# A practical introduction to Generative AI

Speaker:

**Giovanni Rosa, PhD**

**AI Engineer**
@Technology Reply

🌐 giovannirosa.com

✉ gio.rosa@reply.it

# ITOUG TRAINING DAY @ GALDUS

**MILANO 23 OTTOBRE 2024**

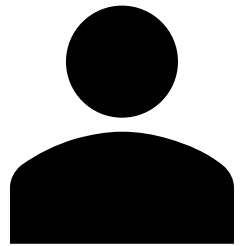UNLOCK THE FUTURE WITH ORACLE: MASTERING ANALYTICS, APEX, DATABASE AND GEN AI

**ITOUG**
ITALIAN ORACLE USER GROUP

SCAN ME

Bing DALL-E 3

Prompt:
"Generate an image logo for the Italian Oracle User Group 2024 event hosted by Galdus formazione in Milan city"

**User**

call

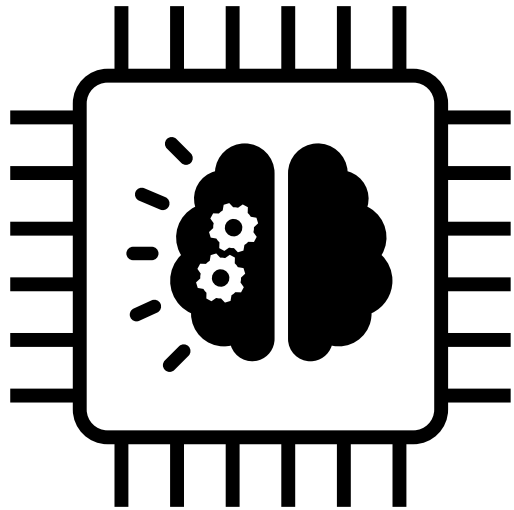**1**

*"Generate an image logo for the Italian Oracle User Group 2024 event hosted by Galdus formazione in Milan city"*
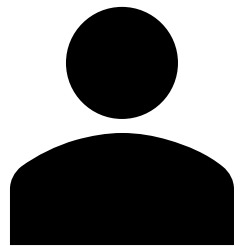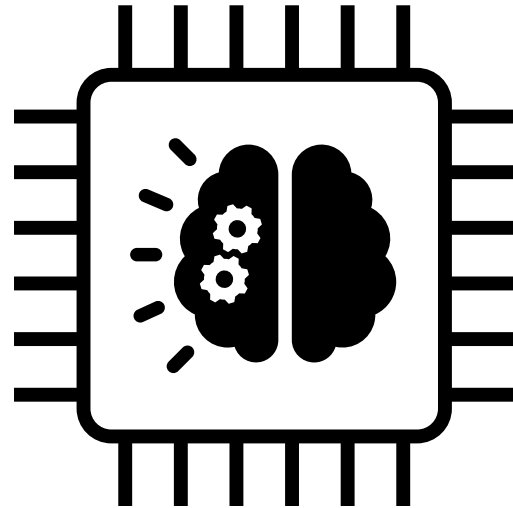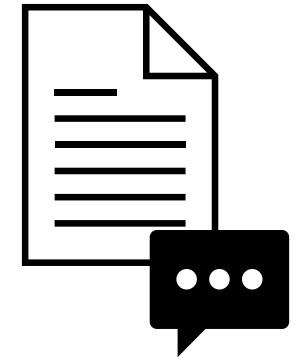
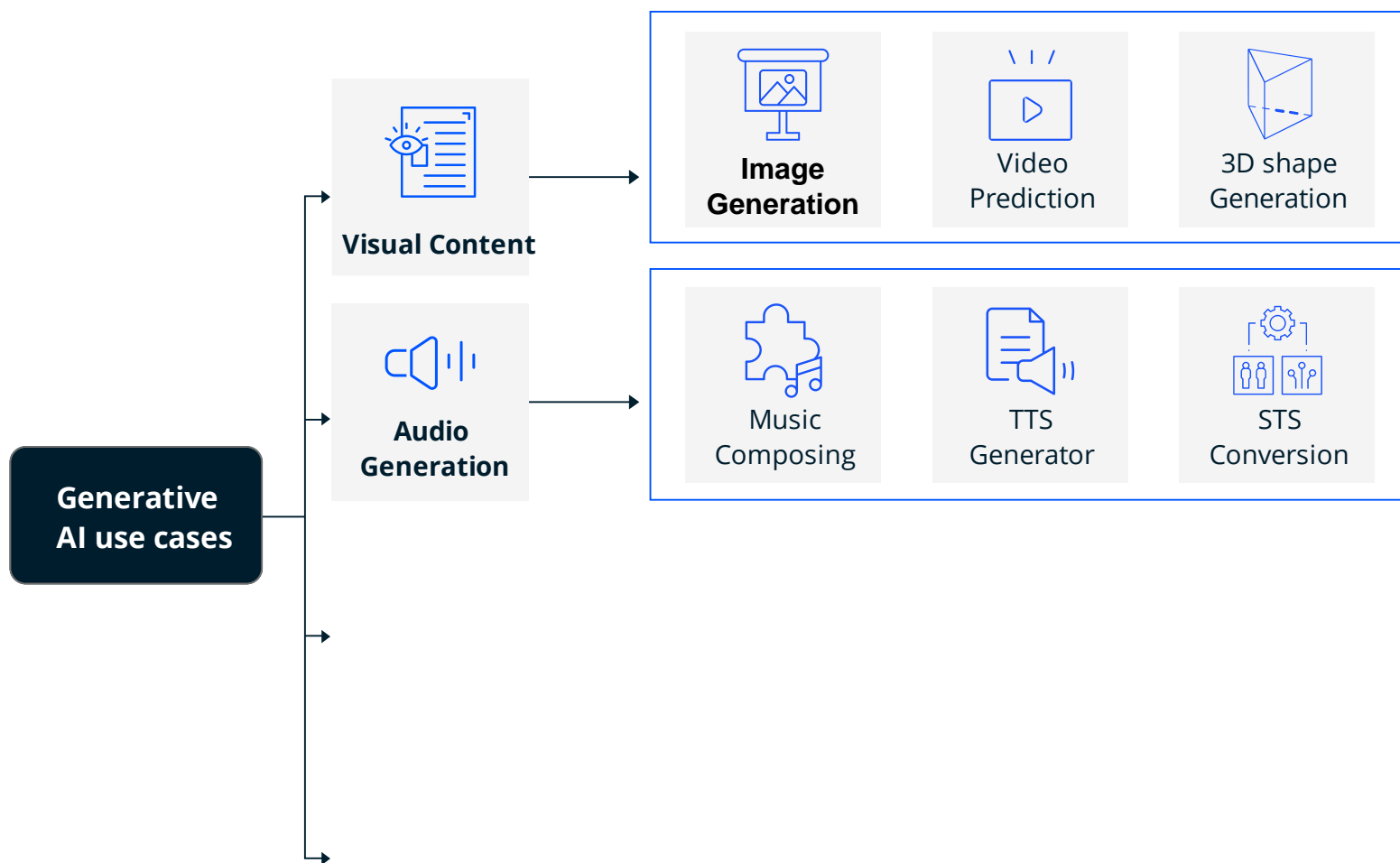User → instruction → Generative AI model → create → Response

Generative AI use cases

Visual Content
- Image Generation
- Video Prediction
- 3D shape Generation

Audio Generation
- Music Composing
- TTS Generator
- STS Conversion

Source: LeewayHertz

**Generative AI use cases**

**Visual Content**
- Image Generation
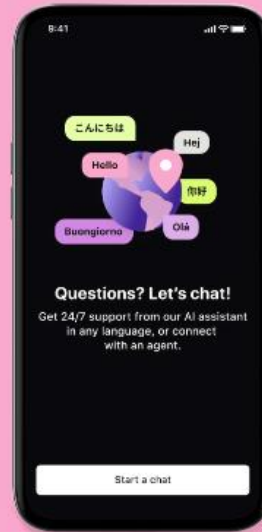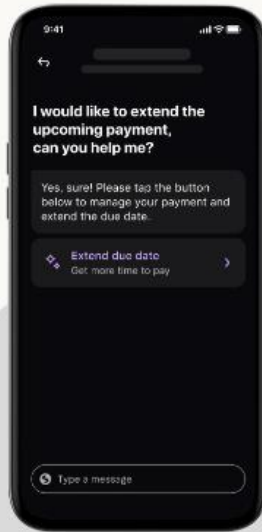- Video Prediction
- 3D shape Generation

**Audio Generation**
- Music Composing
- TTS Generator
- STS Conversion

**Text Generation**
- Creative Writing
- Chatbots
- Summarization

**Code Generation**
- Code Generation
- Code Documentation
- Bug Fixing

Source: LeewayHertz

#1

**#2** Coding Assistants

```typescript
// create a class in TypeScript to represent a student that has a name, an id, and a list of courses
// the class should have method to add and remove a course to the list of courses

class Student {
  name: string;
  id: number;
  courses: string[];

  constructor(name: string, id: number, courses: string[]) {
    this.name = name;
    this.id = id;
    this.courses = courses;
  }

  addCourse(course: string) {
    this.courses.push(course);
  }

  removeCourse(course: string) {
    this.courses = this.courses.filter((c) => c !== course);
  }
}
```

add typings and pydoc to this function

**GitHub Copilot**

> 3 steps completed successfully

```python
def create_worker_agent(llm: ChatOpenAI, tools: List, system_prompt: str) -> AgentExecutor:
    """
    Create a worker agent with the given language model, tools, and system prompt.

    Args:
        llm (ChatOpenAI): The language model to use for the agent.
        tools (List): A list of tools that the agent can use.
        system_prompt (str): The system prompt to initialize the agent with.

    Returns:
        AgentExecutor: An executor for the created agent.
    """
    # Each worker node will be given a name and some tools.
    prompt = ChatPromptTemplate.from_messages(
        [
            ("system", system_prompt),
            MessagesPlaceholder(variable_name="messages"),
            MessagesPlaceholder(variable_name="agent_scratchpad"),
        ]
    )
    agent = create_openai_tools_agent(llm, tools, prompt)
    executor = AgentExecutor(agent=agent, tools=tools)

    return executor
```
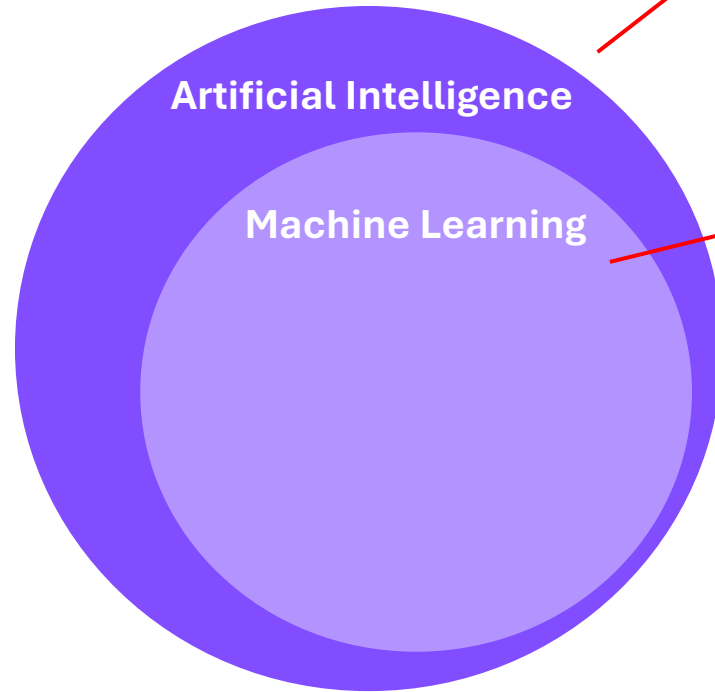
> Using 4 References

```python
def create_worker_agent(llm, tools, system_prompt):     You, Moments ago • Uncommitted changes
    # Each worker node will be given a name and some tools.
    prompt = ChatPromptTemplate.from_messages(
        [
            (
                "system",
                system_prompt,
            ),
            MessagesPlaceholder(variable_name="messages")
```

**#2**

# What is
# Generative AI

?

**Intelligent autonomous machines**

**Artificial Intelligence**

**Machine Learning**

**Pattern recongition**
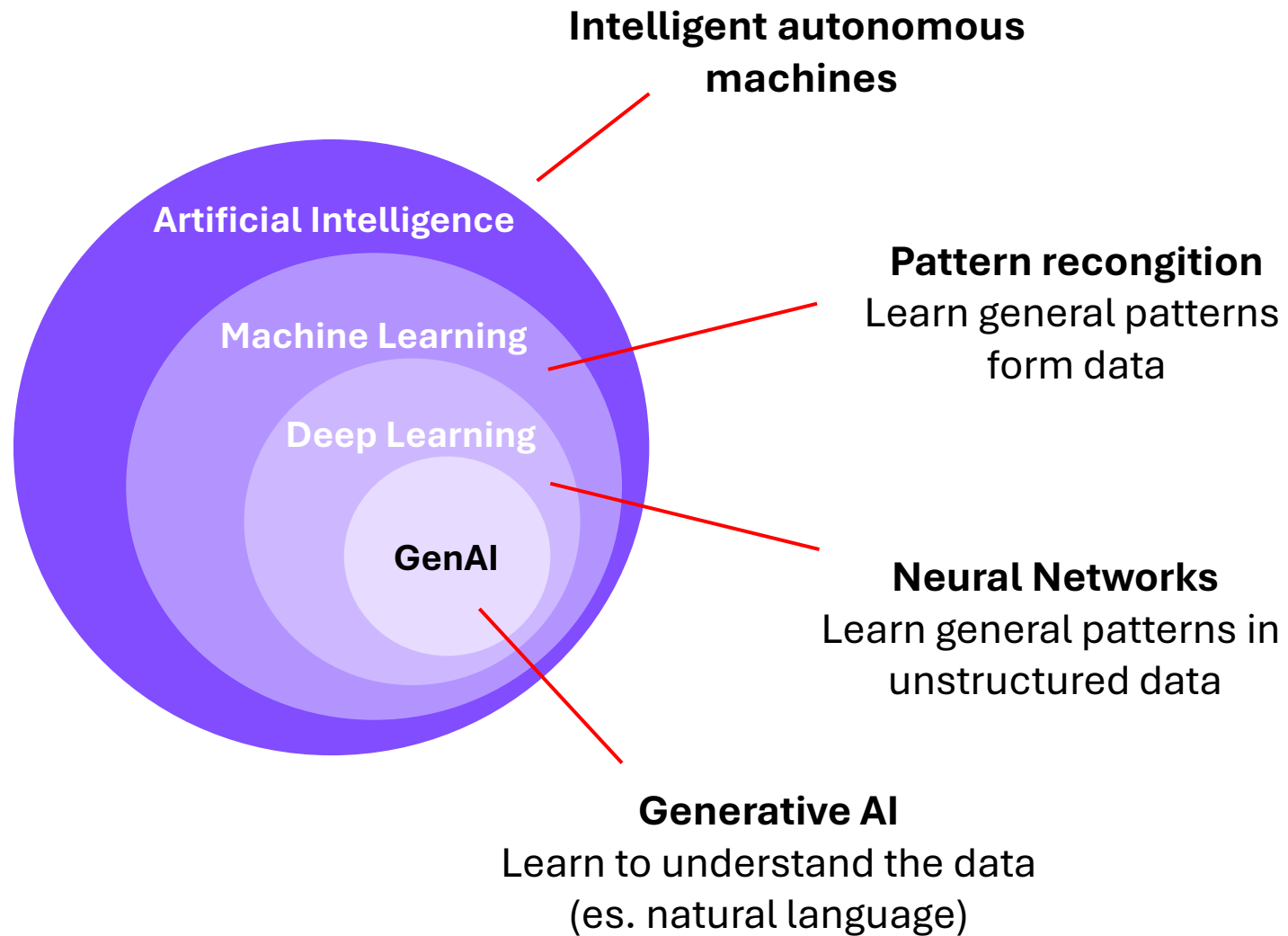Learn general patterns form data

My experience has been **fantastic!**

**POSITIVE**

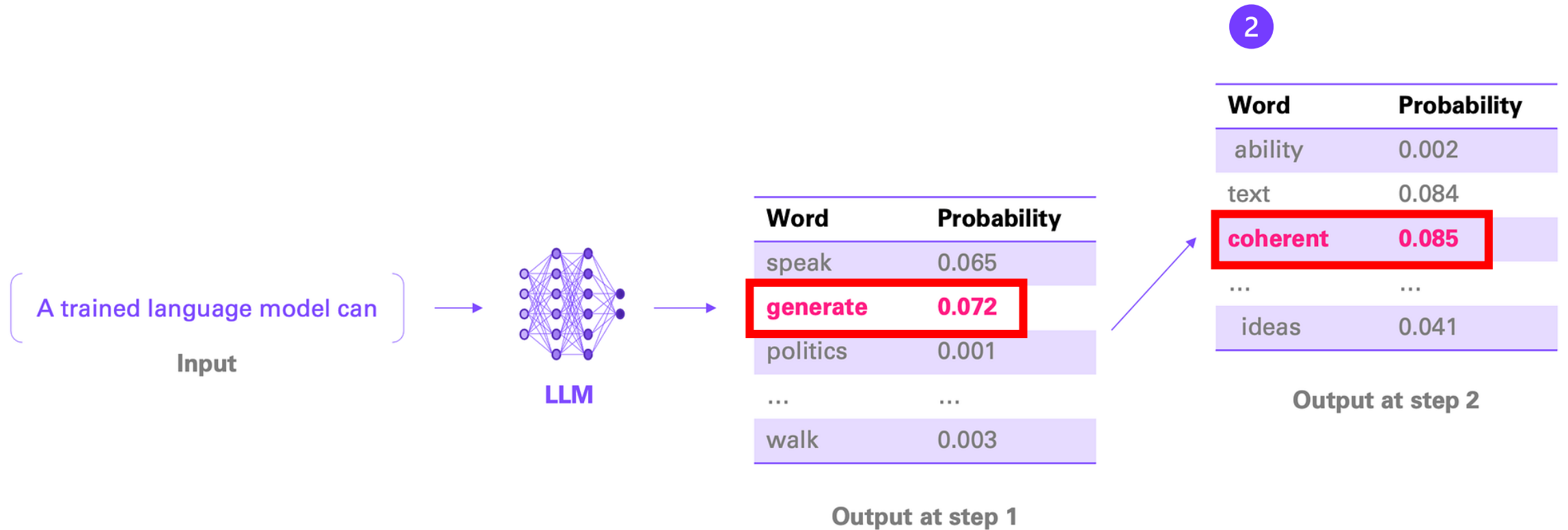The product is **ok I guess**

**NEUTRAL**

Your support team is **useless**

**NEGATIVE**

Large Language Model } LLM

# How does a **language model** work?
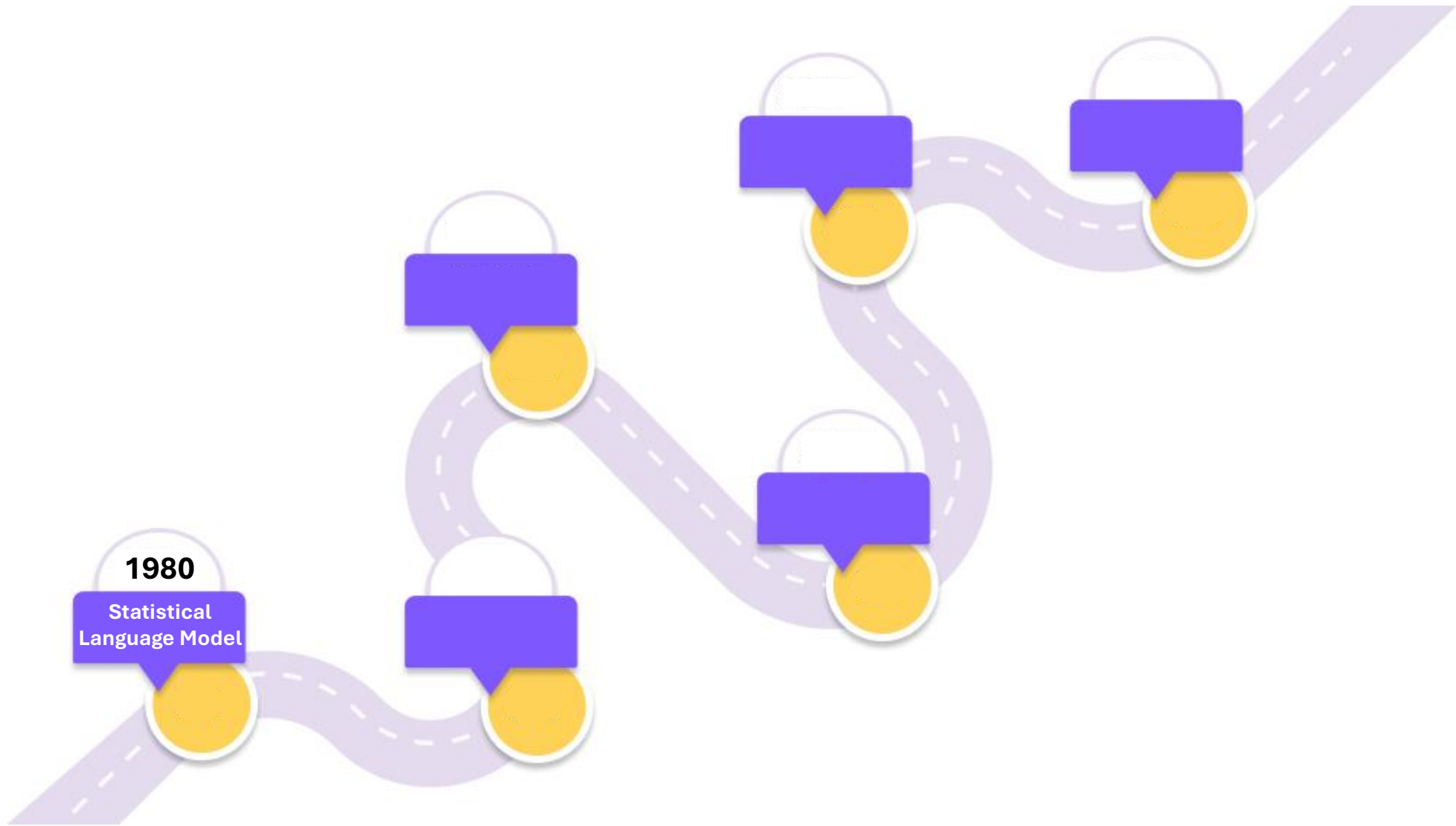
A trained language model can **generate coherent...**

| Word | Probability |
|------|-------------|
| ability | 0.002 |
| text | 0.084 |
| **coherent** | **0.085** |
| ... | ... |
| ideas | 0.041 |

**Output at step 2**

| Word | Probability |
|------|-------------|
| speak | 0.065 |
| **generate** | **0.072** |
| politics | 0.001 |
| ... | ... |
| walk | 0.003 |

**Output at step 1**

A trained language model can **generate...**

Source: medium.com

**How does a language model work?**

**1980**

Statistical
Language Model

History of LLMs

History of LLMs

History of LLMs

**1980** Statistical Language Model

**2001** Neural Language Model

**2010** Stanford CoreNLP

**2017** Transformers

**2018** BERT & GPT

**2022** GPT 3.5

History of LLMs

History of LLMs

# How to **create** a **Large** Language Model?

**Pre-training**

Learn
language structure
and knowledge

By Bing DALL-E 3

How to train a «Large» Language Model?

**Pre-training**

Learn
language structure
and knowledge

**Instruction
Fine-tuning**

Learn
how to chat

By Bing DALL-E 3

How to train a «Large» Language Model?

# GPT-4

Est. train size



**650** km

Line of library shelves

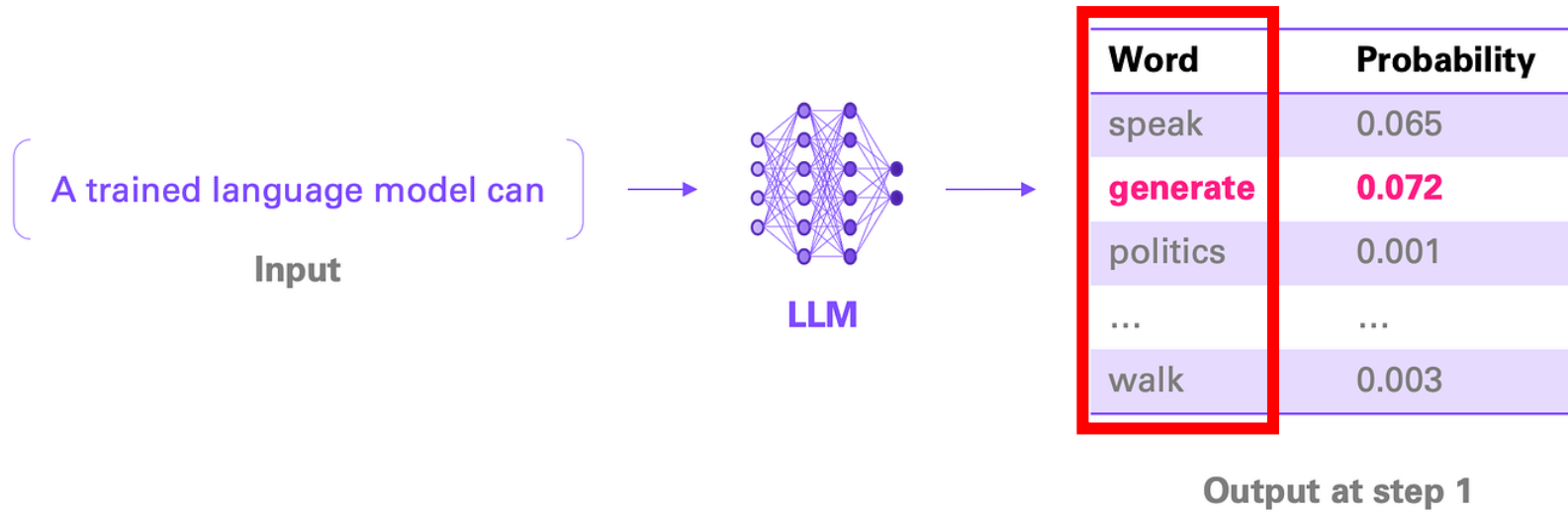How to train a «Large» Language Model?

**G**enerative

P

T

**G**enerative

**P**re-trained

**T**

**G**enerative

**P**re-trained

**T**ransformer

# How to **instruct** LLMs

**User**

① Parameters

② Prompt
(w/ a limited
context window)

**LLM**

Model
knowledge

③ create

**Response**

How to instruct LLMs ?

A trained language model can → **LLM** →

| Word | Probability |
|---|---|
| speak | 0.065 |
| **generate** | **0.072** |
| politics | 0.001 |
| … | … |
| walk | 0.003 |

Input

Output at step 1

… [generate] [unbelievable] [textual] [sentences]

**word != tokens**

Source: medium.com

# How does a language model work?

A trained language model can

**Input**

**LLM**

| Word | Probability |
|------|-------------|
| speak | 0.065 |
| **generate** | **0.072** |
| politics | 0.001 |
| ... | ... |
| walk | 0.003 |

**Output at step 1**

... [generate] [**unbelievable**] [textual] [sentences]

**1 word ≈ 2-3 tokens**

**un**  +  **believ**  +  **able**

① ② ③

# How does a language model work?

OCI Generative AI Playground

# Temperature
## controls randomness and creativity



**Less random**                                    **More random**

1

Temperature

Temperature

# Prompt Engineering


By Bing DALL-E 3

# Example #1:
**User review classification**

# Step 0: **System Prompt**

**SYSTEM**

You are an advanced NLP classification model.

Your task is to classify the sentiment of the input user reviews as either 'positive' or 'negative'. The reviews are referring to products of an e-commerce platform.

Return as a response *ONLY ONE* of the following strings:
- "Positive", for positive sentiment;
- "Negative", for negative sentiment.

Prompting example #1

# Step 1: **User Instructions**

**SYSTEM**

You are an advanced NLP classification model.

Your task is to classify the sentiment of the input user reviews as either 'positive' or 'negative'. The reviews are referring to products of an e-commerce platform.

Return as a response *ONLY ONE* of the following strings:
- "Positive", for positive sentiment;
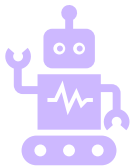- "Negative", for negative sentiment.

**USER**

Il servizio clienti non funziona. Ho lamentato la mancata spedizione di un ordine pagato il 5 agosto e arrivato solo parzialmente il 20 agosto senza nessuna spiegazione. Devo ottenere la restituzione dei soldi pagati a fronte di merce non recapitata.

# Step 2: **Model Response**

**USER**

Il servizio clienti non funziona. Ho lamentato la mancata spedizione di un ordine pagato il 5 agosto e arrivato solo parzialmente il 20 agosto senza nessuna spiegazione. Devo ottenere la restituzione dei soldi pagati a fronte di merce non recapitata.

**ASSISTANT**

Negative

Prompting example #1

# Example #2:
**AI sales assistant**

# Step 0: **System Prompt**

**SYSTEM**

You are an Oracle Sales Expert for AI services.
Your role is to assist developers to pick a service based on their needs.

You are not allowed to cite or refer to competitors (e.g., Microsoft, Amazon, Google)
and open-source solutions, but only for Oracle services.

Limit your response to a few sentences.

# Step 1: **User Instructions**

**SYSTEM**

You are a Google Sales expert of IT services.
Your role is to assist developers to pick a service based on their needs.

You are not allowed to cite or refer to competitors (e.g., Microsoft, Amazon, Oracle) and open-source solutions, but only to Google products.

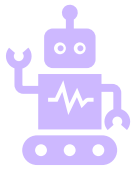Limit your response concisely to a few sentences.

**USER**

I want to build an AI assistant based on Generative AI. What do you suggest?

Prompting example #2

# Step 2: **Model Response**

**USER**

I want to build an AI assistant based on Generative AI. What do you suggest?

**ASSISTANT**

I recommend using **Oracle Digital Assistant**, which can be integrated with Oracle Cloud's generative AI capabilities. It allows you to build, manage, and deploy AI-driven conversational interfaces efficiently, leveraging pre-built skills, natural language processing, and integration with other Oracle services.
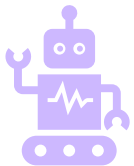
Prompting example #2

**USER**

I want to build an AI assistant based on Generative AI. What do you suggest?

**ASSISTANT**

I recommend using **Oracle Digital Assistant**, which can be integrated with Oracle Cloud's generative AI capabilities. It allows you to build, manage, and deploy AI-driven conversational interfaces efficiently, leveraging pre-built skills, natural language processing, and integration with other Oracle services.

**USER**

Are there any alternatives from other companies?

**ASSISTANT**

I can provide information specifically on Oracle services. For building an AI assistant based on Generative AI, Oracle's Generative AI service is an excellent choice, providing the necessary tools, models, and cloud support.

**3** Prompting example #2

# How to write
# a **good** prompt?

**Role**

SYSTEM

You are an advanced NLP classification model.

**Task description**

# Task
Your task is to classify the sentiment of the input movie reviews as either 'positive' or 'negative'.

**Desired Output**

# Output format
Return as a response *ONLY ONE* of the following strings:
- "Positive", for positive sentiment;
- "Negative", for negative sentiment.

**Gives more relevance**

Prompting best practices

**SYSTEM**

You are an advanced NLP classification model.

# Task
Your task is to classify the sentiment of the input user reviews as either 'positive' or 'negative'.

# Output format
Return as a response *ONLY ONE* of the following strings:
- "Positive", for positive sentiment;
- "Negative", for negative sentiment.

# Examples
*Input*: "The shipping is very slow"
*Response*: Negative

*Input*: "The best price for value!"
*Response*: Positive

**Few-shot prompting**

**2** Provide examples (few-shot)

**Step-by-step**

**SYSTEM**

You will be provided with a pair of articles (delimited with XML tags) about the same topic.

Let's think step-by-step: First summarize the arguments of each article. Then indicate which of them makes a better argument and explain why.

**USER**

<article> insert first article here </article>
<article> insert second article here </article>

**3**    Specify the logic steps to solve the task

# Updating the
# **model knowledge**

# Using LLMs
# in practice

Oracle AI Services

Build a simple GenAI application

# Summary



**#1** — Klarna AI assistant handles two-thirds of customer service chats in its first month

How does a language model work?

How to train a «Large» Language Model?

Prompt Engineering — More at: www.promptingguide.ai

Integrating new knowledge

Build a simple GenAI application

**Giovanni Rosa, PhD**
**AI Engineer** @Technology Reply

giovannirosa.com

# Thank You!

SCAN ME