Mike Ferguson

Chief Architect
**teradata.**
1st massively parallel relational DBMS

Founded 1992 — INTELLIGENT BUSINESS STRATEGIES

C J Date

An Introduction to Database Systems

RELATIONAL DATABASE SELECTED WRITINGS — C. J. DATE

RELATIONAL DATABASE WRITINGS 1985-1989 — C. J. DATE

MASTER CLASS — Normal Forms & All That Jazz — C.J. Date — VIDEO

3

## Mike Ferguson Is A Leading Industry Analyst / Consultant And Conference Chairman Of Big Data LDN - Leading The Industry In Data Management And Analytics

Big Data LDN is the largest data & analytics conference in Europe
- 20000 delegates
- 180 vendors
- 15 theatres
- 300+ speakers

4

## About Intelligent Business Strategies

- A UK-based independent IT analyst and consulting firm founded 1992 specialising in data management and analytics

- Mike Ferguson is an independent IT Industry Analyst and consultant, Conference Chairman of Big Data LDN and a member of the EDM Council CDMC Executive Advisory Board

- Three main lines of business in the areas of **Data Management, BI / ML / AI Analytics and Intelligent Business**

| Research | Education | Consulting |
|---|---|---|
| • Market research<br>  • 4<sup>th</sup> Industrial Revolution Survey<br><br>• D&A product research<br>  • Data Catalogs<br>  • Data Governance<br>  • Data Fabric<br>  • Data Science Workbenches<br>  • Analytical Databases | • How to Govern Data Across a Distributed Data Landscape<br>• Practical Guidelines for Implementing a Data Mesh<br>• Building a Competitive Data Strategy for a Data Driven Enterprise<br>• Data Catalogs<br>• Data Warehouse Modernisation<br>• Data Warehouse Migration to the Cloud<br>• Embedded Analytics, Intelligent Apps & AI Automation<br><br>• Public classes (anyone)<br>• On-site classes (single client)<br>  • Customers, vendors, systems integrators<br>• On-line (public & on-site classes) | • Customer consulting services<br>  • D&A Strategy, Data Architecture<br>  • D&A Technology selection<br>  • D&A Reviews, Data Governance<br>  • Project implementation advisory<br>• Vendor advisory services<br>  • Product strategy<br>  • Product positioning & go to market<br>  • Marketing support<br>    • Speaking at vendor events<br>    • White papers<br>    • Webinars<br>• Venture Capitalists<br>  • Due-diligence, Asset advisory |

www.intelligentbusiness.biz

5

---

## Topics

- What is generative AI?

- What are the business benefits of generative AI?

- How is generative AI being used in data management?

- What does this mean for business going forward?

- What should you do to get started?

6

## Topics – Where Are We?

➤ What is generative AI?

▪ What are the business benefits of generative AI?

▪ How is generative AI being used in data management?

▪ What does this mean for business going forward?

▪ What should you do to get started?

7

## Top Ten Key Trends In Data Management And Analytics (D&A)

1. Generative AI – innovation in every area of business
2. Hybrid multi-cloud computing is now the norm
3. Architecture modernisation - integration of data warehouses, lakehouses, data lakes & streaming
4. Rationalisation towards a common data & analytics software stack for development & governance
   • Do more with less, integration across tools, accelerated development and shared metadata
5. FinOps - CFOs demand visibility of the full cost of the D&A ecosystem with consumption-based pricing
6. Data governance remains very high priority with AI governance now also on the agenda
   • More sources, distributed data complexity, data quality, security, privacy, usage, observability, sharing, retention
   • Poor data culture and weak data governance are barriers to decentralised development
7. Increasing demand for lower and lower latency data
8. Democratisation and acceleration of data and analytics development
   • E.g., Data Mesh, citizen data engineering, DataOps, autoML, MLOps, CI/CD, code Vs low code / no code
9. Compliant sharing and reuse of data and analytical products in data marketplaces / exchanges
10. Growth in intelligent applications, decision intelligence and AI-Automation

8

## What Is Generative AI?

### What is Generative AI?

A subset of deep learning where multi-layer neural network models **generate new content** such as text, images, audio, video, code, and synthetic data in response to prompts based on what the models have learned from patterns in the content they were trained on

- General use cases
  - Text generation
  - Virtual assistant – chat
  - Conversational search
  - Summarisation – text extraction
  - Code generation
  - Synthetic data generation
  - Image generation and classification
  - Video generation
- Benefits
  - Improve customer and employee experience
  - Productivity, automation
  - Ease of use, lower skills bar
  - Democratisation of D&A development

9

## Generative AI – What Are Large Language Models (LLMs)?

- Neural network models with billions of connections

- Pre-trained using unsupervised learning on very large bodies of public content and code to:
  - Predict the next token / word and generate it → repeat
  - Append the token to content generated to date

- Foundation LLM examples

| Vendor | Foundation Models |
|---|---|
| Amazon | Titan & CodeWhisperer, Bloom |
| Databricks | Dolly |
| Google | Codey, LaMDA & PaLM2<br>• Gemini is the chat interface to PaLM2 |
| Open AI / Microsoft | Generative Pre-Trained Transformer (GPT4)<br>• ChatGPT is a chat interface to GPT |

Google Pathways Language Model (PaLM2)

QUESTION ANSWERING
SEMANTIC PARSING
PROVERBS
ARITHMETIC
CODE COMPLETION
GENERAL KNOWLEDGE
READING COMPREHENSION
SUMMARIZATION

LOGICAL INFERENCE CHAINS
COMMON-SENSE REASONING
PATTERN RECOGNITION
TRANSLATION
DIALOGUE
JOKE EXPLANATIONS
PHYSICS QA
LANGUAGE UNDERSTANDING

540 billion parameters

*As the scale of the model increases, the performance improves across tasks while also unlocking new capabilities.*

Image source: Google

- Foundation LLMs can be fine-tuned for specific tasks
  - Using supervised learning on smaller labelled data

10

## Generative AI Example – Code Generation Using LLMs, e.g., PaLM

| Prompt | Model Response |
|---|---|
| # Given a string representing musical notes in a special ASCII format, # parse this string and compute a total length of the music in beats. # Legend: #   'o' – whole note, last four beats #   'o\|' – half note, last two beats #   '.\|' – quarter note, last one beat  notes = 'o o\| .\| o\| o\| .\| .\| .\| o o' | # 1. Split the string into a list of notes. # 2. Iterate over the list and compute the total length of the music. # 3. Print the result.  # 1. notes = notes.split()  # 2. total_length = 0 For note in notes:     if note == 'o':         total_length += 4     if note == 'o\|':         total_length += 2     if note == '.\|':         total_length += 1  # 3. print(total_length) |

*Examples of a fine-tuned PaLM 540B model on text-to-code tasks, such as GSM8K-(
code-to-code tasks, such as Transcoder.*

Image source: Google

**Text-to-code**
Automatically generate code from a prompt

**PaLM-Coder DeepFix Code Repair**
(trained on Python)
Auto debug code so it compiles successfully

```
 1 #include <stdio.h>
 2 int main() {
 3     int a[i];
 4     int i;
 5     scanf("%d", &n);
 6     scanf("%d", &a[i]);
 7     int count = 0;
 8     for (i = 0; i < n; i++) {
 9         for (j = 0; j < n; j++) {
10             if (a[i] == a[j])
11                 count = count + 1;
12         }
13     }
14     if (count == a[i])
15         printf("Yes");
16     else
17         printf("No");
18     return 0;
19 }
```

```
 1 #include <stdio.h>
 2 int main() {
 3     int a[100], i, n, count = 0;
 4     scanf("%d", &n);
 5     for (i = 0; i < n; i++)
 6         scanf("%d", &a[i]);
 7     for (i = 0; i < n; i++) {
 8         for (int j = 0; j < n; j++) {
 9             if (a[i] == a[j])
10                 count = count + 1;
11         }
12     }
13     if (count == n)
14         printf("Yes");
15     else
16         printf("No");
17     return 0;
18 }
```

*An example from the DeepFix Code Repair task. The fine-tuned PaLM-Coder 540B fixes compilation errors (**left**, in red) to a version of code that compiles (**right**).*

Image source: Google

**11**

## Generative AI Concerns / Issues

| Concern/Issue | More Detail |
|---|---|
| Hallucinations | • LLM is unreliable because it can make up wrong answers when it doesn't know, e.g., <br> • The wrong calculation or filter on a SQL statement <br> • The wrong transformation <br> • May be caused by a model not being given enough data, cab be reduced by fine tuning |
| LLMs are trained on very large bodies of public data | • No domain knowledge <br> • No understanding of a company's data <br> • No understanding of a company's metadata, business terms, metrics <br> • No understanding of different user roles or business context |
| Governance for AI | • Safe use - How to train LLMs on data that is safe to use? <br> • Compliance – use of regulated sensitive data may expose that data and/or violate privacy <br> • Prompt governance - Avoiding sensitive and commercially confidential data in prompts <br> • Keeping a prompt log |
| Infringement of IP and copyright | • Reputation or legal implications |
| Ethics | • Ensuring no bias, ensure use that will not cause harm to people or your business |
| Criminal / Social use | • Deepfake, abuse (e.g., bad code), fraud, phishing… |

**12**

## How Does Generative AI Work?

1. Select the Foundation LLM you want to use
   e.g., Open AI GPT4

2. Use it as is OR fine tune it to create a customised LLM for your needs
   - A combination of supervised learning and re-enforcement learning with human feedback
   - A task which can be minutes to hours

Fine tuning LLM using supervised learning
Training data → e.g., Labeled prompts → Foundation LLM → Eval-uate → N / Y → Deploy → Tuned LLM

3. Utilise the customised fine-tuned LLM by invoking it from within applications and tools

Validate response & use ← Tool or Application → Prompt / Respond → GPT 4 / Tuned LLM → AI-generated content (text, code, …) → Analytical data store

**13**

## Tuning A LLM Can Be Done In A Data Science WorkBench
## _ E.g., Fine Tuning A LLM In Google Vertex AI And Deploying It To An Endpoint

**14**

7

## Supplementing And Training Open AI ChatGPT With Your Own Knowledge Base – Creating Embeddings And Indexes

- ChatGPT has been trained on general knowledge
  - It has no or limited domain specific knowledge
  - Hallucinations can occur – when a LLM makes up information

- What are embeddings?
  - They supplement ChatGPT's knowledge base
  - Provide additional information from your own knowledge base that meets your needs
  - Helps provide more relevant, reliable responses
  - Help enable faster retrieval and similarity search
  - Provides long-term memory for LLMs

- Vector indexes
  - Created on vector databases
  - Enable faster search

Internal or external sources
Text documents, Web content, CSV files
Split into chunks
Embedding ML model
vectors
Embeddings vector database
An embedding is a vector of floating-point numbers

Tool or Application
Prompt
Response
AI-generated content (text, code, …)
Trained LLM

S3 directory or files
Google BigQuery
Tool or catalog metadata
CSV files, Text documents
Internal or external sources
Split into chunks
GPT LLM (SimpleVectorIndex)
index
JSON file

Reference article https://beebom.com/how-train-ai-chatbot-custom-knowledge-base-chatgpt-api/

15

## Generating Embeddings In A Vector Database

Creating embeddings addresses the hallucination problem associated with LLM responses
Vectors augment the prompt with enterprise-specific content to produce better responses

Text documents
Natural Language Embedding Model

Images
Image Embedding Model

**Text Vector Table**

| id | vector | text |
|---|---|---|
| 1 | [0.8, 0.5, 1.6, -2.5, …] | "It was the best of times, it was the worst of times, it was.." |
| 2 | [1.1, 0.3, 0.6, -1.3, …] | "It is a truth universally acknowledged, that a single man.." |
| 3 | [1.3, 0.1, 0.2, -1.1, …] | "It was a bright cold day in April, and the clocks were striking.." |
| … | … | … |

**Image Vector Table**

| id | vector | Image |
|---|---|---|
| 1 | [0.5, 1.5, 2.6, -1.1, …] | |
| 2 | [1.0, 0.9, 1.6, -1.3, …] | |
| 3 | [0.6, 1.1, 1.3, -0.9, …] | |
| … | … | … |

Image source and Copyright ©: Oracle - LRN1412 "Enabling Generative AI with AI Vector Search in Oracle Database"

16

## LLM API Frameworks Are Also Now Available To Help Make Use Of LLMs In Tools And Applications

- Several new frameworks have emerged to develop LLM powered Gen AI applications and tools
  - E.g., Amazon Bedrock, Dust, LangChain, Steamship…
- LangChain is a framework for developing applications and tools powered by LLMs
  - LangChain is often used inside data management and BI tools

**17**

## AI Development Using LLM API Frameworks - LangChain Is Used To Connect LLMs To Your Own Data And Also Use The LLM To Help You Take Actions



```
template = """\
You are a regional sales manager for a pharmaceutical company.
What are the total sales for {product} in {my_region} last {time_period}
"""
```

You can have different prompt template libraries for different technical and business use cases
E.g. prompts for data engineering, data science, marketing insights, financial insights, sales insights, HR insights, …

**18**

## Development Tools Now Exist To Build Generative AI Applications – E.g. Google Vertex AI Search And Conversation (formerly Generative AI App Builder)

Google Vertex AI Search and Vertex AI Conversation
- Brings together Google's foundation AI models and search
- Includes a model garden of Google's foundation and partner AI models
- Build apps without writing any code that make use of generative AI
- Transact during conversation with the bots you build

Conversational AI using LLMs is likely to appear in ALL major applications in the enterprise

Gen Vertex AI Model Garden

Foundation models can be tuned and tested in Vertex AI

**19**

## Business Wants Generative AI And ML Models 'Wired Into' Every Application To Improve Productivity, Reduce Costs, And Shorten Time To Action For Better Profitability

Get employees to decide which tasks should use generative AI

AI-Assistant Productivity Automation

Use generative AI to improve ways of working that help yield better returns

Employees

**Intelligent Application** (HR)

**Intelligent Application** (Marketing)

**Intelligent Application** (Finance)

**Intelligent Application** (Sales)

**Intelligent Application** (Procure-ment)

**Intelligent Application** (Service)

**Intelligent Application** (Operations)

**Intelligent Application** (Distribution)

Trained Generative AI models + predictive / classification / and clustering ML models

Suppliers

Partners Customers

Things

Intelligent applications

Governed access to Generative AI LLMs and trained ML models as services across the business is critical to success

Generative AI does not replace other ML models in use in the enterprise

**20**

## Business Prompt Examples
### – Getting Things Done With Generative AI In Customer Facing Business Functions

**Marketing**
- Create a lead gen page for a webinar with a sign-up form
- Create personalised emails to all customer likely to buy our new load product
- Create tweets to promote our new product
- Create a blog introducing our new product

**Customer Service**
- What is the best reply for this customer?
- Create personalised emails to all customers about our new loan product

**Sales**
- Who are the top 10 contacts likely to buy the new product?
- Recommend the best contact to go see next?
- Write an informal introduction email
- What insights are available on Acme Corporation?

Use generative AI to change customer and employee experiences and simplify complex tasks
- AI powered voice and text conversations with customers
- Build generative AI based personalised marketing into customer experiences
- Help employees get things done quicker

21

## Topics – Where Are We?

- What is generative AI?
- ➤ What are the business benefits of generative AI?
- How is generative AI being used in data management?
- How is generative AI being used in data science and BI
- What does this mean for business going forward?
- What should you do to get started?

22

### Business Benefits Of Generative AI In Data Management And Analytics

| Benefit | Examples |
|---|---|
| Improved productivity | AI-powered – conversational data search, assisted metadata curation |
| | AI-automation of tasks – e.g., code generation |
| | AI-automated actions e.g., governance actions |
| Natural language-based user interfaces | Make tools much easier to use |
| | Broadens the use of tools to lesser skilled users |
| Guidance | AI-powered recommendations |
| Enriched answers | AI-powered natural language explanation of insights and their business impact |
| | AI-powered natural language explanation of data, data pipelines and data products |
| Continuous learning | AI continues to learn as you use |

23

### Topics – Where Are We?

- What is generative AI?
- What are the business benefits of generative AI?
- ➢ How is generative AI being used in data management?
- What does this mean for business going forward?
- What should you do to get started?

24

## Mega Trends
### – Generative AI Has Already Emerged In Data Management And Analytics Software

| Data & Analytics Category | Example Data & AI Products Already Supporting Generative AI |
|---|---|
| Data Catalogs | Atlan |
| | Collibra |
| | Data.world |
| Data Integration | SnapLogic SnapGPT |
| | Denodo |
| Data Governance | Informatica CLAIRE GPT (glossaries, policies, rules…) SodaGPT |
| Database Management Systems | ArangoDB |
| | Google Duet AI in BigQuery and BigQueryML |
| | Snowflake, Oracle |
| | Neo4j, Kinetica |
| Master Data Management | ViaMedici PIM |
| Self-service analytics / BI | ThoughtSpot Sage |
| | Tableau GPT & Einstein CoPilot |
| Data Science Workbench | Databricks LakehouseIQ & English API, LakehouseAI |
| | IBM Watsonx.ai |
| | Google Vertex AI |
| Decision Intelligence | Aera |

**25**

## Generative AI Has Emerged In Almost Every Aspect Of Data Management



Data engineering

Data Modelling

Data Virtualisation

Generative AI

Customised LLMs

Data catalog

DBMS / Lakehouse

Business Glossary

Database Migration

Data Marketplace

MDM

Data Governance

**26**

## Generative AI In Data Catalogs
### – The Promise Of Using Data Catalogs To Train LLMs

Data catalog

**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)**

Image source: Unifying Large Language Models and Knowledge Graphs: A Roadmap

Shirui Pan, *Senior Member, IEEE*, Linhao Luo,
Yufei Wang, Chen Chen, Jiapu Wang, Xindong Wu, *Fellow, IEEE*

Training LLMs with your own knowledge graph

Metadata knowledge graph

Data catalog

Data catalogs will provide much needed context to LLMs

27

## Generative AI In Data Catalogs – What's Possible?

Data catalog

- **Conversational data search**
  - Find data and other artefacts more rapidly
  - Research question generation based on data collections

- **Automated metadata enrichment to accelerate curation**
  - Metadata infused prompts
  - Automated AI generated metadata enrichment at scale during data discovery
  - Auto-generation and recommendation of business terms, synonyms and term descriptions in the business glossary
  - Auto generation and recommendation of editable data asset descriptions
  - Provision auto-generated metadata into other tools via catalog APIs
  - AI generated catalog update notifications

- **Automated query generation and summarisation**

- **Automated synthetic sample data generation**

Other tools

API

LLM

Data Catalog

metadata infused prompts
response

notify people of changes
slack

GPT LLM indexed with data catalog metadata

Data catalog metadata

Data discovery & classification

Data sources

28

Using Generative AI To Simplifying The Data Catalog User Interface
– Conversational Data Search Using Natural Language And Voice To Find Data

Copyright © Intelligent Business Strategies 1992-2024                    29



Prompt-Based Search And AI-Assisted Question Generation In Data Catalogs
– Product Example: Data.world

data.world now has generative AI bots called Eureka bots

Copyright © Intelligent Business Strategies 1992-2024                    30

Using Generative AI For Data Catalog Metadata Enrichment
– Auto-Generation Of Business Term Descriptions In Atlan Business Glossary

Business Glossary   Data catalog

Source: Atlan

Add a Readme, pick from a template and auto generate the Readme text using a generative LLM

AI generate

Copyright © Intelligent Business Strategies 1992-2024   **31**



Generative AI In Data Catalogs
– Automated Metadata Enrichment Using Atlan AI To Generate A Description Of An Asset

FCT_ORDERS table currently has no description

Request Atlan AI to generate one

Data catalog

Source: Atlan

Copyright © Intelligent Business Strategies 1992-2024   **32**

## Atlan Generated Contextual Description
### – E.g. The Description Includes The Database The Table Is In, The Schema



Data catalog

Generated description

User can accept this by clicking "Apply"

The user can also edit the description before accepting it

Source: Atlan

Copyright © Intelligent Business Strategies 1992-2024

33

## Atlan AI Activity Shows A Log Of The Changes Made By A User Using Atlan AI



Data catalog

'Atlan keeps an activity' log that shows the description was updated by a user using Atlan AI

Source: Atlan

Copyright © Intelligent Business Strategies 1992-2024

34

Generative AI And Data Catalogs - Automated Metadata Enrichment Using Atlan AI To Auto-Generate Column Descriptions And Business Data Names

Data catalog

Source: Atlan

Copyright © Intelligent Business Strategies 1992-2024

35



Atlan AI Activity Shows A Log Of The Column Description Changes Made By A User Using Atlan AI

Data catalog

Accepted AI-generated column descriptions are immediately visible

The activity log shows they were updated by a user using Atlan AI

Source: Atlan

Copyright © Intelligent Business Strategies 1992-2024

36

## Generative AI In Data Catalogs - Automated Metadata Enrichment Using AI-Generated Column Descriptions In Collibra



Source: Collibra

37

## Generative AI In Data Catalogs – Automatic AI-Generated Metadata Enrichment At Scale During Source Data Discovery In Atlan



Select what metadata you want automatically generated

As Snowflake data assets are automatically discovered in Atlan, generative AI can be invoked to auto-enrich the metadata

Data catalog

Source: Atlan

38

## Synthetic Sample Data Generation In Atlan Data Catalog Using Generative AI

Data catalog

| # | CustomerName | CCN | BillingID | SSN |
|---|---|---|---|---|
| 1 | John Doe | 1234-5678-9012-3456 | 1000001 | 123-45-6789 |
| 2 | Jane Smith | 2345-6789-0123-4567 | 1000002 | 234-56-7890 |
| 3 | Bob Johnson | 3456-7890-1234-5678 | 1000003 | 345-67-8901 |
| 4 | Samantha Lee | 4567-8901-2345-6789 | 1000004 | 456-78-9012 |
| 5 | David Kim | 5678-9012-3456-7890 | 1000005 | 567-89-0123 |
| 6 | Emily Chen | 6789-0123-4567-8901 | 1000006 | 678-90-1234 |
| 7 | Michael Brown | 7890-1234-5678-9012 | 1000007 | 789-01-2345 |
| 8 | Karen Davis | 8901-2345-6789-0123 | 1000008 | 890-12-3456 |
| 9 | Tom Wilson | 9012-3456-7890-1234 | 1000009 | 901-23-4567 |
| 10 | Lisa Nguyen | 0123-4567-8901-2345 | 1000010 | 012-34-5678 |

Source: Atlan

Copyright © Intelligent Business Strategies 1992-2024

**39**

## Generative AI And Data Modelling - Ellie.ai Is a Popular Data Modelling Tool That Can Use Generative AI to Help Build Data Models

Data Modelling

Use natural language prompt-based data modelling

Copyright © Intelligent Business Strategies 1992-2024

**40**

**Generative AI In Data Engineering – SnapLogic SnapGPT Prompt Based Data Engineering**

Data engineering

Generated pipeline preview

Import the pipeline

41



**Generative AI Prompt Based Data Engineering – SnapLogic SnapGPT Pipeline Configuration Wizard**

Data engineering

Source: SnapLogic

42

**Prompt Based Data Engineering Example**
**– Parsing Text, Extracting Data And Populating A Graph Database Using Gen AI**

Text documents → LLM → JSON → Code gen LLM

e.g., CVs | Parse documents and extract structured data from documents | | Generate Cypher code to write the data into the database

Neo4J Graph DB

**Prompt to extract data from text for the Person node**

person_prompt_tpl="""From the Resume text for a job aspirant below, extract Entities strictly as instructed below
1. First, look for the Person Entity type in the text and extract the needed information defined below: `id` property of each entity must be alphanumeric and must be unique among the entities. You will be referring this property to define the relationship between entities. NEVER create new entity types that aren't mentioned below. Document must be summarized and stored inside Person entity under `description` property
   Entity Types:
   label:'Person',id:string,role:string,description:string //Person Node
2. Description property should be a crisp text summary and MUST NOT be more than 100 characters
3. If you cannot find any information on the entities & relationships above, it is okay to return empty value. DO NOT create fictious data
4. Do NOT create duplicate entities
5. Restrict yourself to extract only Person information. No Position, Company, Education or Skill information should be focussed.
6. NEVER Impute missing values Example Output JSON: {"entities": [{"label":"Person","id":"person1","role":"Prompt Developer","description":"Prompt Developer with more than 30 years of LLM experience"}]}

Question: Now, extract the Person for the text below –
$ctext
Answer:
"""

Source: https://github.com/neo4j-partners/neo4j-generative-ai-google-cloud

43

---

**Generative AI In Data Engineering - Synthetic Data Generation Using LLMs**
**e.g., Tinybird Mockingbird – Generate A Schema And Synthetic Streaming Test Data**

Data engineering

tinybird

**Choose your source of data** ✕

Use one of our streaming data samples

**Web Analytics Events**
20 events second / 10 minutes
Sample user navigation events based on our starter kit

**Log Analytics Events**
20 events second / 10 minutes
Sample log events based on our starter kit

**Create your own**
20 events second / 10 minutes
Define your schema from a prompt and send sample data

Connect to other sources using our integration guides

**Google Pub/Sub**      **Google Storage**

Source: Tinybird

Prompt based test data schema generation

← **Create your own data sample** ✕

Describe the type of Data Source you want to generate. We will propose a schema from your description. Confirm it to start appending events from the browser and use it in a pipe.

Schema description

Generate a schema of financial transactions to be used for fraud detection that can detect potential fraud by location, time of day, and purchase amount among other commonly used columns used for financial transactions.

Schema preview

`location` String       `card_number` String      `user_id` String
`time` DateTime         `merchant` String         `device_id` String
`amount` Float32        `transaction_type` String `ip_address` String
`card_type` String      `transaction_id` String

Confirm and ingest

Source: Tinybird       44

---

## Generative AI In Data Virtualisation
– Product Example: Denodo NLP Queries And LLMs



Source: Denodo

45

## Generative AI In The Database
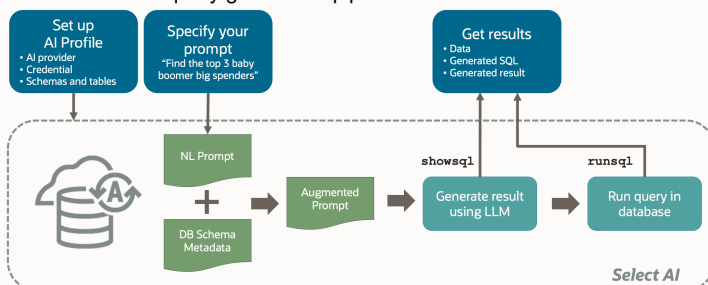– SQL Generation Using LLMs In Oracle Autonomous Database

DBMS

Enabled using the new **DBMS_CLOUD_AI** PL/SQL package

New AI keyword

```
SELECT AI
What are our top 10 streamed movies that were released after 2018
```

Automated SQL query generation pipeline

**Set up AI Profile**
• AI provider
• Credential
• Schemas and tables

**Specify your prompt**
"Find the top 3 baby boomer big spenders"

**Get results**
• Data
• Generated SQL
• Generated result

**Actions**
**runsql** - return the SQL result set (default)
**showsql** – return the generated query
**narrate** – return a conversational result
**chat** - general AI chat

NL Prompt

+

DB Schema Metadata

Augmented Prompt

showsql

Generate result using LLM

runsql

Run query in database

*Select AI*

Image source and copyright © Oracle

Use generative AI models in combination with your database data
You can also combine relational and vector data in the same query

46

## Generative AI In The Database – Oracle AI Vector Search

Native support for generating vectors – New SQL EMBEDDINGS Function



Images in database are embedded into vectors and stored:
**EMBEDDING(**resnet_50 **USING** data_img**)**

Input image is embedded into vector:
**EMBEDDING(**resnet_50 **USING** input_img**)**

Oracle Database

Input Image (client)

Load image as BLOB

**Image Embedding Generation**

Vector embedding

Search for similar image vectors

Return Top Matches

Output Matches

Image source and Copyright © Oracle

You can also import Open Neural Network Exchange (ONNX) embeddings models into the DBMS from object storage

47

---

## Generative AI In The Database
## – Oracle Has Added A New VECTOR Datatype To Store Vectors Directly In The DBMS

- New VECTOR datatype (with underlying BLOB storage for long-term extensibility)
  - VECTOR (<optional # of dimensions>, <optional format for dimension values>)

```
CREATE TABLE My_Images (id number, image BLOB, img_vec VECTOR(768, FLOAT32))
```

- Oracle 23c clients for Javascript and Python support VECTOR type, and so can **insert** Vectors directly
- Vector DML Operations
  - Insert ……. TO_VECTOR() converts a string representing an array of vector dimensions into VECTOR:

```
CREATE TABLE vec_tab(id number, dataVec VECTOR(3, 'FLOAT32'))
INSERT into vec_tab values (1, TO_VECTOR('[1.1, 2.2, 3.3]')
```

- Several new vector functions are supported including Vector_Distance, Vector_Avg, Vector_Norm
  - E.g., Get the Top-5 Nearest Vectors to a given query

```
SELECT id from tab order by VECTOR_DISTANCE(data_vec, :queryVec)
fetch first 5 rows only;
```
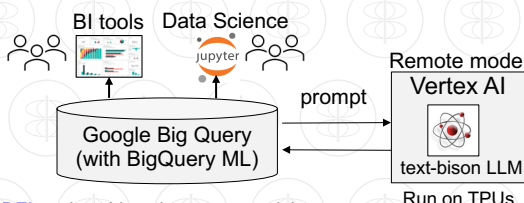
48

## Generative AI In The Database – Google Duet AI In BigQuery AND BigQuery ML Which Has Been Extended To Use LLMs In SQL Queries

DBMS

**Duet AI** in Big Query enables you to:
- **Generate** a SQL query
- **Complete** a SQL query
- **Explain** a SQL query

BI tools   Data Science

prompt → Remote model
Vertex AI

Google Big Query
(with BigQuery ML)

text-bison LLM
Run on TPUs

Uses
- Classification
- Sentiment Analysis
- Entity extraction
- Extractive Question Answering
- Summarization
- Re-writing text in a different style
- Ad copy generation
- Concept ideation

CREATE MODEL project_id.mydataset.mymodel
REMOTE WITH CONNECTION `myproject.us.test_connection`
OPTIONS(REMOTE_SERVICE_TYPE="CLOUD_AI_LARGE_LANGUAGE_MODEL_V1")

SELECT * FROM ML.GENERATE_TEXT( MODEL mydataset.llm_model TABLE mydataset.prompt_table,
STRUCT( 0.2 AS temperature, 75 AS max_output_tokens, 0.3 AS top_p, 15 AS top_k, TRUE AS flatten_json_output));
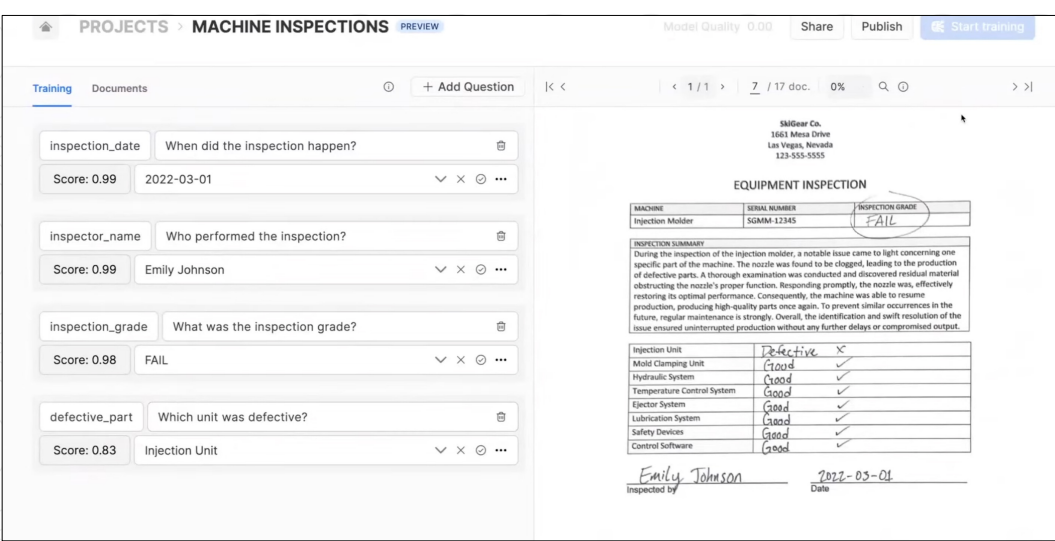
- Provides prompt data from a table column that's named prompt
- Returns a shorter generated text response
- Returns a more probable generated text response
- Flattens the JSON response into separate columns

49

## Generative AI In The Database – Use Natural Language Queries On Unstructured Documents Via Snowflake Document AI

DBMS

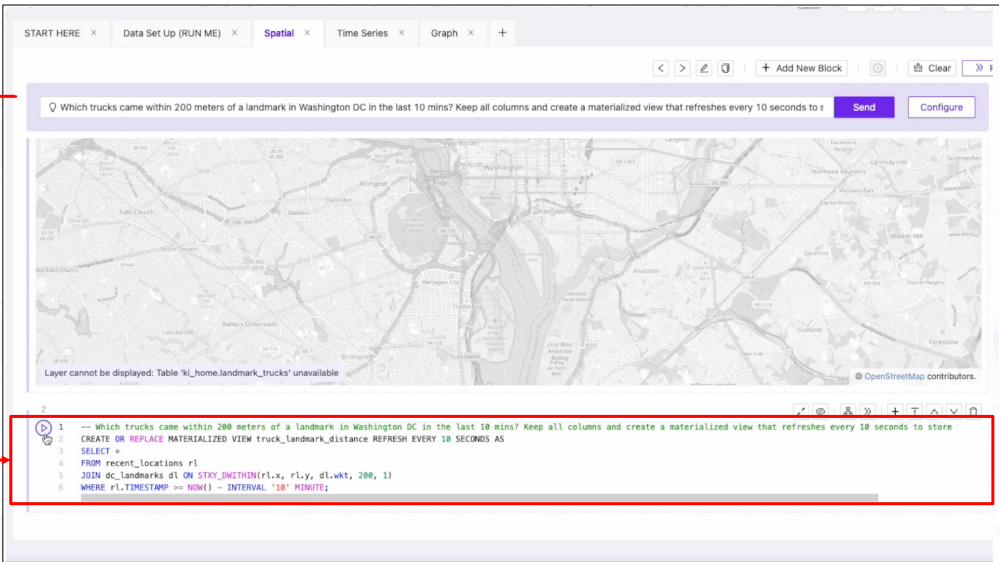PROJECTS > **MACHINE INSPECTIONS** PREVIEW

Model Quality 0.00   Share   Publish   Start training

Training   Documents   + Add Question   |< <   < 1/1 >   7 / 17 doc.   0%   Q

| inspection_date | When did the inspection happen? |
| Score: 0.99 | 2022-03-01 |

| inspector_name | Who performed the inspection? |
| Score: 0.99 | Emily Johnson |

| inspection_grade | What was the inspection grade? |
| Score: 0.98 | FAIL |

| defective_part | Which unit was defective? |
| Score: 0.83 | Injection Unit |

SkiGear Co.
1661 Mesa Drive
Las Vegas, Nevada
123-555-5555

**EQUIPMENT INSPECTION**

| MACHINE | SERIAL NUMBER | INSPECTION GRADE |
| Injection Molder | SGMM-12345 | FAIL |

INSPECTION SUMMARY
During the inspection of the injection molder, a notable issue came to light concerning one specific part of the machine. The nozzle was found to be clogged, leading to the production of defective parts. A thorough examination was conducted and discovered residual material obstructing the nozzle's proper function. Responding promptly, the nozzle was, effectively restoring its optimal performance. Consequently, the machine was able to resume production, producing high-quality parts once again. To prevent similar occurrences in the future, regular maintenance is strongly. Overall, the identification and swift resolution of the issue ensured uninterrupted production without any further delays or compromised output.

| Injection Unit | Defective | X |
| Mold Clamping System | Good | ✓ |
| Hydraulic System | Good | ✓ |
| Temperature Control System | Good | ✓ |
| Ejector System | Good | ✓ |
| Lubrication System | Good | ✓ |
| Safety Devices | Good | ✓ |
| Control Software | Good | ✓ |

Emily Johnson        2022-03-01
Inspected by              Date

Source: Snowflake

50

## Generative AI In The DBMS – Natural Language Queries In Kinetica DBMS Become AI-Generated Geospatial SQL Queries



**DBMS**

Note that LLMs can run in Kinetica because this is a MPP SQL DBMS running on GPUs
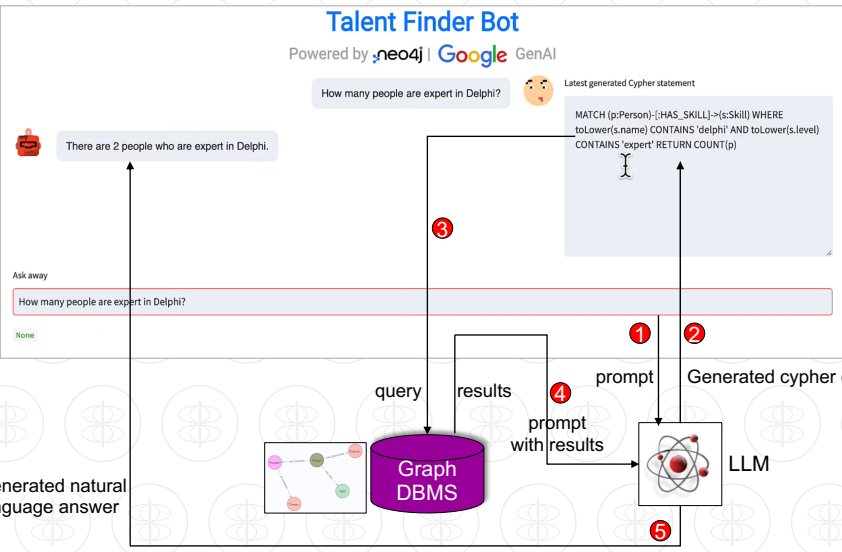
Source: Kinetica

Copyright © Intelligent Business Strategies 1992-2024

51

## Generative AI And Graph Databases – Natural Language Queries And Natural Language Responses In Neo4j - Generates Cypher Queries From Text
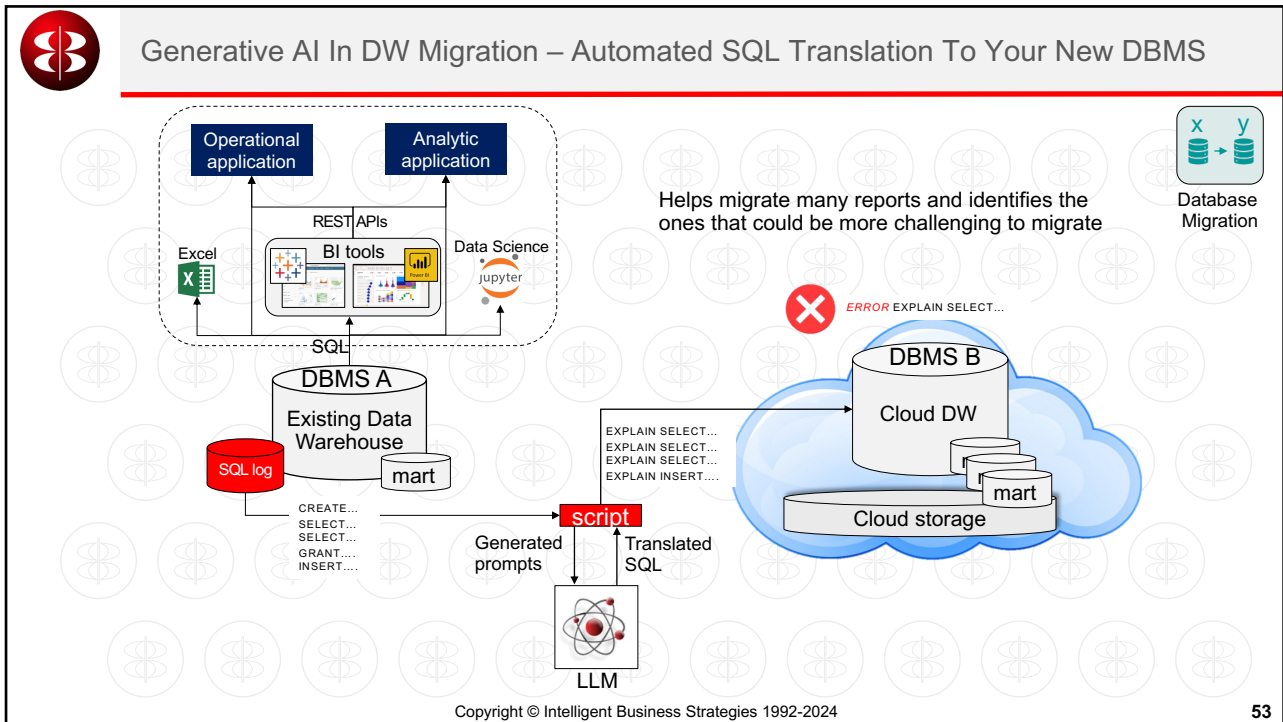


**DBMS**

Copyright © Intelligent Business Strategies 1992-2024

52

## Generative AI In DW Migration – Automated SQL Translation To Your New DBMS



Helps migrate many reports and identifies the ones that could be more challenging to migrate

Database Migration

53

## Generative AI In Data Governance

Data Governance

- Provide a generative AI-Assistant to data stewards

| Data Governance Discipline | Auto AI-Generation Of |
|---|---|
| Data Quality | Data quality validation check rules |
| | Data quality cleansing rules and code |
| | Synthetic data to replace missing values in data to improve data completeness |
| Master Data Management | Data matching rules |
| | Data survivorship rules |
| | Master data descriptions |
| | Product information e.g., digital assets - catalogs, product brief, marketing content |
| Data privacy | Data privacy policies |
| | Synthetic data that excludes PII data to avoid risks |
| Data access security | Data access security policies |
| | Data loss prevention policies |
| Data sharing | Data sharing policies |
| Data retention | Data retention policies |

54

## Generative AI In Data Governance – Prompt Based Discovery Of Policies In Source Systems Scanned By The Catalog In Data.world



Data Governance

Source: Data.world

55

## Generative AI In Data Governance – It Is Also Possible To Automatically Extract And Generate Data Governance Rules From Documents, E.g., IBM Watson Knowledge Catalog



Data Governance

56

## Generative AI In Data Governance – SodaGPT Lets Business Analysts Create Data Quality Checks In Natural Language That Can Be Shared With Data Engineers To Fix Data



Data Governance

This automatically generates SODA Check Language (SodaCL) from prompts to check for data quality problems

SodaCL is a low code abstract language built on SQL and Spark

Copyright © Intelligent Business Strategies 1992-2024

57

## Generative AI In Master Data Management (MDM)
## – Product Example: ViaMedici Product MDM & PIM AI-Powered Text Generation



MDM

ViaMedici infuses prompts with master data and passes them to an LLM to generate briefing text, a marketing poem and marketing tweets for a product

Copyright © Intelligent Business Strategies 1992-2024

58

Copyright © Intelligent Business Strategies 1992-2024

## Generative AI In Data Marketplaces



Data Marketplace

59

## Topics – Where Are We?

- What is generative AI?
- What are the business benefits of generative AI in data management and analytics?
- How is generative AI being used in data management?
- ➢ How is generative AI being used in data science and BI
- What does this mean for business going forward?
- What should you do to get started?

60

## Generative AI And Business Intelligence And Data Science

- **Gen-AI in business intelligence**
  - AI recommended questions from auto-analysis of your data to quick start query & analysis on your data
  - AI-generated queries from voice
  - AI-generated queries from natural language text questions
  - AI generated answers
  - Auto-predict and generate charts from your data
  - Natural language generation to explain insights
  - Metadata enrichment, e.g., generate synonyms for data

- **Gen-AI in data science**
  - AI-generated code
  - Prompt based querying of data
  - Foundation models and model tuning
  - Building generative AI applications

BI tool

BI tool

Generative AI

Customised LLMs

jupyter

Data science workbench

61

## Generative AI In Business Intelligence – Natural Language Questions – E.g., ThoughtSpot Sage Auto Analyses The Data First And Recommends Questions
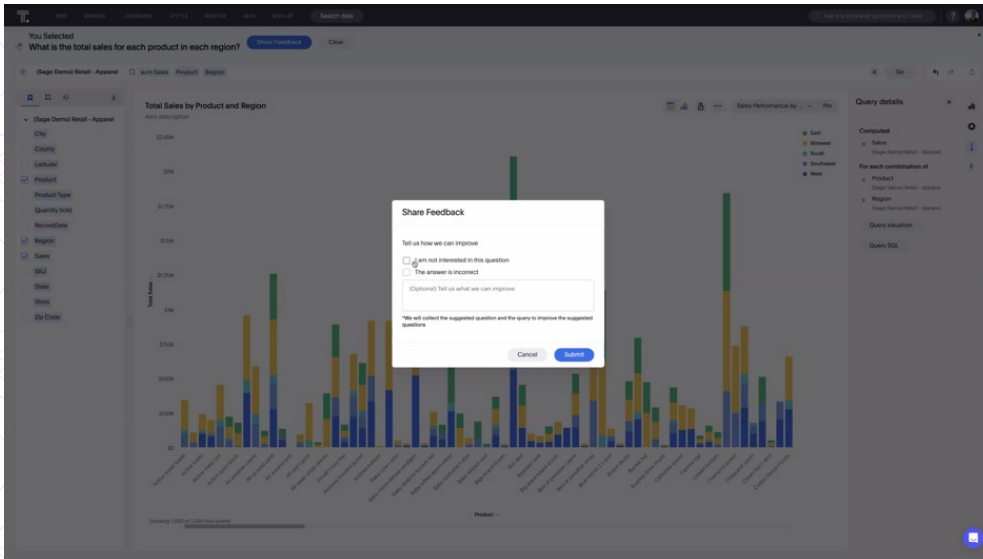
Recommended questions

Source: ThoughtSpot

ThoughtSpot will make use of multiple LLMs

62

## Generative AI In BI Tools - ThoughtSpot Sage Lets You Provide Feedback To Enable The AI Engine To Improve Question Recommendations It Creates



You can say you are not interested in a question so that it learns to get better at generating questions you are interested in
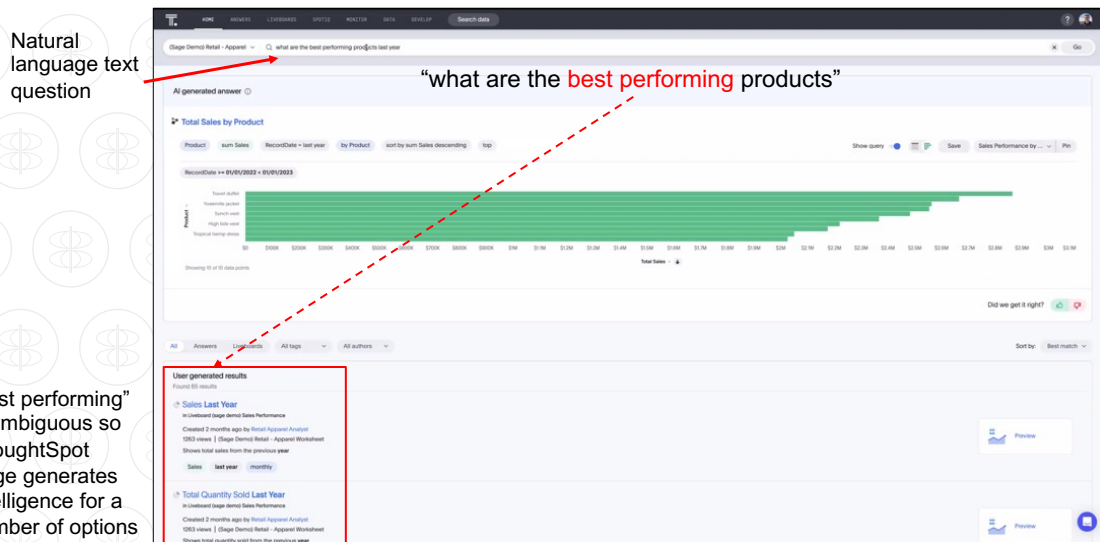
Source: ThoughtSpot

63

## Generative AI In Business Intelligence – ThoughtSpot Sage Generates Search Queries From Natural Language Text Questions While Dealing With Ambiguity

Natural language text question

"what are the best performing products"

"best performing" is ambiguous so ThoughtSpot Sage generates intelligence for a number of options



Source: ThoughtSpot

64

## Generative AI In Business Intelligence – GPT Generates Synonyms So ThoughtSpot Sage Search Engine Understands Synonyms For The Same Data



ThoughtSpot Sage queries the schema and then uses GPT to recommend synonyms to the terms in the schema

This allows users to query using synonyms and still get the same answer

Source: ThoughtSpot

Copyright © Intelligent Business Strategies 1992-2024

65

## Generative AI In Data Science – You Can Now Build Custom LLMs In Data Science Workbenches – E.g., Google Vertex AI Now Includes A Generative AI Studio



Source: Google

Copyright © Intelligent Business Strategies 1992-2024

66

## Generative AI In Data Science – IBM Watsonx.ai Data Science Workbench Allow You To Build, And Tune LLMs And Make Use Of A Prompt Lab To Work With LLMs



Source: IBM

67

## Generative AI In Data Science – Databricks LakehouseIQ Prompt Based Artificially Intelligent Query Generation



Source: Databricks

Lakehouse

Prompt

LakehouseIQ knows the company has two European sales regions and so automatically adds them to the query.

It has also learned from other queries, dashboards and notebooks that used this dataset and so automatically adds a filter to exclude internal usage

68

## Generative AI In Data Science – Databricks English SDK

Source code          Compiler          Byte code

English → Generative AI → PySpark

Generative AI's expert knowledge about Spark is built into the English SDK

| Data ingestion | spark_ai = SparkAI() auto_df = spark_ai.create_df("2022 UK national auto sales by brand") |
|---|---|
| Dataframes | spark_ai.activate()  # Activate the df.ai methods<br>auto_df.ai.plot("pie chart for UK sales market shares, show the top 5 brands and the sum of others")<br><br>Market Share of Top 5 Brands and Others<br><br>Toyota 13.9%, Ford 13.3%, Chevrolet 11.3%, Honda 6.07%, Hyundai 5.45%, Others 49.4%<br>Others, Toyota, Ford, Chevrolet, Honda, Hyundai<br><br>auto_top_growth_df=auto_df.ai.transform("top brand with the highest growth")<br>auto_top_growth_df.show() |

| brand | uk_sales_2022 | sales_change_vs_2021 |
|---|---|---|
| Toyota | 134726 | 327 |

69

## Topics – Where Are We?

- What is generative AI?
- What are the business benefits of generative AI?
- How is generative AI being used in data management?
➢ What does this mean for business going forward?
- What should you do to get started?

70

## What Does Generative AI In Data And Analytics Mean For Business?

| Data Producers |
| --- |
| Conversational data search in data catalogs |
| Automated metadata enrichment at scale |
| Automated physical schema generation |
| Prompt-based data engineering |
| Synthetic data generation |
| AI-assisted generation of data quality validation rules, master data matching rules and data governance policies |
| AI-assisted generation of digital product information |
| AI-assisted generation of product information in MDM (Digital web pages, brochures, marketing content) |

| Data Consumers |
| --- |
| Conversational data search in data marketplaces |
| Automatic exploration of data, data pipelines and lineage |
| Natural language queries |
| AI-generated natural language answers |
| Auto generated virtual views on data |
| |
| |
| |

- Generative AI in data and analytics:
  - Lowers the skills bar and broadens inclusion
  - Accelerates development
  - Explains business insights and the business impact
  - Shortens the time to value
  - Shortens the time to act

**Everyone will get a self-learning AI-Assistant!**

71

---

## Topics – Where Are We?

- What is generative AI?
- What are the business benefits of generative AI?
- How is generative AI being used in data management?
- What does this mean for business going forward?
- ➢ What should you do to get started?

72

## Conclusion - Getting Started

- Generative AI is rapidly finding its way into many different data management, data governance, BI and data science tools with more to come

- Fine tune foundation LLMs to understand your data and your metadata

- Exploit generative AI in data management to:
  - Improve data catalogs through conversational data search and automated data curation
  - Accelerate and start automating data governance tasks
  - Increase the number of citizen data engineers
  - Accelerate data engineering to produce data products more rapidly
  - Accelerate the build of master data
  - Auto generate product information and marketing content from product master data

73

## About Mike Ferguson

www.intelligentbusiness.biz

mferguson@intelligentbusiness.biz

@mikeferguson1

(+44) 1625 520700

Mike Ferguson is Managing Director of Intelligent Business Strategies Limited. As an independent IT industry analyst and consultant, he specialises in BI / analytics and data management. With over 40 years of IT experience, Mike has consulted for dozens of companies on BI/Analytics, data strategy, technology selection, enterprise architecture, and data management. Mike is also conference chairman of Big Data LDN, the largest data and analytics conference in Europe and a member of the EDM Council CDMC Executive Advisory Board. He has spoken at events all over the world and written numerous articles. Formerly he was a principal and co-founder of Codd and Date – the inventors of the Relational Model (which caused the birth of relational databases and the SQL language), and Chief Architect at Teradata on the Teradata DBMS He teaches popular master classes in Data Strategy, Data Warehouse Modernisation, Practical Guidelines for Implementing a Data Mesh, Big Data, How to Govern Data across a Distributed Data Landscape, Machine Learning and Advanced Analytics, and Embedded Analytics, Intelligent Apps and AI Automation

### Thank You!

74