



UAI

**Universidad Abierta
Interamericana**

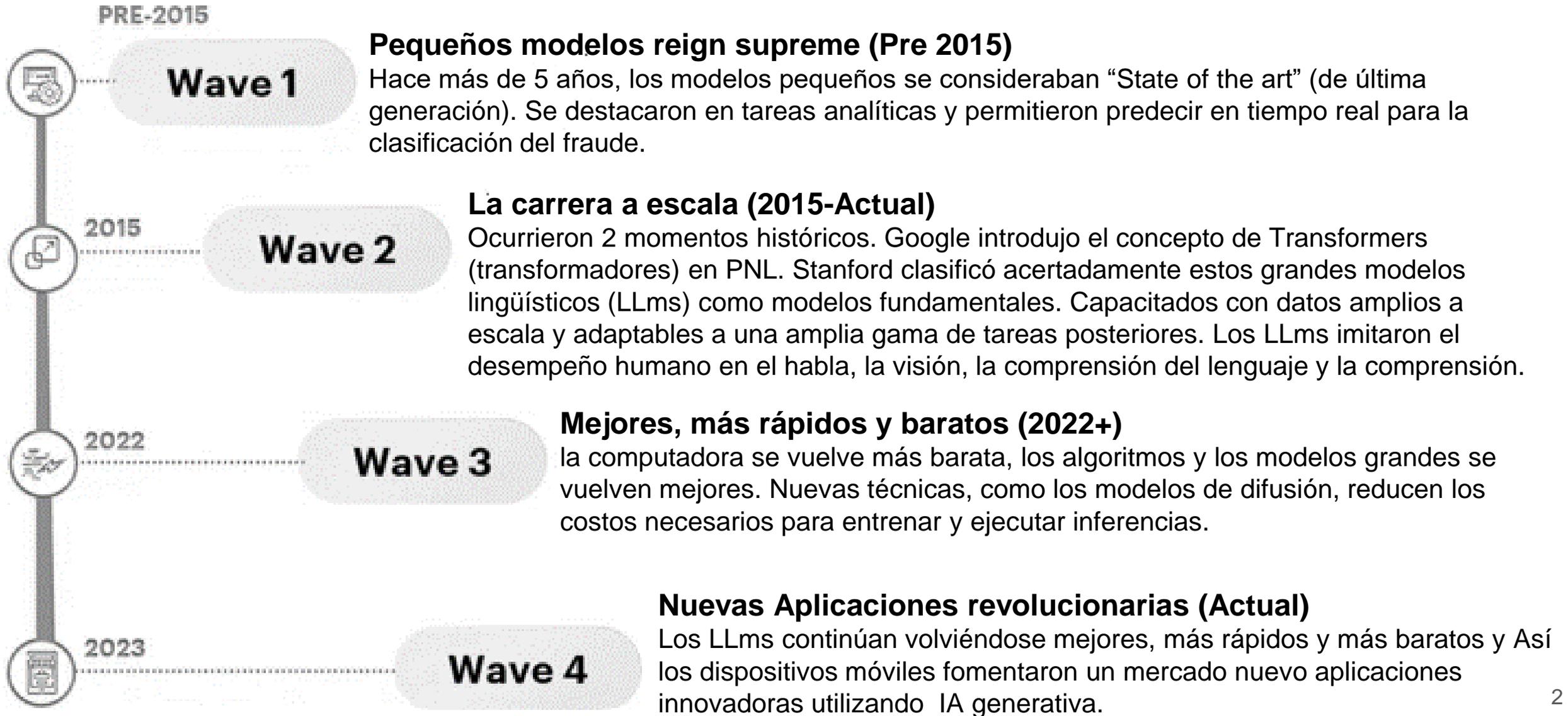
IA generativa en la Industria

CITII 24

Ing. Maximiliano Bonaccorsi

LLMs (large language models)

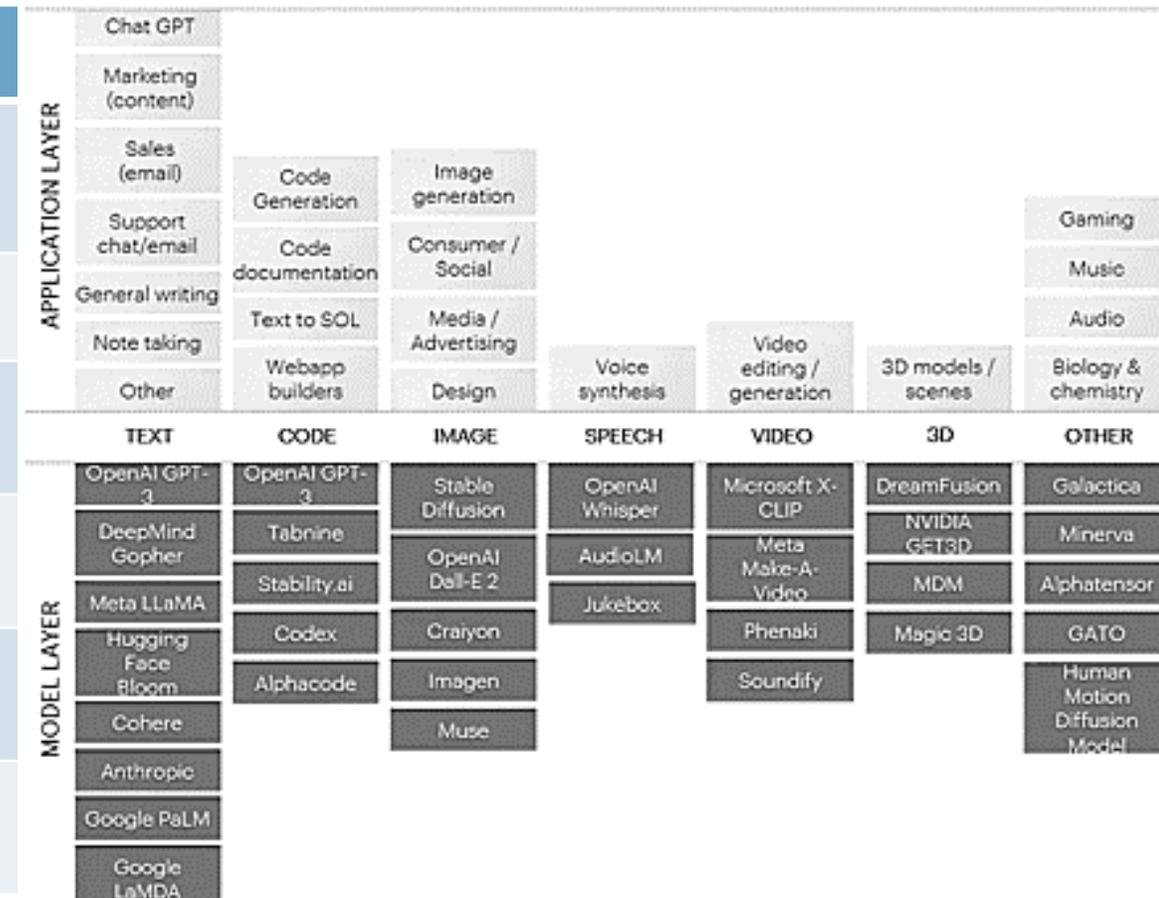
Generación de modelos: Better models, more data, more computer



LLMs (large language models)

Modelos y apps de IA generativa

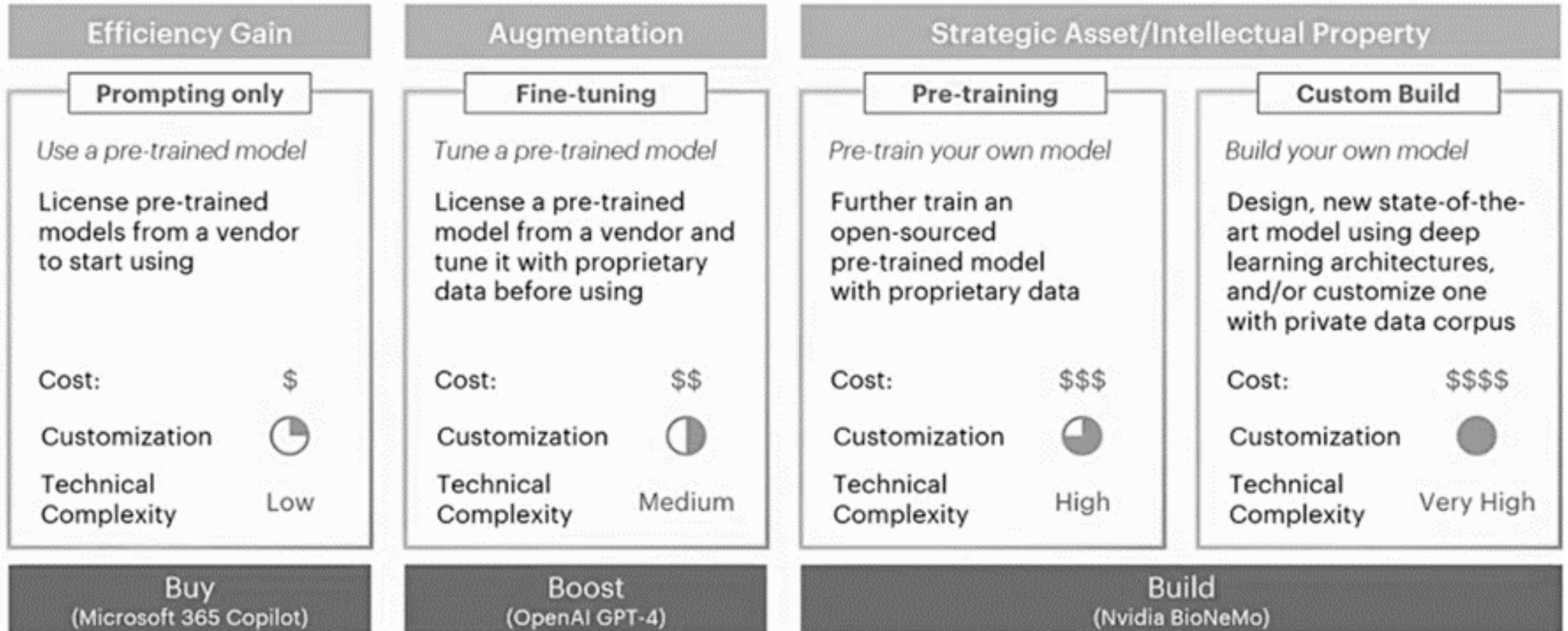
Tipo	Descripción
Modelos de texto	Decente en escritura corta/forma media. Necesidad de resultados de mayor calidad, contenido más extenso y un mejor ajuste específico del dominio
Generación Código	Gran impacto sobre la producción del desarrollo (ej. Github Copilot)
Imágenes	Advenimiento de diferentes estilos estéticos y técnicas de edición y modificación.
Speech synthesis	por un tiempo (¡hola siri!), pero las aplicaciones empresariales y de consumo están mejorando
Video y 3d models	Potencial para destrabar mercados creativos como cine, juegos, VR, arquitectura y diseño de productos físicos
Otros Dominios	Modelos Fundamentales (LLMs) ocurren a través de muchos campos desde la música hasta las farmacéuticas



IA generativa

Enfoques para habilitar modelos de base (foundation models)

Hay 3 enfoques principales para habilitar modelos de base y su complejidad técnica y costos aumentan con el nivel de personalización requerido.



IA generativa

Ejemplo: Enfoques de habilitación de modelos para codificación asistida por IA

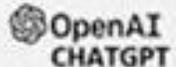
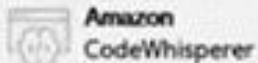
Se debe encontrar la combinación correcta de enfoques para impulsar el valor empresarial y crear su propia llm para tareas específicas

Comprar

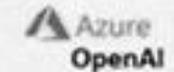
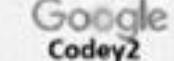
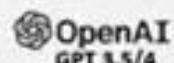
Purchasing or licensing a pre-existing model from a vendor or service provider.

"I need a code assistant that provides my development team with recommendations based on industry best practices"

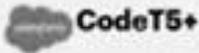
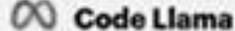
GenAI-enabled Coding Tools



Code Model Providers



Open-source Code Models

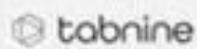


Mejorar

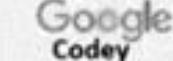
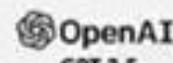
Using an existing AI model as a foundation and then customizing it to fit specific use cases

"I need a code assistant that utilizes industry best practices but also takes into account my organization's unique standards"

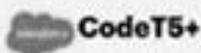
GenAI-enabled Coding Tools



Code Model Providers



Open-source Code Models



Crear

Designing and developing the model from the ground up, resulting in a custom, private model

"I need a code assistant that is exclusively trained on my in-house coding language or to reduce risk"

GenAI-enabled Code Tools



Code Model Providers

Open-source Code Models

IA generativa

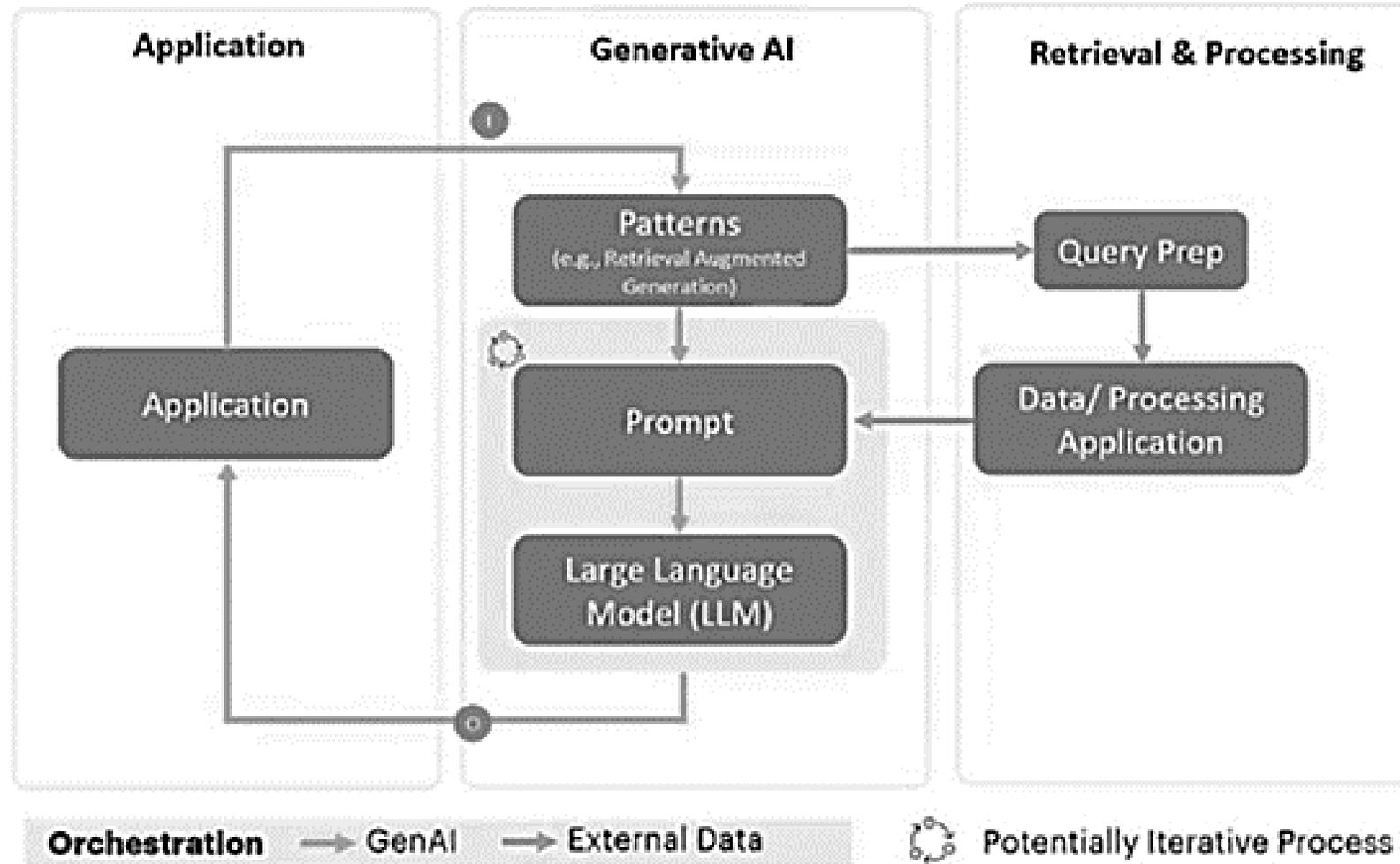
Mejorar

La siguiente fase agrega datos para habilitar la asistencia por copilot.

Ejemplo: patrón de generación aumentada de recuperación aplicado a la entrega de tecnología.

Outputs

- Code
- Documentation
- Tests
- Ticket Resolution
- Requirements
- Workflows
- Scripts
- Pipelines
- Q/A driven insights

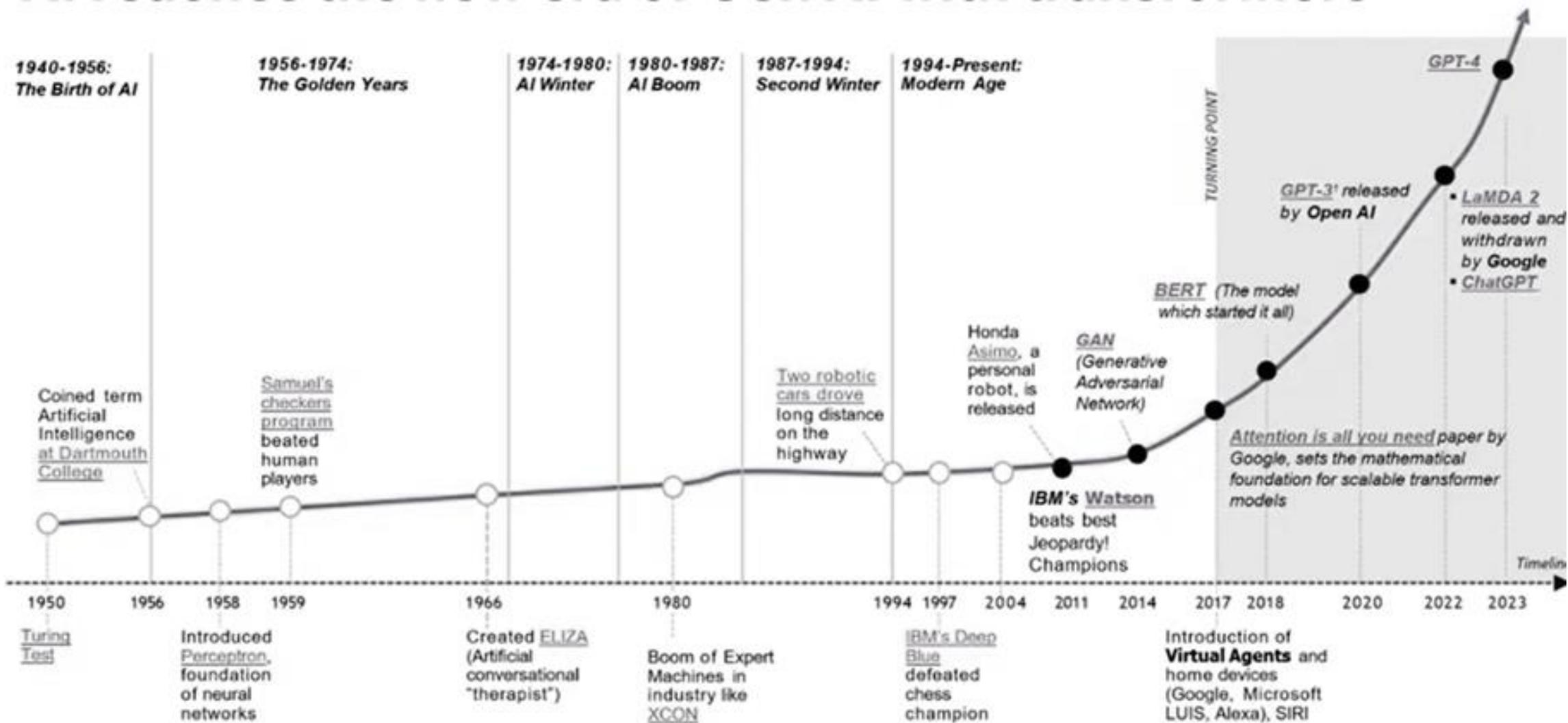


Knowledge Base

- Software Code & Configuration Files
- Infrastructure as Code Templates
- Architecture
- Design Specification
- Vendor manuals
- Requirements
- Workflows
- System Logs
- Standard Operating Procedures
- Business Process

IA generativa

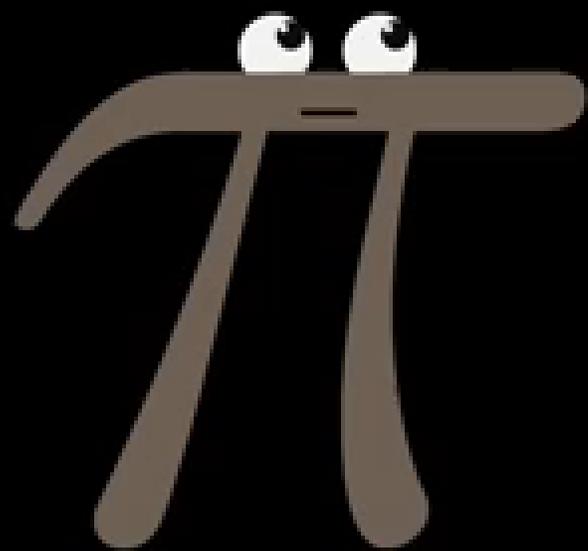
AI llega a la nueva era de gen ai con transformers



IA generativa

**la nueva era de IA Generativa surge con
Transformers**

Veamos su funcionamiento



IA generativa

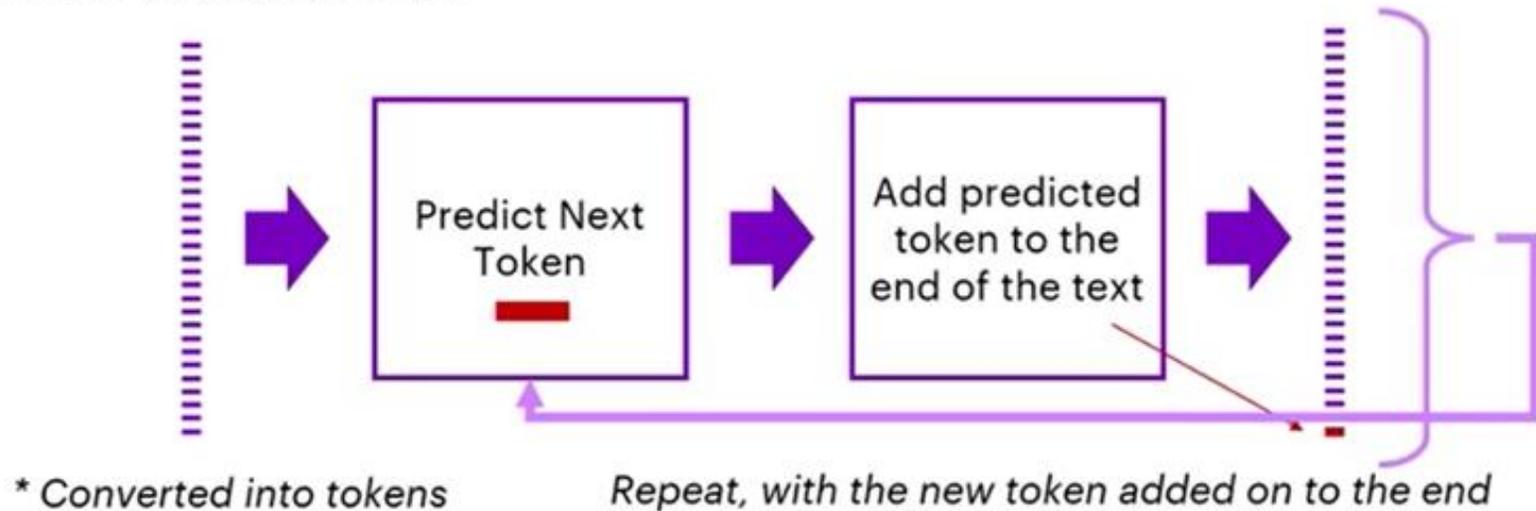
¿Como funciona Gen AI?

Parte 1- Predicción

Ejemplo de Token:

is a superlyawesomelycoolness IT services company.

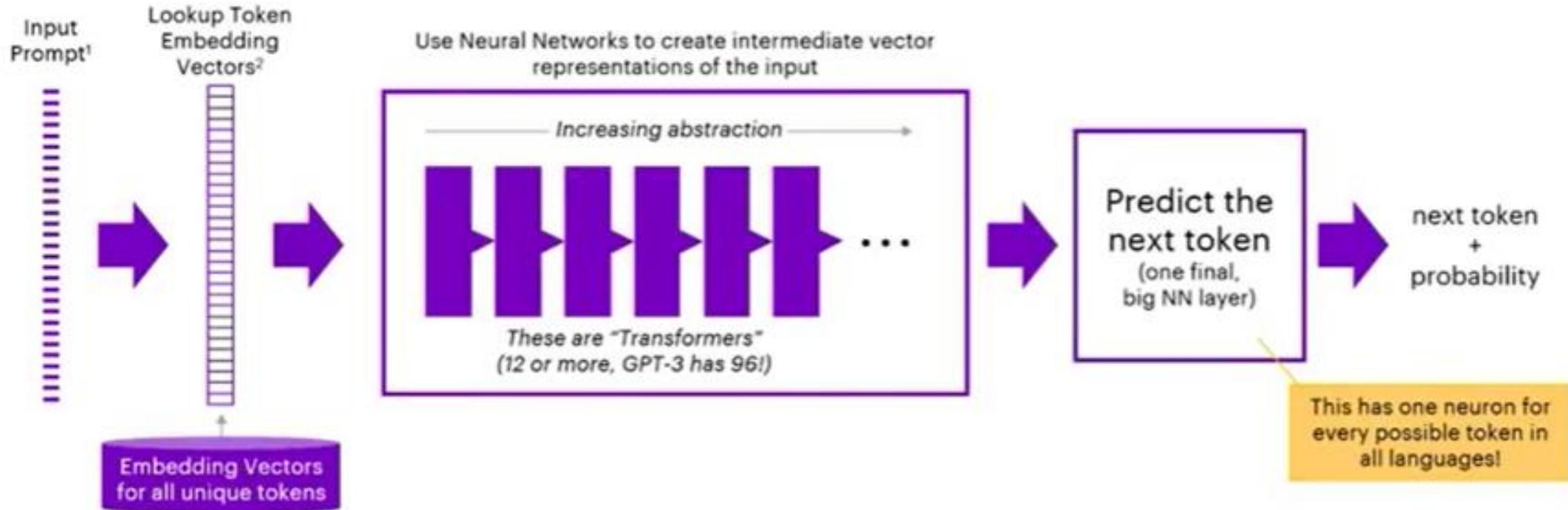
Input Text* (aka, the prompt)



IA generativa

¿Como funciona Gen AI?

Parte2- Transformers



Nota:

1. Convertido a Token
2. Cada token único tiene un vector, cada vector es una lista de 700 a 13000 números de punto flotante (el valor del vector actual para cada token único está entrenado por el proceso de entrenamiento NN).
3. El vector del token se modifica según la posición dentro del prompt (mensaje).

Cada Transformer:

1. combina todos los vectores simbólicos en el input (entrada) entre sí (esta es llamada la "self attention layer" y este complejo mashup permite que la red neuronal aprenda la sintaxis)
2. tiene múltiples (al menos 2) capas de red neuronal en su interior (con millones de neuronas y billones de pesos de conexiones de red neuronal interconectadas)

IA generativa

¿Como funciona Gen AI?

Parte3- Training

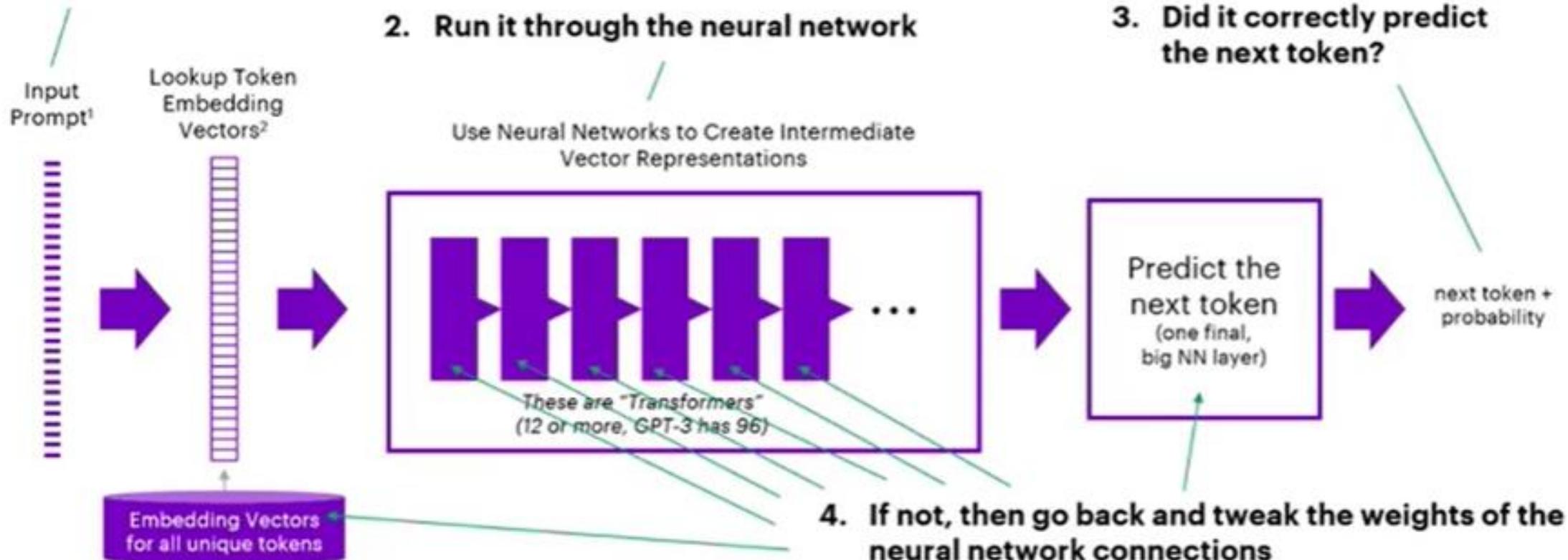
1. Take some text from the internet

2. Run it through the neural network

3. Did it correctly predict the next token?

4. If not, then go back and tweak the weights of the neural network connections

5. Repeat for months and months until it learns

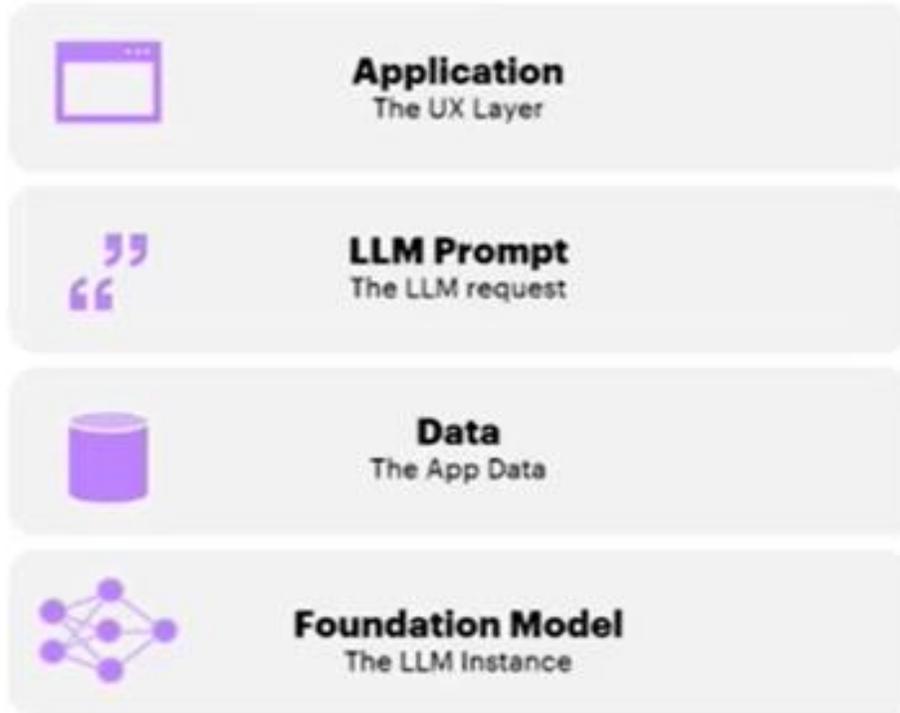


IA generativa

Utilización de modelos MMLs en aplicaciones

El Flujo operacional para una aplicación no difiere en
Diferentes tipos de aplicaciones

The Layers of Generative AI Application



EXAMPLES

ChatGPT

ChatBot

- the UI/UX for user input, which includes instruction and data

The LLM Prompts

- managed by the Chatbot
- filters user input
- Includes user input, chat history, and more

Application Data

- Chatbot determines if data needs to be retrieved via Plugins (new)

Foundation Model

- OpenAI GPT 3.5 / 4.0

Github CoPilot

IDE Extension & Service

- the inside IDE UI/UX for code completion and generation
- The backend service for processing IDE requests.

The LLM Prompts

- managed by Github CoPilot
- Includes developer input, and other code collected from dev workspace

Application Data

- Does not apply

Foundation Model

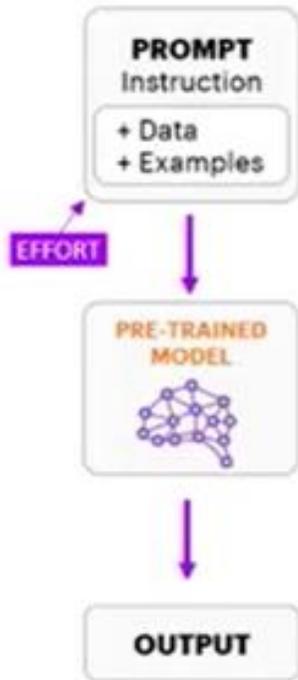
- OpenAI GPT 3.5

IA generativa

3 enfoques clave para adaptar los LLMs a necesidades específicas

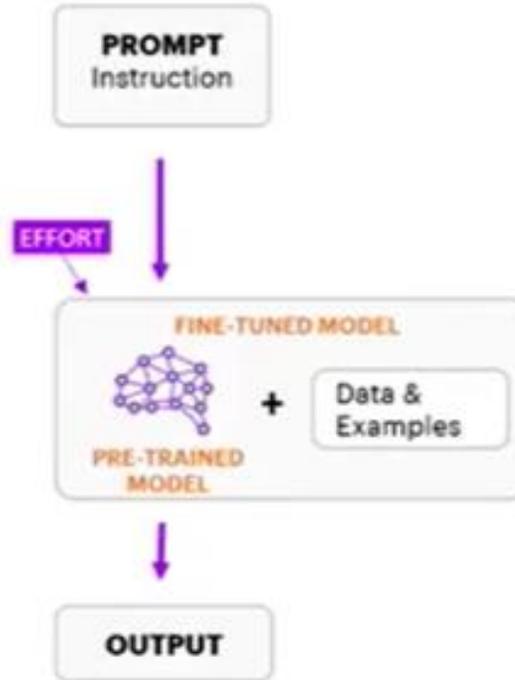
Prompt Engineering

Tailor the prompt to a task



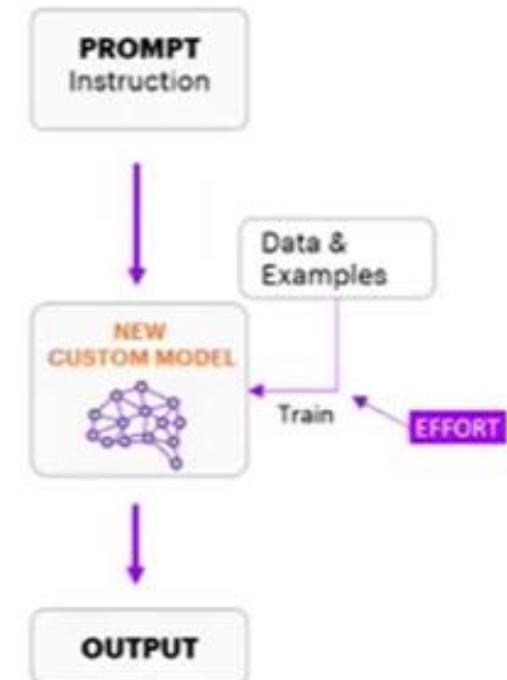
Fine-tuning

Adapt a pre-trained model for a task



Pre-train / Custom

Build a new the model for a task



La complejidad técnica y los costos escalan hacia arriba según el nivel de customización que sea requerida

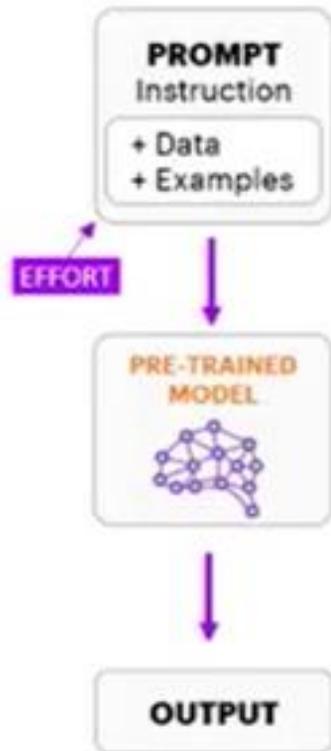
IA generativa

3 enfoques clave para adaptar los LLMs a necesidades específicas

adaptar el mensaje a una pregunta

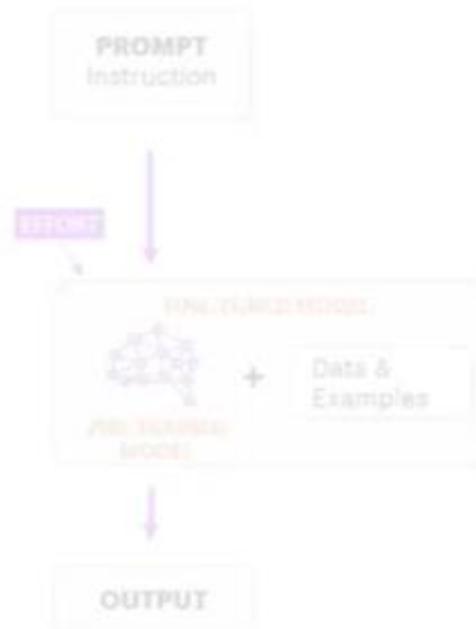
Prompt Engineering

Tailor the prompt to a task



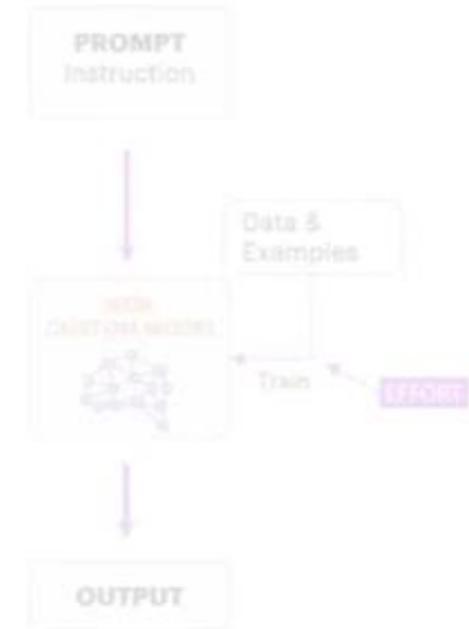
Fine-tuning

Adapt a pre-trained model for a task



Pre-train / Custom

Build a new the model for a task



La complejidad técnica y los costos escalan hacia arriba según el nivel de customización que sea requerida

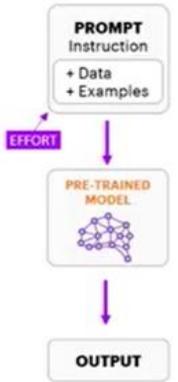
IA generativa

Prompt engineering es la forma de optimizar los resultados Gen IA

Prompt engineering es una técnica utilizada en LLMs para guiar o dirigir su producción proporcionándoles un mensaje específico, un conjunto de instrucciones o datos, donde el objetivo es mejorar la precisión y relevancia de la salida del modelo.

Prompt Engineering

Tailor the prompt to a task



01. Content

Give it content to process

Often, the machine will need information to perform a task. Include this in the prompt.

- Document text
 - E.g., Documentation, policies, procedures, etc.
- Tabular data (e.g., CSV)
- Lists of things to choose from (e.g., APIs or data sets)

02. Instructions

Tell it what to do

This is where you tell it what you want to accomplish.

- Answer question
- Create a plan
- Summarize some content
- Process some content
- And, most importantly, what sort of output you want

03. Examples

Provide examples, if needed

If it is unable to follow instructions, you may need to provide examples.

- Content understanding
- Aspects or dimensions of interest
- What to extract
- Output formatting

IA generativa

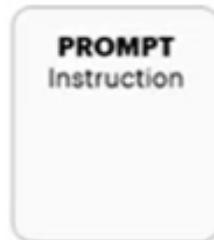
Hay 3 técnicas de aprendizaje en contexto.

Entender las técnicas de aprendizaje en contexto provee instrucciones adicionales para ayudar la interpretación

Técnicas de aprendizaje en contexto: definiendo los datos limitados.

Zero shot learning

Zero-shot learning aims to correctly interpret tasks without any specific examples of those tasks during training, purely relying on the ability to understand instructions and context.



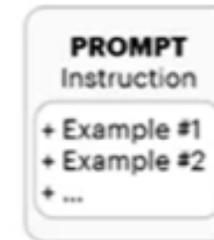
One-shot learning

One-shot learning focuses on the ability of a model to accurately make predictions after seeing just a single training example.



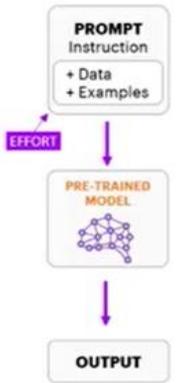
Few shot learning

Few-shot learning allows a model to generalize or make predictions after seeing a few examples.



Prompt Engineering

Tailor the prompt to a task



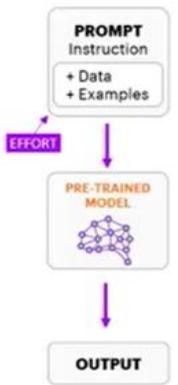
IA generativa

Adaptando LLMs con prompt engineering.

Un modelo aprende a través de ejemplos analizando los datos de entrada (inputs), identificando patrones y realizando predicciones o decisiones, mientras que el razonamiento en cadena de pensamiento permite que el modelo se base en conocimientos previos, conectando diferentes conceptos para generar conocimientos integrales.

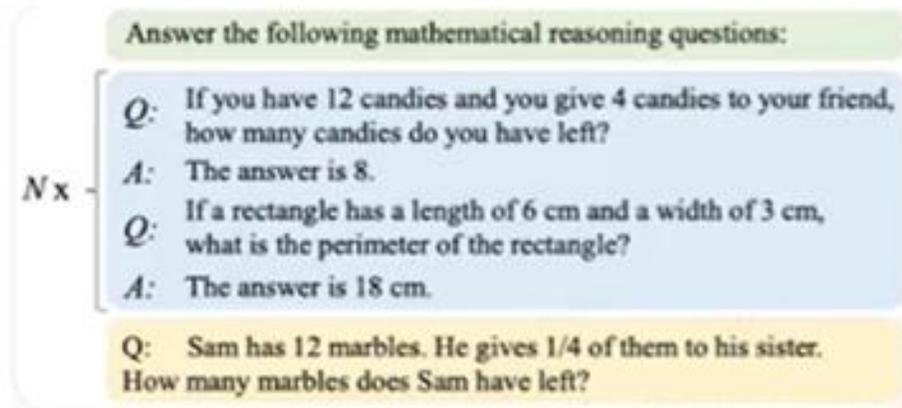
Prompt Engineering

Tailor the prompt to a task



aprendizaje en contexto

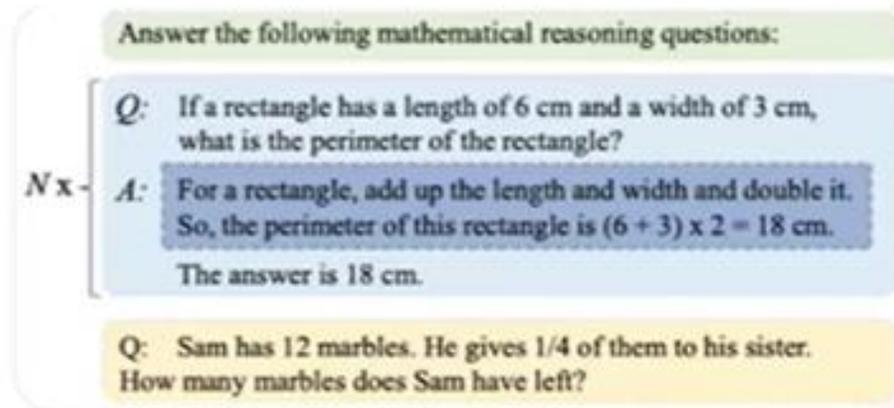
In-Context Learning



A: The answer is 9.

cadena de pensamiento

Chain-of-Thought Prompting



A: He gives $(1/4) \times 12 = 3$ marbles. So Sam is left with $12 - 3 = 9$ marbles. The answer is 9.

: Task description

: Demonstration

: Chain-of-Thought

: Query

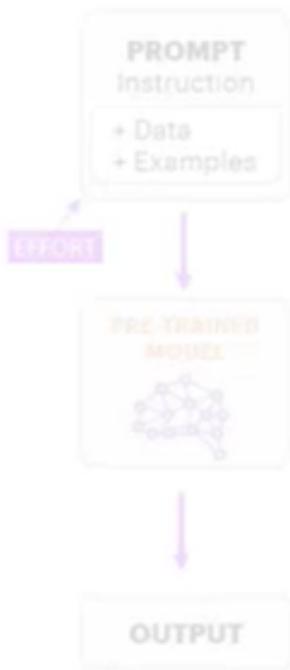
IA generativa

3 enfoques clave para adaptar los LLMs a necesidades específicas

adaptar el modelo pre-entrenado para una tarea

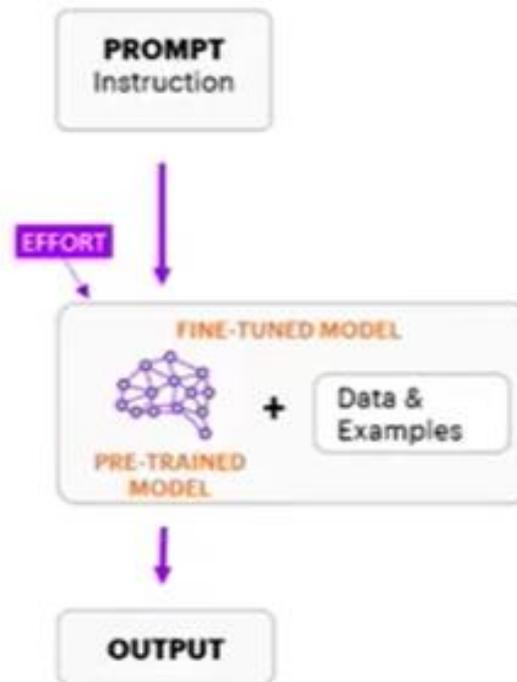
Prompt Engineering

Tailor the prompt to a task



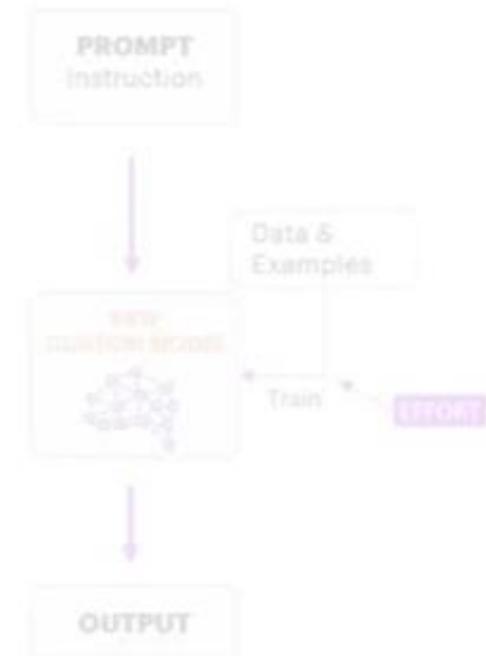
Fine-tuning

Adapt a pre-trained model for a task



Pre-train / Custom

Build a new the model for a task

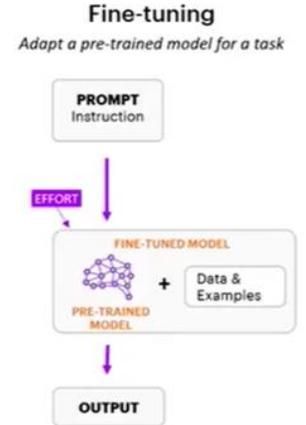


La complejidad técnica y los costos escalan hacia arriba según el nivel de customización que sea requerida

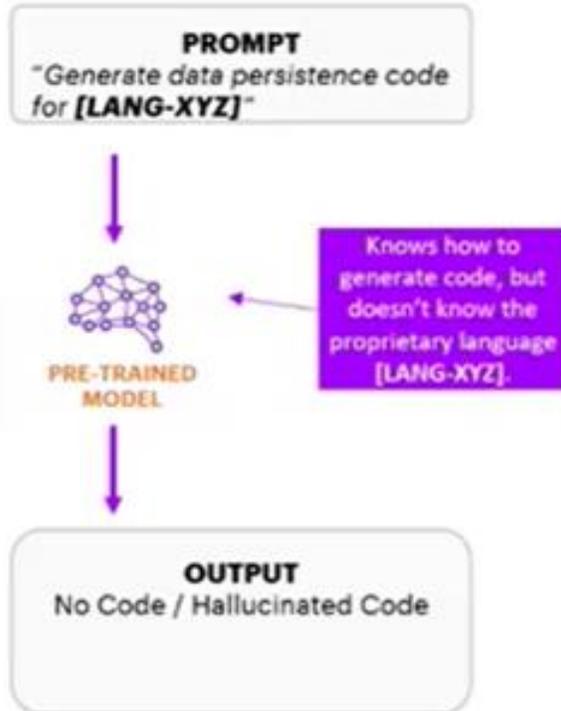
IA generativa

Adaptar LLms a necesidades específicas con Fine-tune

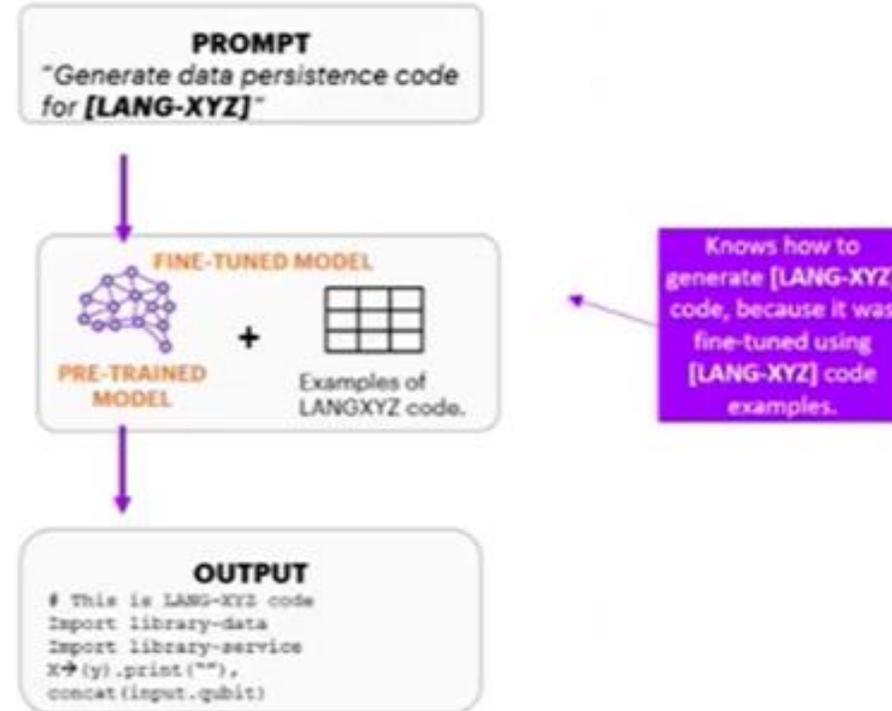
Las habilidades de los modelos se pueden adaptar aún más (o fine-tune) para adaptarse a casos de uso únicos después del entrenamiento.



Code Generation w/ Pre-trained Model Only



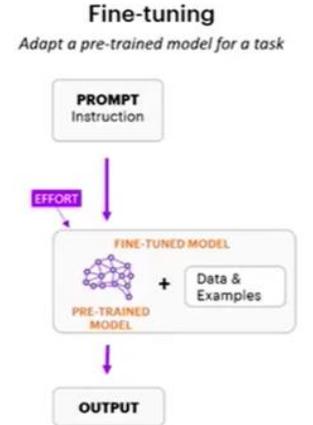
Code Generation w/ Fine-tuned Model



IA generativa

Adaptar LLMs a necesidades específicas con Fine-tune

Las habilidades de los modelos se pueden adaptar aún más (o fine-tune) para adaptarse a casos de uso únicos después del entrenamiento.



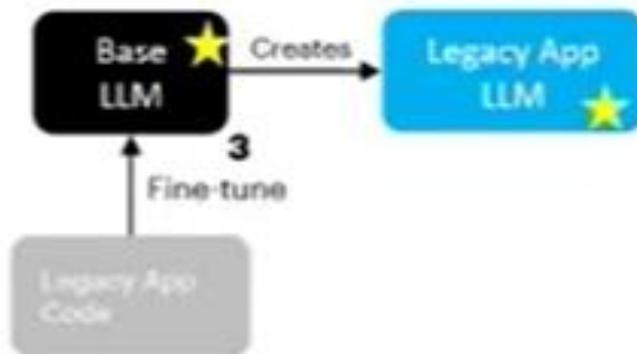
Fine-tune LLM with source code

To be able to contextually query **source code**.



Query examples:

- What is the purpose of the overall application?
- Generate a sequence diagram of the checkout process.
- Generate (reverse engineer) the requirements for payment processing module.



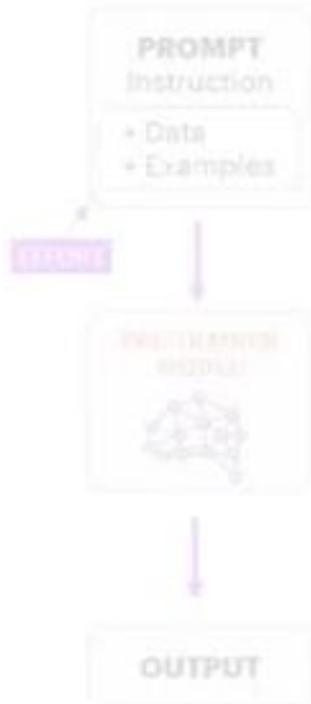
- 1 Standard embedding models
- 2 Pre-trained commercial or open-source model that performs well generating and interpreting (legacy language; potentially fine-tuned model for better language understanding)
- 3 Pre-trained commercial or open-source model that that can be fine-tuned with legacy source code (to increase model's knowledge of private source code)

IA generativa

3 enfoques clave para adaptar los LLMs a necesidades específicas

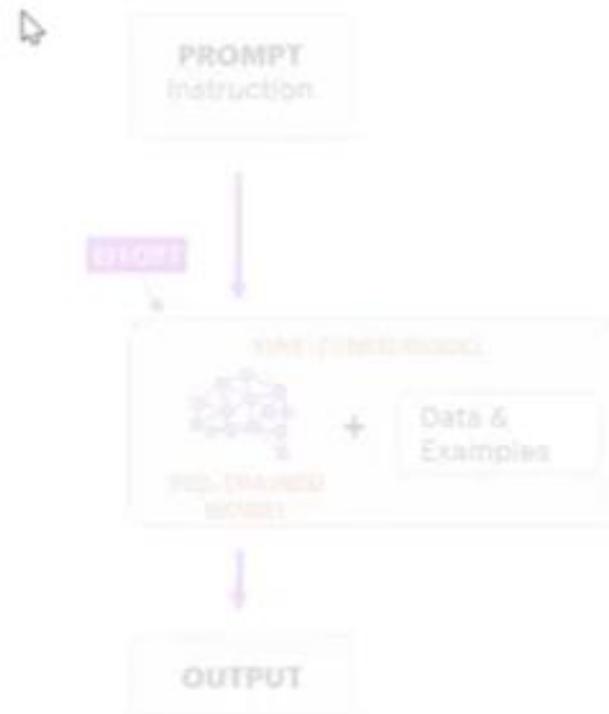
Prompt Engineering

Tailor the prompt to a task



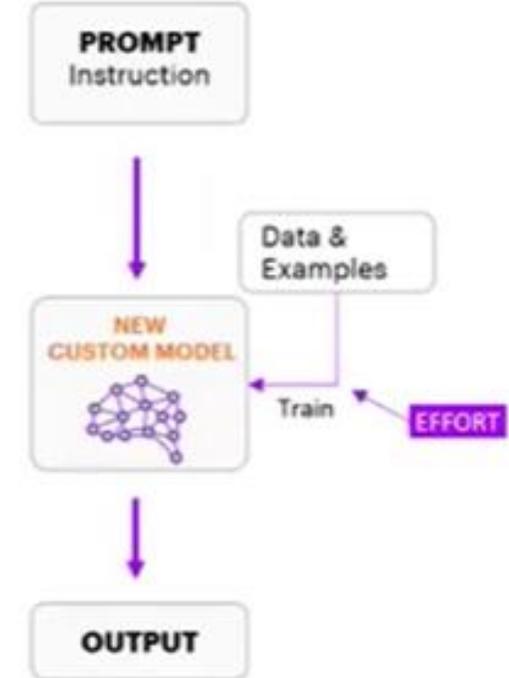
Fine-tuning

Adapt a pre-trained model for a task



Pre-train / Custom

Build a new the model for a task



La complejidad técnica y los costos escalan hacia arriba según el nivel de customización que sea requerida

IA generativa

Componentes arquitectónicos clave de la aplicación Gen AI.

Una aplicación gen IA es una implementación de uno o más patrones de uso gen ia. Prompt engineering se utiliza para construir la indicación con las instrucciones, datos y ejemplos necesarios para que LLMs complete el objetivo o la tarea.

Application

- The interface for end-users to access features enabled by GenAI

Patterns

- The generative AI pattern applied.

Query Prep

- The preparation of query to enterprise data store.

Data

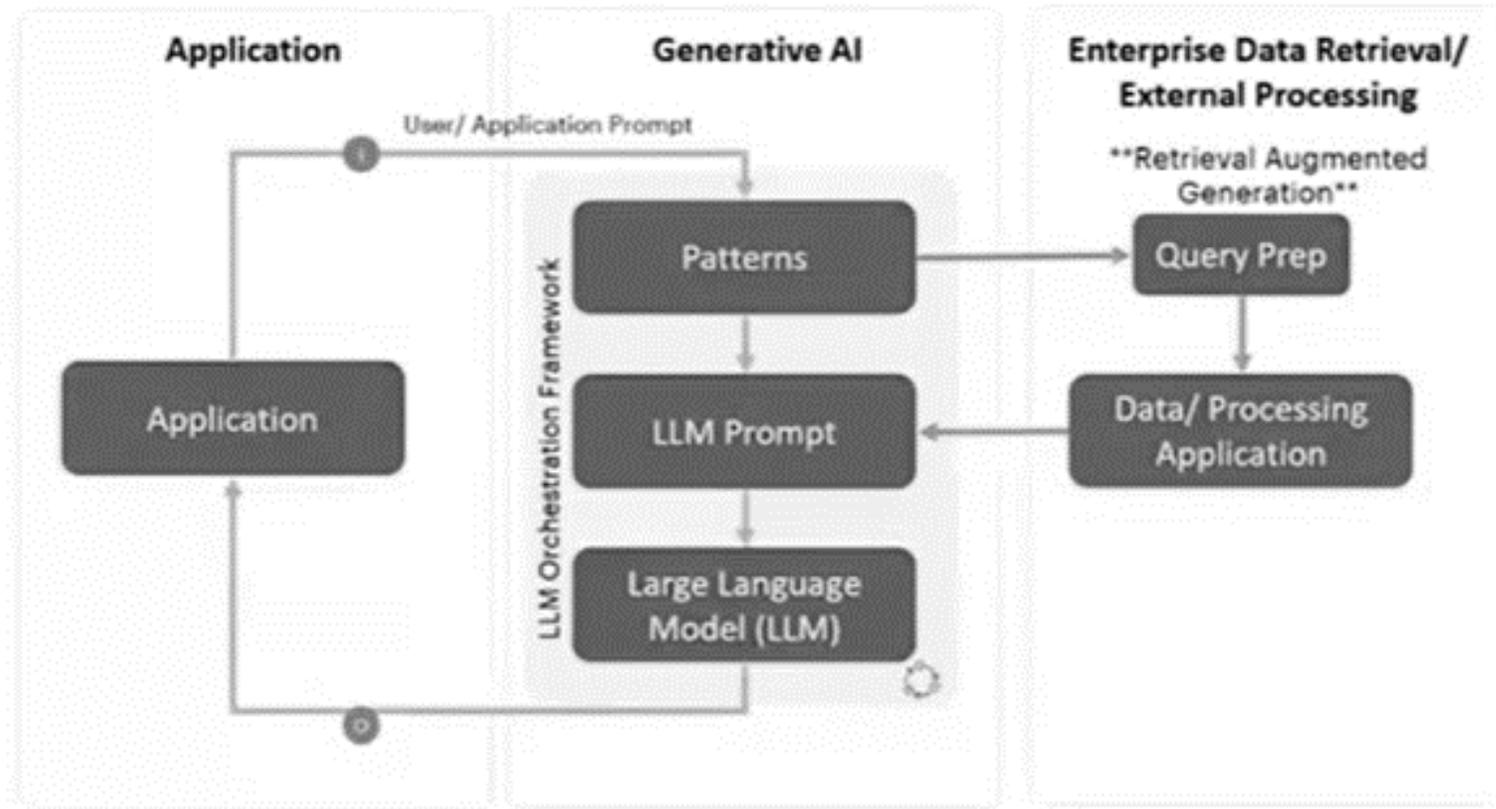
- The querying the enterprise application or data store.

Prompt

- What is sent to the LLM
- An instantiation of a pattern
- Includes original query, the history, the retrieved data.

Large Language Model (LLM)

- Generates (infers) an answer to question from data embedded in the LLM
- The LLM could be universal model, fine-tuned, or custom



Orchestration



GenAI



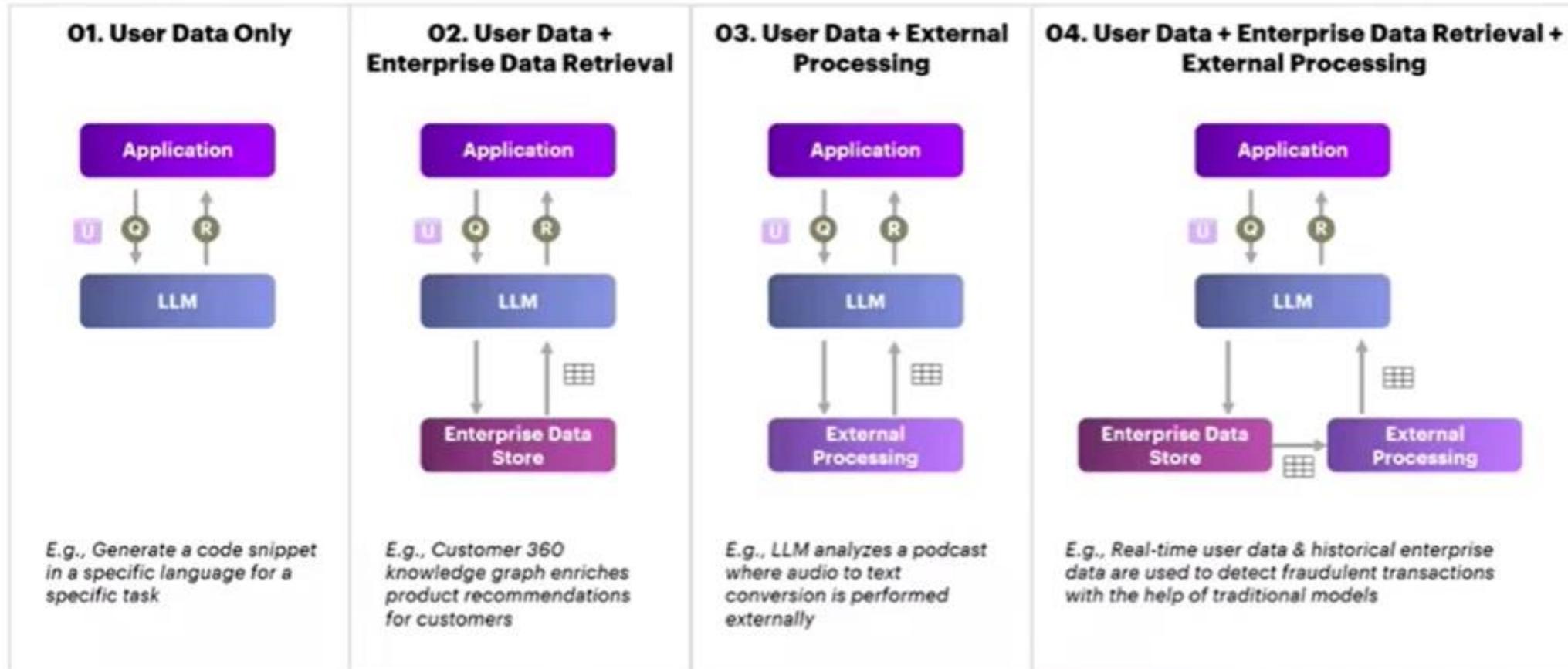
Enterprise Data



Potentially Iterative Process

IA generativa

4 arquetipos centrales para la Gen AI



Legends

U - User data

Q - Query

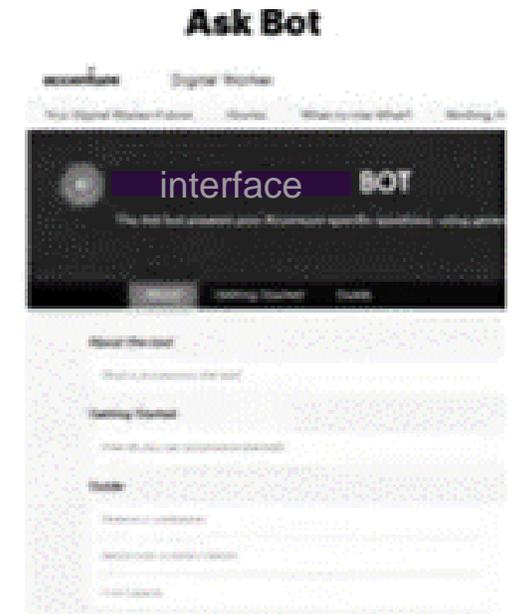
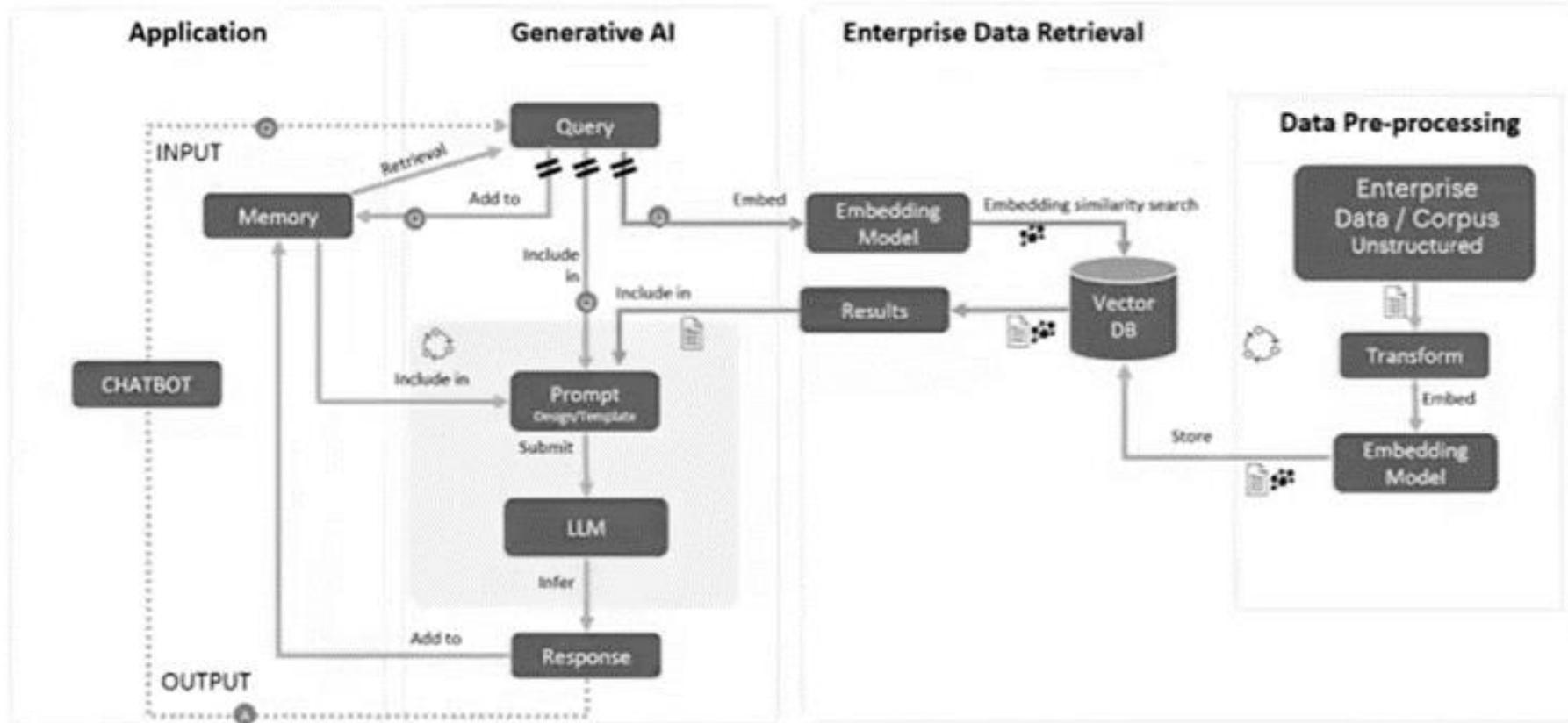
R - Response

Grid icon - Enterprise data | Externally processed data/insights

IA generativa

Búsqueda Semántica / Vector de búsqueda

La recuperación de datos empresariales incorpora la consulta para buscar semánticamente en la base de datos vectorial almacenes de artículos que contienen información relevante, que se agregará como contexto al mensaje.



Legends

Orchestration



GenAI



Enterprise Data



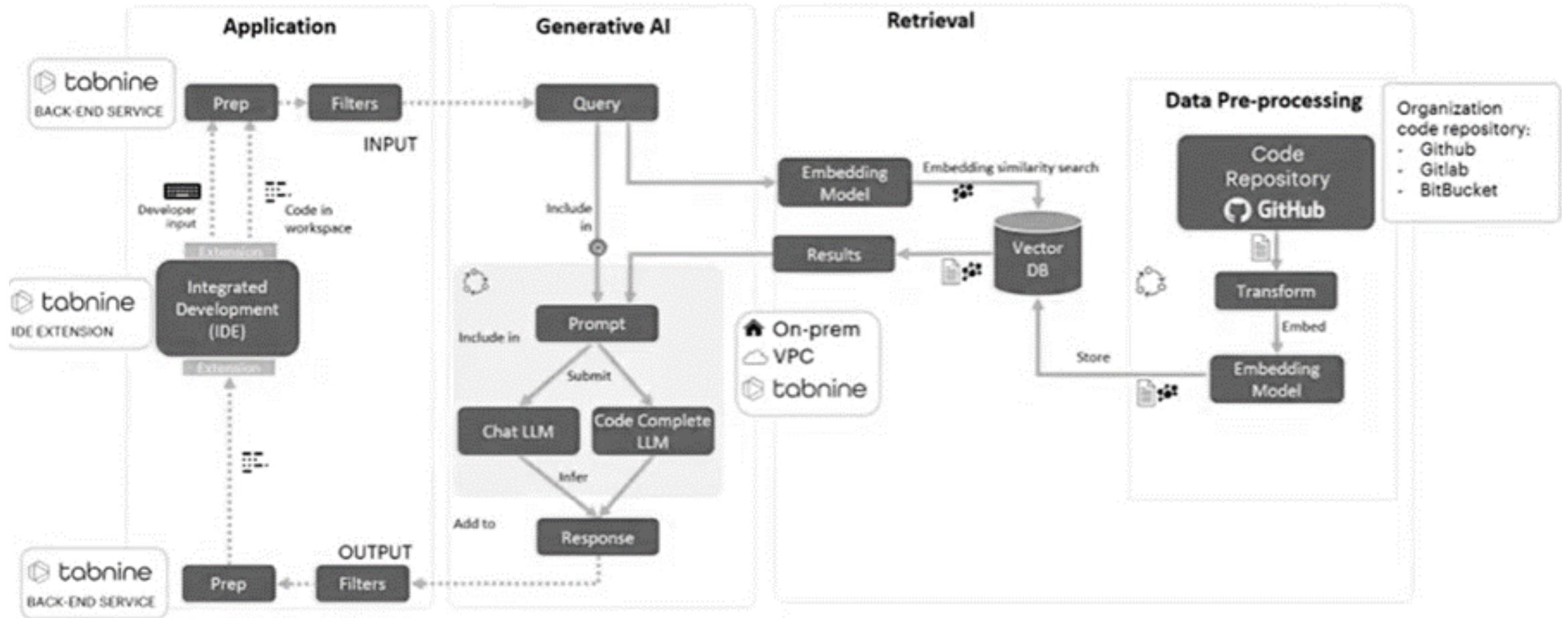
Potentially Iterative Process



Parallel

IA generativa

Ejemplo: Asistente de código (generación y completar código)



Legends

Orchestration

→ GenAI

→ Retrieval

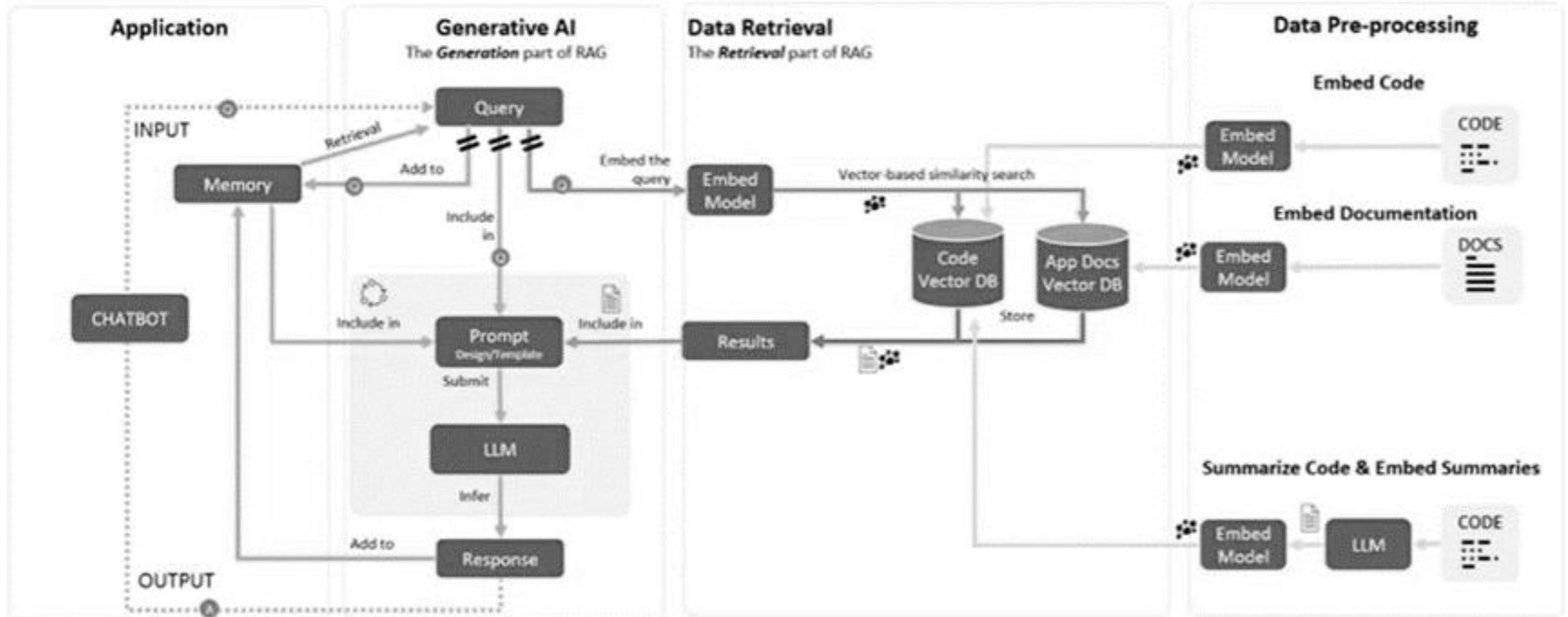


Potentially Iterative Process

≡ Parallel

IA generativa

Ejemplo: Asesor de Código (conoce tu código)



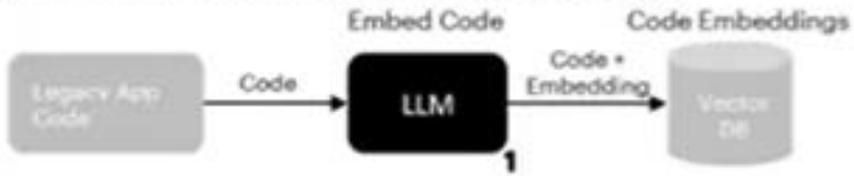
IA generativa

Preparación de datos para Asesor de Código



Embed source code

To be able to semantically search **source code**.



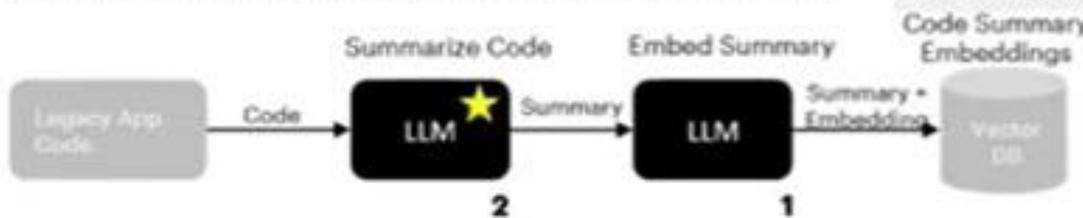
Query examples:

- Find me code that uses this library.
- Find me code that is similar to this [code].



Summarize code and embed summaries

To be able semantically search **source code summaries**.



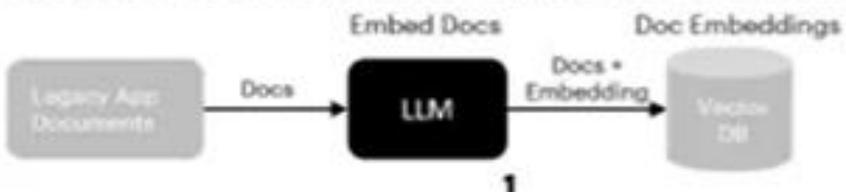
Query examples:

- Find me code that enable end-users to set application preferences



Embed application documentation

To be able to semantically search **application documentation**.



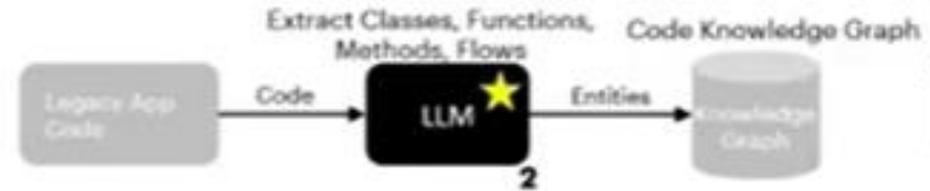
Query examples:

- Find me text(s) that describes how end-user will use this application feature.



Create knowledge graph of source code

To be able to contextually query **source code**.



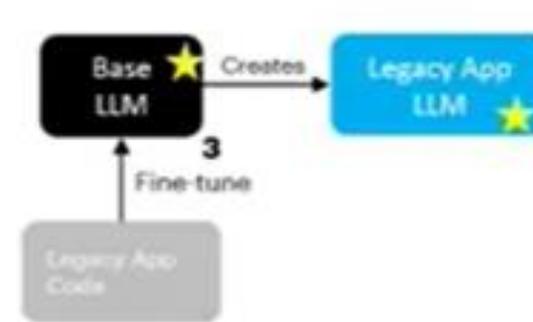
Query examples:

- How does a transaction flow through the system?
- What are all the software components involved in a data flow?
- What components would be impacted by a change to [entity]?



Fine-tune LLM with source code

To be able to contextually query **source code**.



Query examples:

- What is the purpose of the overall application?
- Generate a sequence diagram of the checkout process.
- Generate [reverse engineer] the requirements for payment processing module.

- 1 Standard embedding models
- 2 Pre-trained commercial or open-source model that performs well generating and interpreting (legacy language; potentially fine-tuned model for better language understanding)
- 3 Pre-trained commercial or open-source model that that can be fine-tuned with legacy source code (to increase model's knowledge of private source code)

IA generativa

Consideraciones para adoptar un LLM en la industria

Consideración	Descripción
Selección Modelo LLM	El Modelo LLm Aplicable <ol style="list-style-type: none">1. Puro LLm2. Valor agregado al LLm3. Open Source LLm
Accesibilidad y Desarrollo	Opciones de desarrollo que hagan accesible al LLm para uso de negocios
Enfoque de Adaptacion	Adaptación de modelos para trabajar con datos de organizaciones para lograr objetivos de precisión y rendimiento.
Preparación empresarial	Garantizar que los modelos cumplan con las exigencias de seguridad, confiabilidad, interoperabilidad y responsabilidad de la empresa.
LLm Ops	modelo operativo, procesos, marcos para producir y monitorear LLms para la empresa

IA generativa

Ejemplo de creación LLMs comprimido OpenSource con RAG

Utilización de Modelo Open Source, del tipo meta Llama pre-entrenado como LLM y backend de integración para un asistente de investigación local basado en RAG.

Utilizando fuente de información web (wikipedia) y datos locales como Pds, Txt, etc



```
from llama_cpp import Llama

# Cargar el modelo preentrenado
modelo = Llama.load("https://spanishchat-7b.Q4_K_M.gguf")

# Cargar los documentos y las páginas web
documents = ["Documento1.txt", "Documento2.txt", "Documento3.txt"]
web_pages = ["https://en.wikipedia.org/wiki/Fender_Stratocaster",
             "https://en.wikipedia.org/wiki/Gibson_Les_Paul",
             "https://en.wikipedia.org/wiki/Guitar_amplifier" ]

# Procesar los documentos y las páginas web
processed_docs = modelo.process_docs(documents)
processed_web_pages = modelo.process_web_pages(web_pages)

# Realizar una búsqueda de texto en los documentos y las páginas web
query = "texto de búsqueda"
results = modelo.search(processed_docs, processed_web_pages, query)

# Imprimir los resultados
for result in results:
    print(result)
```

IA generativa en la industria



¡¡Muchas Gracias!!

¿Preguntas?

Ing. Maximiliano Bonaccorsi