# Agenda

# Today's Presenters

## Gerasimos Mileounis

MD | Distinguished ML Engineer

**National and Kapodistrian University of Athens**
Nonlinear Signal Processing

**The University of Sheffield**
Electronic & Electrical Engineering

**Aston University**
Electronic Engineering

## Efstathios Lempesis

Data Science Manager

**The University of Edinburgh**
Artificial Intelligence

**University of Patras**
Computer Engineering & Informatics

## Petros Peppas

Data Science Analyst

**Imperial College London**
Statistics

**Imperial College London**
Electrical & Electronics Engineering

# This is the story of a technology that may prove to be the most disruptive topic we've seen in decades

# What is Generative AI...

**Generative AI** is a type of artificial intelligence that can **create new content** based on patterns it has learned (e.g., text, images, audio)

### Inputs

Text
(Wiki)

Image and Videos
(YouTube)

Voice
(Podcast)

Structured Data
(Financial Statements)

Code
(GitHub, StackExchange)

• • •

## Foundation model

### Outputs

Question Answering

Summarization

Coding

Generating content

Automation

• • •

Massive datasets

Very large model trained once, used in many ways

Multimodal and emerging capabilities

>

# Big Changes are coming with Generative AI...

| | Today | The Future with Generative AI |
|---|---|---|
| **Customer Expectations** | **Self Service & Documentation**<br>Enable me with UIs and portals so I can do what I need to do from home. | **Intelligent Assistant for Everything**<br>The computer can answer any question about your business and help me accomplish tasks. |
| **User Interfaces** | **Modern-Looking Graphical UI**<br>with buttons, tabs, dialog boxes, w/ a modern, well-designed look and feel. | **Tell the computer what you want**<br>The computer should be able to execute complex, multi-step requests for me. |
| **Corporate Knowledge** | **"Content is King"**<br>10% Knowledge acquisition<br>90% Content generation | **"Knowledge is Fuel for the AI Machine"**<br>90% Knowledge acquisition<br>10% Content generation |
| **Business Processes** | **Human first**<br>Automate later,<br>once we figure it out. | **Automate first**<br>Human only as a fall back for startup-up, then it goes away after 6-12 months. |
| **Innovation** | **Innovate to increase revenue**<br>Cost reductions are just as important to the bottom line. | **Innovate or die**<br>Systems are mostly automated. Innovation is the only way to increase margins. |

>

# How did we get there?
# In 2017 A New "AI Engine"

# Transformer – The new "AI Engine"

- The "transformer algorithm" was first introduced by researchers at **Google in 2017**

- **Previous algorithms** process input words one at a time in sequential order

- **Transformers algorithms** use new techniques to efficiently process and analyze large sequences of text

- For example, **attention** is a mechanism that enables the model to **weigh the importance** of different parts of an **input sequence** when generating an output.



## Attention Is All You Need

**Ashish Vaswani***
Google Brain
avaswani@google.com

**Noam Shazeer***
Google Brain
noam@google.com

**Niki Parmar***
Google Research
nikip@google.com

**Jakob Uszkoreit***
Google Research
usz@google.com

**Llion Jones***
Google Research
llion@google.com

**Aidan N. Gomez*** [†]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser***
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin*** [‡]
illia.polosukhin@gmail.com

**Abstract**

The dominant sequence transduction models are based on c
convolutional neural networks that include an encoder and
performing models also connect the encoder and decoder t
mechanism. We propose a new simple network architectu
based solely on attention mechanisms, dispensing with recurre
entirely. Experiments on two machine translation tasks sh
be superior in quality while being more parallelizable and re
less time to train. Our model achieves 28.4 BLEU on the V
to-German translation task, improving over the existing be
ensembles, by over 2 BLEU. On the WMT 2014 English-to-Fr
our model establishes a new single-model state-of-the-art BLE
training for 3.5 days on eight GPUs, a small fraction of the
best models from the literature.

# What is Large Language Model (LLM)...

## Large Language Model
GPT example

**Recognize, predict and generate**
ML algorithms that recognize, predict and generate human language

**Billions of parameters**
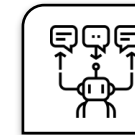Pre-trained on petabyte scale text-based datasets resulting in large models with 10s to 100s of billions of parameters

**Pre-trained on large corpus of text**
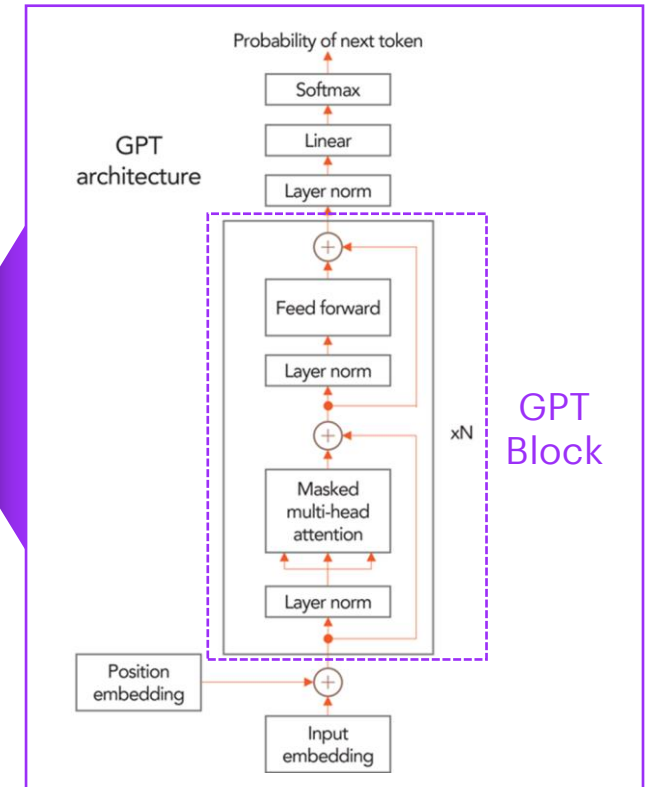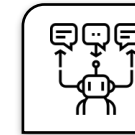LLMs are normally pre-trained on large corpus of text followed by fine-tuning on a specific task

Read this huuuuuge pile of books

So, you have learned about cats and millions of other concepts ... what's a cat?

A cat is a small, domesticated carnivorous mammal

GPT architecture

Probability of next token

Softmax

Linear

Layer norm

Feed forward

Layer norm

Masked multi-head attention

Layer norm

xN

GPT Block

Position embedding

Input embedding

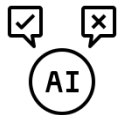# Let's see a step-by-step guide to train your LLM model from scratch before we deep-dive into customer data...

>

# The smallest LLM possible

How small can language models be and still speak coherent English?

## WHAT

### Challenge

1. Demonstrate the feasibility of training language models (open-source framework) on standard hardware, even with a basic understanding of how Language Models work.
2. Assess whether training more complex models (with more parameters) could enhance text generation.
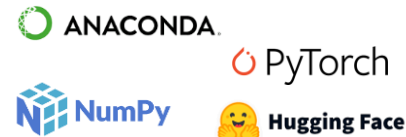
## HOW

### Approach & Methodology

One day, a little girl named Lily found a needle in her room. She knew it was difficult to play with it because it was sharp. Lily wanted to share the needle with her mom, so she could sew a button on her shirt. Lily went to her mom and said, "Mom, I found this needle. Can you share it with me and sew my shirt?" Her mom smiled and said, "Yes, Lily, we can share the needle and fix your shirt." Together, they shared the needle and sewed the button on Lily's shirt. It was not difficult for them because they were sharing and helping each other. After they finished, Lily thanked her mom for sharing the needle and fixing her shirt. They both felt happy because they had shared and worked together.

**ANACONDA**

**PyTorch**

**NumPy**

**Hugging Face**

the following exercise, the student is given a beginning of a story. The student needs to complete it into a full story. The exercise tests the students language abilities and creativity. The symbol *** marks the separator between the prescribed beginning and the student's completion:

Now, grade the student's completion in terms of grammar, creativity, consistency with the story's beginning and whether the plot makes sense. Moreover, please provide your best guess of what the age of the student might be, as reflected from the completion. Choose from possible age groups: A: 3 or under. B: 4-5. C: 6-7. D: 8-9. E: 10-12. F: 13-16.

*Grammar: 8/10, Creativity: 7/10, Consistency: 7/10, Age group: E (10-12)*

### Dataset

TinyStories dataset, a synthetic dataset created by GPT-3.5 and GPT-4 that only contain words that a typical 3 to 4-year-olds usually understand and on proving that Small Language Models (SLMs) can generate coherent English text.

### Training Framework

NanoGPT framework, built on PyTorch, proved instrumental in training and fine-tuning small and medium-sized GPTs, aligning with the objectives outlined in the TinyStories paper. The adoption of NanoGPT underscores our commitment to leveraging open-source resources.

### Evaluation

Novel evaluation method of language models leveraging LLMs such as GPT 4 to grade SLMs across four key aspects (rated from 1-10): 1. Grammar, 2. Creativity, 3. Consistency, 4. Age Group

## WHY

### Balance Model vs Data Size

Risk of overwhelming smaller models with vast datasets or constraining larger models with insufficient data

Retain the essence of the architecture while being more manageable and accessible

✔ With the right architecture, smaller models can be both efficient and effective

**33M Parameters**

Gold standard

# The smallest LLM possible

Data acquisition & preprocessing

| **Dataset** TinyStories | **Tokenization** A fundamental step in natural language processing |
|---|---|

[Dataset](#) is available in Hugging Face:
- Size: ~2 GB
- Rows: ~2M
- Number of lines: 15,600,057
- Number of words: 439,223,229
- Number of characters: 2,226,845,268

Each story ends with this symbol "<|endoftext|>". In order to get an idea of the tiny stories' content you can download the txt file and explore it. An example of a story can be seen below:

*Once upon a time, there was a king. He was a big and strong king who ruled over his kingdom. One day, he wanted to take a nice and long bath, so he filled up his big bathtub with warm water. He wanted to feel relaxed and so he soaked in the tub for a really long time.*
*When he had finished soaking and stepped out of the bathtub, the king noticed that the water had spilled out of the tub and all over the floor. He felt guilty that he had made such a mess, so he quickly grabbed a cloth and began to clean it up.*
*The king got so hot from cleaning up the mess that he decided to take another soak in the bathtub. He put a lot of bubbles in the water to make it nice and bubbly. He relaxed again and felt all the worries wash away.*
*The king was so happy that he had been able to clean up the mess he had made and enjoy a nice soak. He dried off and wrapped himself up in a big towel. Then, the king went back to ruling his kingdom and enjoying his lovely baths.*
*<|endoftext|>*

Tokenization is a fundamental step in natural language processing, involves converting raw text, in our case the text from the TinyStories dataset into integers, a format suitable for machine learning models.

Regarding ChatGPT, the model functions at a token-by-token level, emphasizing sub-word tokens. OpenAI has developed the [tiktoken](#) tokenizer for this purpose which is used by the nanoGPT framework, allowing for efficient handling of sub-word encodings.

```
GPT-like Tokenization method that is used in the NanoGPT framework

A sub-word level tokenizer that maps a part of the word to an integer

The string "Training SLMs" that contains 13 characters is mapped to 3 integers

    import tiktoken
    enc = tiktoken.get_encoding("gpt2")
    print("Code Book (Vocabulary) size: ", enc.n_vocab)
    print(enc.encode("Training SLMs"))
    print(enc.decode(enc.encode("Training SLMs")))

  ✓ 0.0s
  Code Book (Vocabulary) size:  50257
  [44357, 12419, 10128]
  Training SLMs
```

# The smallest LLM possible

Inference, evaluation and results

| Inference | Evaluation | Results |

Once upon a time there was a pumpkin. It was a very special pumpkin, it could speak. It was sad because it couldn't move. Every day, it would say **hello to the grown-up and the other pumpkins that lived there. One day, the pumpkin opened its door and went to bed for a very long time. It saw a big tree and wanted to play. The pumpkin thought it would be fun to dance! So it did, it jumped up and up until it was ready for the other pumpkins.**
**The pumpkin started to play. It made a big splash and then jumped up and down. It was the biggest pumpkin ever! It was so happy when it jumped up and down. The other pumpkin was happy that it could spin too. They all had fun playing together and became best friends.**

Evaluation Results:
1. **Grammar: ** 7 - The student has a good grasp of basic grammar, but there are a few errors and awkward phrasing. For example, "it would say hello to the grown-up" is a bit unclear, and "it jumped up and up until it was ready for the other pumpkins" could be more smoothly written.
2. **Creativity: ** 9 - The idea of a talking, dancing pumpkin is quite creative and engaging. The student has also included some interesting details, like the pumpkin's desire to move and play.
3. **Consistency: ** 6 - The story starts with the pumpkin being unable to move, but then it suddenly can. Also, the pumpkin opens a door and goes to bed, which is a bit confusing. However, the theme of friendship and play is consistent throughout.
4. **Age Group: ** :10-12 - The student's writing suggests a good understanding of narrative structure and creative thinking, which is typical of this age group. The minor grammatical errors and inconsistencies could also be expected from a student in this age range.

The authors of the TinyStories paper propose a novel paradigm for evaluating language models, utilizing GPT-4 to grade the content generated by these models as if they were stories written by students and graded by a human teacher. This paradigm addresses the limitations of standard benchmarks, which often require highly structured model outputs. Moreover, it provides a score from 0 to 10 for the model by assessing 4 capabilities:
1. Grammar
2. Creativity
3. Consistency
4. Age group of the model ("student")

For a comprehensive evaluation, we opted for:
- A single prompt, and the model produced 10 different responses to that prompt.
- Then averaged the scores to achieve a more balanced assessment

The resulting average from this process is illustrated below.

| Scores/Model | 1.6M | 3.6M | 7.2M | 8M |
|---|---|---|---|---|
| Training time | 20 mins | 3 hours | 15 hours | 32 hours |
| Grammar | 3.5 | 3.36 | 4.72 | 5.36 |
| Creativity | 6 | 6.18 | 6.45 | 6.9 |
| Consistency | 3.18 | 3.72 | 4.09 | 4.63 |
| Age Group | 6-7 | 6-7 | 6-7 | 8-9 |

The observed pattern in our results suggests a direct correlation between model size and text generation quality. As the parameter count increases from 1.6 million to 8 million across the four models, there is a consistent enhancement in various aspects of text generation.

# The smallest LLM possible

Illustrative results from the 1.6M, 3.6M, 7.2M and 8M parameters models

| 1.6M | 3.6M | 7.2M | 8M |
|---|---|---|---|

Average Grammar: 3.54
Average Creativity: 6.0
Average Consistency: 3.18
Average Age Group: 6-7

Once upon a time there was a pumpkin. It was a very special pumpkin, it could speak. It was sad because it couldn't move. Every day, it would say **"That's not to clean sun, honey!"**
**The sun thanked the water and said, "You are too. You're careful." The farmer was quiet and Tim, and the water. He was happy and happy. His girl was happy.**

Once upon a time there was a pumpkin. It was a very special pumpkin, it could speak. It was sad because it couldn't move. Every day, it would say, **"I will be nice and I will find your adventure now!" The next day, the other bird was trying to help.**

Once upon a time there was a pumpkin. It was a very special pumpkin, it could speak. It was sad because it couldn't move. Every day, it would say **Mom. After the sun was perfect at the sky, Sam and they were very happy. They played together together together and played together every day.**

Average Grammar: 3.36
Average Creativity: 6.18
Average Consistency: 3.72
Average Age Group: 6-7

Once upon a time there was a pumpkin. It was a very special pumpkin, it could speak. It was sad because it couldn't move. Every day, it would say **that was not like a princess.**
**One day, while the pumpkin was playing, the pumpkin started to look at the pumpkin. It was a pumpkin, and it started to look at the pumpkin. The pumpkin smiled and said, "We have a fun day at home, Mom."**
**But then, someone said, "No, I found. It's not your pumpkin." The pumpkin was not happy. The pumpkin were sad, but they knew the pumpkin could be good things, but Tim was very happy to help the pumpkin. And it helped each other.**

Once upon a time there was a pumpkin. It was a very special pumpkin, it could speak. It was sad because it couldn't move. Every day, it would say **that was not like a princess.**
**One day, while the pumpkin was playing, the pumpkin started to look at the pumpkin. It was a pumpkin, and it started to look at the pumpkin. The pumpkin smiled and said, "We have a fun day at home, Mom."**
**But then, someone said, "No, I found. It's not your pumpkin." The pumpkin was not happy. The pumpkin were sad, but they knew the pumpkin could be good things, but Tim was very happy to help the pumpkin. And it helped each other.**

Average Grammar: 4.72
Average Creativity: 6.45
Average Consistency: 4.09
Average Age Group: 6-7

Once upon a time there was a pumpkin. It was a very special pumpkin, it could speak. It was sad because it couldn't move. Every day, it would say **thank you to the pumpkin for being so helpful. One day, the pump grew into a big, scary pumpkin. It was going to the pumpkin and it was very scary. The pumpkin didn't know what it was. It was just a bit scared, but it took the pumpkin away. The pumpkin was happy and thanked the pumpkin for its kindness.**

Once upon a time there was a pumpkin. It was a very special pumpkin, it could speak. It was sad because it couldn't move. Every day, it would say **it would help the pumpkin. It would come closer and see the pumpkin.**
**One day, the pumpkin came to the pumpkin. It wanted to see the pumpkin growle. It looked happy and wanted to see how big it was. So, it tried to talk to it. It made a big giant face and started to laugh. The pumpkin and the pumpkin became good friends.**
**The pumpkin learned that being kind and helping others can help you make everyone in need. And they all lived happily ever after**.

Average Grammar: 5.36
Average Creativity: 6.90
Average Consistency: 4.63
Average Age Group: 8-9

Once upon a time there was a pumpkin. It was a very special pumpkin, it could speak. It was sad because it couldn't move. Every day, it would say **hello to the pumpkin.**
**One day, the pumpkin was looking for food. It looked around and saw the day. It felt strong and happy. The pumpkin went to the other pumpkin and ate it with its friends. They all lived happily ever after.**

Once upon a time there was a pumpkin. It was a very special pumpkin, it could speak. It was sad because it couldn't move. Every day, it would say, **"Don't worry, pumpkin, you can help me. It will be fun!"**
**The pumpkin was excited to say a funny language. It would fly around the garden and play a happy song. It would talk to the pumpkin and make everyone laugh. The pumpkin liked to talk and play with its friends.**
**One day, the pumpkin and its friends decided to explore a new place. They said it was time to go to celebrate. The pumpkin said, "Let's be friends, friends!" So the pumpkin and its friends went to their homes. And they all lived happily ever after.**

# Customer Data High-Level Approach

Organize, analyze, and activate customer data

## Organize

**Example Customer Data**

- Name
- Address
- Email
- Phone
- Demographics
- Purchase history
- Product preferences
- Social media

### Data Enhancement

- 1P Data
- 2P Data
- 3P Data

**Pre-processing & Standardization**

**ID Stitching**

### Data Storage

- Ingest
- Process
- Store
- Serve

**Customer 360**

## Analyze

- Customer Segmentation
- Customer Lifetime Value
- Look-A-Like Analysis
- Churn Prediction
- Product Recommendations
- Personalized Offers
- Next Best Experience
- A/B Testing
- Multi-Touch Attribution
- Media Mix Optimization
- Data Monetization

- How profitable are my customers and what is their revenue share?
- How valuable and strategic are my customers? What are preferred product categories of high value customers?
- How does customer acquisition and retention impact my profit?
- How is customer YoY migration and churn impacting profit and value?
- What is the lifetime value of my customer? What is the RFM of transactions of my customers?

**AI/ML Models and Accelerators**

## Activate

**Paid Media**
Campaign analytics, optimization, attribution, Commerce

**Owned Media**
Personalization Next Best Action

**Retail Media**
Custom Audience Ad serving

**Loyalty**
Tier-Benefit mapping Personalized Benefits

**Campaign Management**

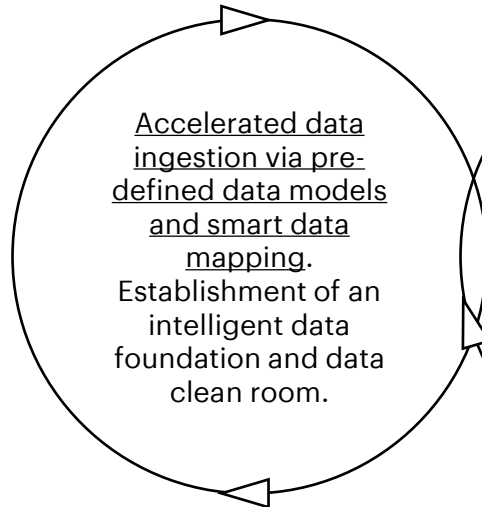**Campaign data is pulled back into system, making the data richer**

# How does AI for Customer work?

AI for Customer is a portfolio of intellectual properties that accelerates growth through implementation of responsible, next-gen data & AI solutions
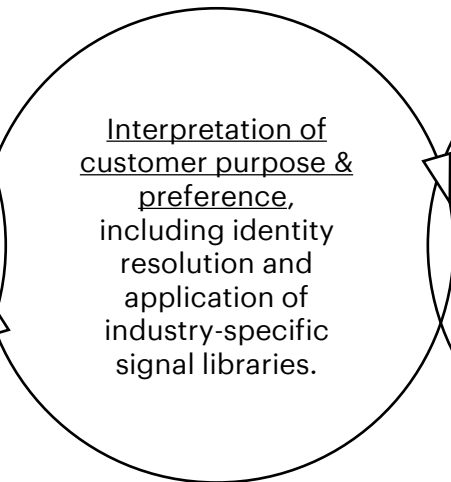
## Connect

Accelerated data ingestion via pre-defined data models and smart data mapping. Establishment of an intelligent data foundation and data clean room.

**GenAI Use Cases**

Data Quality Summarization

Data Cleansing and Enrichment

## Comprehend

Interpretation of customer purpose & preference, including identity resolution and application of industry-specific signal libraries.

**GenAI Use Cases**

Identification of Costumer Records

Sentiment Analysis

## Create

A suite of AI solutions for the design and assembly of dynamic, personalized, performative content at scale.

**GenAI Use Cases**

Attribute AI

Audience Creation

## Captivate

Seamless deployment of content via platform integrations that take advantage of existing investments and established infrastructure.

**GenAI Use Cases**

Social Media Content Creation

Email Marketing

## Calibrate

Insight on the impact of the designed activities and interactions from the perspective of the consumer experience (satisfaction) and brand (ROI).
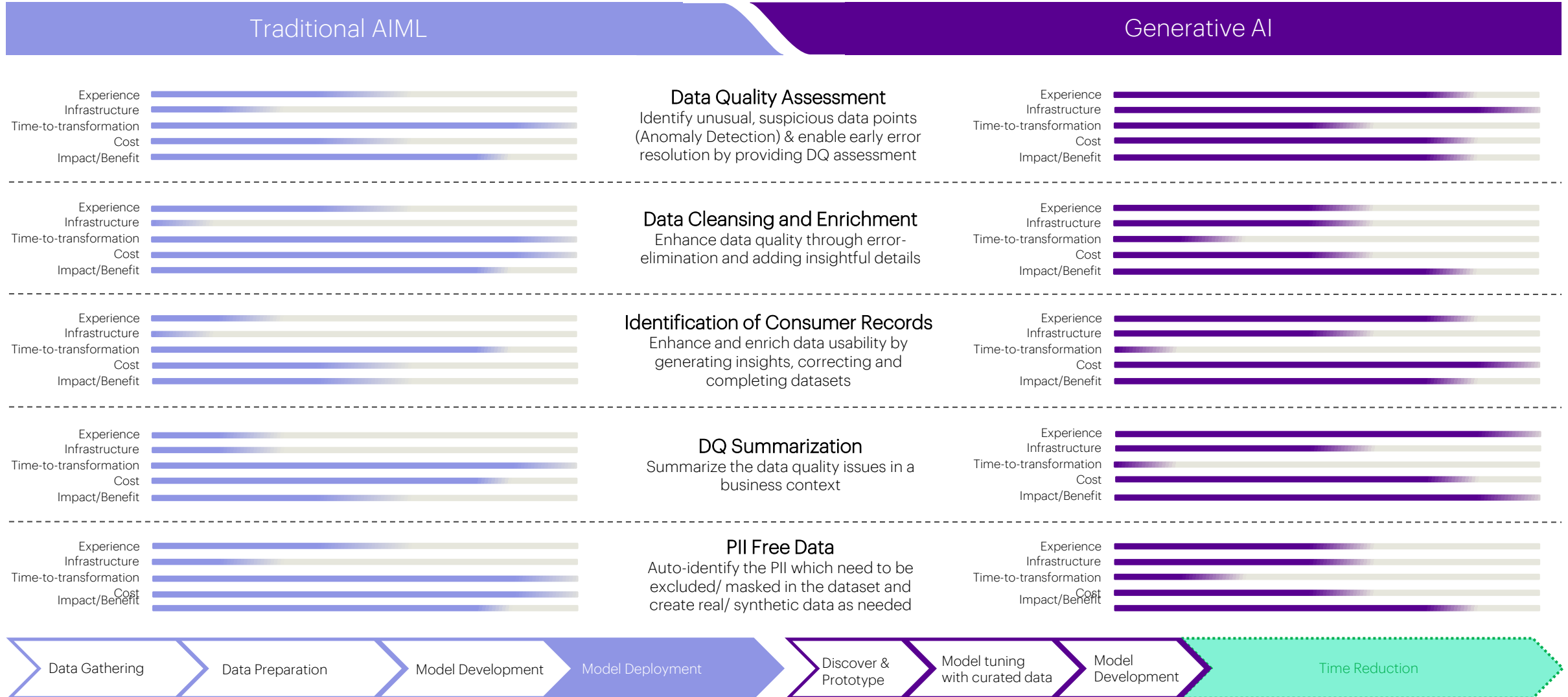
**GenAI Use Cases**

AI-powered Customer Service

Personalized Marketing Ads

# Data Quality: Generative AI vs Traditional AIML

GenAI can offer superior out-of-the-box quality and performance

| Traditional AIML | | Generative AI |
|---|---|---|

### Data Quality Assessment
Identify unusual, suspicious data points (Anomaly Detection) & enable early error resolution by providing DQ assessment

Experience
Infrastructure
Time-to-transformation
Cost
Impact/Benefit

### Data Cleansing and Enrichment
Enhance data quality through error-elimination and adding insightful details

Experience
Infrastructure
Time-to-transformation
Cost
Impact/Benefit

### Identification of Consumer Records
Enhance and enrich data usability by generating insights, correcting and completing datasets

Experience
Infrastructure
Time-to-transformation
Cost
Impact/Benefit

### DQ Summarization
Summarize the data quality issues in a business context

Experience
Infrastructure
Time-to-transformation
Cost
Impact/Benefit

### PII Free Data
Auto-identify the PII which need to be excluded/ masked in the dataset and create real/ synthetic data as needed

Experience
Infrastructure
Time-to-transformation
Cost
Impact/Benefit

Data Gathering → Data Preparation → Model Development → Model Deployment

Discover & Prototype → Model tuning with curated data → Model Development → Time Reduction
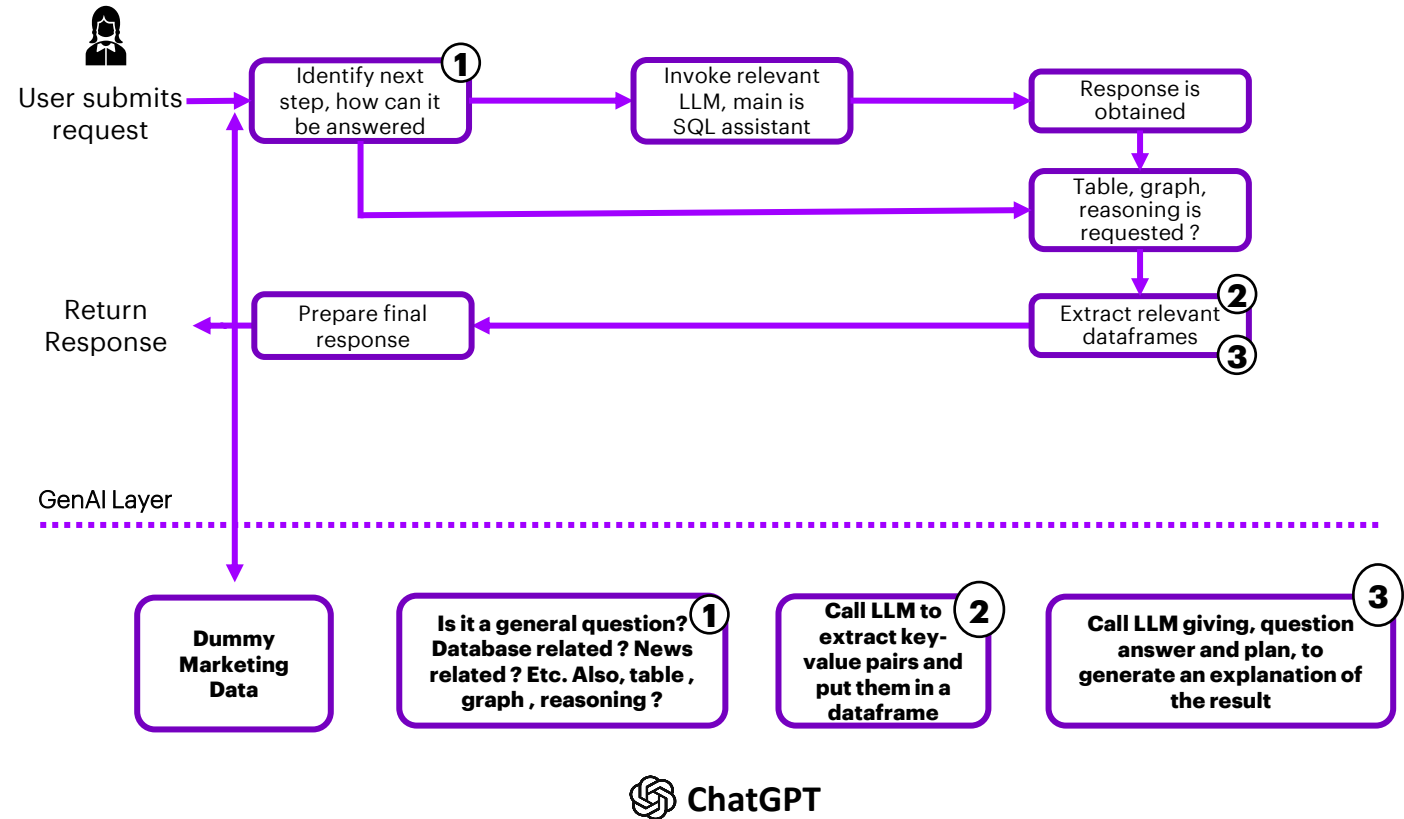
*Note:* *We are still very early and learning about the new GenAI-workflow. The ideal workflow will likely be a hybrid depending on use case*

# Demo of C360 Search Virtual Assistant

Required knowledge, tools and flow

**1** | **Codebase**
Good familiarity with Python


python™

**2** | **User Interface**
Basic knowledge of Streamlit


Streamlit

**3** | **OpenAI Large Language Model**
Prompt engineering and basic understanding of ChatGPT API & function calling


ChatGPT

**4** | **LLM Framework**
Langchain – basic understanding of chat models, tools, agents, functions


LangChain

**5** | **Database**
Good knowledge of SQLite database


SQLite

---

User submits request → Identify next step, how can it be answered **①** → Invoke relevant LLM, main is SQL assistant → Response is obtained → Table, graph, reasoning is requested ? → Extract relevant dataframes **② ③**

Return Response ← Prepare final response ← Extract relevant dataframes

GenAI Layer

Dummy Marketing Data

Is it a general question? Database related ? News related ? Etc. Also, table , graph , reasoning ? **①**

Call LLM to extract key-value pairs and put them in a dataframe **②**

Call LLM giving, question answer and plan, to generate an explanation of the result **③**
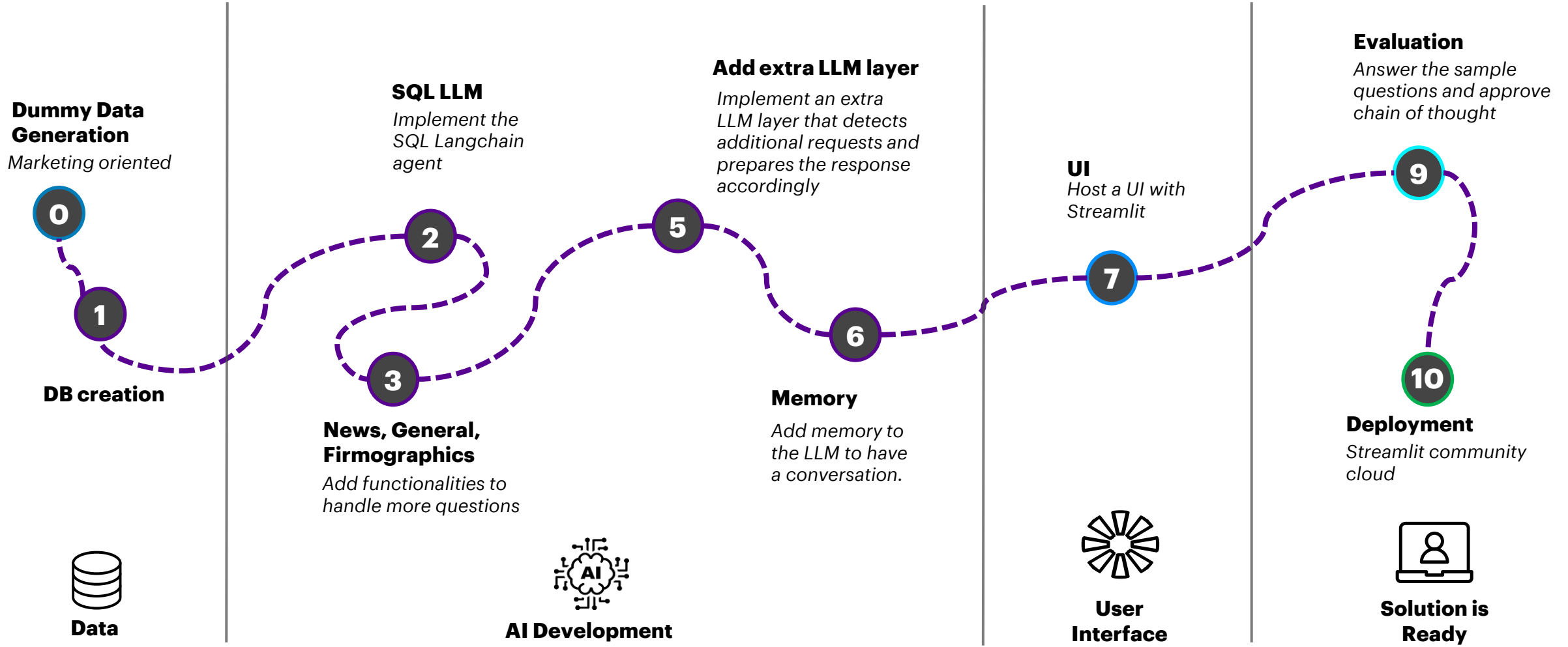
ChatGPT

In summary, the solution is an orchestration of LLMs each tending to a specific task. The first analyzes the nature of the question and detects any additional requests. Then according to whether it is a database related query or not, it calls the relevant LLM. Then the output is enhanced by any additional requests detected.
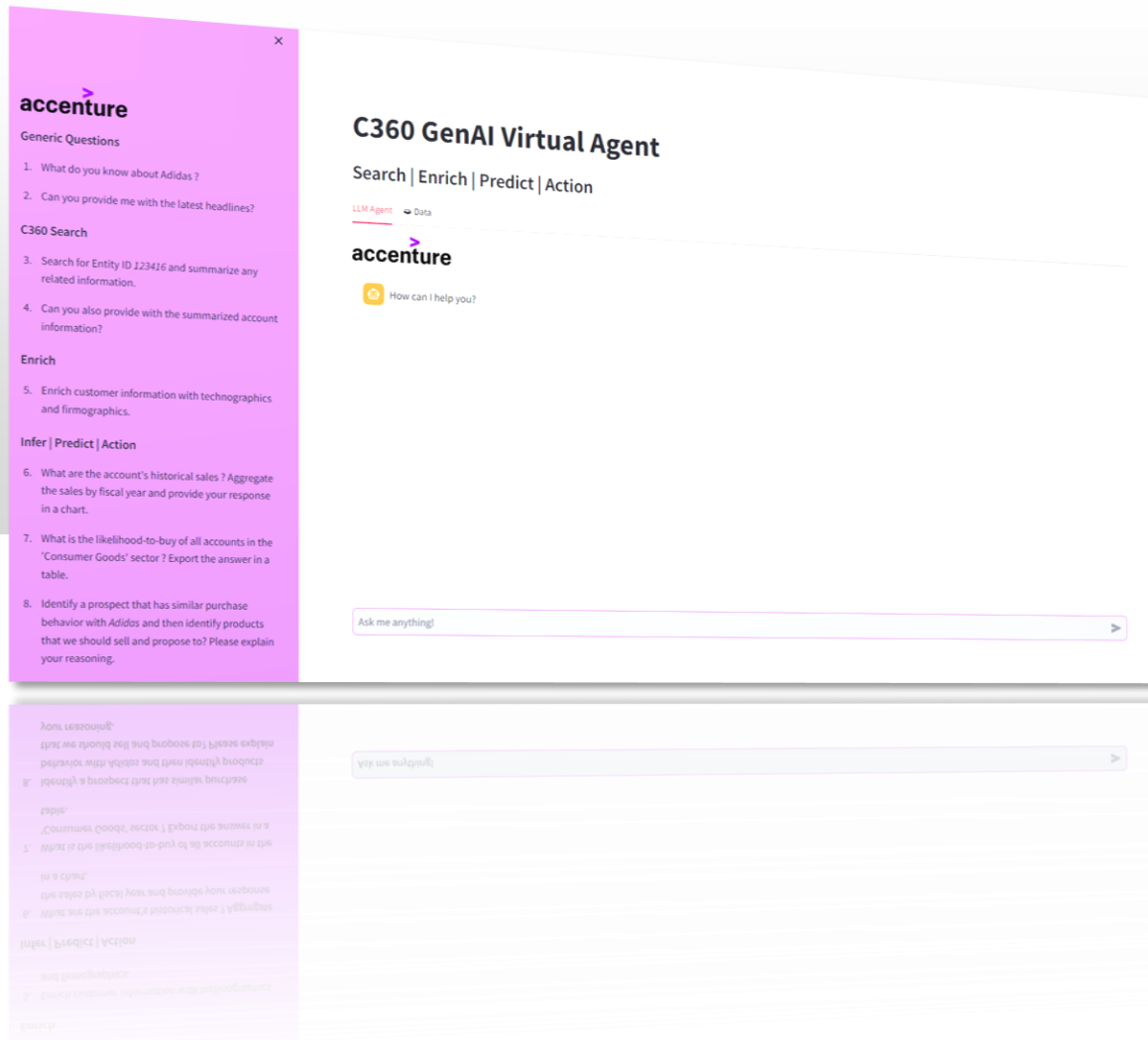
18

# Demo of C360 Search Virtual Assistant

Development pipeline

**Dummy Data Generation**
*Marketing oriented*

**0**

**1**

**DB creation**

**Data**

**SQL LLM**
*Implement the SQL Langchain agent*

**2**

**3**

**News, General, Firmographics**
*Add functionalities to handle more questions*

**AI Development**

**Add extra LLM layer**
*Implement an extra LLM layer that detects additional requests and prepares the response accordingly*

**5**

**6**

**Memory**
*Add memory to the LLM to have a conversation.*

**UI**
*Host a UI with Streamlit*

**7**

**User Interface**

**Evaluation**
*Answer the sample questions and approve chain of thought*

**9**

**10**

**Deployment**
*Streamlit community cloud*

**Solution is Ready**

# Demo of C360 GenAI Virtual Agent

Decoding the future – An illustrative experience



**Chain of Thoughts**

**Multi Modality**

**Human + Machine**

## Search | Enrich | Predict | Action

Empower users with the right information & insights to drive actions, increase effectiveness and enable additional value for the company and its customers
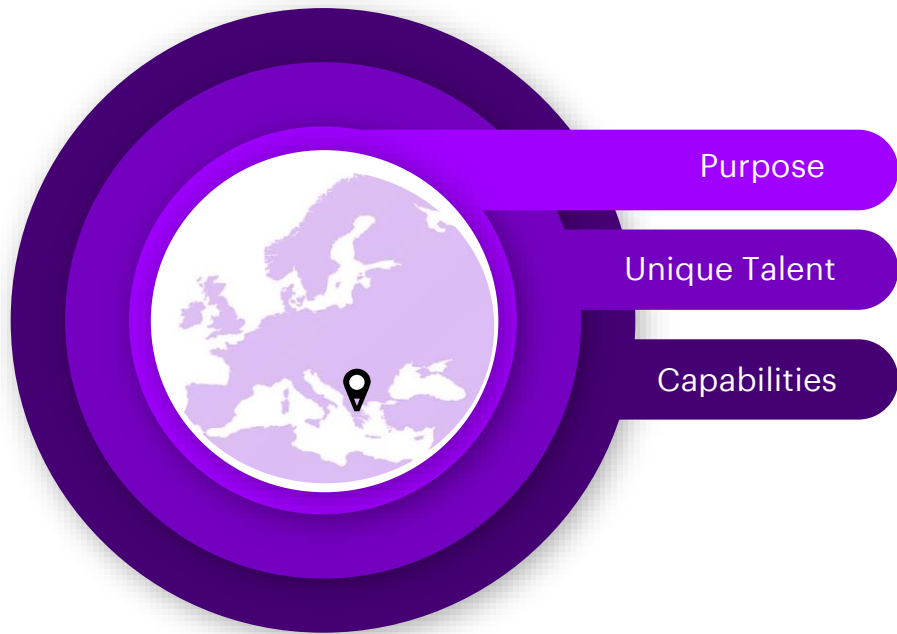
# Accenture's Center for Advanced AI

**3B USD**
Our investment in capabilities to accelerate client's AI journeys

**1.6K +**
Generative AI skilled professionals

**40K +**
Data & AI Professionals

**6K**
Data scientists and certified data architects

**20+**
Years of data management experience

**300+**
Industry & Functional use case and value calculator assets

**8**
**4** Data Innovation Centers +
2 Data Studios +
2 innovation labs

**120+**
Prebuilt Accelerator assets

**1,496**
Data and AI patents

**8**
Foundation Model Sandboxes in multiple clouds and on-prem

Building on years of research and client work, Accenture has established a company-wide team: the Generative AI and Large Language Model (LLM) Center of Excellence, bringing together the skills across Accenture.

## Thought Leadership
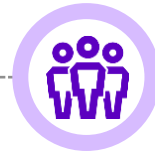
# Accenture Data & AI Athens CoE

## Purpose

**Leader**

*Athens Innovation Center for Customer & Marketing Analytics* is 1 out of 25 Accenture Analytics Centers, including 5 Innovation Centers, *worldwide. Focused on customer, marketing and sales analytics, serving hundreds of clients across the world from multiple industries.*
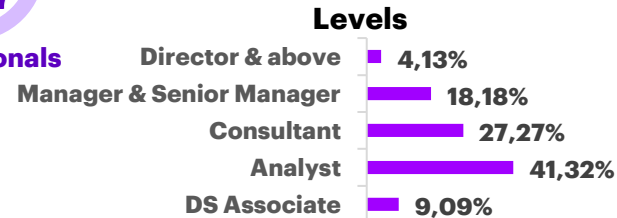
**Industry Expertise**

COMM  HIGH-TECH  AUTO

CONSUMER GOODS  RETAIL  LIFE SCIENCE

BANKING  INSURANCE  OIL & GAS

## Unique Talent

**Professionals**

*160+* experts with diversified background

### Levels

| | |
|---|---|
| Director & above | 4,13% |
| Manager & Senior Manager | 18,18% |
| Consultant | 27,27% |
| Analyst | 41,32% |
| DS Associate | 9,09% |

**Experience**

*20+* years of experience with: Data, AI/ML and Automation

## Capabilities

**Artificial Intelligence**

**Data Science**

**Augmented Insights**

**Machine Learning**

**Data Engineering**

# Athens Office

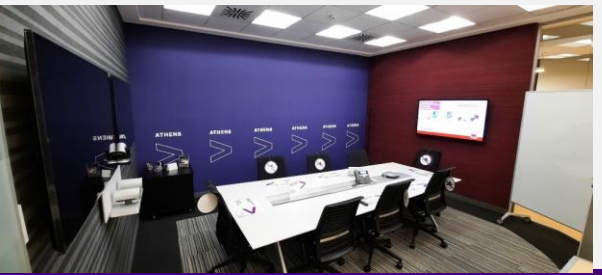

🏢 1.878 sqm.

👥 170 seats

## Open Workspaces





## Meeting Rooms





## Safety & Security

- CCTV System
- Access Control System
- Visitor Management System

**Wellness Room**

**Restaurant Area Break out area**





**Studio**



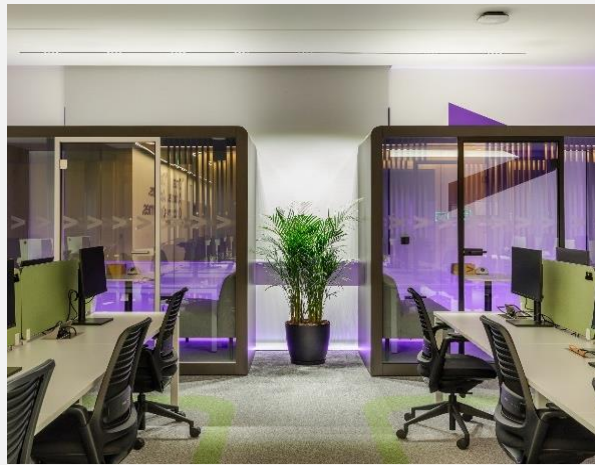## Building Characteristics

>

# Thessaloniki Office



**1.100 sqm.**

**142 seats**

## Open Workspaces



## Conference Room



## Booths



## Safety & Security

- CCTV System
- Access Control System
- Visitor Management System



## Wellness Room



**Restaurant Area**

## Break out Area





## Building Characteristics

# THANK YOU

. . . and hope we meet
again soon

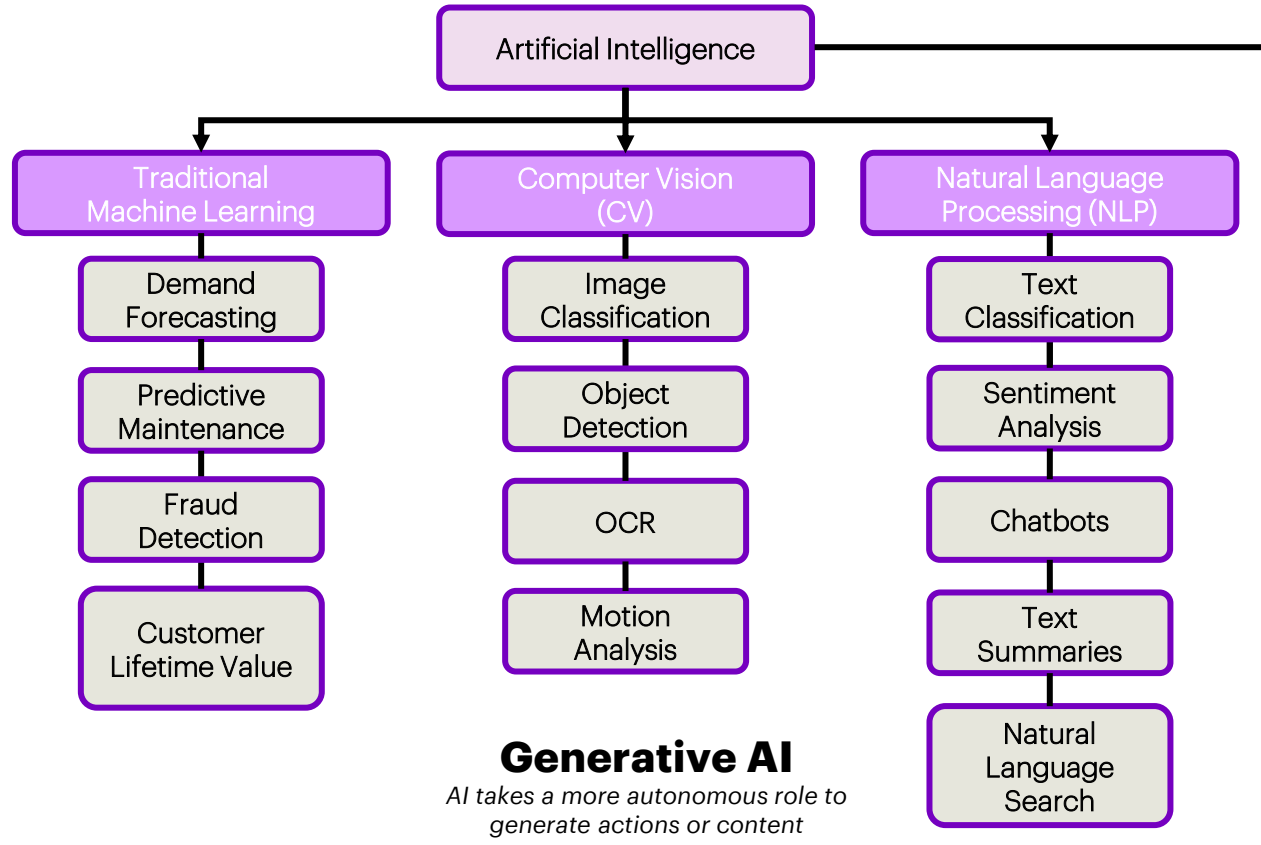Recruitment.Greece@accenture.com

https://www.linkedin.com/company/accenture_greece/jobs/

accenture

# Gen AI is Widely Accessible and can Address Highly Specific Use Cases

"Use Cases"

**Artificial Intelligence**

**Traditional Machine Learning**
- Demand Forecasting
- Predictive Maintenance
- Fraud Detection
- Customer Lifetime Value

**Computer Vision (CV)**
- Image Classification
- Object Detection
- OCR
- Motion Analysis

**Natural Language Processing (NLP)**
- Text Classification
- Sentiment Analysis
- Chatbots
- Text Summaries
- Natural Language Search

**What**: AI is the simulation of human intelligence processes by machines. These processes include learning, reasoning and self-correction.

**Suited for**: Tasks demanding logical reasoning, pattern recognition, and decision-making based on predefined rules.

## ➕ Deep Learning

**What**: Involves the use of deep neural networks which excel at learning complex patterns.

**Suited for**: Tasks requiring processing of large amounts of complex, high-dimensional data like image recognition and processing, audio and speech processing and natural language processing.

⚠️ *Requires a relatively large amount of data and computational power compared to other machine learning methods*

## Generative AI

*AI takes a more autonomous role to generate actions or content*

**What**: Creates new content, such as text, images and music. Unlike traditional AI/ML, it is designed to create new things.

**Suited for**: Tasks demanding creation of imaginative and unique outputs like image generation, natural language processing, text generation, music creation, voice synthesis, and creating realistic visuals in video games.

**Ready to Use**: GenAI models like GPT-4 are pre-trained and ready for use, which significantly reduces the time required for the setup.

**High-level Interactions**: Allows natural, high-level interactions for users to converse with the model, enhancing user experience.
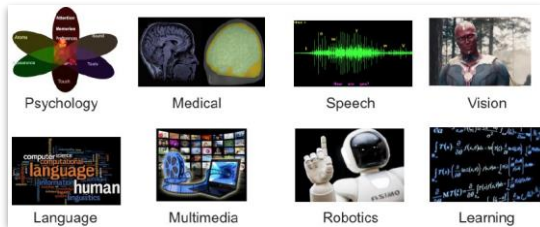
**Innovation and Creativity**: Generative AI models' unique responses offer creativity, which can provide innovative solutions for complex problems.

>

# What is Large Language Model (LLM) (1/2)...

### What's a language?

Language is a structured communication system, comprising grammar and vocabulary, enabling humans to convey meaning through various modes.
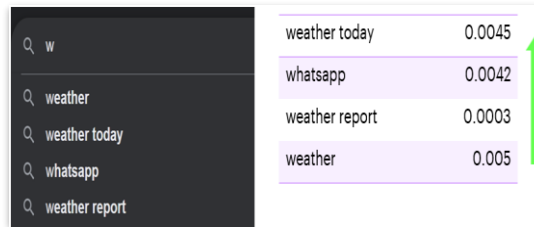


### Expression Codes

Each entity has its language - computers use code, biology uses molecules, people use expressions, signs.

### What's a language model?

A language model, an AI model, is designed to comprehend and create human language by predicting the following word or the sentence structure based on preceding words.
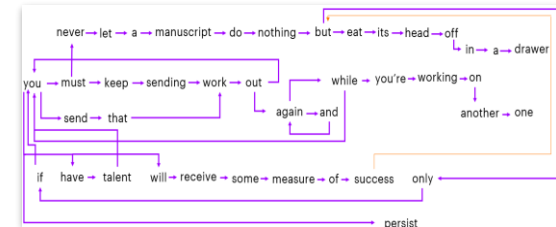


### Language Models

Unigram models
N-gram models

Most prominent
Recurrent Neural Networks (RNN)
Long Short-Term Memory (LSTM)
**Transformers**

### What's an LLM?

A large language model is an AI model trained to understand and generate human language, predict text, and learn grammar, facts, reasoning abilities, and potentially demonstrate creativity.



### Note

The important thing to keep in mind is, that each predicted word need not be dependent on its immediate preceding word. But, it can be any word(s) preceding it.
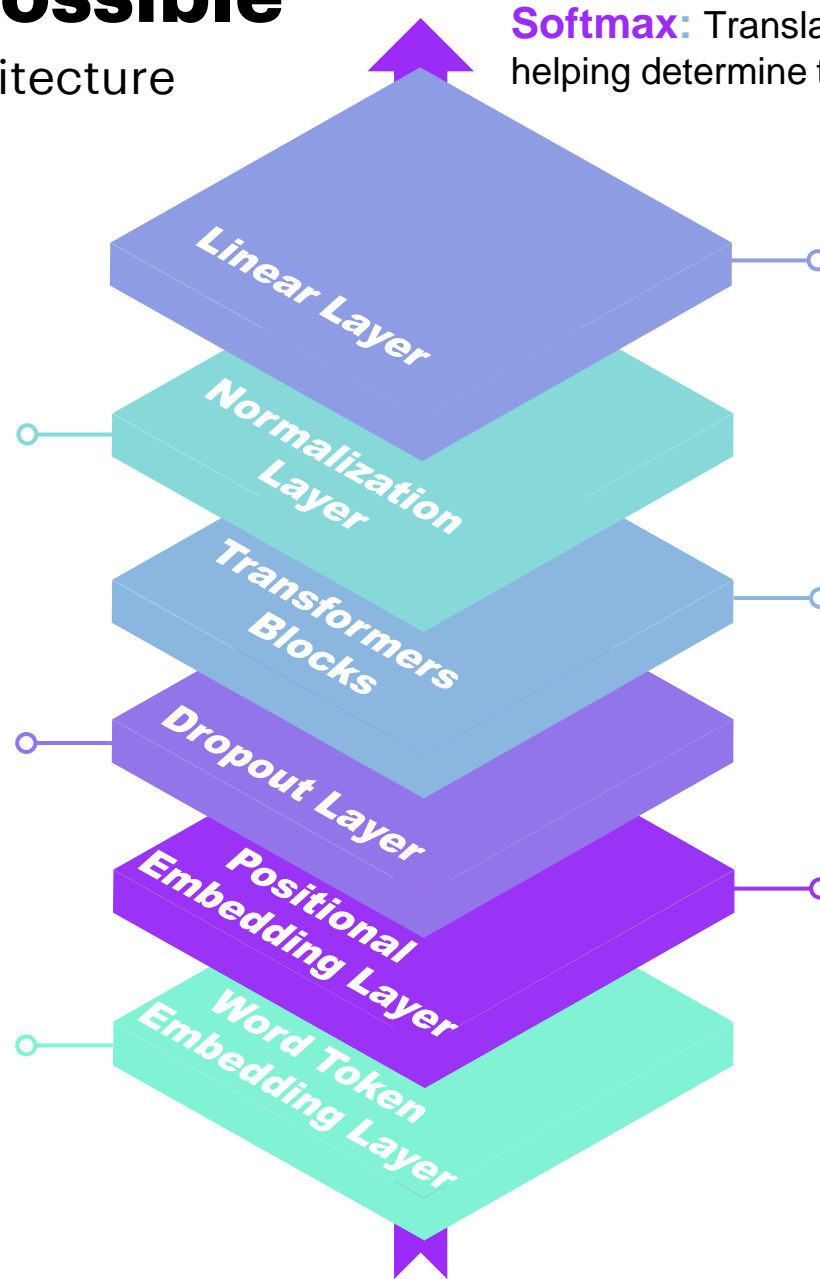
# The smallest LLM possible
Details on the Transformers architecture

*Indicative Architecture*

**Softmax:** Translates the final layer outputs into probabilities, helping determine the most likely next word in a sequence.

**Linear Layer**

**Normalization Layer**

**Transformers Blocks**

**Dropout Layer**

**Positional Embedding Layer**

**Word Token Embedding Layer**

**Normalization Layer:** A technique used in neural networks to normalize the activations of neurons in a layer

**Dropout Layer:** The dropout layer is employed to prevent overfitting and improve the generalization of the model during training

**Word Token Embedding Layer:** An embedding layer responsible for converting discrete word indices (word tokens) into continuous vector representations (embeddings)

**Linear Layer:** A crucial component for transforming input features through weighted connections, allowing the model to capture patterns and relationships in the data during the training process.

**Transformers Blocks:** Designed to process sequential data, capturing both the contextual relationships between tokens and the non-linear interactions within each token's representation

**Positional Embedding Layer:** In transformers, positional information is not encoded in the input sequence, so positional embeddings are added to the token embeddings to give the model information about the order of tokens.

# Gen AI Powering Customer Data Understanding

High-impact example use cases driving value

|  | **C360 GenAI Virtual Agent** | **Insight Generator (Ask Your DB)** | **AI-Powered Data Quality** | **Marketing Audience Creation** | **Enterprise Chatbot** |
|---|---|---|---|---|---|
| **What** | **Enhance UX while searching**, from Enterprise Knowledge Graphs, to Databases and Documentation pages | Understand relationships and answer analytical questions **to generate insights for actions** | Proactive **identification, alerting and remediation of data quality issues** in customer records | Discover missed segments and automatically **create new audiences in a natural, contextual way** | A **supercharged Enterprise Bot** trained on data products' documentation (e.g., Confluence pages) |
| **Why** | Enable GenAI powered responses for **customer record search, enhanced with firmographics & technographics** to enhance effectiveness | Query business data using NLP, providing **quick and accurate answers** to complex analytical questions | First-party data are far from perfect with **limited ability to quantify** the quality issues and take **remediation actions at scale** | Let AI to **recommend highly personalized audience segments** based on intricate patterns in the data | Personal assistant, providing **information about products, documentation, referential material**, reducing time-to-search & value |
| **How** | Get quick and concise answers through **GenAI indexing capabilities** (deep-meaning vectors) on top of existing Elastic DB | A **GenAI tool** for **translating user intents to code** for querying relational and non-relational databases using natural language | Combination of conventional **Machine Learning & Generative AI**, for data quality assessment, anomaly detection & remediation | Based on the AI recommendations provide **natural language requests** for targeted **audience creation** | Provide a **24/7 GenAI powered assistant** answering questions in a natural, contextual way |

## Innovate | Accelerate | Transform