# AI 3-in-1: Agents, RAG, and Local Models

Presented by Brent Laster &

Tech Skills Transformations LLC
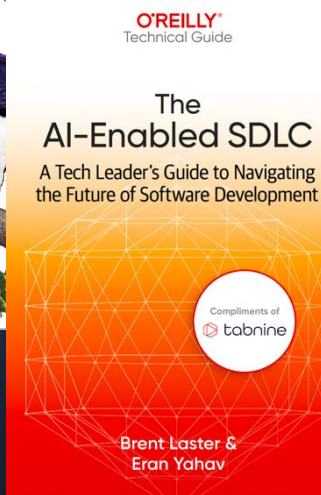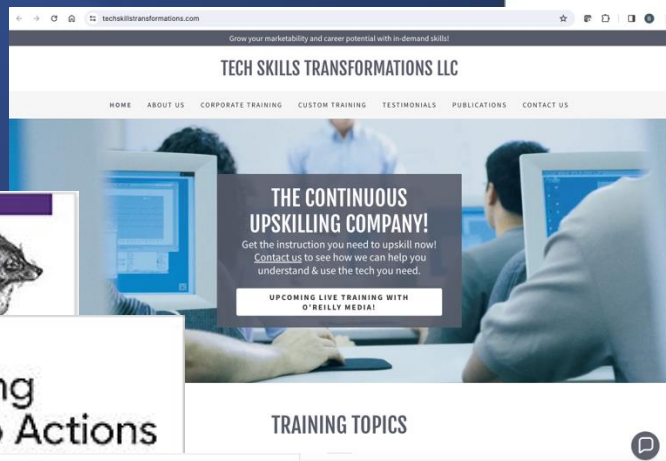
TECHUPSKILLS ®

# About me

- Founder, Tech Skills Transformations LLC
- https://getskillsnow.com
- info@getskillsnow.com

- Long career in corporate as dev, manager, and director in DevOps and other areas
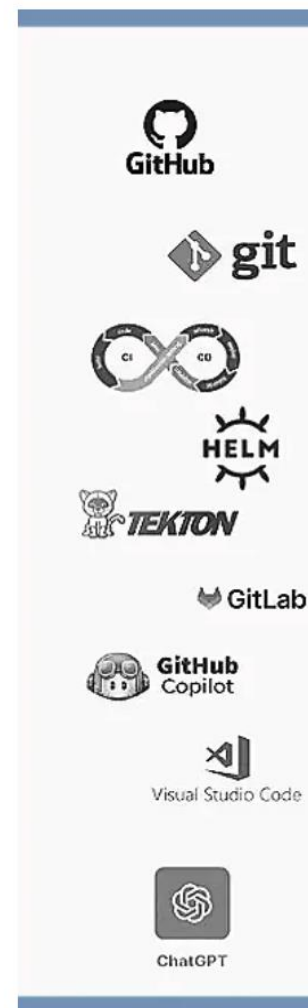
- Author
  - O'Reilly "reports"
  - Books
    - Professional Git
    - Jenkins 2 – Up and Running
    - Learning GitHub Actions
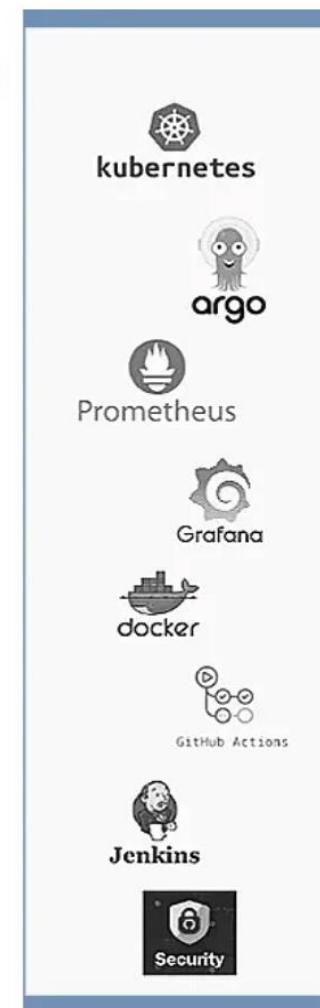    - Learning GitHub Copilot
    - AI-Enabled SDLC

- Speaker

- Social media
  - ❑ **LinkedIn: brentlaster**
  - ❑ **X: @BrentCLaster**
  - ❑ **Bluesky: brentclaster.bsky.social**
  - ❑ **GitHub: brentlaster**

LISTED BELOW ARE A FEW OF THE TECHNOLOGIES FOR WHICH WE OFFER HANDS-ON TRAINING. THESE CAN BE CUSTOMIZED FOR ANY SIZE TEAM FROM 1-100 AND FOR ANY LEVEL FROM BEGINNER TO ADVANCED.

- ArgoCD
- Containers
- Docker
- Gerrit
- Git
- GitHub Actions
- GitHub Codespaces
- GitHub Copilot
- GitHub Foundations
- GitHub Security
- GitLab
- GitOps
- Gradle
- Grafana
- Helm
- Jenkins
- Kubernetes
- Kustomize
- LLMs
- Prometheus
- Tekton
- VS Code

# Running models locally
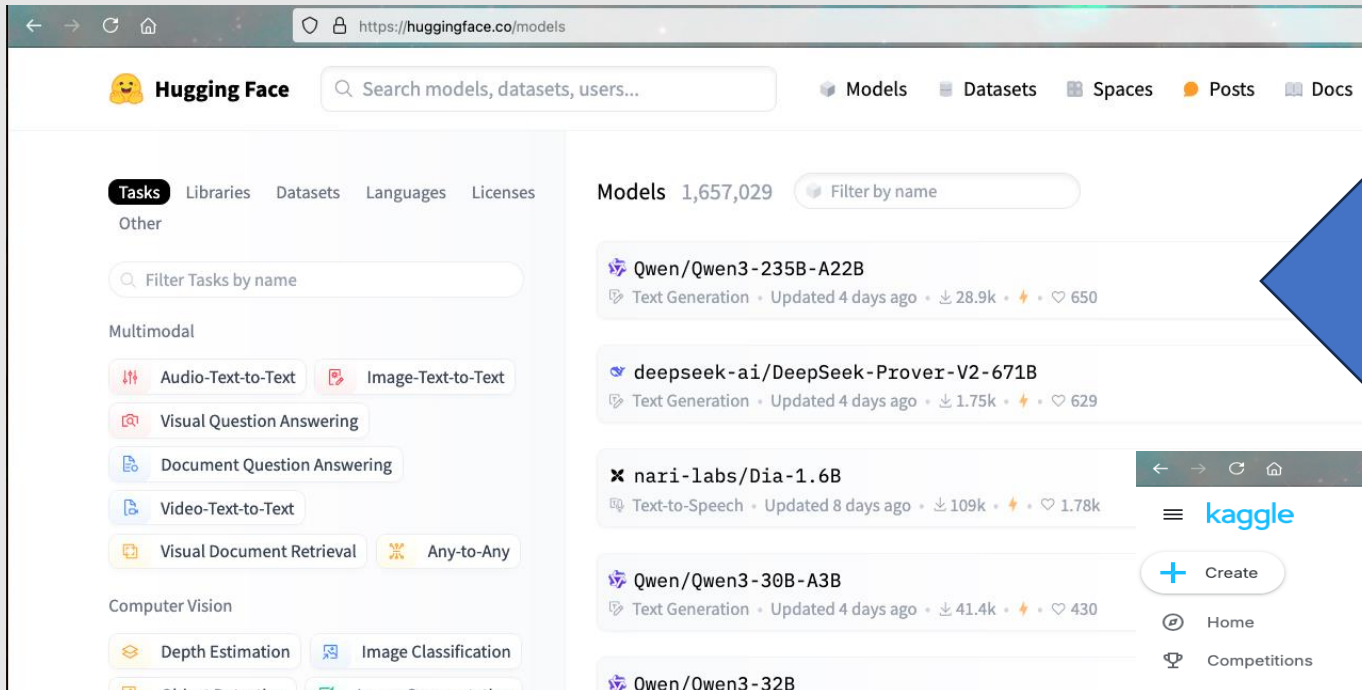
# Why run models locally?

- Privacy - no need to share data

- Gives you control over setup, configuration, and customization options
  - Can tailor LLM to your needs, experiment with settings, integrate into your infra

- Can easily swap between different models for different tasks

- Work in offline mode

- Cost savings
  - No charges for subscriptions or API calls
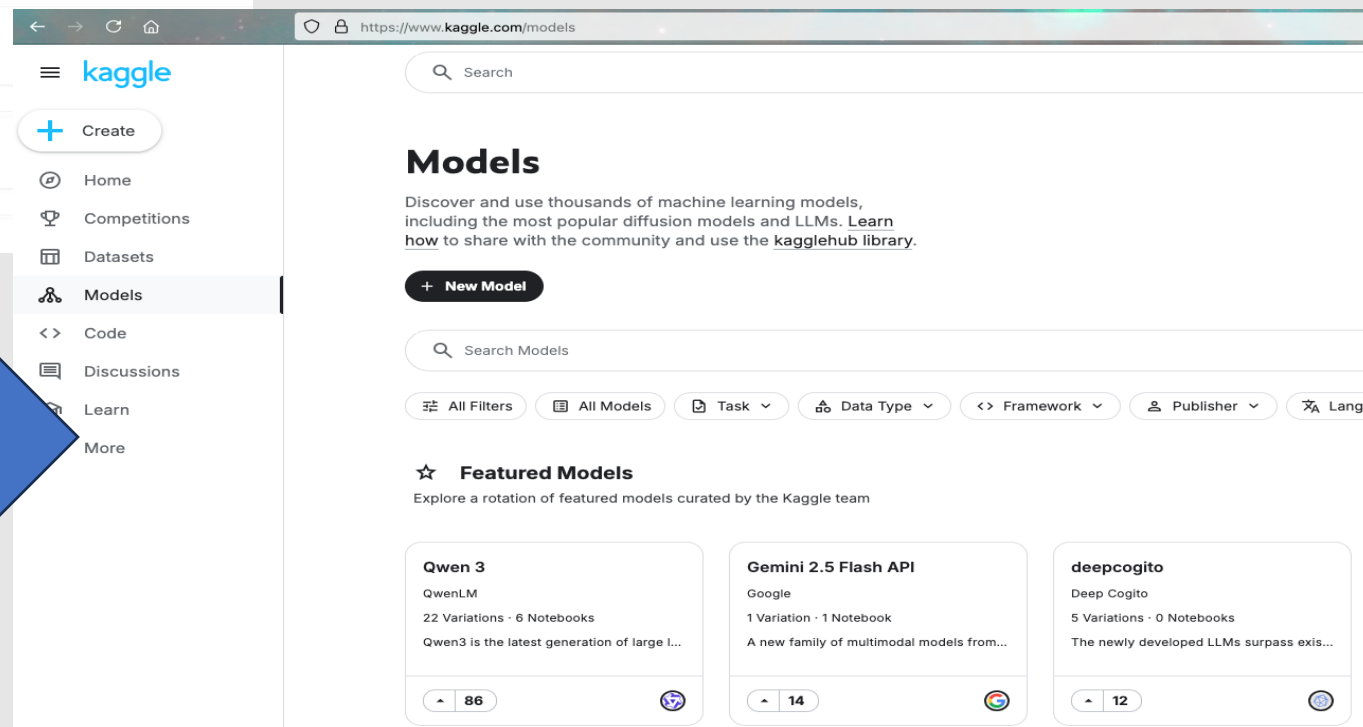
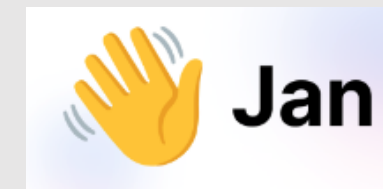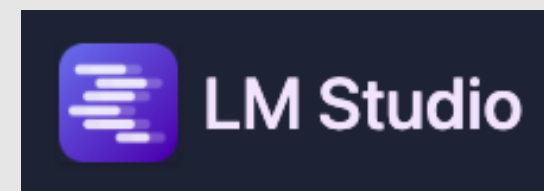- No censoring of results

# Where to get models +

http://huggingface.co/models

http://kaggle.com/models

# Options for running LLMs locally

- GPT4All - https://github.com/nomic-ai/gpt4all

- LM Studio - https://lmstudio.ai

- Jan AI - https://jan.ai

- llama.cpp - https://github.com/ggerganov/llama.cpp

- LlamaFile - https://github.com/Mozilla-Ocho/llamafile

- Ollama - https://ollama.com/

- HuggingFace Transformers - https://huggingface.co/docs/transformers

- More!

# Ollama

- Command line tool for downloading, exploring and using LLMs on local machine

- open source

- supports most of Hugging Face's popular models

- allows uploading new ones

- Links:
  - main site:  https://ollama.com
  - GitHub: https://github.com/ollama/

- Advantages
  - speeds up and simplifies
    - » model selection and download
    - » configuring endpoints
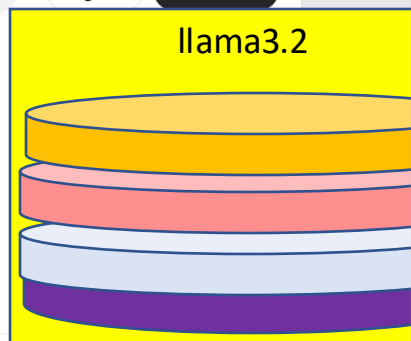    - » integration with Python or JavaScript codebase

# Working with Ollama #1



(py_env) @gwstudent2 → /workspaces/3in1 (main) $

Ollama Search — ollama.com/search

Discord   GitHub   Models

Embedding   Vision   Tools

Llama 3.2

gemma3

The current, most capable model that runs on a single GPU.

vision   1b   4b   12b   27b

4.1M Pulls   21 Tags   Updated 2 weeks ago

llama3.2

*ollama pull*

Ollama

# Working with Ollama #2

PROBLEMS    TERMINAL    ⋯         🎲 bash - 3in1  ＋ ∨  ▯  🗑  ⋯  ∧  ✕

○ (py_env) @gwstudent2 → /workspaces/3in1 (main) $ ▮

```
>>> Briefly explain what
an AI model is
```

*ollama run*

```
>>> query
```

llama3.2

# Working with Ollama #3

```
○ (py_env) @gwstudent2 → /workspaces/3in1 (main) $ ▯
```
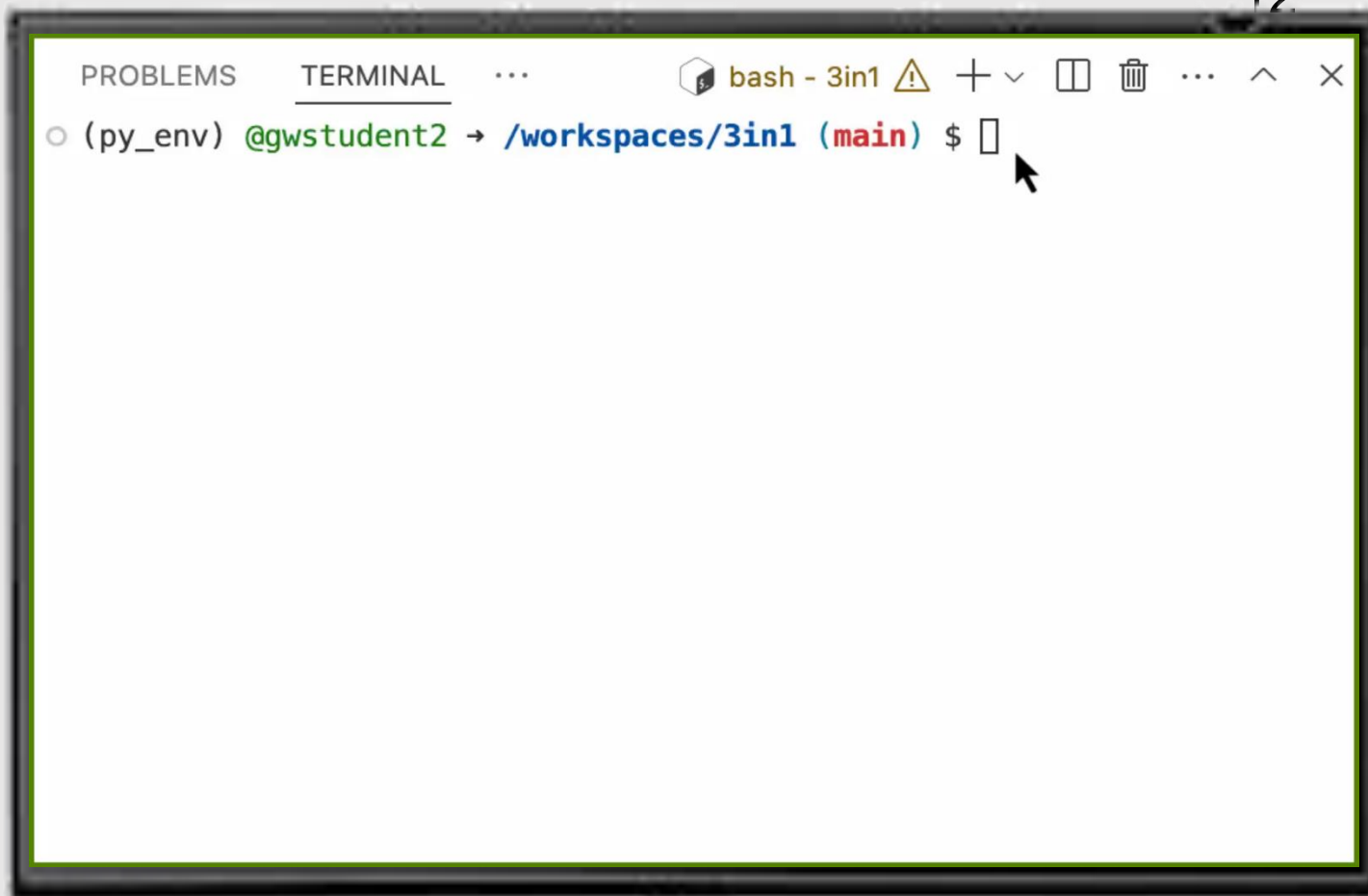
```python
from openai import OpenAI

# Initialize the client to talk to your local Ollama server
client = OpenAI(
    base_url="http://localhost:11434/v1",
    api_key="n-a"
)

# Send a chat completion request
completion = client.chat.completions.create(
    model="llama3.2",
    messages=[
        {"role": "system", "content": "Always answer in 3 bullet points."},
        {"role": "user",   "content": "Tell me what AI is."}
    ],
    temperature=0.7,
)

# Print out the assistant's reply
print(completion.choices[0].message.content)
```

*ollama serve*

Ollama

*http://localhost:11434/v1*

llama3.2

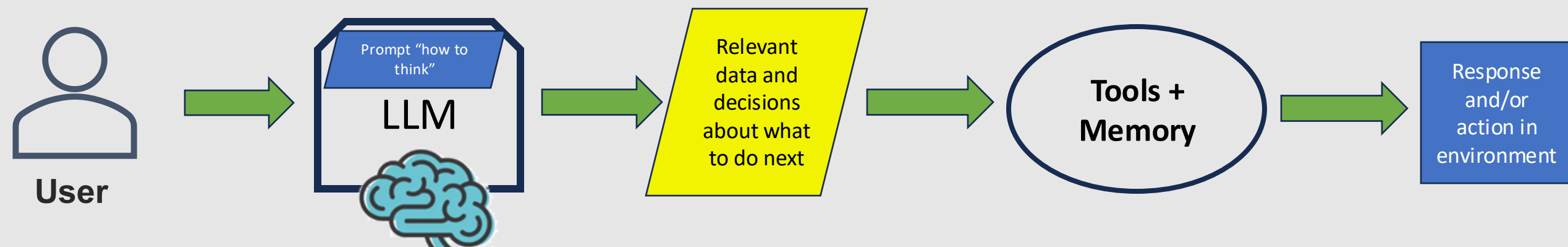# Demo #1 – Simple program to work with local model

# Agents

# What is an AI Agent?

- A **system** that operates within an **environment** by using **sensors** to **perceive** information, a **decision-making mechanism to process and reason about the data**, and **actuators to take actions that influence or update/respond to the environment**

- This interaction enables the agent to achieve specific goals autonomously while continuously learning and adapting over time

- Agents use LLMs to identify key data, drive decisions, and communicate naturally

# Architectural Features of AI Agents

**Planning**

- AI autonomously outlines and executes a logical series of steps for accomplishing a given objective.
- Provides the AI with a way to dynamically adapt its approach based on real-time data and feedback..
- Might employ reflection to evaluate and improve responses
- Example: A research agent plans search → summarize → generate report.

**Tool Use**

- AI agents interact with external APIs, databases, and functions.
- Enhances LLMs by providing access to real-world knowledge.
- Reduces hallucinations by using retrieval-augmented generation (RAG).
- Example: Calling a Python function to perform complex calculations.

**Memory**

- Short-term handles tasks; long term stores knowledge and experience
- Memory ensures consistency and efficiency in multi-step decisions
- Memory recalls preferences to enhance personalization and user experience
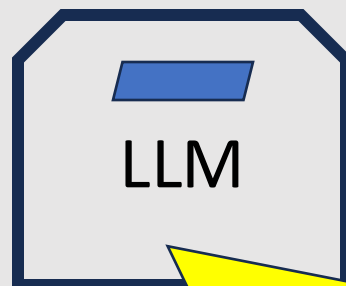- Example: Storing user preferences for future reference or personalized responses

**system_message="""You are an AI assistant designed to help users accurately and efficiently. Your primary goal is to provide precise, helpful, and clear responses.**

**You have access to the following tools:**

**Tool Name: find_weather, Description: Get weather for a location., Arguments: latitude: float, longitude: float, Outputs: string**

**You should think step by step in order to fulfill the objective with a reasoning process divided into Thought/Action/Observation. This cycle can repeat multiple times if needed.**

**You should first reflect with "Thought: {your_thoughts}" on the current query, then (if necessary), call a tool with the proper JSON formatting "Action: {JSON_BLOB}", or else print your final answer starting with the prefix "Final Answer:""""**

# Agent Example
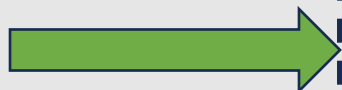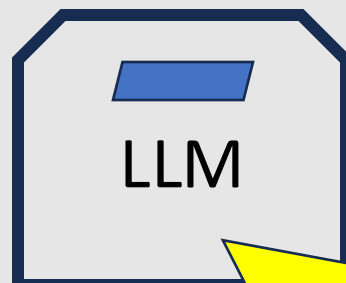
system_message="""You are an AI assistant designed to help users accurately and efficiently. Your primary goal is to provide precise, helpful, and clear responses.
You have access to the following tools:
Tool Name: find_weather, Description: Get weather for a location., Arguments: latitude: float, longitude: float, Outputs: string
You should think step by step in order to fulfill the objective with a reasoning process divided into Thought/Action/Observation. This cycle can repeat multiple times if needed.
You should first reflect with "Thought: {your_thoughts}" on the current query, then (if necessary), call a tool with the proper JSON formatting "Action: {JSON_BLOB}", or else print your final answer starting with the prefix "Final Answer:" """

**What's the weather in Paris?**

**User**

LLM

**Chain of Thought – Step 2: Decide to use tool**
Thought: "I need real-time data, so I will call the '**find_weather**' tool. First, I need to get the latitude and longitude for the tool call.

```
AIResponse(
    tool_calls=[{
        name:
"find_weather"
        parameters:  {
            latitude:
"48.8566",
            longitude:
"2.3522",
        },
        id: "call_tool123",
        type: "tool_invoke"
    }]
)
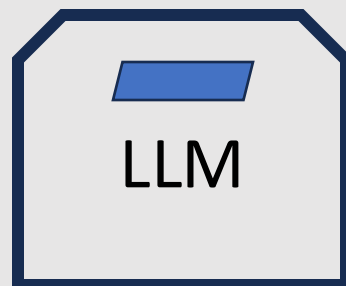```

**Weather Search Tool**

**AI Agent**

# Agent Example

**system_message=**"""You are an AI assistant designed to help users accurately and efficiently. Your primary goal is to provide precise, helpful, and clear responses.
You have access to the following tools:
Tool Name: find_weather, Description: Get weather for a location., Arguments: latitude: float, longitude: float, Outputs: string
You should think step by step in order to fulfill the objective with a reasoning process divided into Thought/Action/Observation. This cycle can repeat multiple times if needed.
You should first reflect with "Thought: {your_thoughts}" on the current query, then (if necessary), call a tool with the proper JSON formatting "Action: {JSON_BLOB}", or else print your final answer starting with the prefix "Final Answer:" """

What's the weather in Paris?

**User**

**LLM**

```
AIResponse(
    tool_calls=[{
        name:
"find_weather"
        parameters: {
            latitude:
"48.8566",
            longitude:
"2.3522",
        },
        id: "call_tool123",
        type: "tool_invoke"
    }]
)
```

Agent parses LLM output identifies JSON tool call, parses it, forms it into actual tool call

```
{
    name:
"find_weather"
    parameters: {
        latitude:
"48.8566",
        longitude:
"2.3522",
    },
    id: "call_tool123",
    type: "tool_invoke"
}
```
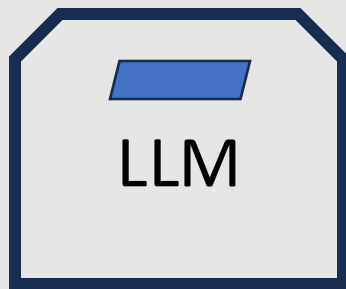
**Weather Search Tool**

**AI Agent**

# Agent Example

system_message="""You are an AI assistant designed to help users accurately and efficiently. Your primary goal is to provide precise, helpful, and clear responses.
You have access to the following tools:
Tool Name: find_weather, Description: Get weather for a location., Arguments: latitude: float, longitude: float, Outputs: string
You should think step by step in order to fulfill the objective with a reasoning process divided into Thought/Action/Observation. This cycle can repeat multiple times if needed.
You should first reflect with "Thought: {your_thoughts}" on the current query, then (if necessary), call a tool with the proper JSON formatting "Action: {JSON_BLOB}", or else print your final answer starting with the prefix "Final Answer:" """

What's the weather in Paris?

**User**

**LLM**

```
AIResponse(
    tool_calls=[{
        name:
"find_weather"
        parameters: {
            latitude:
"48.8566",
            longitude:
"2.3522",
        },
        id: "call_tool123"
```

**Agent executes tool call**

    }]
)

**Weather Search Tool**

```
{
    name:
"find_weather"
    parameters: {
        latitude:
"48.8566",
        longitude:
"2.3522",
    },
    id: "call_tool123",
    type: "tool_invoke"
```
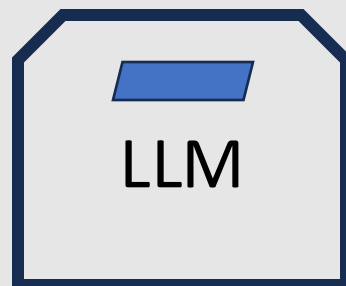
**AI Agent**

# Agent Example

**system_message=**"""You are an AI assistant designed to help users accurately and efficiently. Your primary goal is to provide precise, helpful, and clear responses.
You have access to the following tools:
Tool Name: find_weather, Description: Get weather for a location., Arguments: latitude: float, longitude: float, Outputs: string
You should think step by step in order to fulfill the objective with a reasoning process divided into Thought/Action/Observation. This cycle can repeat multiple times if needed.
You should first reflect with "Thought: {your_thoughts}" on the current query, then (if necessary), call a tool with the proper JSON formatting "Action: {JSON_BLOB}", or else print your final answer starting with the prefix "Final Answer:" """

What's the weather in Paris?

**User**

**LLM**

```
AIResponse(
    tool_calls=[{
        name:
"find_weather"
        parameters:  {
            latitude:
"48.8566",
            longitude:
"2.3522",
        },
        id: "call_tool123",
        type: "tool_invoke"
    }]
)
```

Weather tool returns result

```
ToolResponse(
    content="53 and
rainy",

name="find_weather",
    tool_invoke_id:
"call_tool123"
)
```

**Weather Search Tool**

```
{
    name:
"find_weather"
    parameters:  {
        latitude:
"48.8566",
        longitude:
"2.3522",
    },
    id: "call_tool123",
    type: "tool_invoke"
}
```
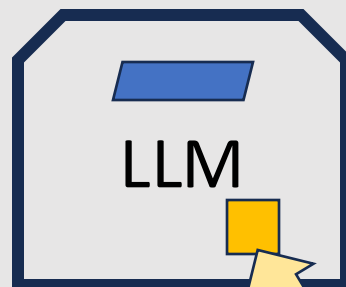
**AI Agent**

# Agent Example
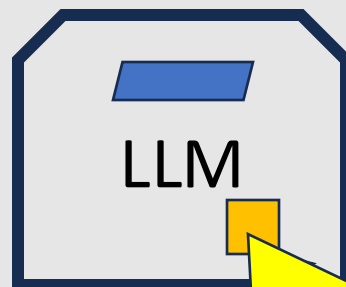
system_message="""You are an AI assistant designed to help users accurately and efficiently. Your primary goal is to provide precise, helpful, and clear responses.
You have access to the following tools:
Tool Name: find_weather, Description: Get weather for a location., Arguments: latitude: float, longitude: float, Outputs: string
You should think step by step in order to fulfill the objective with a reasoning process divided into Thought/Action/Observation. This cycle can repeat multiple times if needed.
You should first reflect with "Thought: {your_thoughts}" on the current query, then (if necessary), call a tool with the proper JSON formatting "Action: {JSON_BLOB}", or else print your final answer starting with the prefix "Final Answer:" """

What's the weather in Paris?

**User**

LLM

AIResponse(
    tool_calls=[{
        name:
"find_weather"
        parameters:  {
            latitude:
"48.8566",
            longitude:
"...",
        },
    d: "call_tool123",
    type: "tool_invoke"
        }]
)

**Agent includes tool output in message/prompt back to model**

ToolResponse(
    content="53 and rainy",

name="find_weather",
    tool_invoke_id:
"call_tool123"
)

**Weather Search Tool**

{
    name:
"find_weather"
    parameters:  {
        latitude:
"48.8566",
        longitude:
"2.3522",
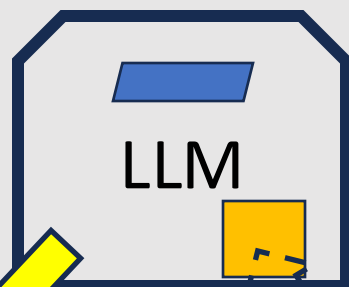    },
    id: "call_tool123",
    type: "tool_invoke"

**AI Agent**

system_message="""You are an AI assistant designed to help users accurately and efficiently. Your primary goal is to provide precise, helpful, and clear responses.
You have access to the following tools:
Tool Name: find_weather, Description: Get weather for a location., Arguments: latitude: float, longitude: float, Outputs: string
You should think step by step in order to fulfill the objective with a reasoning process divided into Thought/Action/Observation. This cycle can repeat multiple times if needed.
You should first reflect with "Thought: {your_thoughts}" on the current query, then (if necessary), call a tool with the proper JSON formatting "Action: {JSON_BLOB}", or else print your final answer starting with the prefix "Final Answer:""""

What's the weather in Paris?

**User**

**LLM**

```
AIResponse(
    tool_calls=[{
        name:
"find_weather"
        parameters: {
            latitude:
"48.8566",
            longitude:
"2.3522",
        },
        id: "call_tool123",
        type: "tool_invoke"
    }]
)
```

```
{
    name:
"find_weather"
    parameters: {
        latitude:
"48.8566",
        longitude:
"2.3522",
    },
    id: "call_tool123",
    type: "tool_invoke"
}
```

**Chain of Thought – Step 3 : Interpret JSON Response**
Thought: "The tool returned weather data for Paris. I will summarize the information concisely.

```
name="find_weather",
    tool_invoke_id:
"call_tool123"
)
```
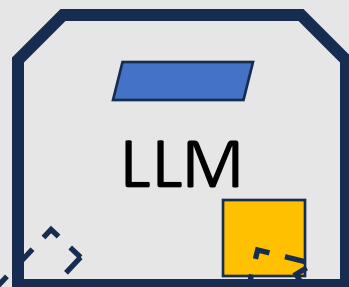
**Weather Search Tool**

**AI Agent**

# Agent Example

system_message="""You are an AI assistant designed to help users accurately and efficiently. Your primary goal is to provide precise, helpful, and clear responses.
You have access to the following tools:
Tool Name: find_weather, Description: Get weather for a location., Arguments: latitude: float, longitude: float, Outputs: string
You should think step by step in order to fulfill the objective with a reasoning process divided into Thought/Action/Observation. This cycle can repeat multiple times if needed.
You should first reflect with "Thought: {your_thoughts}" on the current query, then (if necessary), call a tool with the proper JSON formatting "Action: {JSON_BLOB}", or else print your final answer starting with the prefix "Final Answer:" """

What's the weather in Paris?

**User**

LLM

AIResponse(
    tool_calls=[{
        name:
"find_weather"
        parameters: {
            latitude:
"48.8566",
            longitude:
"2.3522",
        },
        id: "call_tool123",
        type: "tool_invoke"
    }]
)

{
    name:
"find_weather"
    parameters: {
        latitude:
"48.8566",
        longitude:
"2.3522",
    },
    id: "call_tool123",
    type: "tool_invoke"
}

AIFinalResponse(
    content="The current weather in Paris is 53 degrees Celsius with light rain."
)

ToolResponse(
    content="53 and rainy",

name="find_weather",
    tool_invoke_id:
"call_tool123"
)

**Weather Search Tool**

**AI Agent**

# Agent Example

**system_message="""**You are an AI assistant designed to help users accurately and efficiently. Your primary goal is to provide precise, helpful, and clear responses.
You have access to the following tools:
Tool Name: find_weather, Description: Get weather for a location., Arguments: latitude: float, longitude: float, Outputs: string
You should think step by step in order to fulfill the objective with a reasoning process divided into Thought/Action/Observation. This cycle can repeat multiple times if needed.
You should first reflect with "Thought: {your_thoughts}" on the current query, then (if necessary), call a tool with the proper JSON formatting "Action: {JSON_BLOB}", or else print your final answer starting with the prefix "Final Answer:"**"""**

What's the weather in Paris?

**User**

**LLM**

```
AIResponse(
    tool_calls=[{
        name:
"find_weather"
        parameters: {
            location: "Paris",
        },
        id: "call_tool123",
        type: "tool_invoke"
    }]
)
```

```
{
    name:
"find_weather"
    parameters: {
        latitude:
"48.8566",
        longitude:
"2.3522",
    },
    id: "call_tool123",
    type: "tool_invoke"
}
```

```
AIFinalResponse(
    content="The
current weather in Paris
is 53 degrees Celsius
with light rain."
)
```

```
ToolResponse(
    content="53 and
rainy",

name="find_weather",
    tool_invoke_id:
"call_tool123"
)
```

**Weather Search Tool**

**AI Agent**

# Demo #2 – Adding agency to our code

# RAG

# What is RAG and how does it work?

- **Combination of retrieval and generation**: RAG combines information retrieval (like a search engine) with text generation (like a language model).

- **Uses external knowledge**: Instead of relying solely on pre-trained knowledge, RAG retrieves relevant documents or data from an external source (like a database or private knowledge bases) to generate more accurate and up-to-date responses.

- **Improves factual accuracy**: By pulling in real-time data or documents, RAG reduces the risk of generating factually incorrect or outdated information.

- **Two-step process**:
  - **Retrieve**: The model searches for relevant information from a knowledge source.
  - **Generate**: It then uses the retrieved data to create a coherent, contextually accurate answer.



**Retrieval Phase**

User Query → Embed Query → Vector Search → Top-k Documents

**Generation Phase**

Combine Query + Document → LLM Completion → Final Answer

Source: https://blogs.nvidia.com/blog/what-is-retrieval-augmented-generation/

# How is RAG setup?

*Doc Ingestion and Retrieval*



- You provide data sources and point application to them
- Info is retrieved from the data sources and tokenized, embedded and stored in a data store
- For queries/prompts, application gathers results (most relevant ones) from the vector database with your data

# Embeddings

- Embeddings represent text as sets of numeric data - tensors (lots of dimensions)

- Each dimension stores some info about the text's meaning, context, or syntactical aspects

- Words or sentences with similar meanings are stored closer together in the vector space
    - If two pieces of text are similar syntactically, they will have similar embeddings (smaller distance between their vectors)

- During training, models learn to place text with similar meanings closer together in the embedding space

- Common pre-trained models used for generating embeddings include BERT and variants (RoBERTa, DistilBERT)

- Once you have embeddings, you can use them for NLP tasks like semantic search, text classification, sentiment analysis



B
(0,0,255)

Whi
(255,255

R
(255,0,0)

G
(0,255

# Understanding vectors in AI

- Collection of data points that encapsulate an item's relationship to other items

Distance to Raleigh, NC USA

# Understanding vectors in AI

- Collection of data points that encapsulate an item's relationship to other items

Distance to Raleigh, NC USA

Madrid, Spain

{ 3,960 }

# Understanding vectors in AI

- Collection of data points that encapsulate an item's relationship to other items

Distance to Raleigh, NC USA

Kinshasa, DRC

{
3,960
6,609
}

# Understanding vectors in AI

- Collection of data points that encapsulate an item's relationship to other items

Distance to Raleigh, NC USA

San Francisco

3,960
6,609
2,839.4

# Understanding vectors in AI

- Collection of data points that encapsulate an item's relationship to other items



Distance to Raleigh, NC USA

3,960
6,609
2,839.4
6,001

Cairo, Egypt

# Understanding vectors in AI

- Collection of data points that encapsulate an item's relationship to other items



Distance to Raleigh, NC USA

$$\begin{bmatrix} 3{,}960 \\ 6{,}609 \\ 2{,}839.4 \\ 6{,}001 \\ 507.6 \end{bmatrix}$$

# Understanding vectors in AI

- Collection of data points that encapsulate an item's relationship to other items

Distance to Raleigh, NC USA

London, UK

$$\begin{cases} 3,960 \\ 6,609 \\ 2,839.4 \\ 6,001 \\ 507.6 \\ 3,872 \end{cases}$$

# Understanding vectors in AI

- Collection of data points that encapsulate an item's relationship to other items

Distance to Raleigh, NC USA

Nairobi, Kenya

3,960
6,609
2,839.4
6,001
507.6
3,872
7,679

# Understanding vectors in AI

- Collection of data points that encapsulate an item's relationship to other items

Distance to Raleigh, NC USA

$$
\begin{Bmatrix}
3{,}960 \\
6{,}609 \\
2{,}839.4 \\
6{,}001 \\
507.6 \\
3{,}872 \\
7{,}679 \\
2{,}870.1
\end{Bmatrix}
$$

Seattle, WA

# Semantic meaning / relationships

- Suppose we have 3 words
- King and Queen are more similar to each other than they are to lunch
- In order for neural net to understand the relationships, each word needs to be represented as a vector
- Suppose each word is represented by a 2-dimensional vector

King $\begin{bmatrix} -130.16 \\ 89.5 \end{bmatrix}$

Queen $\begin{bmatrix} -115.43 \\ 95.2 \end{bmatrix}$

Lunch $\begin{bmatrix} -89.5 \\ 34.3 \end{bmatrix}$

# Embedding space

- Plotting in 2-dimensional embedding space shows relationships
- Way to let NN understand relationships between words
- We want the NN to learn that King and Queen are more similar to each other than they are to lunch

King $\begin{bmatrix} -130.16 \\ 89.5 \end{bmatrix}$

Queen $\begin{bmatrix} -115.43 \\ 95.2 \end{bmatrix}$

Lunch $\begin{bmatrix} -89.5 \\ 34.3 \end{bmatrix}$



2-dimensional space for word embeddings

# Searching for Vectors - similarity metrics

- 3 metrics commonly used to determine similarity of two vectors (2-dimensional representation)

*imagine 3 vectors - a,b,c*

$$a = \begin{bmatrix} .01 \\ .07 \\ .1 \end{bmatrix} \quad b = \begin{bmatrix} .01 \\ .08 \\ .11 \end{bmatrix} \quad c = \begin{bmatrix} .91 \\ .57 \\ .6 \end{bmatrix}$$

**Cosine similarity** - measure the angle between two vectors; values from -1 to 1; 1 = both point in same direction; -1 point in opposite directions; 0 = orthogonal (perpendicular)

**Cosine similarity**

$$sim(u, v) = \frac{u \cdot v}{\|u\|\|v\|} = \frac{\sum_{i=1}^{n} a_n b_n}{\sqrt{\sum_{i=1}^{n} u_n^2}\sqrt{\sum_{i=1}^{n} v_n^2}}$$

$$sim(a, b) = \frac{(a_1 * b_1) + (a_2 * b_2) + (a_3 * b_3)}{\sqrt{a_1^2 + a_2^2 + a_3^2}\sqrt{b_1^2 + b_2^2 + b_3^2}}$$

$$= \frac{(0.01 * 0.01) + (0.07 * 0.08) + (0.1 * 0.11)}{\sqrt{0.01^2 + 0.07^2 + 0.1^2}\sqrt{0.01^2 + 0.08^2 + 0.11^2}}$$

**0.0141**

**Dot product / inner product** - measures how well 2 vectors align with each other; values from - ∞ to ∞; positive values indicate vectors are in same direction; negative values indicate opposite directions; 0 = orthogonal

**Dot product / inner product**

$$u \cdot v = |u||v|cos\theta = \sum_{i=1}^{n} a_n b_n$$

$$a \cdot b = (a_1 b_1) + (a_2 b_2) + (a_3 b_3)$$
$$= (0.01 * 0.01) + (0.07 * 0.08) + (0.1 * 0.11)$$

Dot product formula

**0.0167**

**Euclidean distance** - measures the distance between two vectors; values from 0 to ∞; 0 = identical; larger numbers farther apart

**Euclidean distance**

$$d(u, v) = \sqrt{\sum_{i=1}^{n} (u_i - v_i)^2}$$

$$d(a, b) = \sqrt{(b_1 - a_1)^2 + (b_2 - a_2)^2 + (b_3 - a_3)^2}$$

$$= \sqrt{(0.01 - 0.01)^2 + (0.08 - 0.07)^2 + (0.11 - 0.1)^2}$$

Euclidean distance formula

**0.9998**

credit: https://towardsdatascience.com/similarity-metrics-in-nlp-acc0777e234c

# Visualizing Embeddings and Vector Similarity



source: https://projector.tensorflow.org/?config=https://gist.githubusercontent.com/martin-labrecque/4483ff5a104f0b56417585c3bc9a12f1/raw/57348e12a70c8d70c2c573d3dbc0122ac077556b/journaux_config.json

# Vectors and relationships example

- Query - what words are related to "dog" in model "English Wikipedia"?



Source: http://vectors.nlpl.eu/explore/embeddings/en/MOD_enwiki_upos_skipgram_300_2_2021/dog_NOUN/

# Vector Databases

# Vector Databases

- Specialized database that index and stores *vector embeddings*
- Useful for
  - fast retrieval
  - similarity search
- Offer comprehensive data management capabilities
  - metadata storage
  - filtering
  - dynamic querying based on associate metadata
- Scalable and can handle large volumes of vector data
- Support real-time updates
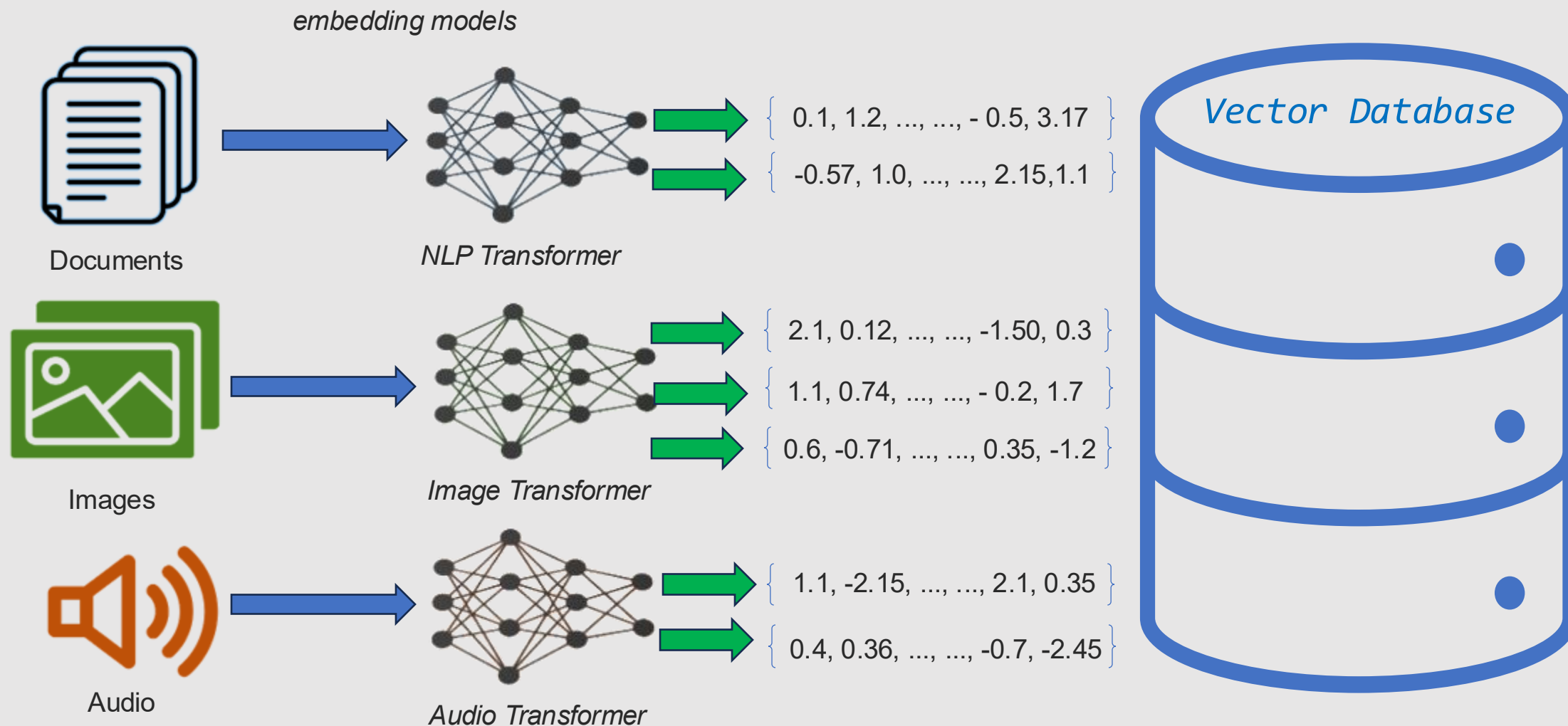- Play key role in AI and ML applications

Vector Database
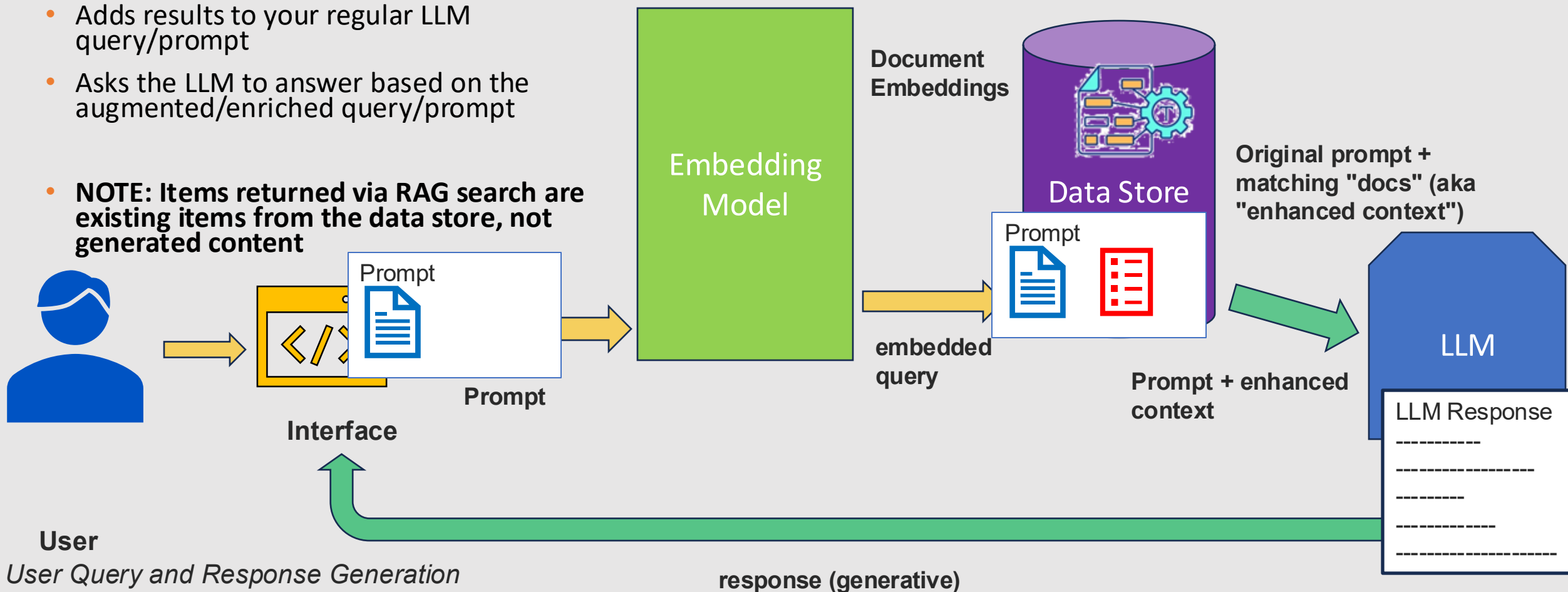
# How data gets into Vector Databases

- Data is input, converted to embeddings (vectors) and stored

- Queries are input, converted to embeddings (vectors) and then **similarity metrics** are used to find results ("nearest neighbors")



*embedding models*

Documents → NLP Transformer → { 0.1, 1.2, ..., ..., - 0.5, 3.17 } { -0.57, 1.0, ..., ..., 2.15,1.1 }

Images → Image Transformer → { 2.1, 0.12, ..., ..., -1.50, 0.3 } { 1.1, 0.74, ..., ..., - 0.2, 1.7 } { 0.6, -0.71, ..., ..., 0.35, -1.2 }

Audio → Audio Transformer → { 1.1, -2.15, ..., ..., 2.1, 0.35 } { 0.4, 0.36, ..., ..., -0.7, -2.45 }

**Vector Database**

# How does RAG work?

- For queries/prompts, application gathers results (most relevant ones) from the vector database with your data

- Adds results to your regular LLM query/prompt

- Asks the LLM to answer based on the augmented/enriched query/prompt

- **NOTE: Items returned via RAG search are existing items from the data store, not generated content**
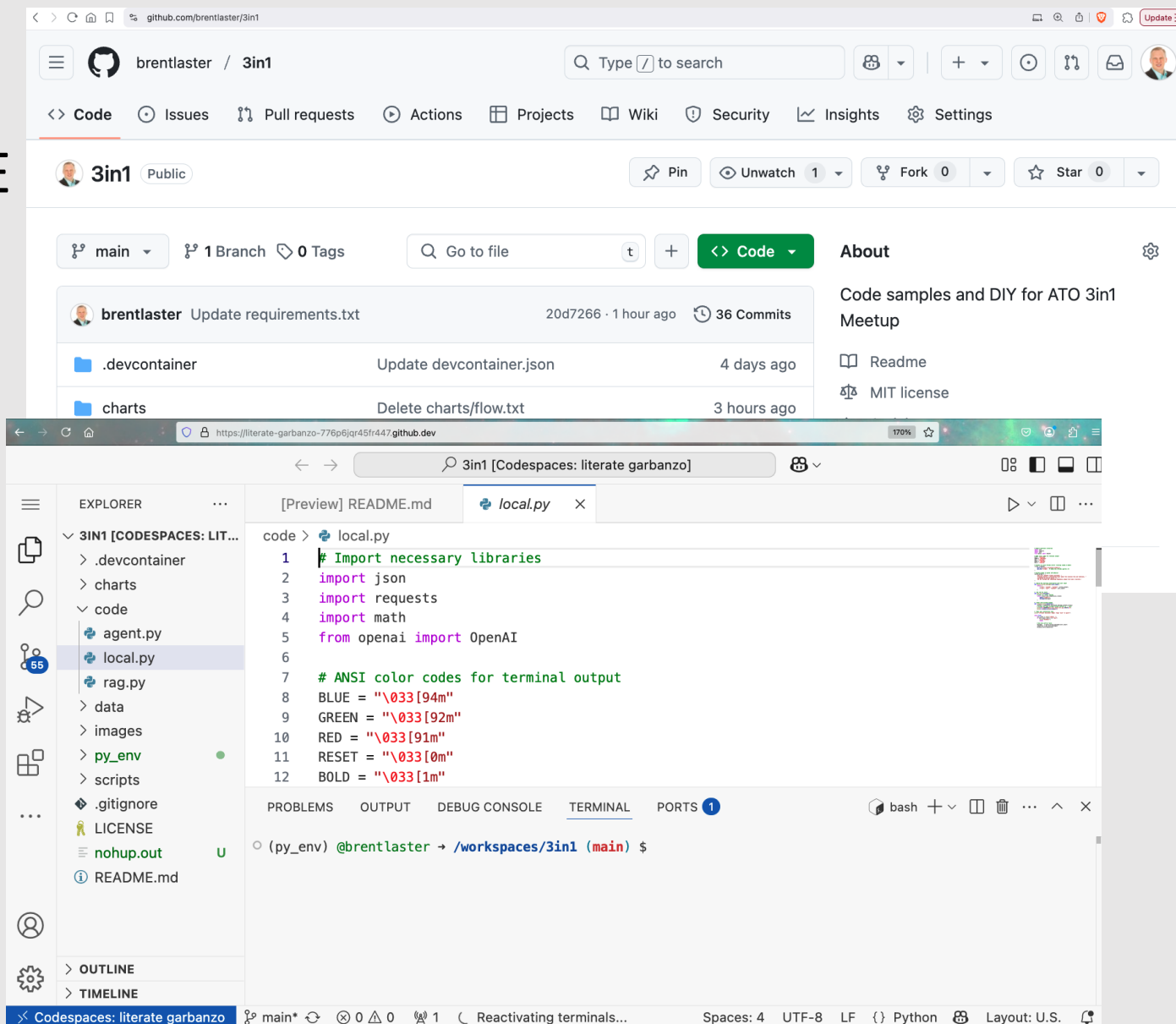


**Document Embeddings**

**Embedding Model**

**Data Store**

Prompt

**Original prompt + matching "docs" (aka "enhanced context")**

**LLM**

LLM Response
----------
----------------
--------
------------
------------------

**embedded query**

**Prompt + enhanced context**

Prompt

**Prompt**

**Interface**

**User**
*User Query and Response Generation*

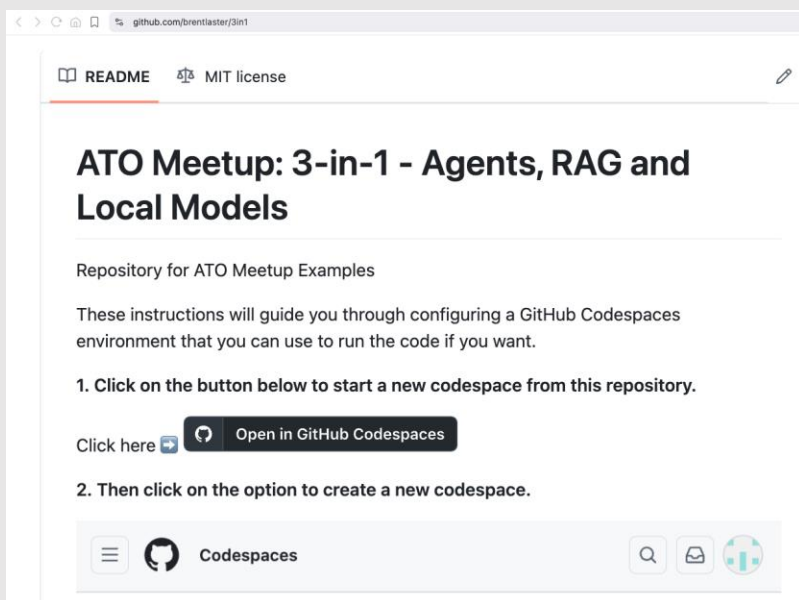**response (generative)**

# Demo #3 – Adding RAG to our code

# DIY – github.com/brentlaster/3in1

- Fork if desired

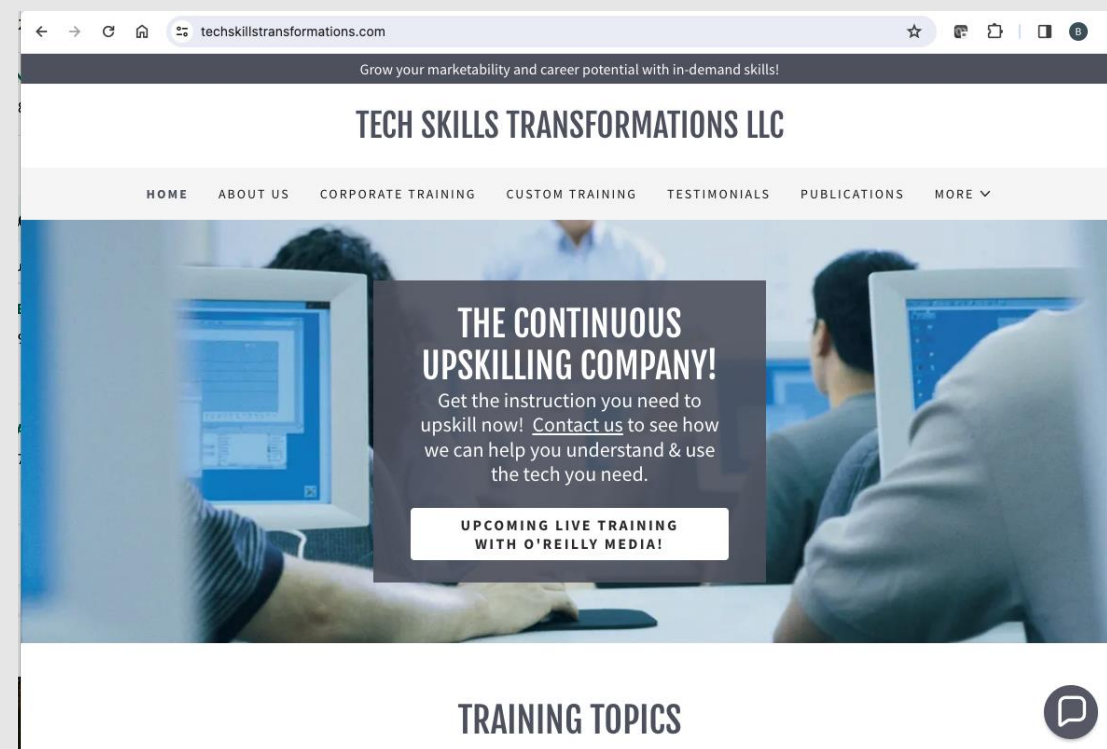- Click on button in README to start codespace

- Follow guide.md

© 2025 Brent C. Laster &
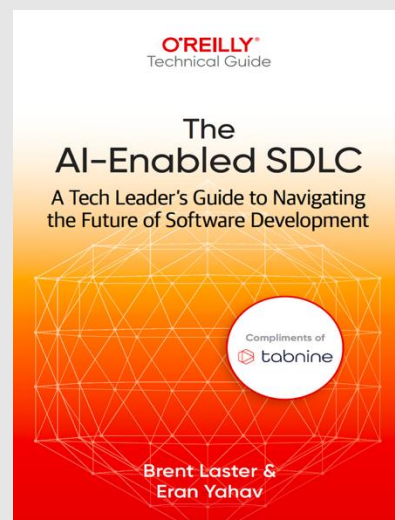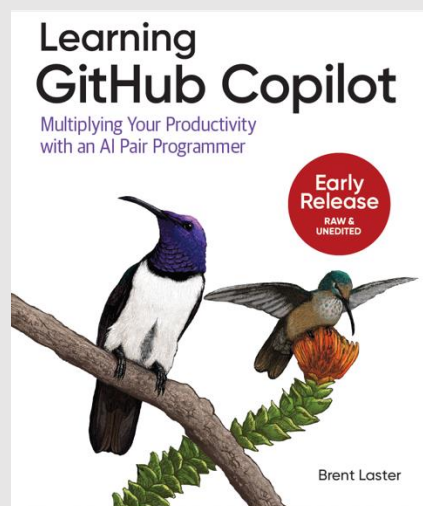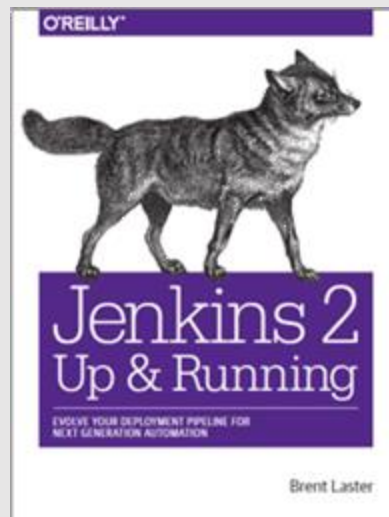Tech Skills Transformations LLC

# That's all - thanks!

**Contact:** training@getskillsnow.com

techskillstransformations.com
getskillsnow.com



□ **LinkedIn: brentlaster**
□ **X: @BrentCLaster**
□ **Bluesky: brentclaster.bsky.social**
□ **GitHub: brentlaster**