**DDVUG**

# Current progress and limitations of GenAI in the field of Data Management

Evaluating GenAI using Willibald

Sebastian Bruns
Hamburg, April 12th 2024

**eraneos**

We prompted ...

"Draw a picture of a Data vault driven by Artificial intelligence"

# Your speaker

**Sebastian Bruns**
*Lead Data Engineer*
*Eraneos Analytics*
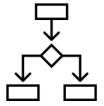
sebastian.bruns@eraneos.com

The speaker is a **technology enthusiast** and loves to discuss DevOps and CI/CD topics

Sebastian's background is **Computer Science** with focus on **Data Engineering**

o   2007 Oracle
o   2010 SAP Business Warehouse
o   2014 SQL Analysis Services
o   2015 NoSQL Databases
o   2016 "Big Data" / ELT
o   2019 Cloud Engineering
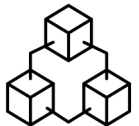
eraneos

# Agenda

### The power of Generative AI (GenAI) in automation

How we use Large Language Models (LLMs) to tackle complex automation tasks

### GenAI applications in Data Warehouse Automation (DWA)

Taking on the Willibald challenge: Evaluation of LLMs in the realm of the DWA benchmark dataset

### Looking Ahead

What we expect from the development in the field of GenAI in 2024 and beyond

eraneos

# We were experts in Machine Learning long before the rise of GenAI. Now we see the opportunity to solve many problems in a faster and more cost efficient way.

**150+**
Data & AI Top Experts

**500+**
Data & AI Projects

**150+**
Customers in various industries

## Electricity Price Forecasting
### Energy Supplier



## Predictive Maintenance
### MRO Service Provider



## AI in Train Disposition
### Railway



## Legacy DWH Migration
### Global Shipping Company

# The power of Generative AI (GenAI) in automation

How we use Large Language Models (LLMs) to tackle

complex automation tasks

# Eraneos has successfully **leveraged GenAI to automate migration and evaluation tasks** across multiple problem domains and industries

**+10** more GenAI projects

## LLM-generated Case Summaries
### Commercial Vehicles OEM



## Legal and Compliance Assistant
### Automotive



## Code Transformation
### Public Sector



Identifying closed tickets that should not have been closed by the service desk

Highlight legal regulations for global rollout of products

We have auto-created 1800 Unit Tests on a legacy C++ codebase with >300 K lines of code

# Eraneos Code Migration Framework as the core building block for a scalable, reliable and cost-efficient automation around LLMs

**1800+**
Unit Test generated

**550€**
GenAI Bill

**80%**
Cost & Time Saving

legacy source code

code indexing

eraneos migration framework

LLM-code creation

shipment to target system

validate and compile

eraneos

# GenAI applications in
# Data Warehouse Automation (DWA)

Taking on the Willibald challenge: Evaluation of LLMs in the realm of the DWA benchmark dataset



## Hello and welcome to Samen- und Pflanzenhandel Willibald

## A company that was only created because of its data.

Normally it's the other way around. But this way we can create the ideal data set and ensure that this data set contains all the typical problems that a data warehouse with Data Vault has to solve. Of course, this also includes requirements for evaluations and the resulting business rules.

**Willibalds Data**   **Github repository**

# The **strength of general-purpose LLMs** is their knowledge and universal applicability to a wide variety of problem statements

## Explain & Understand

- Write comments and documentation
- **DEMO** Create field mappings
- **DEMO** Explain the content data and scripts
- Outline data lineage in legacy systems

## Generate & Propose

- **DEMO** Create data model from raw data
- Generation of test data (synthetic data)
- **DEMO** Create data model from an ERD
- Recommend refactorings (indexes and contraints and performance improvements)

## Assess & Monitor

- Evaluate the functionality of the data model (potentially from business analysts)
- Assess the quality of the data in the database
- Automatically detect anomalies and fraud

## Optimize & Migration

- **DEMO** Migrate the data modeling
- Write data migrate scripts
- Modernize the code base
- Resolve dependencies between (e.g. normalized stacked views)

eraneos

# Creating a basic setup using the markup file and **deducting information about the business context** to build an **initial data vault model**

TASK →

PROMPT →

RESULT

Using the [**readme en.md**] markup file as initial starting point

**The initial data**

The company Willibald is a traditional house and trades with seeds and plants via the internet. In the past, the company sold exclusively via a catalogue. The catalogue was quickly terminated in 2000, a full 4 years after the webshop was opened. Willibald was the first plant supply shop on the internet and is still proud of it today.

When ordering via the w...

**The evaluation**

The desired report for Willibald includes the following key figures on a granular level:

- Quantity The quantity of ordered products per product. This key figure is taken directly from the order item.
- Revenue The amount of money to be paid for the ordered product quantity. For this purpose, the quantity and price from the order item are multiplied and then the discount is deducted.
- Open order A 1 to indicate briefly whether this order is still open.
- Completely delivered order A 1 to indicate briefly whether this order has already been delivered in full.
- Difference between request date and delivery date in days If each item has been delivered, the deviation is calculated from the order.requestdate and the last delivery.deliverydate. If the delivery was made before the requested date, the deviation is negative. Willibald wants to deliver on time, because too early deliveries also cause problems for the customer (care of seedlings).

➕ **customer.csv**

**Text Input**: If we were having larger text input (e.g. documents) we would exceed the token limit at some point and go for RAGs

We prompted ...

"Give me in brief explanation what Willibald is about"

"Create a Data Vault model reflecting Willibalds business processes"

"What challenges do you see modeling the Willibald data?"
"What customers have no phone number?"

This file provides information about the data and data model f... fictitious company called 'Samen und Pflanzenhandel Willibald' that trades seeds and plants online. It explains the company's history, its ordering and delivery process, and its data model. It also details the compa... with horticu... pany's biannu... it outlines ...ures and necess... ncludes directo... ...ta.

Creating a Data Vault model based on the information provided would involve the following steps:

1. **Hubs**: These are the core business keys in your system. In this case, the hubs would likely be Customer, Product, Order, ClubPartner, and Delivery.

2. **Links**: These represent the relationships between business keys. Possible links in this scenario could be Customer–Order, Order–Product, Customer–ClubPartner, Order–Delivery, and Customer–Delivery.

3. **Satellites**: These hold descriptive data (attributes) about the business keys and their relationships. For example, a Customer satellite might hold data like CreditCard, ValidTo, and Credit Card Company. An Order satellite could hold Quantity, Price, Discount, RequestDate, DeliveryDate, and DeliveryStatus.

4. **Business Vault**: This layer would handle the necessary transformations and business rules mentioned, such as unifying the orders of the association partners and assigning roadshow orders to the customers.

"Incomplete Data" and "Data Quality: There are some fields that appear to have missing"

# We want to know **how knowledgeable general-purpose LLMs are about modeling techniques** and how they implement them

TASK ———————————————————→  PROMPT ———————————————————→  RESULT

The [Lieferung] Data from Willibald's "Webshop Testdata 1" contains references to other tables (e.g. [Bestellung]) as well as actual data (e.g. [LieferDatum)]



**This is new**: LLMs understands the *meaning* of the data fields and model them accordingly.

**Non-deterministic**: Column names are often, but not always deducted correctly.

We prompted …

"Model this dataset as **star schema** with fact and dimension tables."

"Model this dataset as **Data Vault** with hubs, links and satellites."

"Model this dataset as **3rd Normal Form.**"

1. Fact Table:
   ◦ Fact_Lieferung (Keys: Lieferung_ID, Produkt_ID, Kunde_ID, Datum_ID, Attributes: Menge)
2. Dimension Tables:
   ◦ Dim_Produkt (Key: Produkt_ID)
   ◦ Dim_Kunde (Key: Kunde_ID)
   ◦ Dim_Datum (Key: Datum_ID, Attributes: Day, Month, Year)

Based on the provided data, here's how it could be modeled in Data Va...

1. Hubs:
   ◦ Hub_Lieferung (Key: Lieferung_ID)
   ◦ Hub_Produkt (Key: Produkt_ID)
   ◦ Hub_Kunde (Key: Kunde_ID)
   ◦ Hub_Datum (Key: Datum)
2. Links:
   ◦ Link_Lieferung_Produkt (Keys: Lieferung_ID, Produkt_ID)
   ◦ Link_Lieferung_Kunde (Keys: Lieferung_ID, Kunde_ID)
   ◦ Link_Lieferung_Datum (Keys: Lieferung_ID, Datum)
3. Satellites:
   ◦ Satellite_Lieferung (Parent: Hub_Lieferung, Attributes: Menge)
   ◦ Satellite_Produkt (Parent: Hub_Produkt, Attributes: None in this case)
   ◦ ...ub_Kunde, Attributes: None in this case)
   ◦ ...ub_Datum, Attributes: None in this case)

1. Data Warehouse:
   ◦ Lieferung (Attributes: Lieferung_ID, Produkt_ID, Kunde_ID, Datum, Menge)
   ◦ Produkt (Attributes: Produkt_ID)
   ◦ Kunde (Attributes: Kunde_ID)
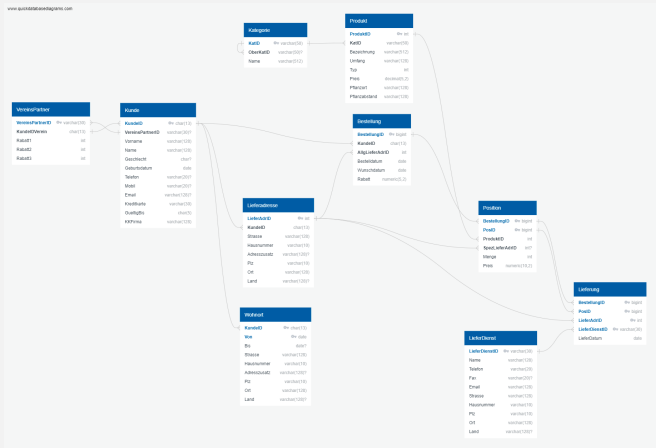   ◦ Datum (Attributes: Datum_ID, Day, Month, Year)
2. Data Marts:
   ◦ Lieferung Data Mart (Attributes: Lieferung_ID, Produkt_Name, Kunde_Name, Datum, Menge)
   ◦ Produkt Data Mart (Attributes: Produkt_ID, Produkt_Name)
   ◦ Kunde Data Mart (Attributes: Kunde_ID, Kunde_Name)
   ◦ Datum Data Mart (Attributes: Datum_ID, Day, Month, Year)

eraneos

# We have been provided the **source table UML diagram** and wish to create the data vault model including build scripts and field mappings

TASK → PROMPT → RESULT

We provide the LLM the file [webshop-willibald.png] without providing any more information than the prompt



Only two years back, this kind of image recognition task would have taken us a full project!

We prompted …

"Create a data vault model from this UML diagram"

"Create the SQL statements that generate the data vault"

"Draw me a mapping of the fields between the source and the target model"



| Original UML Field | Data Vault Entity | Data Vault Attribute |
|---|---|---|
| KundeID | Hub_Kunde | KundeID |
| Name | Sat_Kunde_Details | |
| Vorname | Sat_Kunde_Details | |
| Geschlecht | Sat_Kunde_Details | |
| Geburtsdatum | Sat_Kunde_Details | |
| … | … | |
| ProduktID | Hub_Produkt | |
| Bezeichnung | Sat_Produkt_Details | |
| Preis | Sat_Produkt_Details | |

eraneos

# We are interested in **optimizing an existing data vault model** and wish to get recommendations and explanations

TASK → PROMPT → RESULT

**We give the model an existing but flawed data model:**

```
Hub_Customer
Customer_ID (PK)
Customer_Firstname
Load_Date
Record_Source
```

```
Sat_Customer
Customer_ID (Foreign Key),
Name
Gender
Date_of_Birth
Email
```

**The Firstname doesn't belong into the hub, the satellite is missing meta data fields**

We prompted …

"What can be improved in this data vault model?"

"Create me the SQL statements to apply these changes"

"Why did you split the customer data into personal and contact information?"

"2. Remove Firstname from Hub_Customer"
"3. Add Record Source to Satellite Tables"

```sql
CREATE TABLE Sat_Customer_Personal_Info (
    Customer_ID INT NOT NULL,
    Name VARCHAR(255),
    Gender VARCHAR(10),
    Date_of_Birth DATE,
    Load_Date DATE,
    Record_Source VARCHAR(255),
    FOREIGN KEY (Customer_ID) REFERENCES Hub_Customer(Customer_ID)
);

CREATE TABLE Sat_Customer_Contact_Info (
    Customer_ID INT NOT NULL,
    Phone_Number VARCHAR(20),
    Email VARCHAR(255),
    Load_Date DATE,
    Record_Source VARCHAR(255),
    FOREIGN KEY (Customer_ID) REFERENCES Hub_Customer(Customer_ID)
);
```

"Improved Organization, Enhanced Flexibility, Better Performance, Increased Security, Easier Maintenance"

eraneos

15

# The GenAI automation toolbox offers a **variety of methods** to further **improve the quality and reliability** of the results

| Method | Description | UseCases |
|---|---|---|
| **Prompt Engineering** | Add **specific goals**, commands or output formats to the prompt | <ul><li>Narrow down the options, e.g. what exact model technique</li><li>Prevent / guardrail certain unwanted outputs</li></ul> |
| **Retrieval Augmented Generation** | **Additional data** is retrieved from a (vector) database and handed in along with the query | <ul><li>Internal data (e.g. from your knowledge base) shall be used</li><li>Bigger amounts of text shall be considered</li></ul> |
| **Large Action Models** | Implements a **feedback cycle** provided by another model | <ul><li>From simple compilation checks</li><li>To complex validation via second LLM</li></ul> |
| **Fine-Tuning the Foundation model** | Use additional data to **train the foundation model** new knowledge | <ul><li>Learn on domain knowledge</li><li>Train on specific best-practice patterns</li></ul> |

eraneos

# Large Action Models (LAMs) focus on **understanding actions and orchestrating sequences of actions** to accomplish specific goals without requiring assistance

Engineering Flow of a multi-stage, graph-based AI approach



Query

*Query containing the task description (and the required format)*

**Query Analyzer**

*Analyses which tools are available and when to use them*

*Chooses generator tool*

Docs    Schema

**Generator**

*Gets schema, docs and few-shot examples*

*Iterate until validator approves; this can take multiple rounds*

**Validator**

*Receives the generated output, and validates the proposed solution, e.g., checks the correctness of an SQL statement or the conformity with specific conventions*

**Responder**

*Interprets the results / the information received from the validator and converts it into the required format*

Response

*Response including the results / a successful task execution message*

Tool B

Tool C

Tool D

*Note: This step can be modularized even more and/or connected in additional series, or enriched by additional tool access (this might improve the response quality but comes in with increased complexity and costs)*
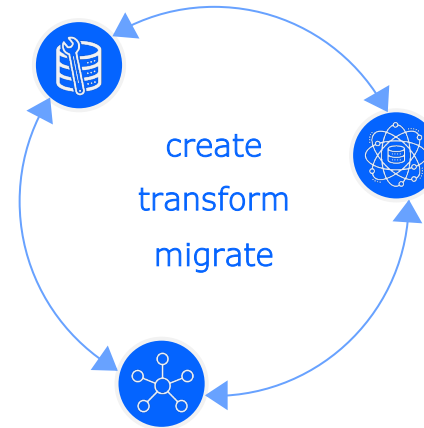
eraneos

# Looking Ahead

What we expect from the development in the field of GenAI in 2024 and beyond



eraneos

Our sample evaluation of general-purpose LLMs shows **impressive out-of-the-box functionality** on data modeling and sheds light on how **next level data modeling automation** might look like

**3rd Normal Form (3NF)**
**according to Inmon**

Write optimized

**Star schema**
**according to Kimball**

Read optimized

create
transform
migrate

**Data Vault**
**according to Linstedt**

Agile approach

We believe that **modeling, optimizing and migration** between data models is soom going to be **largely supported** by specialized tools

eraneos

# We have found **strengths and weaknesses with LLM's** in the context of data management

**Recommendations**

Be **specific about your target** model to get meaningful results

Use tools that **close the gap** between your working environment and the LLM API

Anticipate the **cost aspect** before setting up a scalable solution

Leverage **LLMs to get you started** but do not rely on end-to-end solutions yet

**Limitations**

Don't expect LLMs running in **business critical** production environments soon

**Large-scale automation** is not yet supported by tools – you need to built it yourself

eraneos

We have **6 general hypotheses on GenAI** for the year 2024 and believe we are going to see **productionized AI features in DWA tools by end of this year**
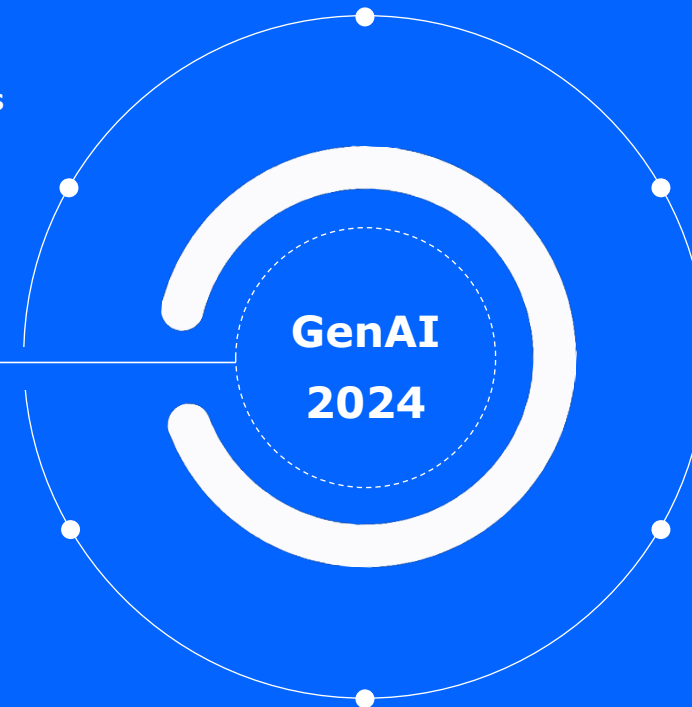
**Smaller, built-in Expert Models**

- Model sizes decrease, while maintaining same performance
- Efficient, specialized models e.g. for data vault modeling in specific domains

**Larger Generalist Models**

- Large foundation models
- Trend to larger context lengths

**Maintained Market Dominance of Big Tech**

- OpenAI, Google, Meta, Anthropic, Mistral, etc.

**GenAI 2024**

**Multimodality will unify**

- Text, image, audio, video, structured data
- Unified and simplified solutions from just reading to seeing, hearing, speaking

**New Skill Requirements**

- Interaction with AI for all
- Not prompt engineer, but standard skill for everyone

**Guardrailing & Security**

- Protection against malicious attacks (e.g., prompt injection)
- Prevention from extraction of raw training data or other confidential information

eraneos