

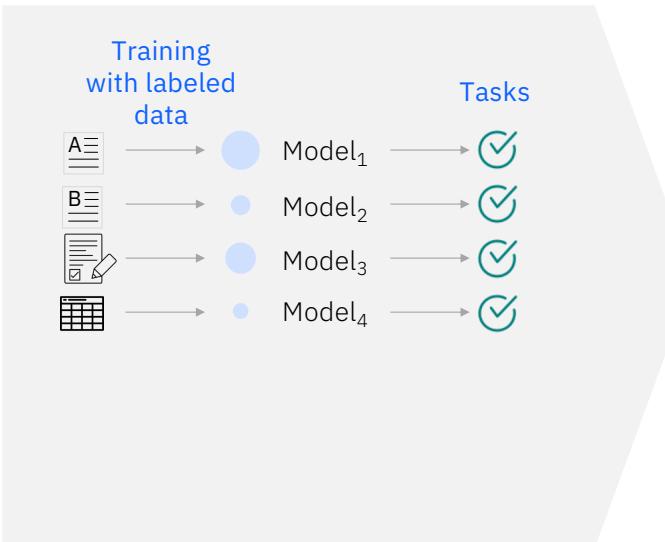
Introduction to key concepts

Agenda

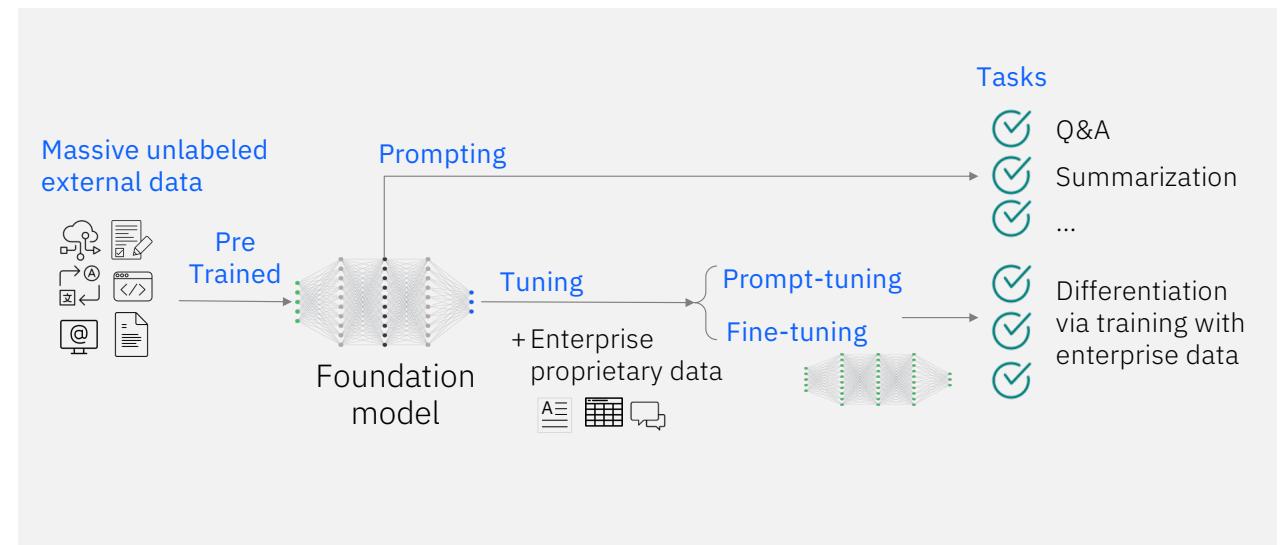
- FM, Prompt Engineering, Prompt Tuning and Fine Tuning
- Use Gen-AI API from python code in Notebooks or VS Code
- Leverage LangChain
- Retrieval Augmented Generation
- Introduction to Embeddings and VectorDB
- Position IBM Watson & other products as relevant

Foundational models enable a new paradigm of data-efficient AI development – generative AI

Traditional AI models



Foundation Models



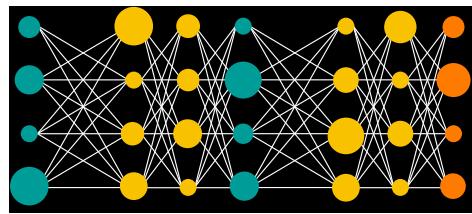
- Individual siloed models
- Require task specific training
- Lots of human supervised training

- Rapid adaptation to multiple tasks with small amounts of task-specific data
- Pre-trained unsupervised learning

Foundation models are ...



Self-supervised
training



Foundation model

Pre-trained

On unlabeled datasets
of different modalities
(e.g., language, time-
series, tabular)

Self-learning

Systems that leverage
self-supervised learning

Multiple applications

Able to learn
generalizable and
adaptable data
representations that can
be effectively used in a
variety of domains and
tasks (code generation,
question answering,
sentiment analysis)

Large language models

A type of foundation
model trained with
language-related data

ChatGPT is based on a
large language model

Foundation models: generalizable and adaptable

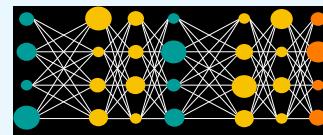
Translation
prompt

Summarization
prompt

Answer finding
prompt

+

Input
text



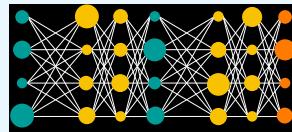
→ Translated input

→ Summarized input

→ Answer to the
input question



Task-specific
fine-tuning



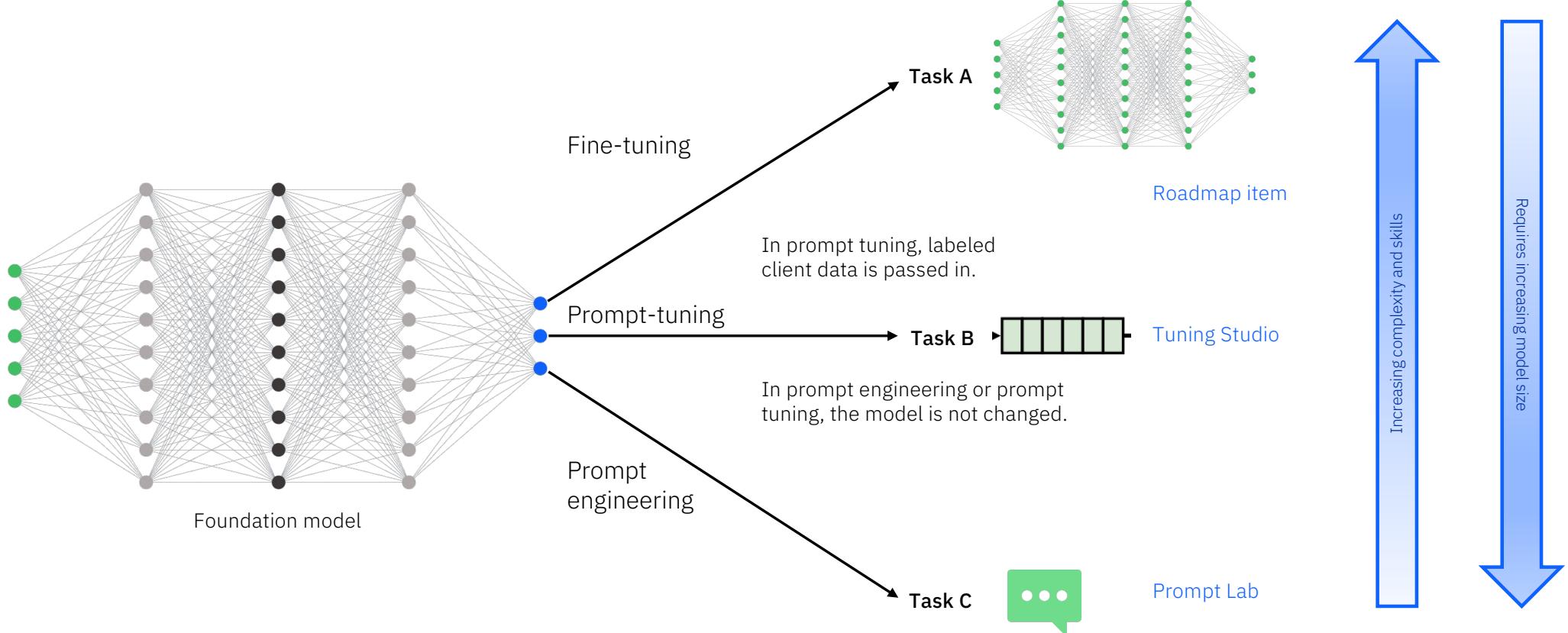
Translation
model

Summarization
model

Answer finding
model

Rapid adaptation to multiple tasks with small amounts of task-specific data

Fine-tuning requires labeled data and more resources to tune the model. When a model is fine-tuned, some of the weights are modified and clients get a private instance of the model.

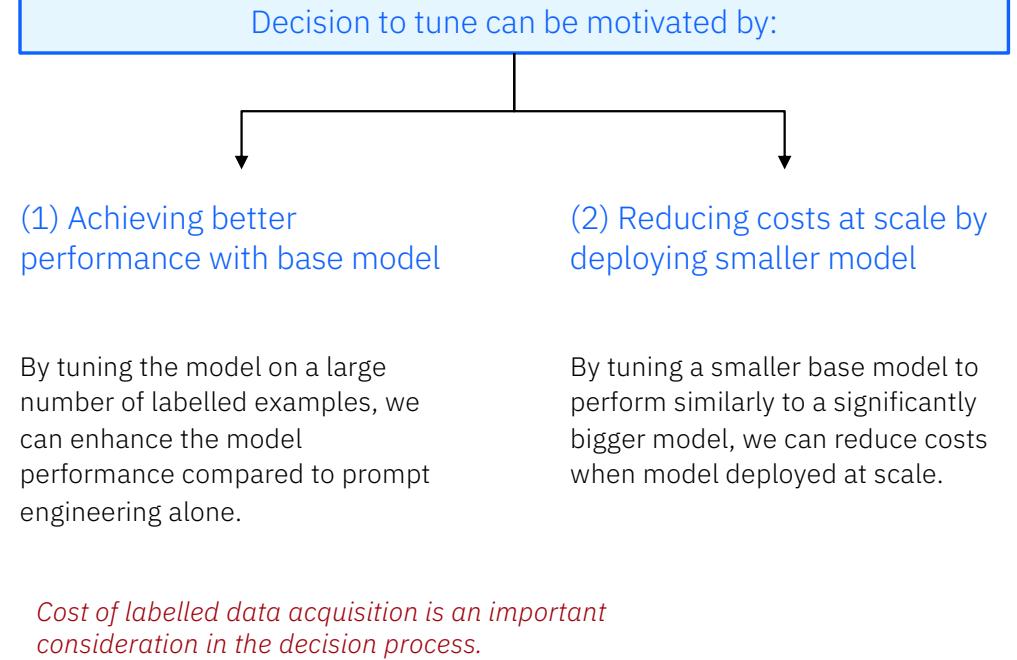


When to tune a model?

Always start with prompt engineering the largest LLM suited for your task.

This should provide some indication that the task is suited to be addressed by LLMs. It is also helpful to experiment with different labelled examples and understand which prompt formats work best on the target task.

12



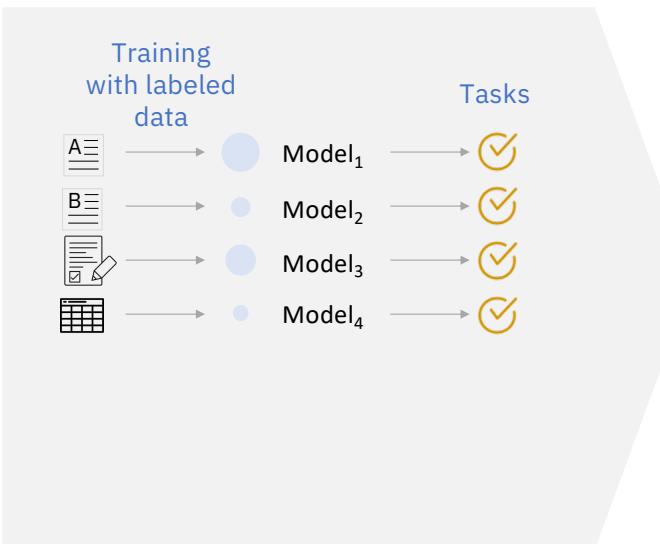
Suggested workflow

	Step 1	Step 2	Step 3
Stage	Initial PoC	Pilot deployment	Deployment at scale
Goal	<i>Prove the use case with minimal effort</i>	<i>Reduce costs as permitted within PoC duration</i>	<i>Maximize ROI</i>
Recommendation	Use the largest model and minimal labeled data Create labeled test dataset to measure model accuracy	Prompt engineer or prompt-tune a medium-size model . You may need to gather additional labeled data	Consider using additional data gathered to fine-tune / prompt-tune a small model . Deploy the tuned model
Inference costs	\$\$\$	\$\$	\$

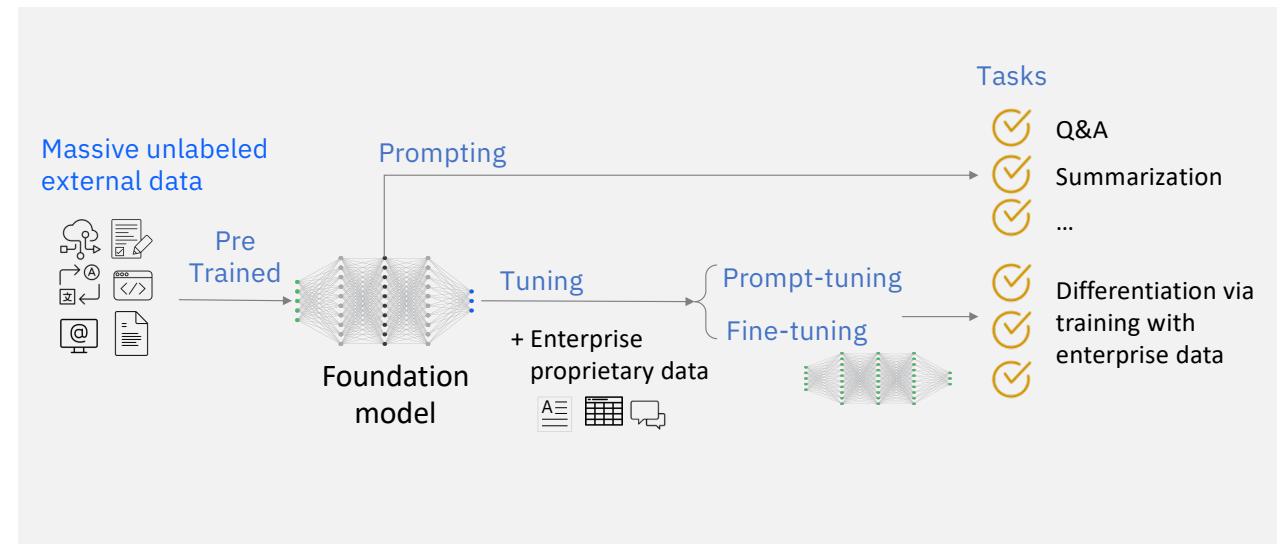
Note: fine-tuning a model requires creating a copy of the model specific to the user. The cost of hosting this model may impact the ROI analysis compared to prompt tuning.

Foundational models enable a new paradigm of data-efficient AI development – generative AI

Traditional AI models



Foundation Models



- Individual siloed models
- Require task specific training
- Lots of human supervised training

- Rapid adaptation to multiple tasks with small amounts of task-specific data
- Pre-trained unsupervised learning

IBM partnership with open-source models provider



HUGGING FACE

- IBM `watsonx.ai` clients have access to the latest and greatest open-source foundation models from Hugging Face.
- The IBM and Hugging Face partnership demonstrates a joint commitment to deliver an open ecosystem to clients, allowing them to find the best foundation models for their business needs.

Most common generative AI tasks implemented today

Summarization

Transform text with domain-specific content into personalized overviews that capture key points.

Conversation summaries, insurance coverage, meeting transcripts, contract information

Classification

Read and classify written input with as few as zero examples.

Sorting of customer complaints, threat and vulnerability classification, sentiment analysis, customer segmentation

Generation

Generate text content for a specific purpose.

Marketing campaigns, job descriptions, blog posts and articles, email drafting support

Extraction

Analyze and extract essential information from unstructured text.

Medical diagnosis support, user research findings

Question-answering

Create a question-answering feature grounded on specific content.

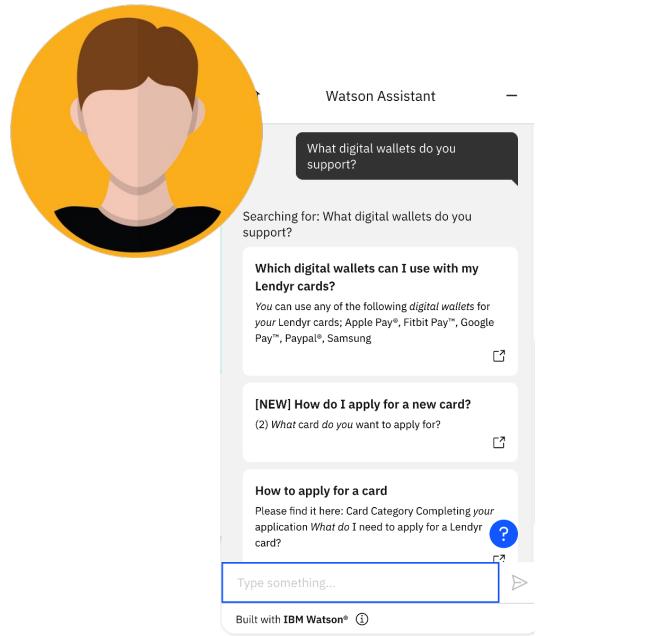
Build a product specific Q&A resource for customer service agents.

LangChain for LLMs

- LangChain is an open-source framework designed to simplify creating applications using LLMs.
- Models
- Prompt Templates
- Parsers
- Chains
- Question Answer

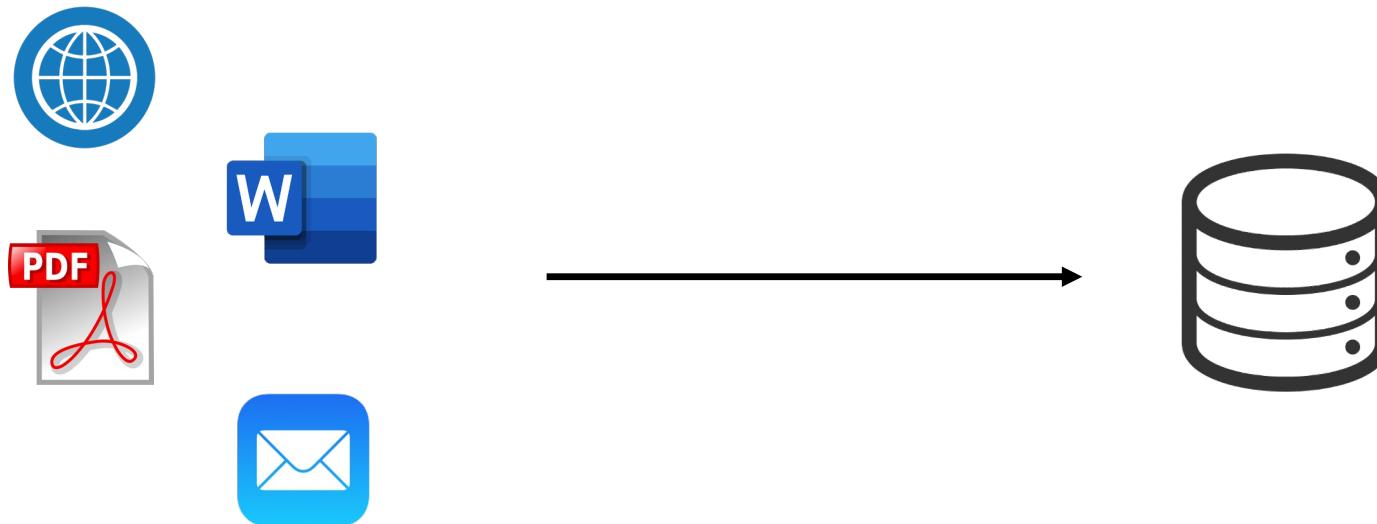
Retrieval Augmented Search - Overview

Conversational search – Q&A for documents



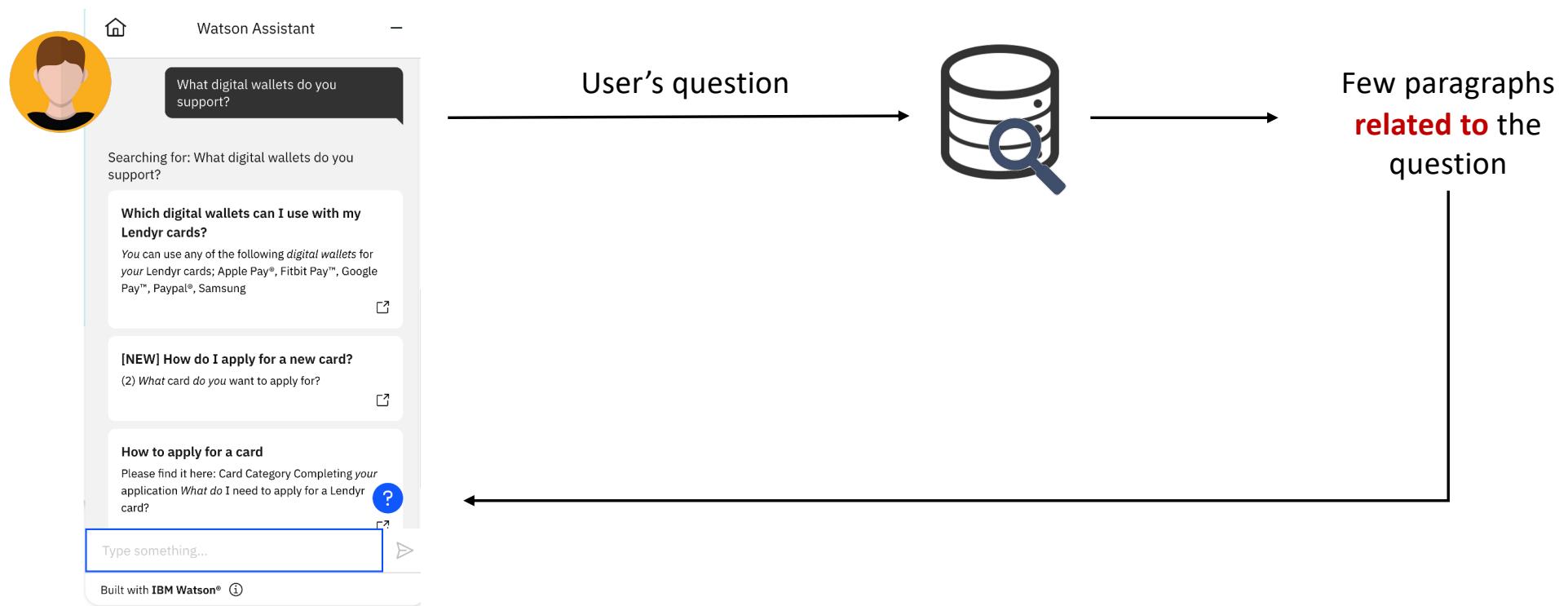
Phase 1

Prepare the data



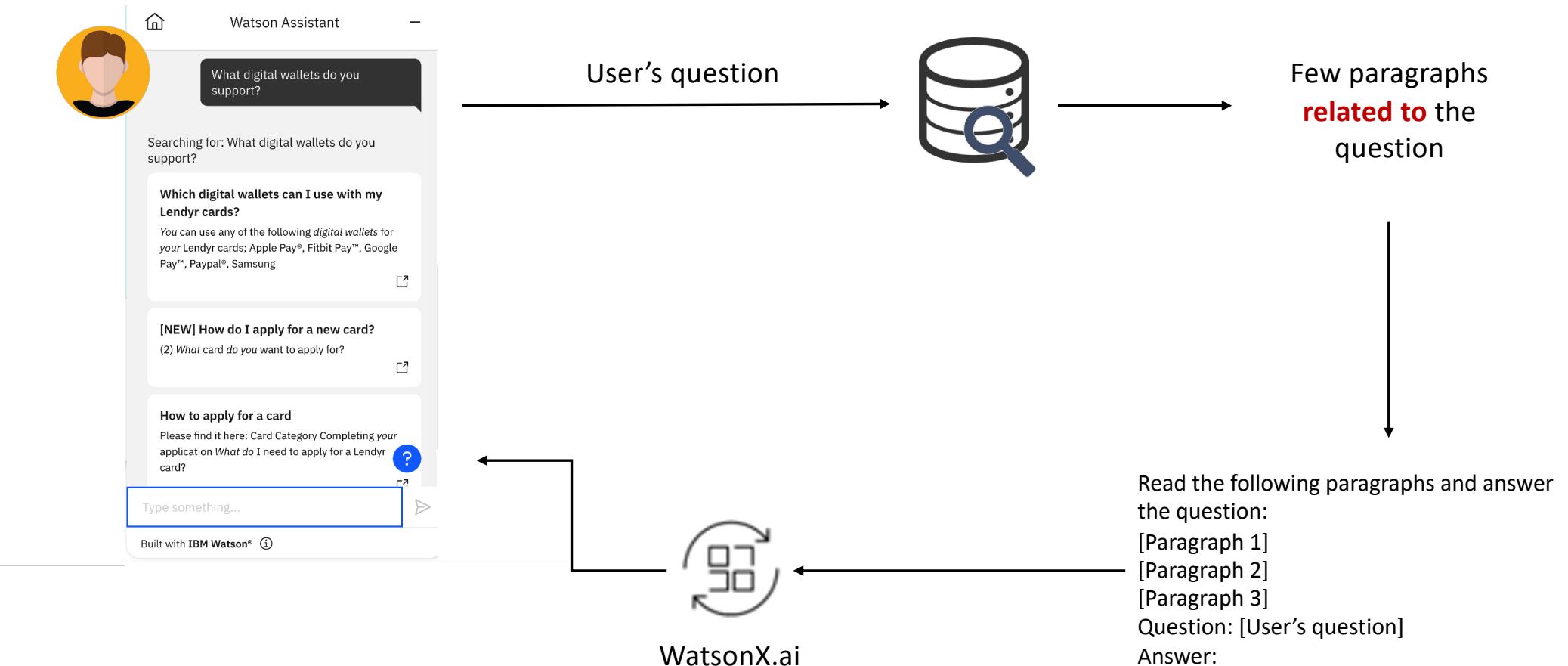
Phase 2

Query the data



Phase 2 (new steps based on LLMs)

Query the data

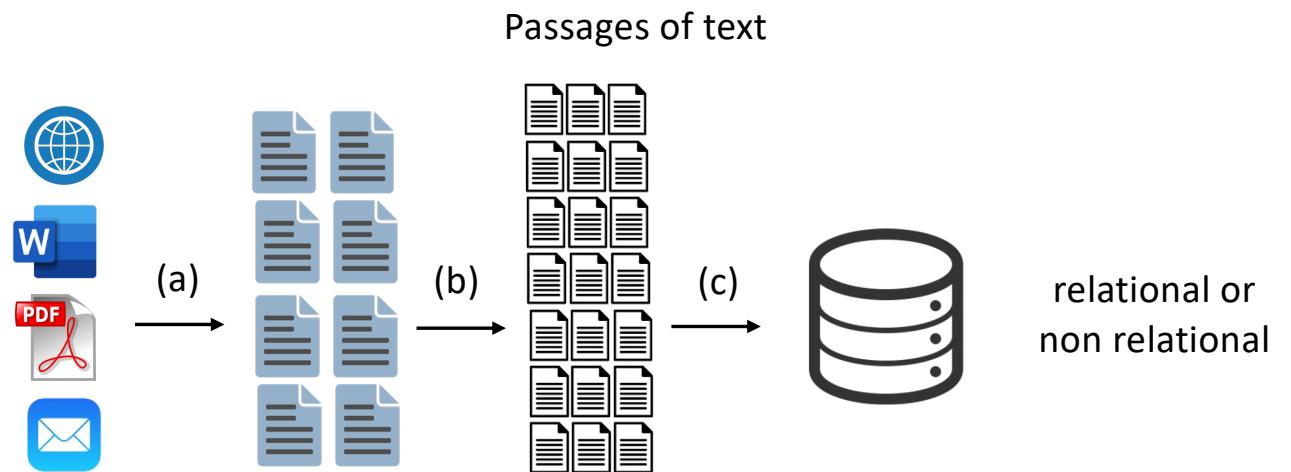


Phase 1: The “traditional” way

Phase 1

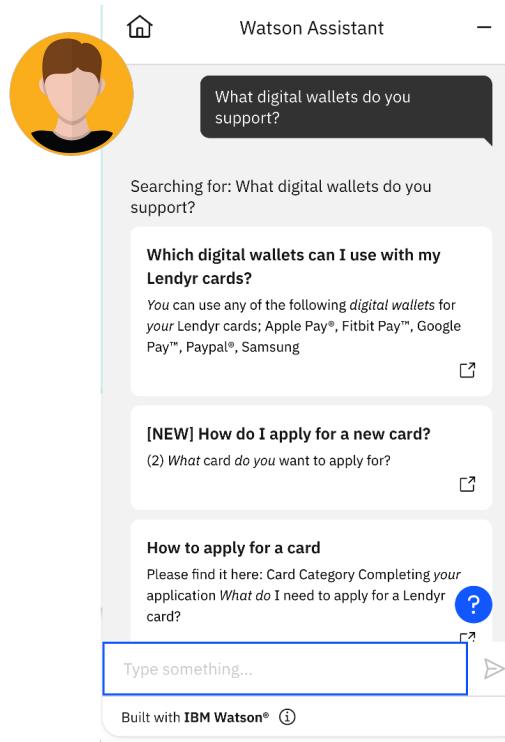
Ingest your data

- (a) Original files to documents
- (b) Documents to chunks
- (c) Chunks to database



Phase 2: Syntactic search

Query the data



How to **apply** for a **card**?



Few paragraphs related to the question

Retrieve all the rows/documents such that the chunk contains either **apply** or **card**



WatsonX.ai

Read the following paragraphs and answer the question:

- [Paragraph 1]
- [Paragraph 2]
- [Paragraph 3]

Question: [User's question]

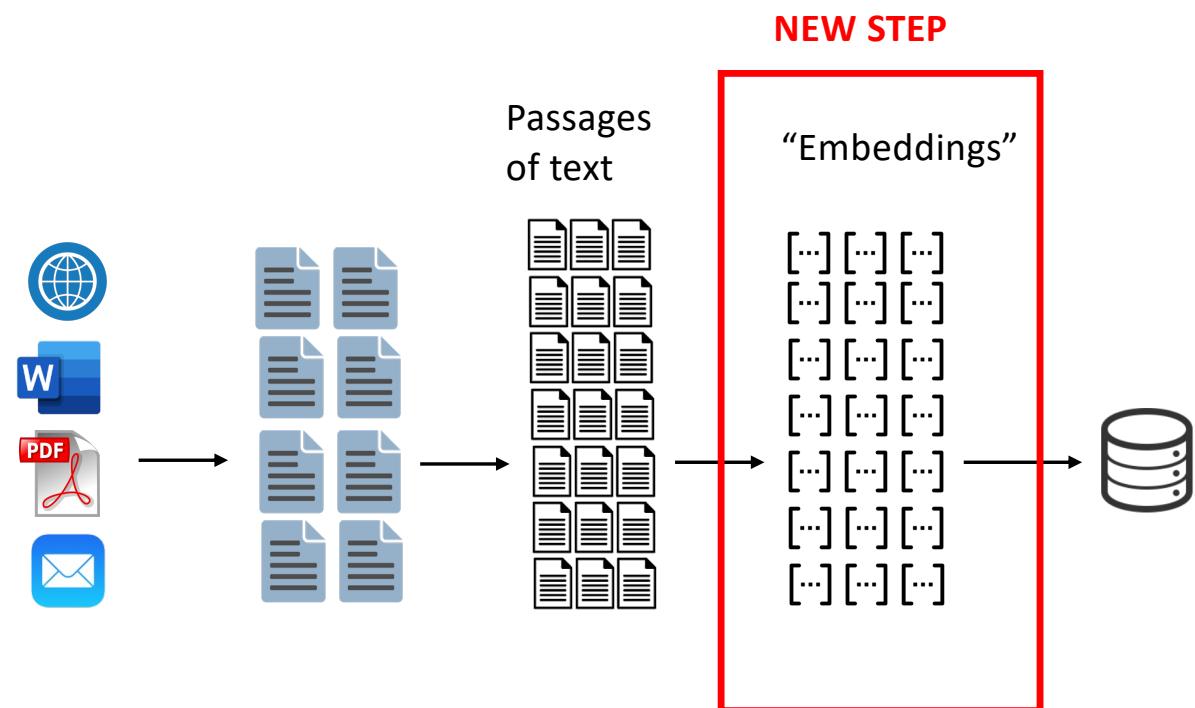
Answer:

Phase 1: The “embeddings” way

Phase 1

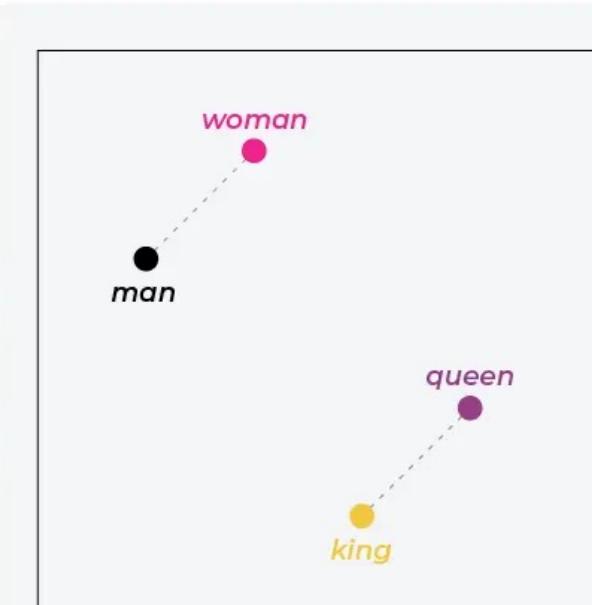
Ingest your data

- (a) Original files to documents
- (b) Documents to chunks
- (c) Chunks to embeddings
- (d) Embeddings to vector store



Phase 1: The “embeddings” way

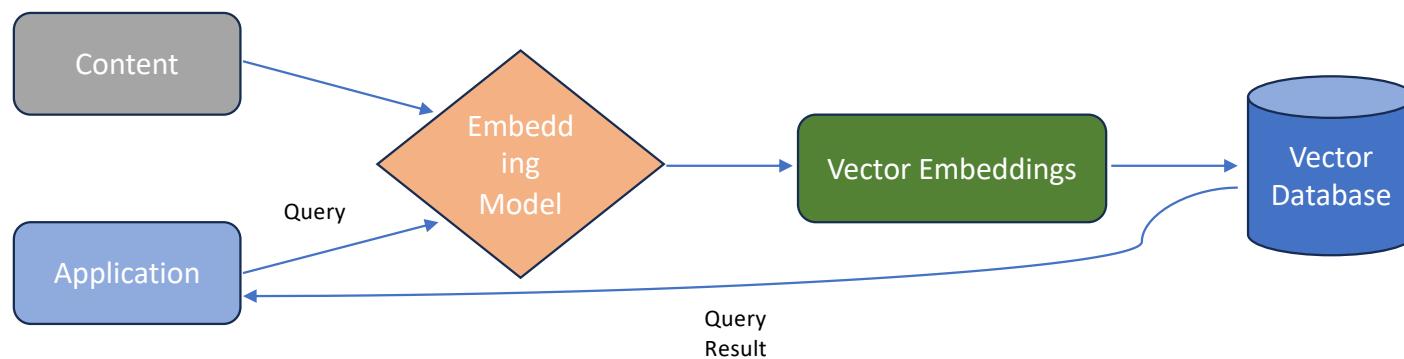
	living being	feline	human	gender	royalty	verb	plural
<i>man</i> →	0.6	-0.2	0.8	0.9	-0.1	-0.9	-0.7
<i>woman</i> →	0.7	0.3	0.8	-0.7	0.1	-0.5	-0.4
<i>king</i> →	0.5	-0.4	0.7	0.8	0.9	-0.7	-0.6
<i>queen</i> →	0.8	-0.1	0.8	-0.9	0.8	-0.5	-0.9
word	Word embedding						



Visualization of word embedding

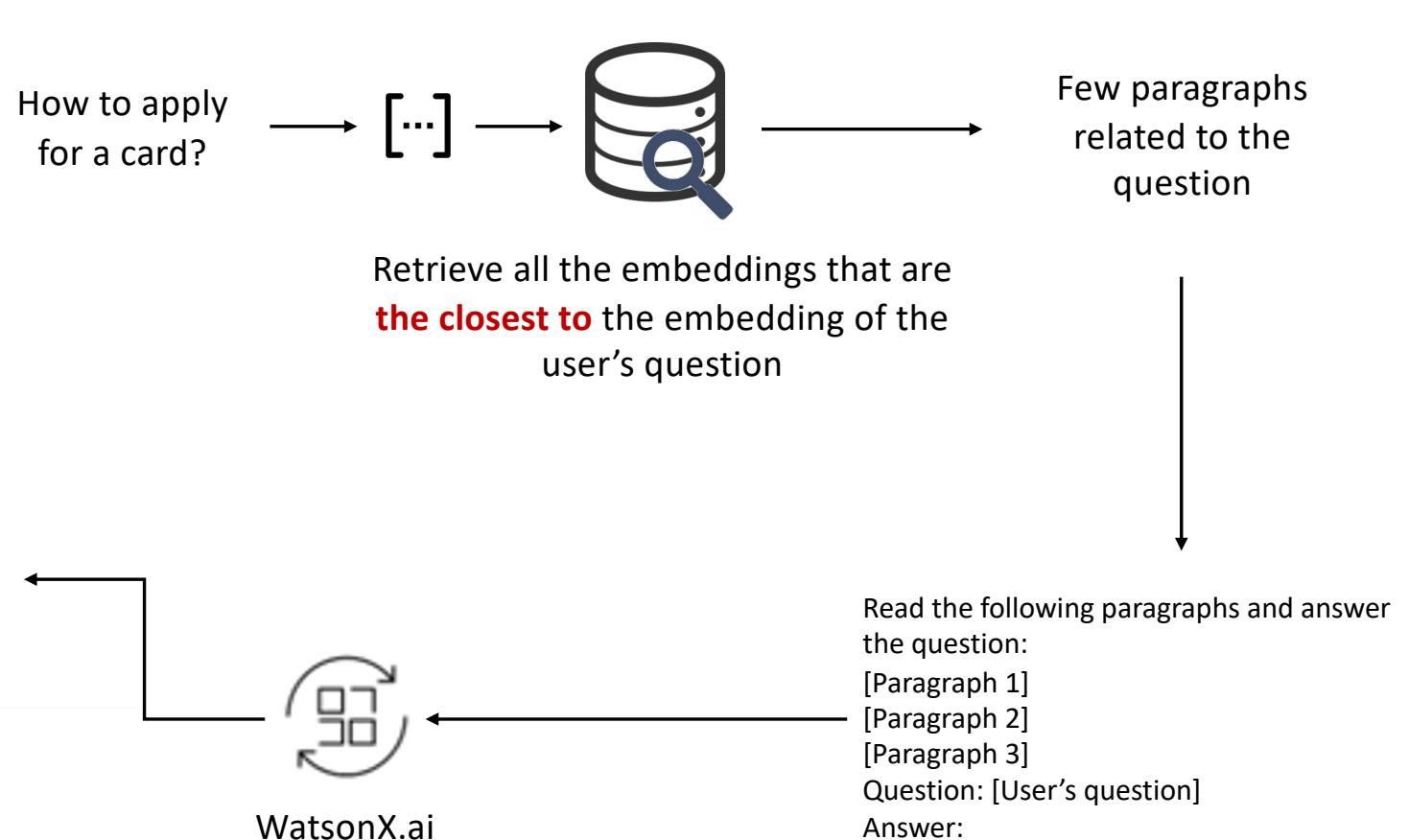
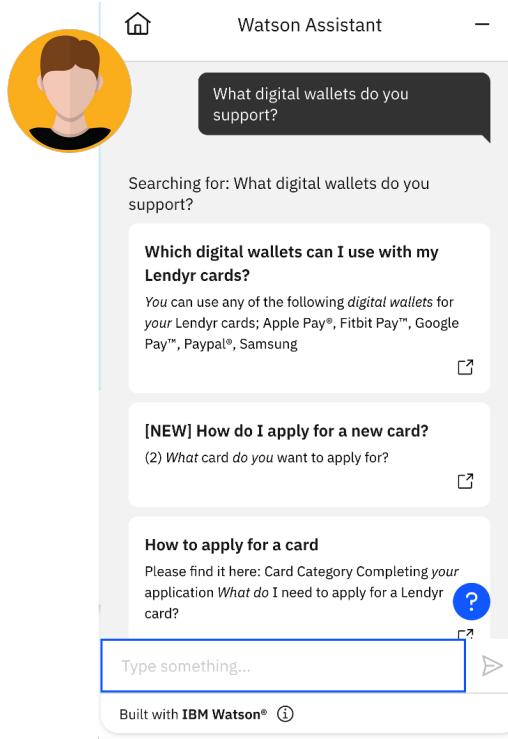
Phase 1: The “embeddings” way – VectorDB

- A vector database is a special type of database that can store high-dimensional vectors which are mathematical representations of the features.
- This data is nothing, but a vector created through embeddings
- Use cases
 - Recommendation systems
 - Anomaly detections
 - NLP



Phase 2: Semantic search

Query the data



Syntactic vs. Semantic search

Why semantic is a “better” way of searching for information

The user expresses himself in his/her own way, whereas the documents usually use “specialized” terms.

Examples:



Paid leave of absence
(IBM HR documents)

Corporate assets
(Bank's code of ethics)

Revenue, profits, benefits
(10K form)



Day off

Company's laptop

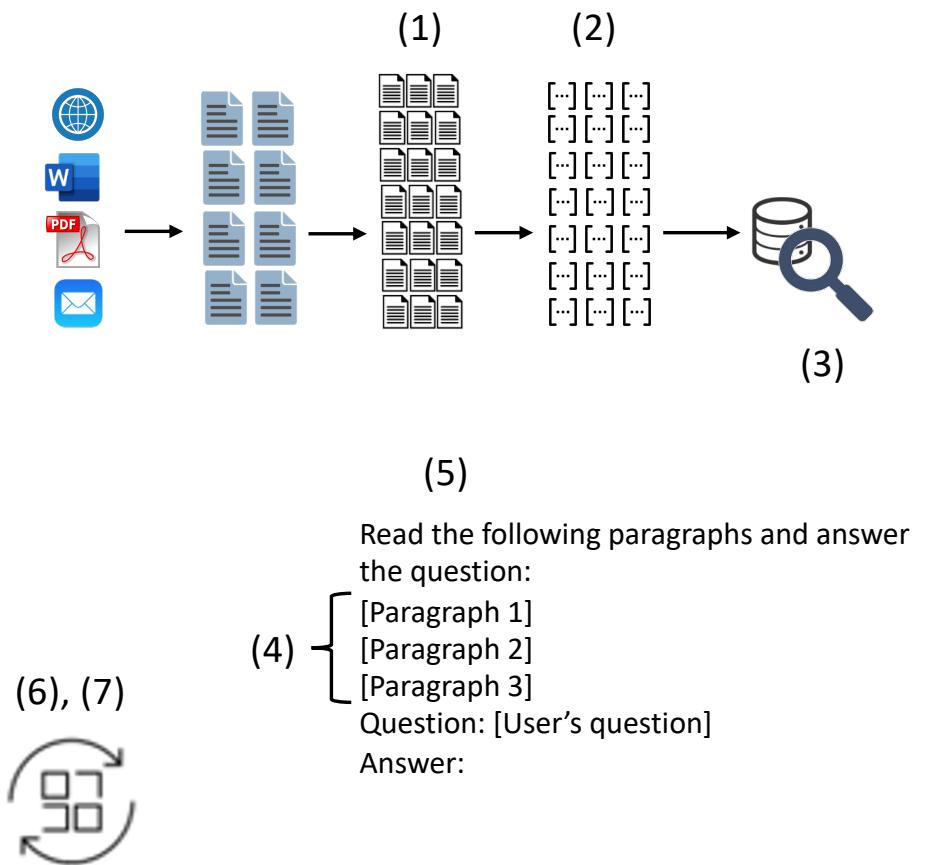
Money

How to improve the accuracy?

Optimize the config at each and every step:

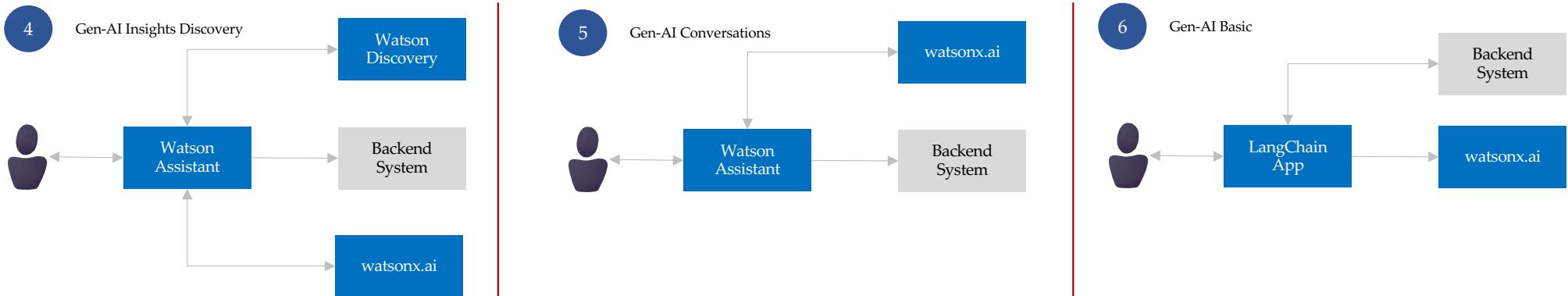
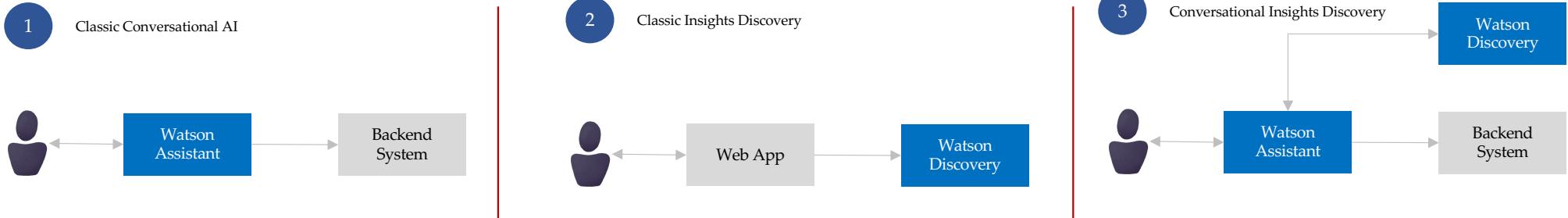
1. Length of the chunks of texts
2. Choice of embeddings library
3. Distance function between embeddings
4. Number of chunks retrieved from the database
5. Prompt
6. LLM parameters (temperature, topK, top, etc)
7. Choice of LLM
8. Etc.

BUT there are more efficient ways

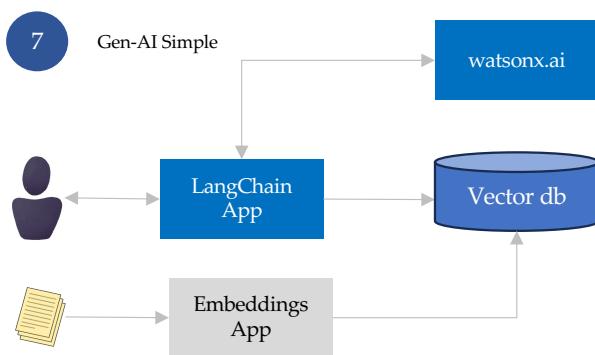


WatsonX.ai

Watson / watsonx.ai Strawman Patterns



- Tasks**
1. Summarization
 2. Q & A
 3. Extraction
 4. Classification
 5. Generation
 6. Transformation
 7. Listing
 8. Comparison

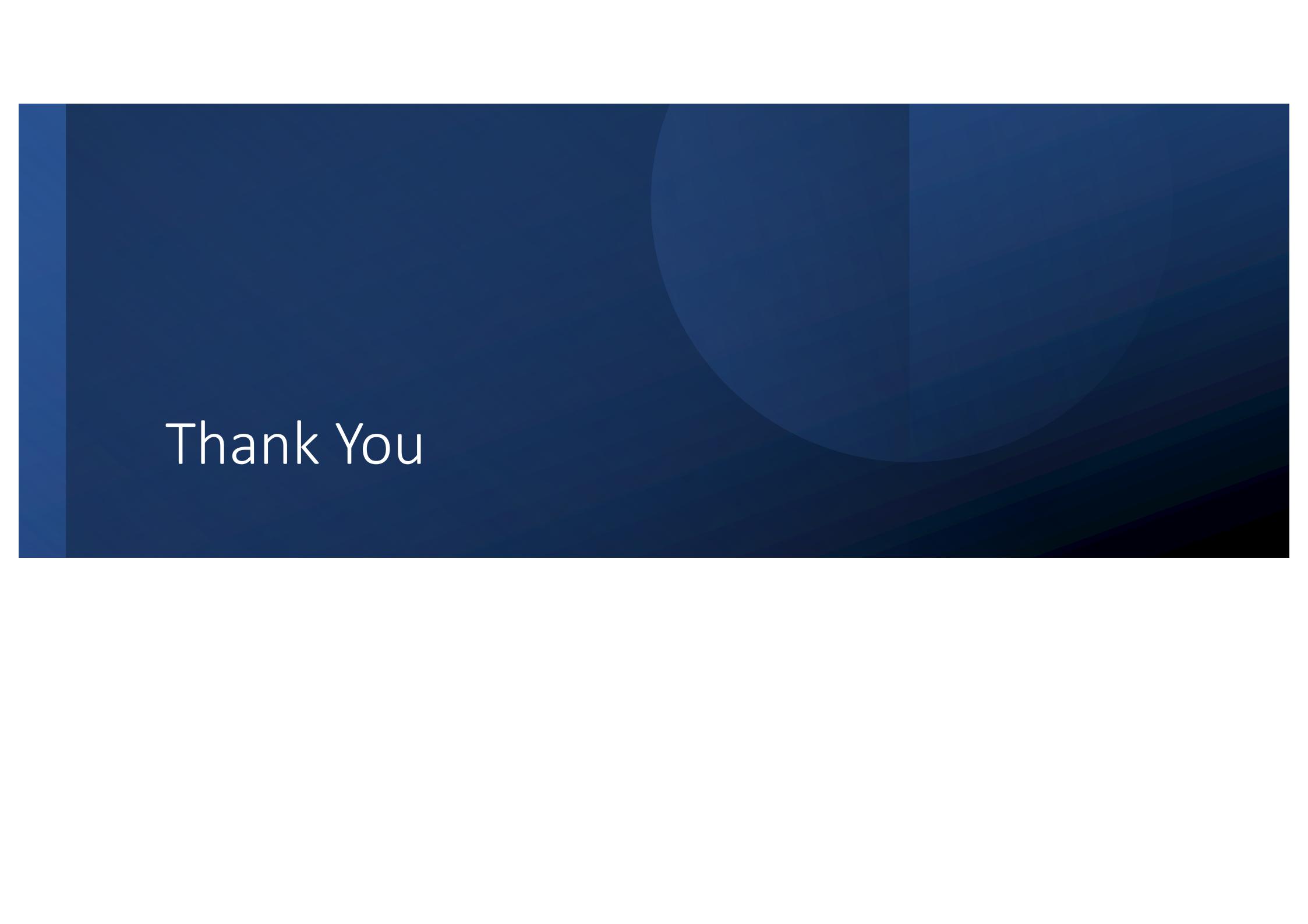


- IBM Tech**
1. Foundation Models
 2. Prompt Lab
 3. Watson Assistant
 4. Watson Discovery
 5. Watson Speech
 6. IBM Cloud

- Sample Technology Stack (non-IBM)**
1. Python
 2. LangChain
 3. ChromaDB, Pinecone...
 4. Flask, FastAPI...
 5. MySQL

Watsonx.ai using Python SDK

Demo



Thank You