



Generative AI for Biomedical Decisions

Sambit Panda, PhD

Christian Cruz

MATRIX AI Consortium

April 29, 2025

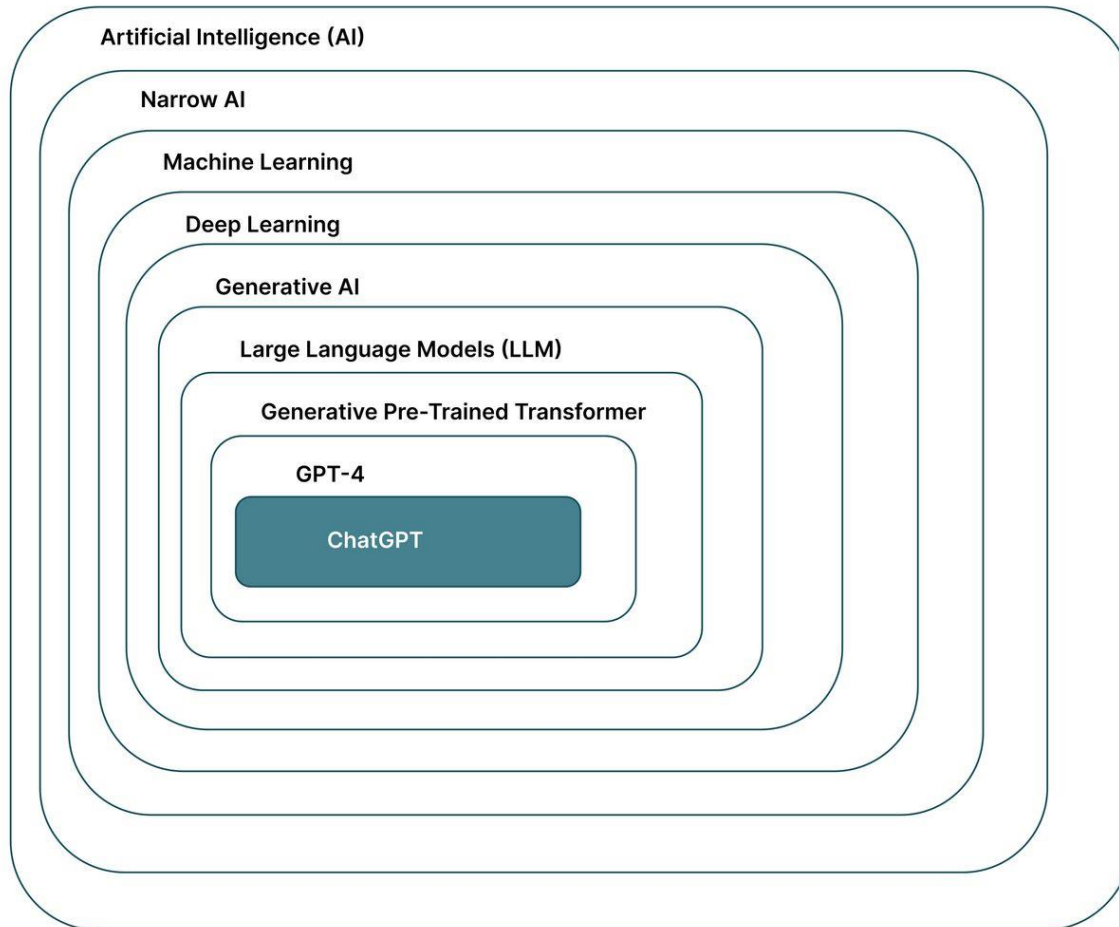
Outline

- Brief overview of Generative AI
- A (recommended) AI workflow
- Customizing performance
- Practical steps to get started
- Pitfalls and considerations
- Conclusion

Overview

What is generative AI?

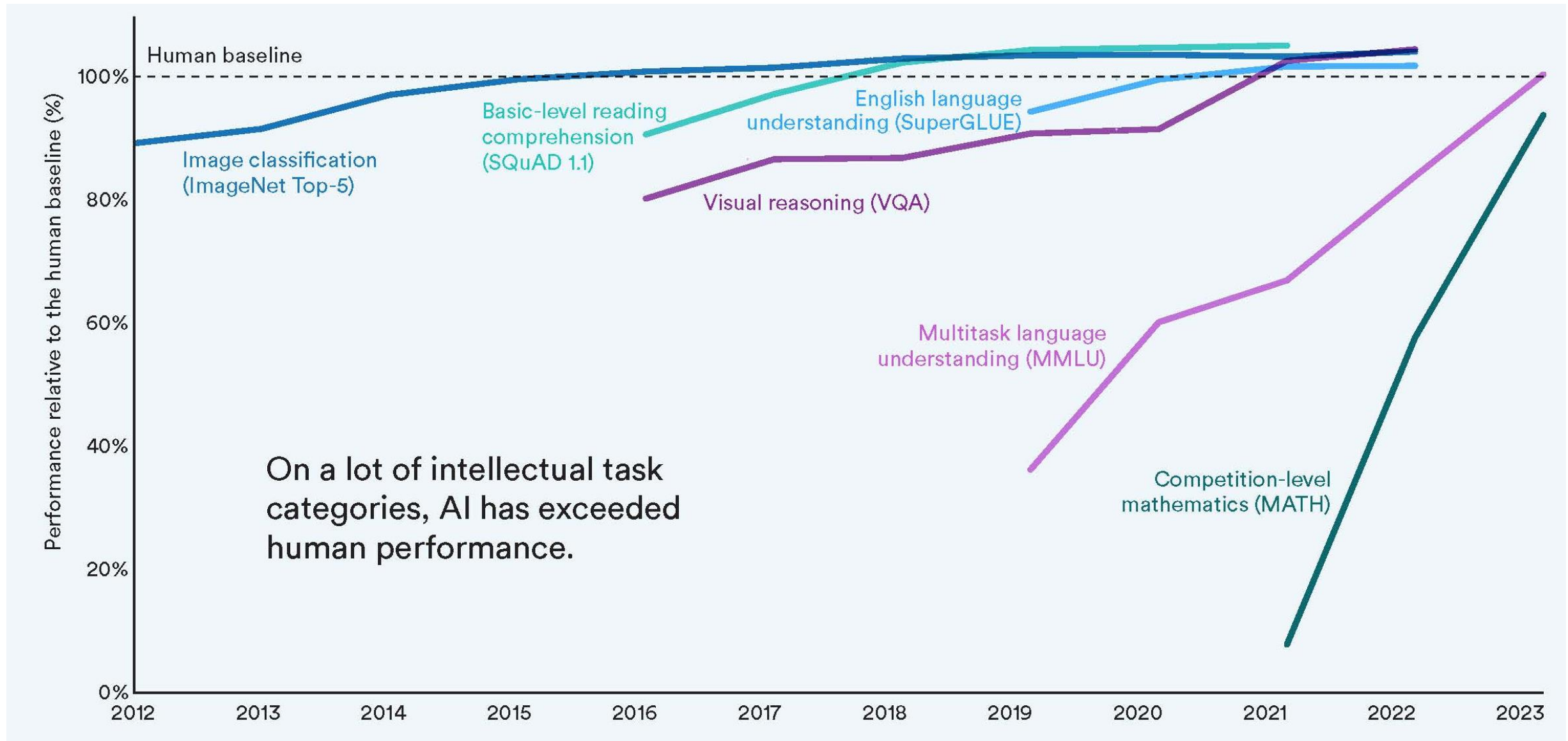
"Generative AI" is a type of artificial intelligence that creates novel content based on patterns learned from existing data.



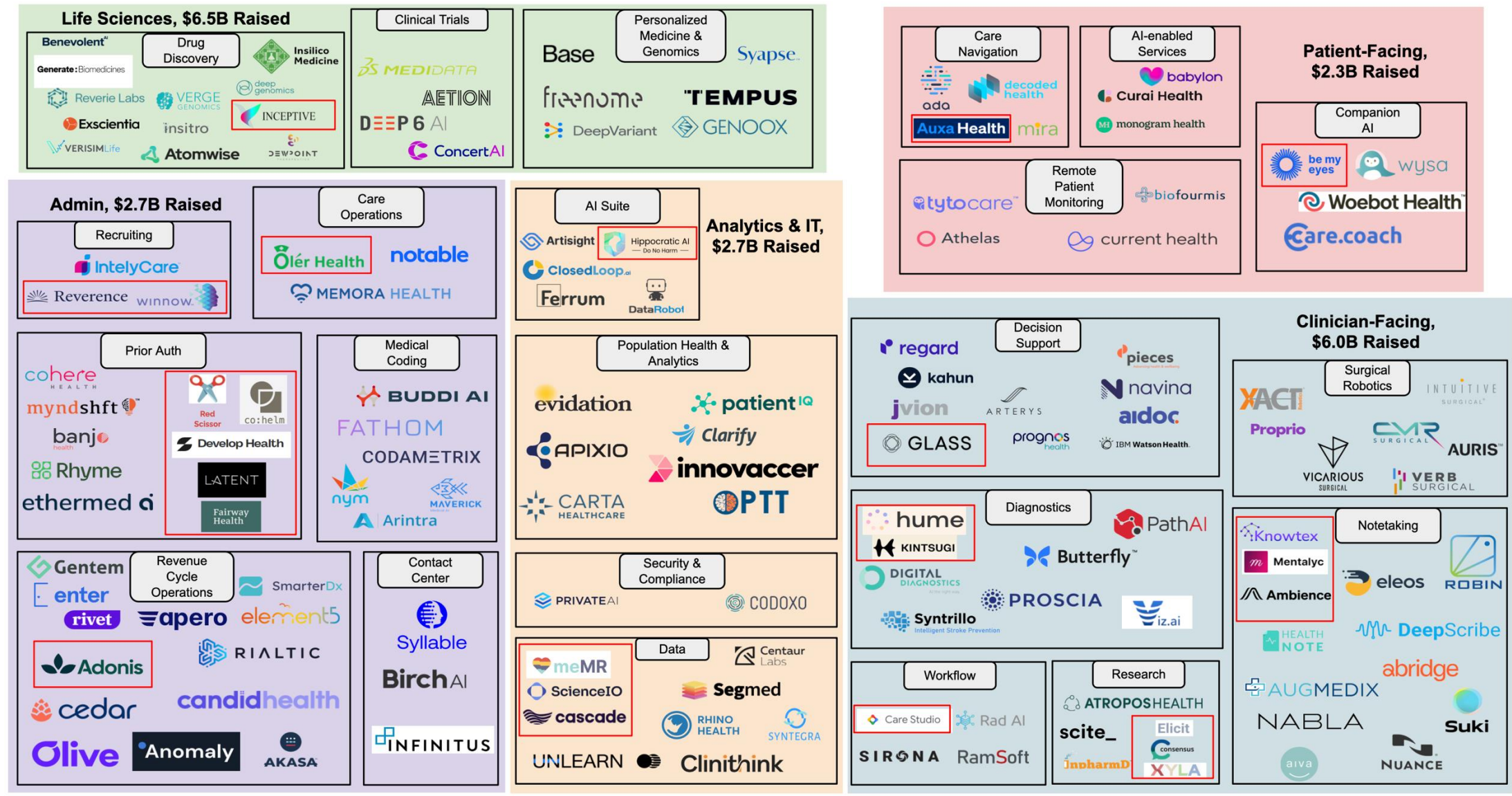
Examples

- **Text:** Summarizing clinical notes, generating trial protocols
- **Biological Sequences:** Designing novel proteins, DNA promoters
- **Images/Videos:** Video and image generation from prompts
- **Data Augmentation:** Synthetic data generation

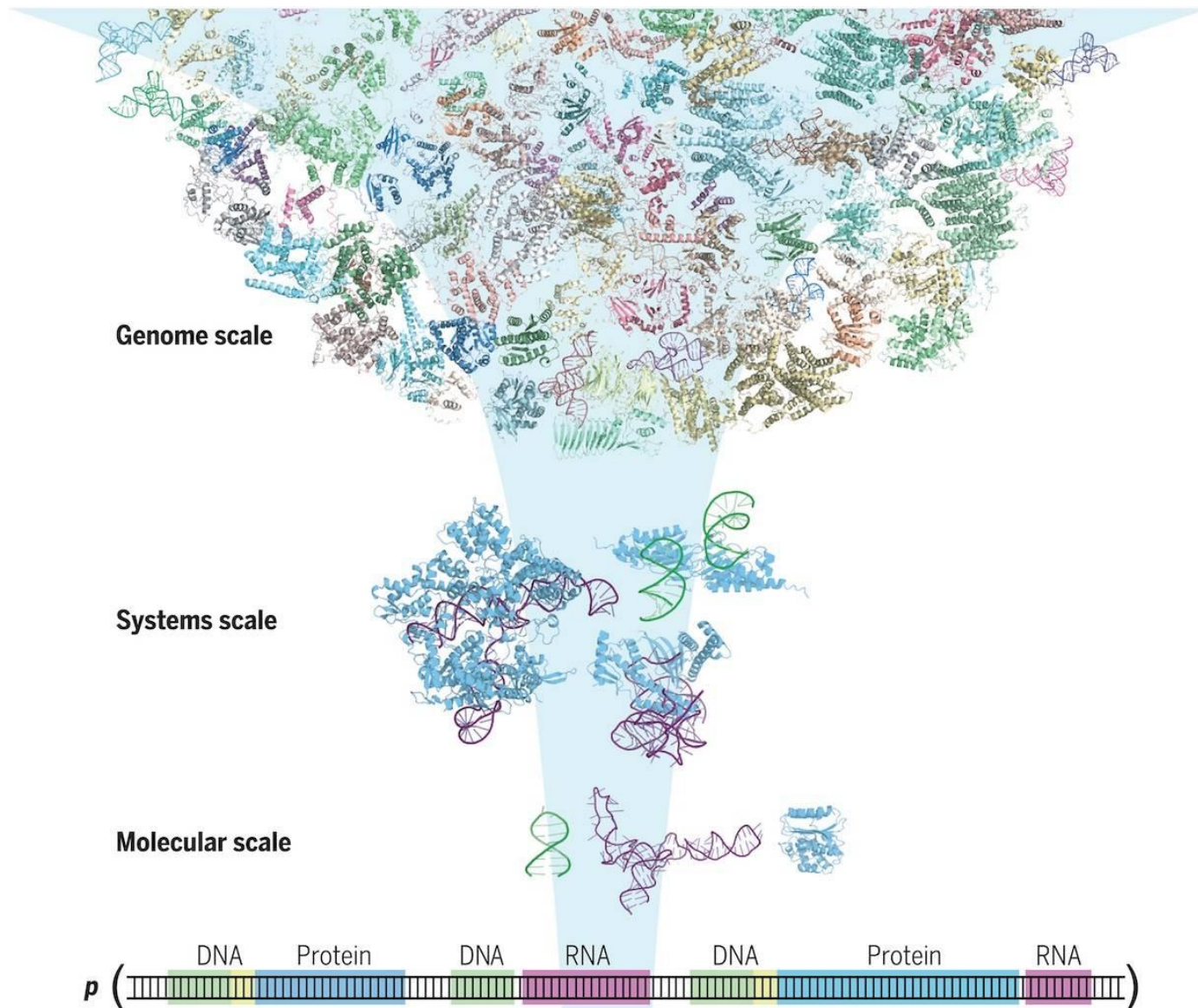
It's hitting (and exceeding) human levels



There are many generative AI medical companies



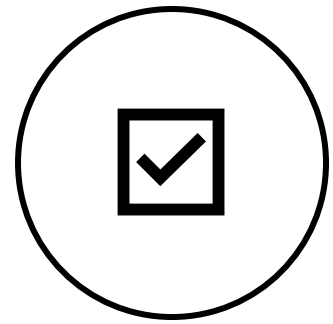
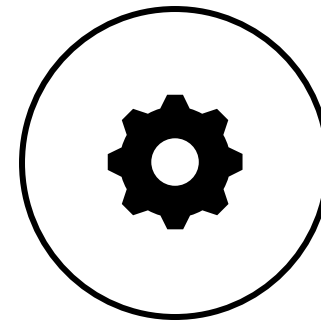
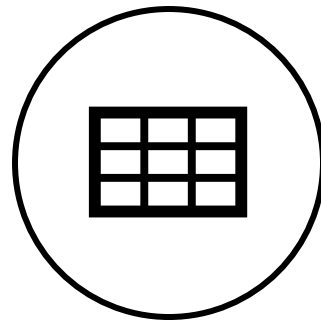
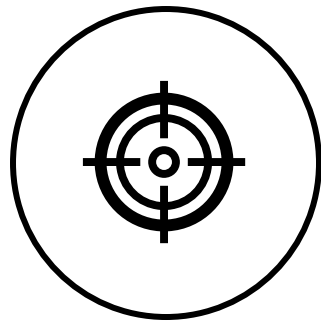
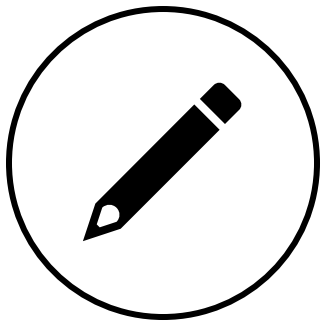
AlphaFold & Evo create protein/DNA sequences



AlphaFold predicts the structure of proteins based on provided sequences

Evo is a generalized model across DNA, RNA, and Proteins

**A (recommended) AI
workflow**



Define the Task

- What (precise) problem are you solving?
- Collect the data you want to solve that problem

Select the Model

- What do you want the model to do?

Clean the Data

- Clean the data
- Split into training and validation sets

Train the Model

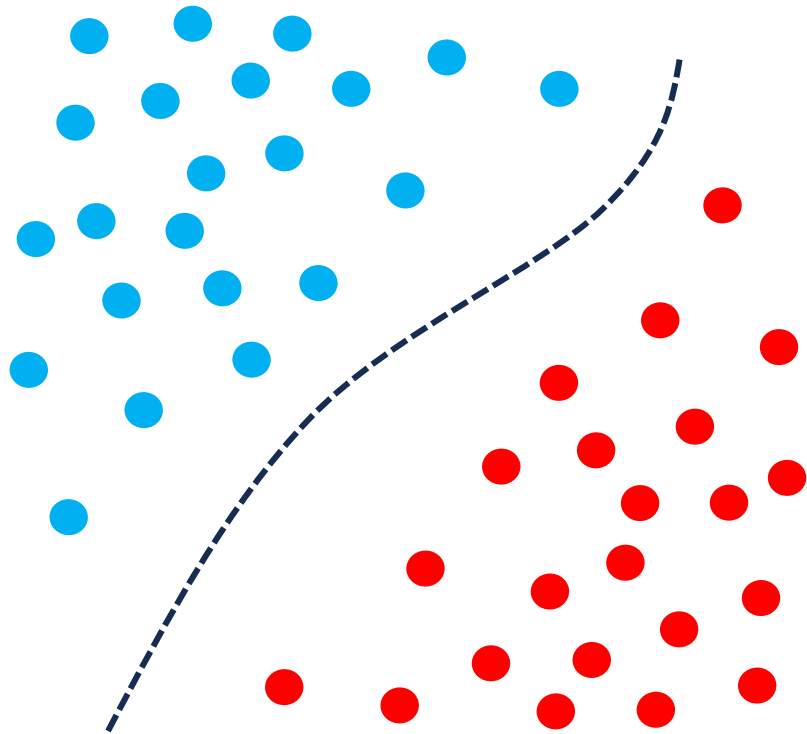
- Choose the model
- Optimize parameters

Evaluate the Model

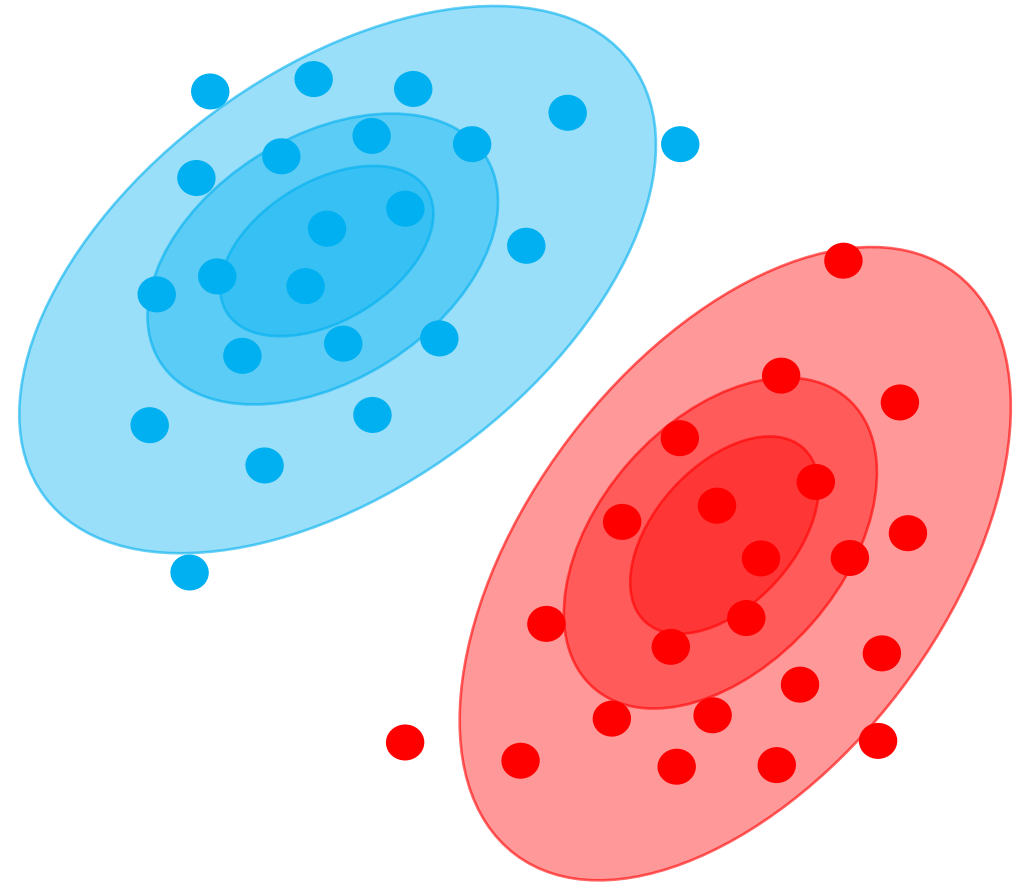
- Choose evaluation metrics
- Assess performance with validation

Define the problem

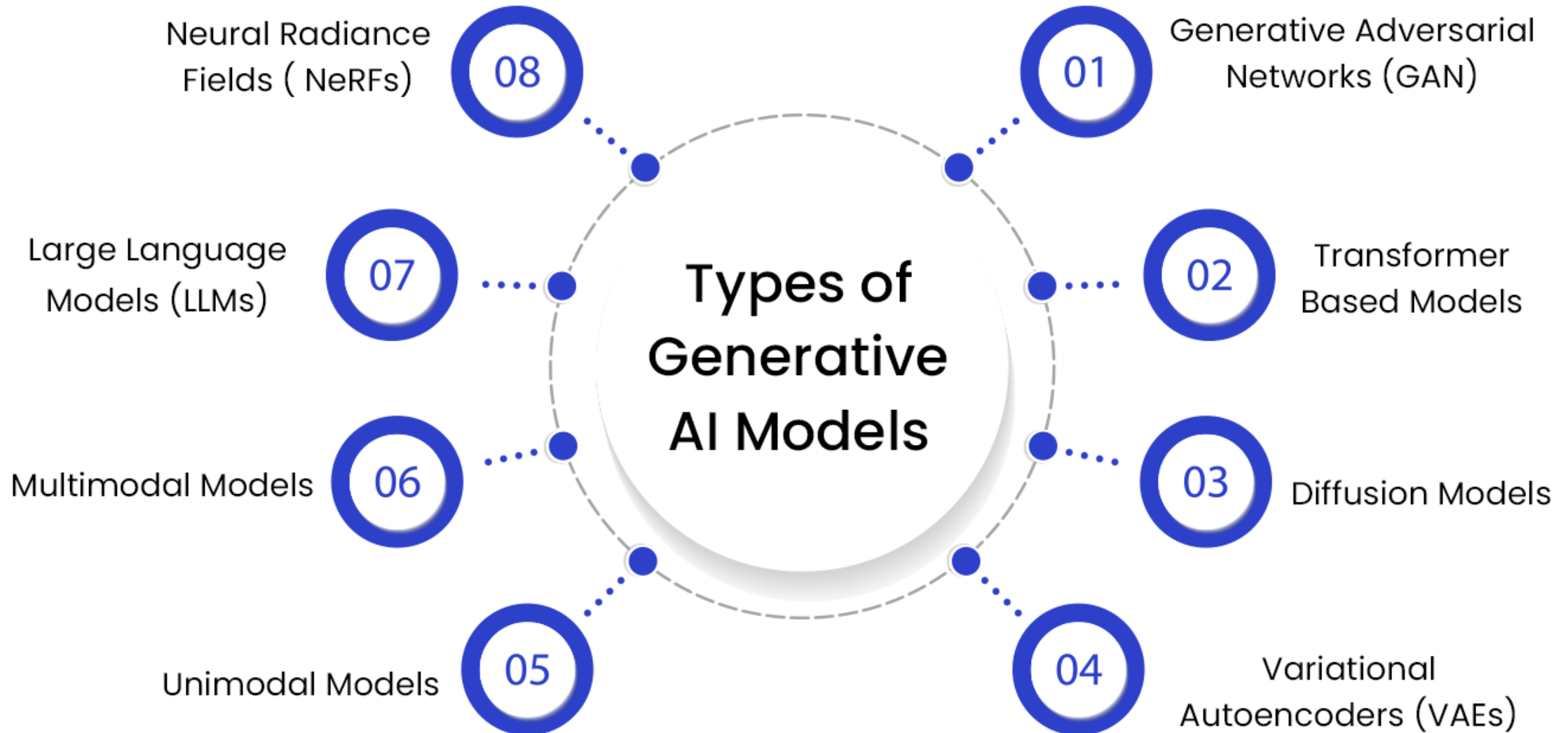
Classical AI Algorithms try to find the optimal boundary to separate data from multiple groups



Generative AI Algorithms try to learn the underlying data generating distribution



Select the model



Clean the data

Up-to-date data



Data should be up-to-date in order to obtain maximum value from the data analysis.



Missing values



Count missing values and analyze where in the data they are missing. Missing values can disrupt some analyses and skew the results.



Duplicates



Duplicate IDs indicate multiple records for one person, e.g. someone holds multiple functions at the same time.



Numerical outliers



Numerical outliers are fairly easy to detect and remove. Define minimum and maximum to spot outliers easily.



Check IDs



Check data labels of all the fields to see whether some categorical values are mislabeled.



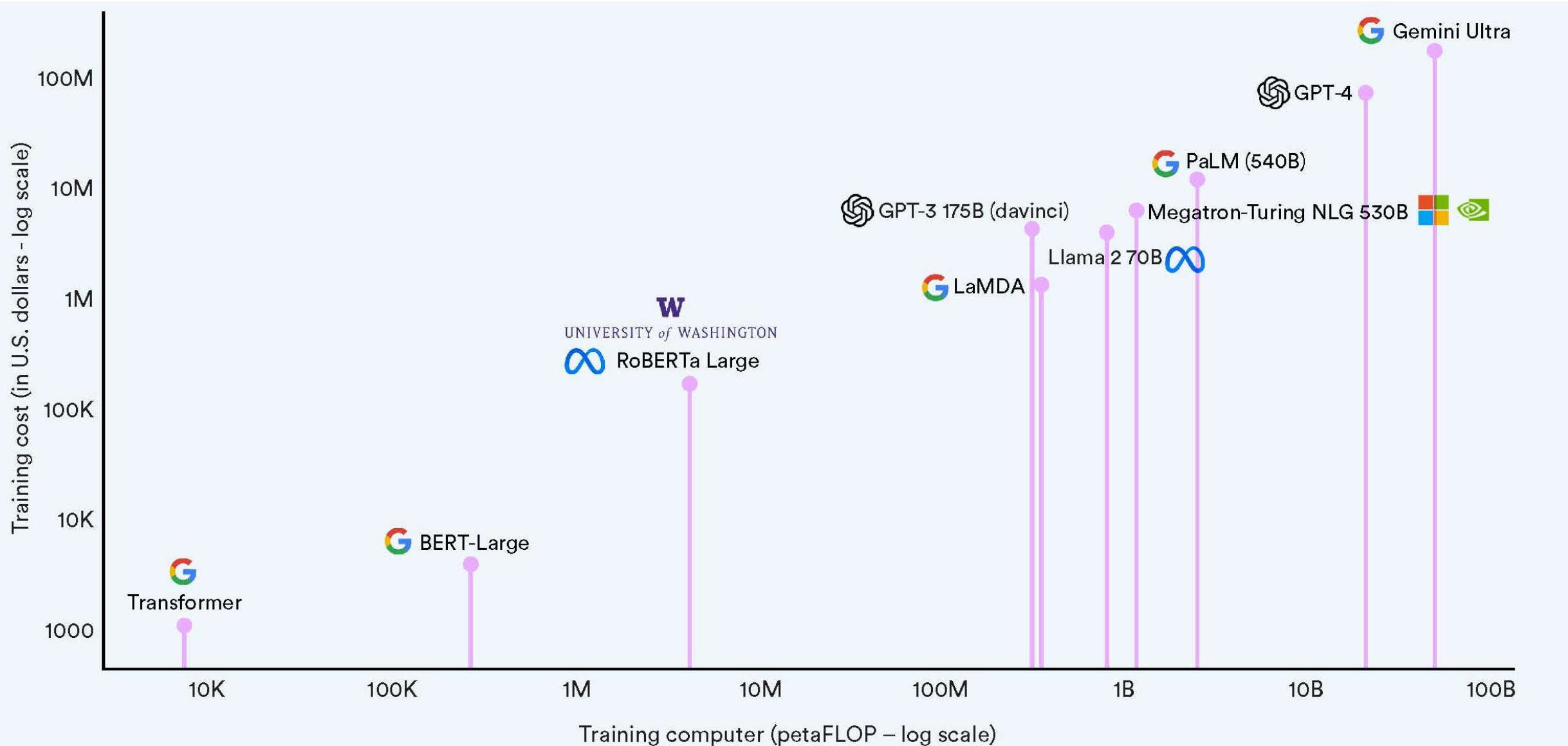
Define valid output



Define valid data labels for categorical data. Define data ranges for numerical variables. Non-matching data is presumably wrong.



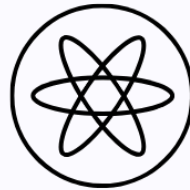
Train the model...



...or just use a pre-trained one

Prompt

Write a positive review for a blender that makes great margaritas and smoothies and is easy to clean.



**Generative
pre-trained
model**



Response

This blender makes amazing margaritas and smoothies, efficiently pulverizing fruits and ice. Its straightforward controls make operation a breeze, perfect for any blending task. Even cleaning up is hassle-free, making it a valuable addition to my kitchen. Reliable and easy to use, this blender is a worthy recommendation for any beverage enthusiast.

...or just use a pre-trained one

This blender makes amazing

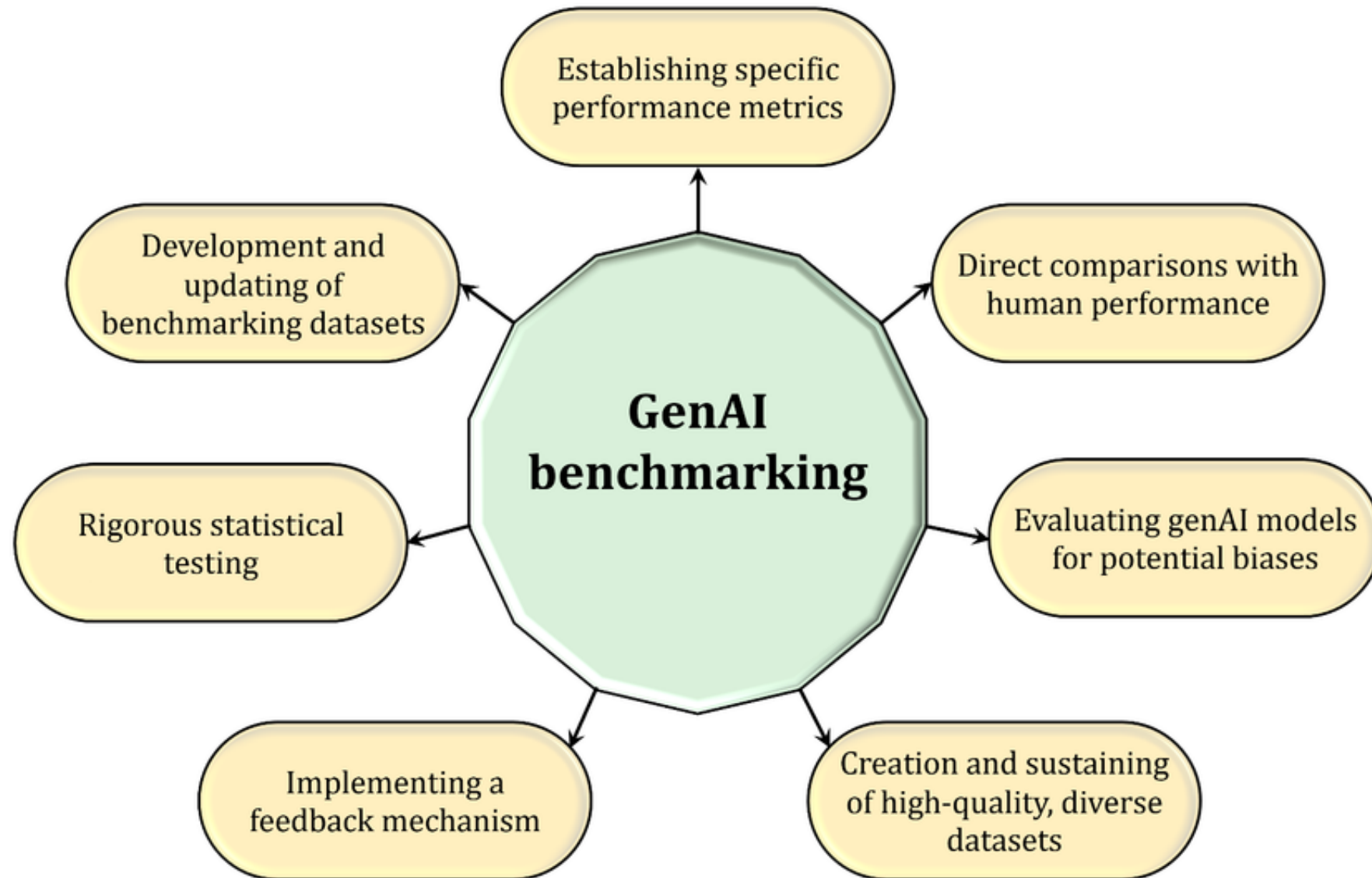
smoothies

shakes

...

bicycles

Evaluate the model



Customizing performance

Strategies to customize performance

- **Prompt engineering:** Design natural language instructions to guide model responses

Strategies to customize performance

- **Prompt engineering:** Design natural language instructions to guide model responses
- **Retrieval Augmented Generation (RAG):** Provide vector store for quick retrieval of relevant information for task

Strategies to customize performance

- **Prompt engineering:** Design natural language instructions to guide model responses
- **Retrieval Augmented Generation (RAG):** Provide vector store for quick retrieval of relevant information for task
- **Fine-tuning:** Augmenting network architecture to improve performance on a small dataset that is specific for the task

Strategies to customize performance

- **Prompt engineering:** Design natural language instructions to guide model responses
- **Retrieval Augmented Generation (RAG):** Provide vector store for quick retrieval of relevant information for task
- **Fine-tuning:** Augmenting network architecture to improve performance on a small dataset that is specific for the task
- **Reinforcement Learning through Human Feedback (RLHF):** Training a model by receiving feedback from humans

Prompt engineering

- **When?** Any-time; best suited for tasks that require a high level of precision and output
- **How?** Systematically construct and iterate instructions to elicit desired response content, structure, and tone

Prompt engineering (example)

A 45-year-old male with a blunt force injury to the head, a Glasgow Coma Scale (GCS) of 8, and loss of consciousness (LOC) for over 10 minutes. What trauma level facility should this patient be transferred to?

Given the description of the patient's injury—a blunt force trauma to the head, a Glasgow Coma Scale (GCS) of 8, and a loss of consciousness (LOC) for over 10 minutes—this patient is likely suffering from a significant head injury and requires immediate and specialized care.

A **GCS of 8** indicates a **moderate to severe traumatic brain injury (TBI)**, and a **LOC for over 10 minutes** suggests a potential for significant neurological compromise. The patient's condition would need close monitoring and likely advanced imaging (such as a CT scan) and surgical intervention if there is any progression or deterioration in neurological status.

For trauma of this severity, the patient should be transferred to a **Level I trauma center**. A Level I trauma center is equipped to handle the most severe trauma cases, offering:

1. **24-hour availability of neurosurgeons** and other specialists (such as orthopedic, trauma, and anesthesiology).
2. **Comprehensive imaging capabilities** (CT, MRI) for rapid diagnosis.
3. **Advanced critical care services** in a dedicated intensive care unit (ICU).
4. **Multidisciplinary teams** for ongoing assessment and management of complex trauma patients.

In summary, the patient should be transferred to a **Level I trauma center** for optimal care.

Prompt engineering (example)

You are designed to help trauma staff make the best decisions for patients at the point of care. You must reference ATLS protocol where applicable. Keep your responses concise and helpful.

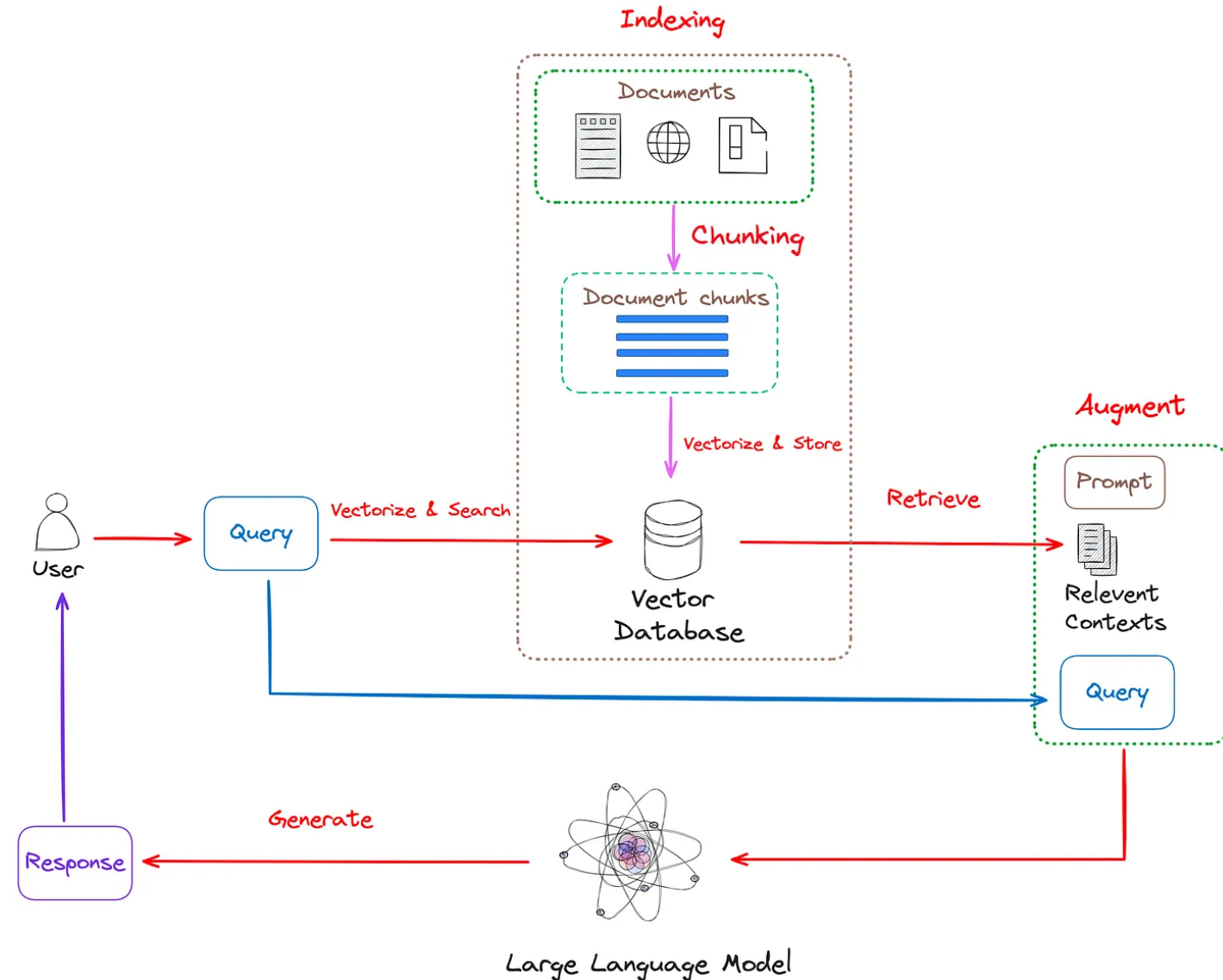
Prompt engineering (example)

A 45-year-old male with a blunt force injury to the head, a Glasgow Coma Scale (GCS) of 8, and loss of consciousness (LOC) for over 10 minutes. What trauma level facility should this patient be transferred to?

Based on the information provided, the patient has a GCS of 8, which indicates a severe head injury. According to the Advanced Trauma Life Support (ATLS) guidelines, patients with a GCS of 8 or less should be transferred to a Level I trauma center for comprehensive evaluation and management. This level of care is equipped to handle severe traumatic brain injuries and provide the necessary surgical interventions and critical care.

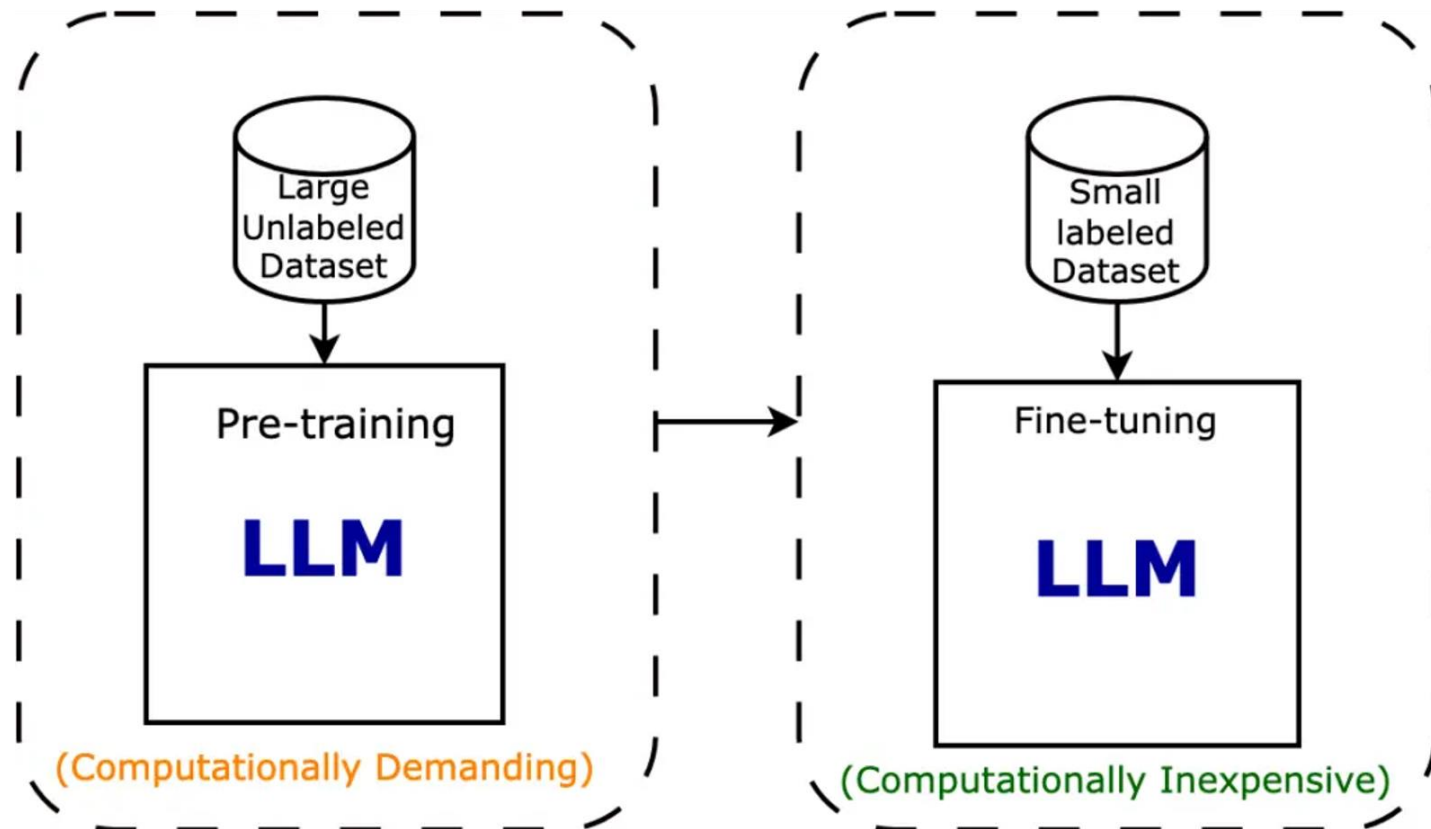
Retrieval Augmented Generation (RAG)

- **When?** Specific context (from files) needs to be given for the desired task
- **How?** Chunk data and vectorize chunks in a database for quick retrieval



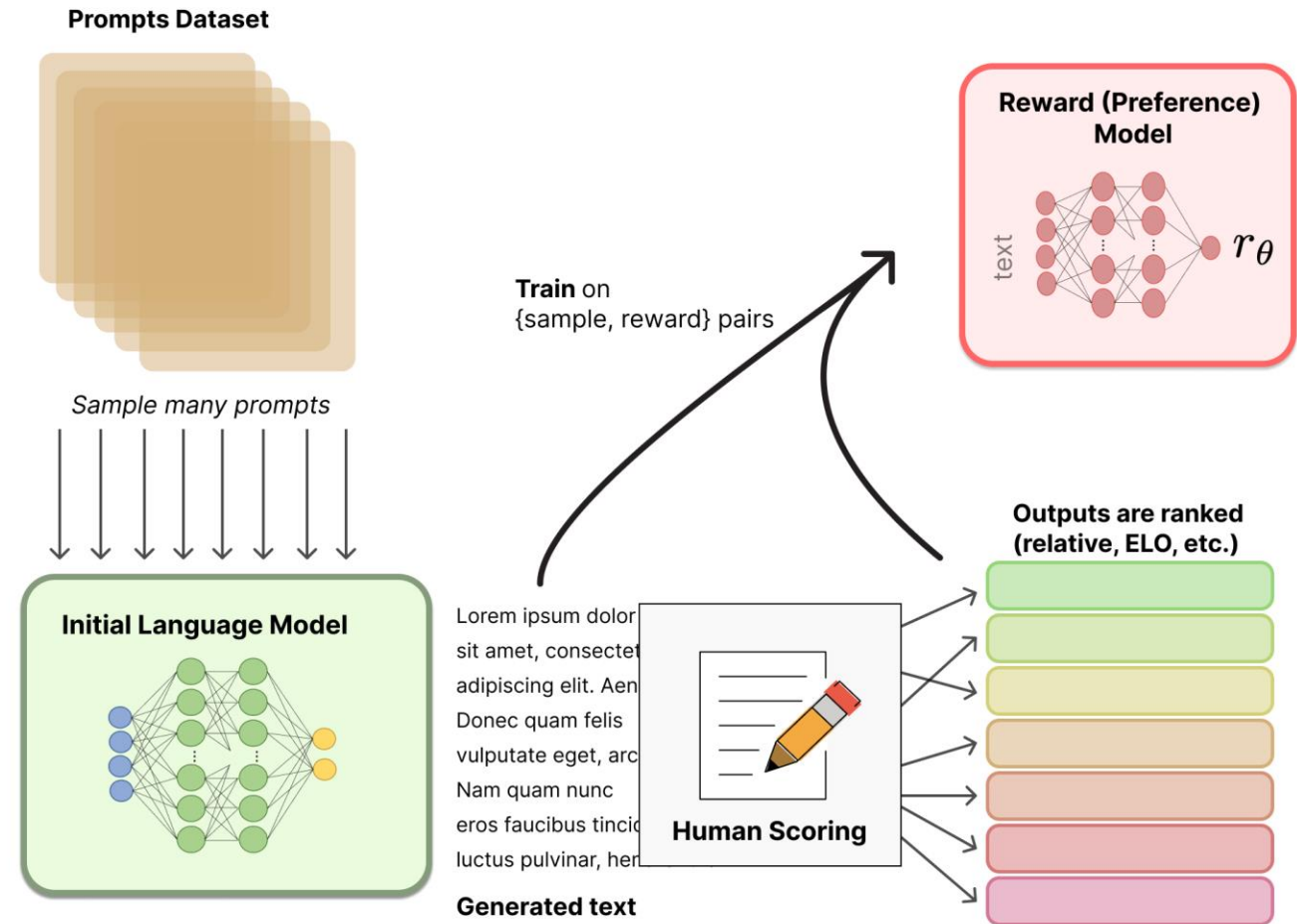
Fine-tuning

- **When?** Large labeled dataset is available and the task is extremely focused
- **How?** Pre-trained model is adjusted by changing weights until desired accuracy is achieved



Reinforcement Learning through Human Feedback (RLHF)

- **When?** High-level of accuracy needed for task and inputs a highly variable
- **How?** Human feedback updates weights of pre-trained model



Practical steps to get started

3 Ways to Start Using Generative AI

1. No Code or Playground

- Pre-trained models or no-code playgrounds through web interfaces
- ChatGPT web interface, Teachable Machine, etc.

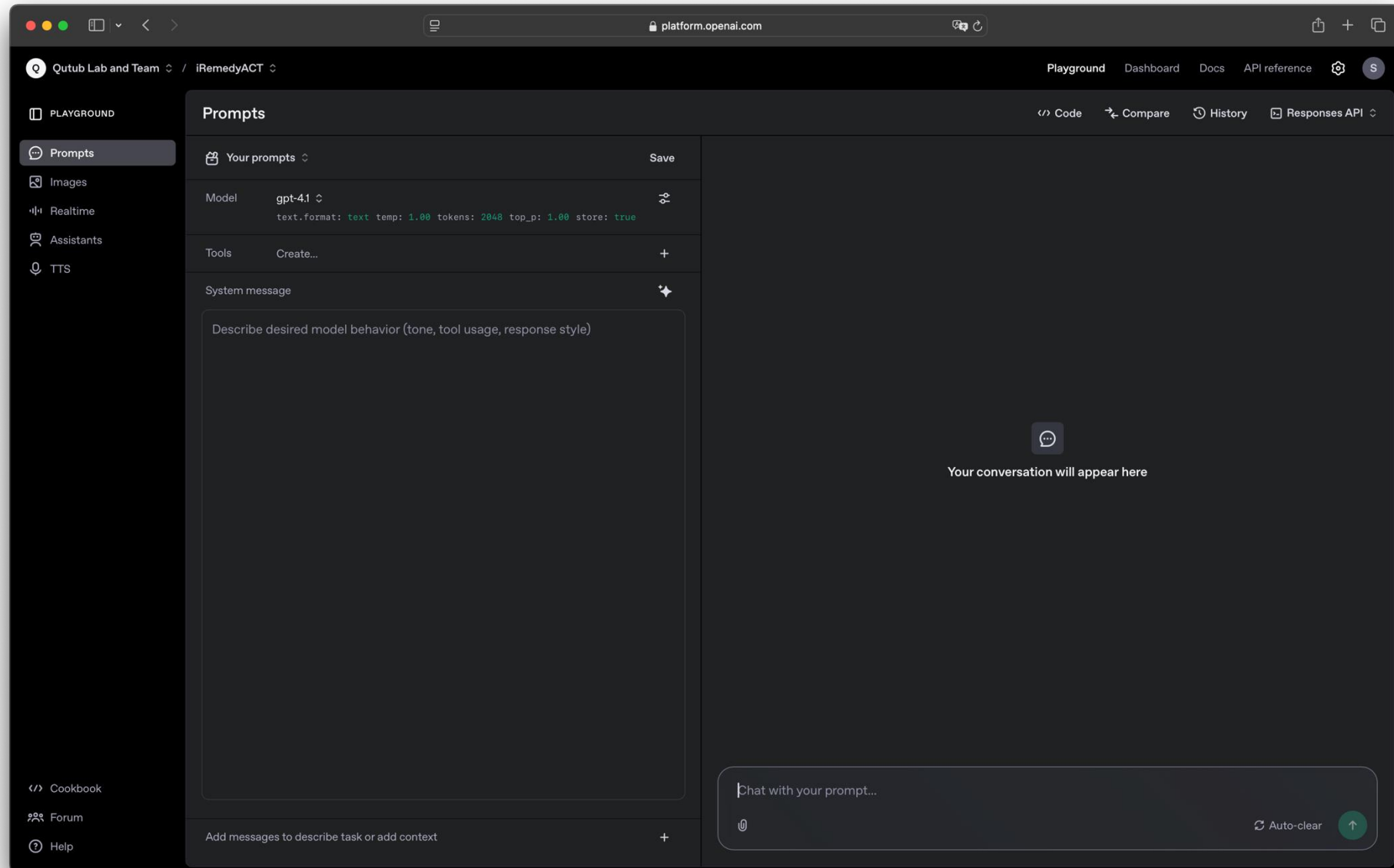
2. Prototypes or Demo Applications

- Python environments using APIs (OpenAI, Gemini, etc.) or open-source models
- Jupyter or Google Colab Notebooks with corresponding Python packages for each model

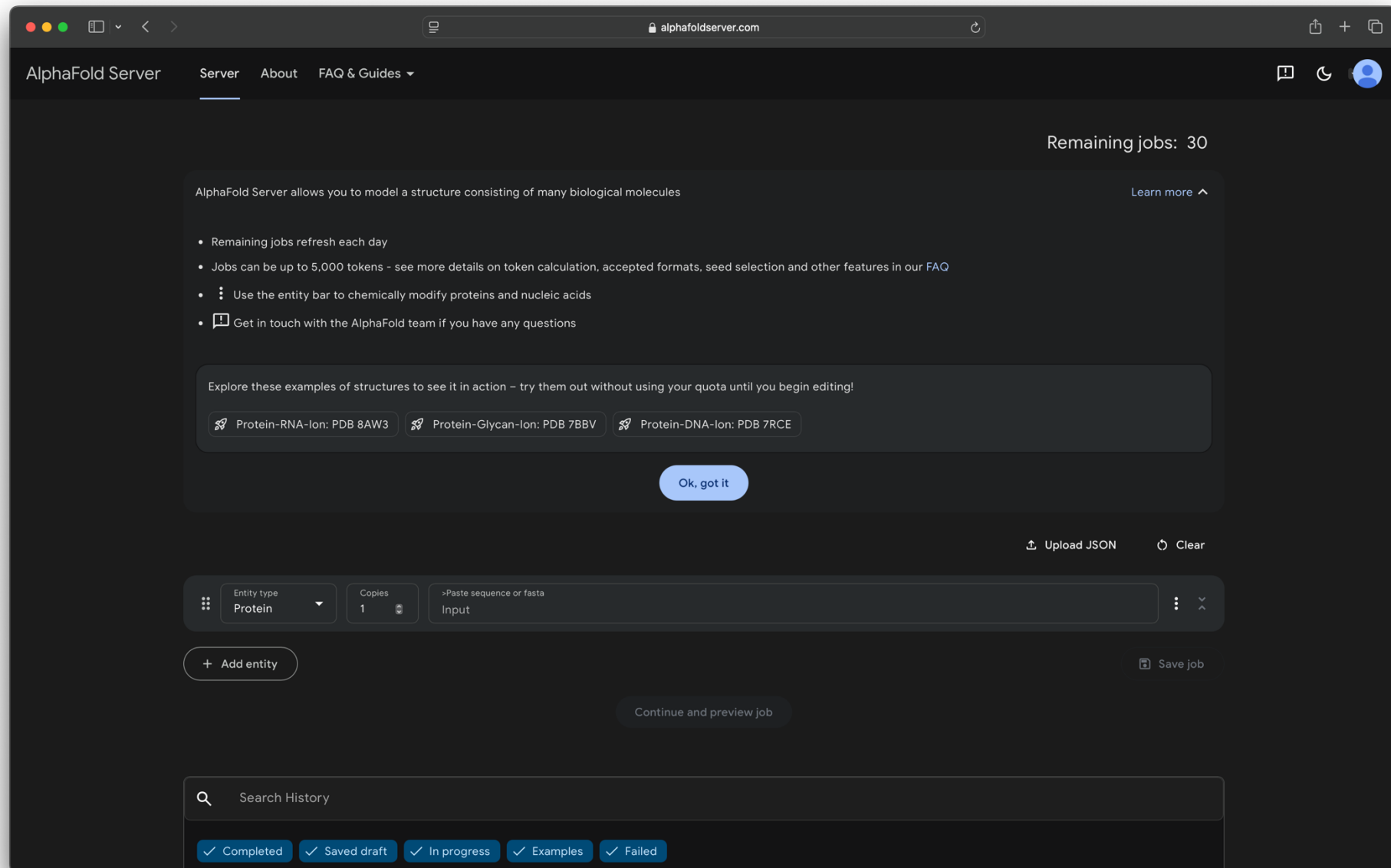
3. Production Level Applications

- Build software pipelines with CI, testing, authentication, etc.
- LangChain, Vercel AI SDK, Auth.js, FastAPI, etc.

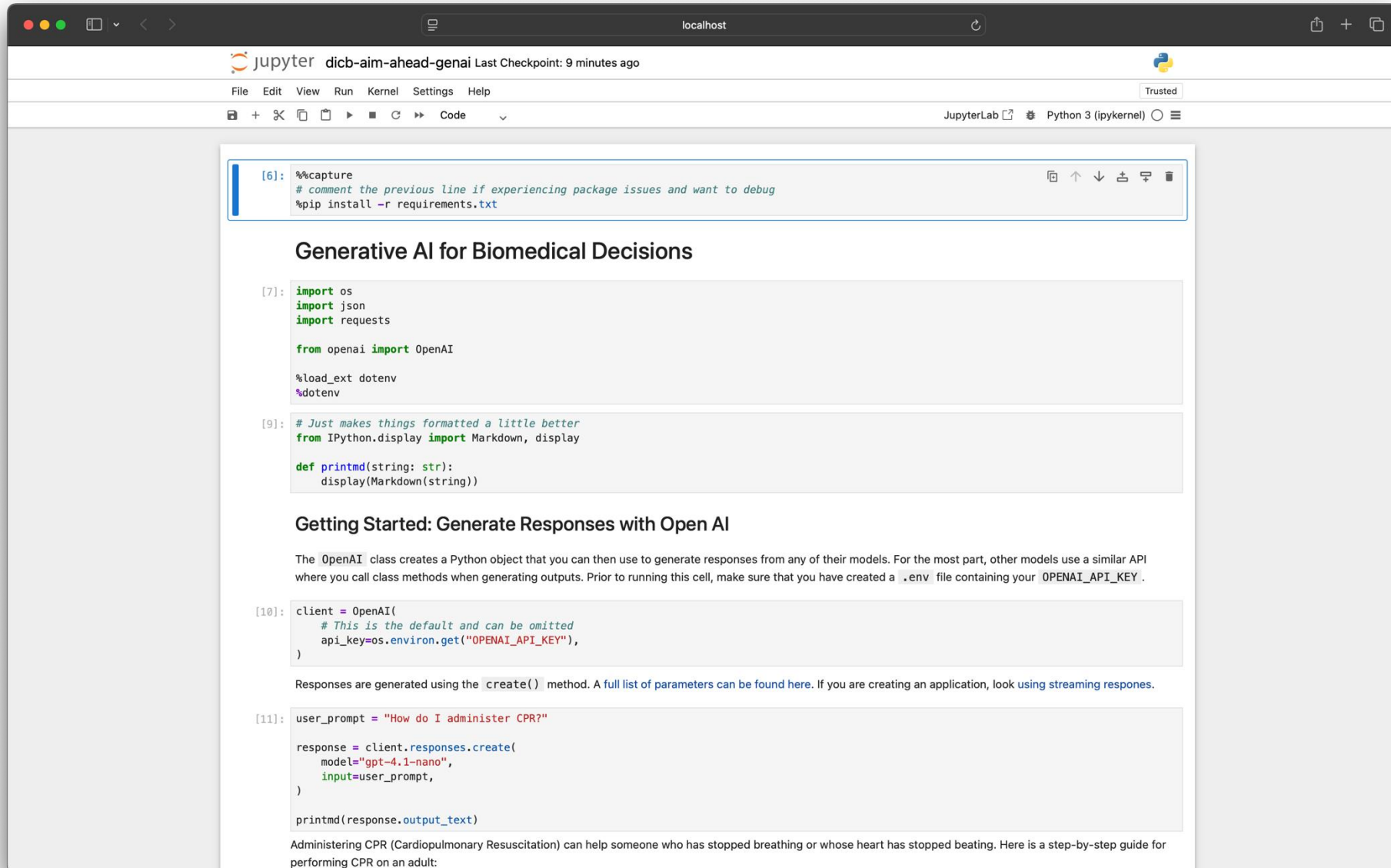
ChatGPT Playground



AlphaFold Server



Python LLM API Example



The screenshot shows a JupyterLab window titled "dicb-aim-ahead-genai" with a "Last Checkpoint: 9 minutes ago" status. The interface includes a menu bar (File, Edit, View, Run, Kernel, Settings, Help) and a toolbar with icons for file operations and execution. The main area displays a Python notebook with the following content:

```
[6]: %capture
# comment the previous line if experiencing package issues and want to debug
%pip install -r requirements.txt
```

Generative AI for Biomedical Decisions

```
[7]: import os
import json
import requests

from openai import OpenAI

%load_ext dotenv
%dotenv

[9]: # Just makes things formatted a little better
from IPython.display import Markdown, display

def printmd(string: str):
    display(Markdown(string))
```

Getting Started: Generate Responses with Open AI

The `OpenAI` class creates a Python object that you can then use to generate responses from any of their models. For the most part, other models use a similar API where you call class methods when generating outputs. Prior to running this cell, make sure that you have created a `.env` file containing your `OPENAI_API_KEY`.

```
[10]: client = OpenAI(
# This is the default and can be omitted
api_key=os.environ.get("OPENAI_API_KEY"),
)
```

Responses are generated using the `create()` method. A [full list of parameters can be found here](#). If you are creating an application, look using [streaming responses](#).

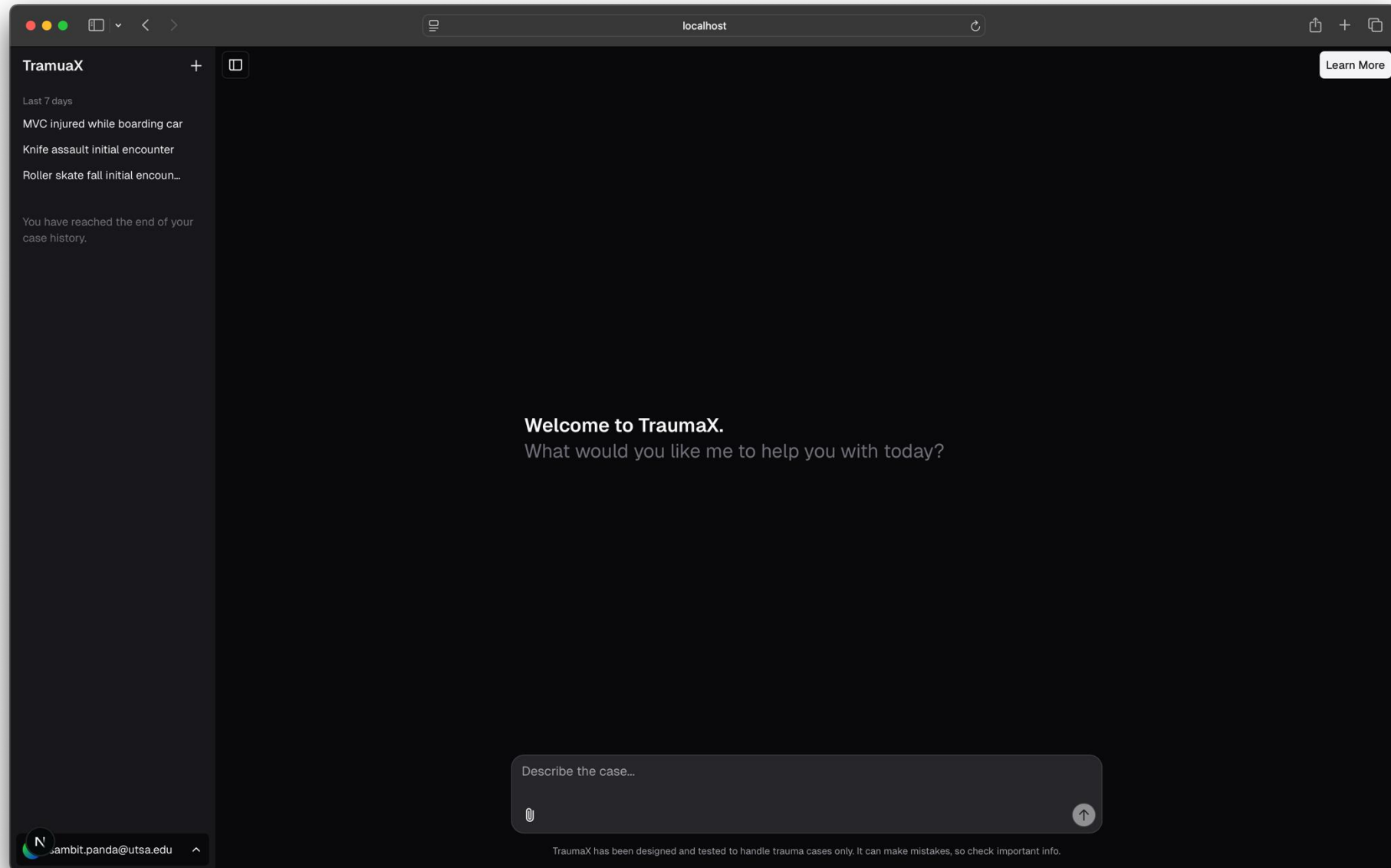
```
[11]: user_prompt = "How do I administer CPR?"

response = client.responses.create(
    model="gpt-4.1-nano",
    input=user_prompt,
)

printmd(response.output_text)
```

Administering CPR (Cardiopulmonary Resuscitation) can help someone who has stopped breathing or whose heart has stopped beating. Here is a step-by-step guide for performing CPR on an adult:

TraumaX Example



Pitfalls and considerations

Common Pitfalls

- **Using generative models when you just need classification:** Generative model adds complexity (cost) and risk without improving outcome

Common Pitfalls

- **Using generative models when you just need classification:** Generative model adds complexity (cost) and risk without improving outcome
- **Replacing real-world validation with synthetic data:** Synthetic examples cannot fully replace real clinical/experimental data

Common Pitfalls

- **Using generative models when you just need classification:** Generative model adds complexity (cost) and risk without improving outcome
- **Replacing real-world validation with synthetic data:** Synthetic examples cannot fully replace real clinical/experimental data
- **Assuming generated outputs are reliable (copy/pasting a source from ChatGPT):** Models can "hallucinate", producing convincing but incorrect outputs

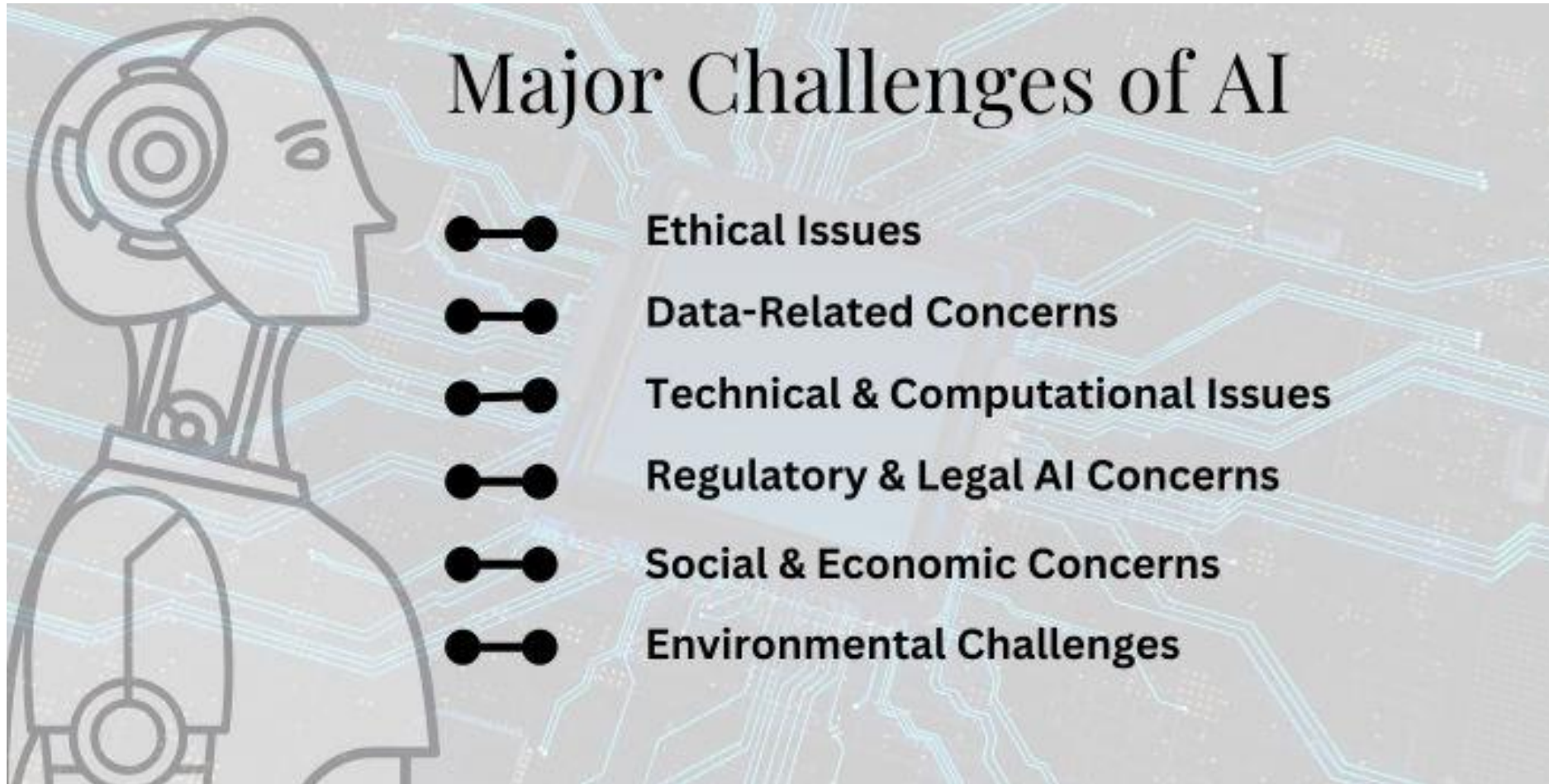
Common Pitfalls

- **Using generative models when you just need classification:** Generative model adds complexity (cost) and risk without improving outcome
- **Replacing real-world validation with synthetic data:** Synthetic examples cannot fully replace real clinical/experimental data
- **Assuming generated outputs are reliable (copy/pasting a source from ChatGPT):** Models can "hallucinate", producing convincing but incorrect outputs
- **Using generative models to clean data instead of improving data collection:** Better to fix experimental design; Imputation methods

Common Pitfalls

- **Using generative models when you just need classification:** Generative model adds complexity (cost) and risk without improving outcome
- **Replacing real-world validation with synthetic data:** Synthetic examples cannot fully replace real clinical/experimental data
- **Assuming generated outputs are reliable (copy/pasting a source from ChatGPT):** Models can "hallucinate", producing convincing but incorrect outputs
- **Using generative models to clean data instead of improving data collection:** Better to fix experimental design; Imputation methods
- **Applying models outside their training domain:** Model will hallucinate (make stuff up) (e.g., asking ChatGPT to draft a niche protocol for molecular biology)

Things to keep in mind



Conclusion

- Introduced generative AI and the landscape of tools
- Provided a recommended workflow for doing any biomedical discovery using AI
- Learned how to customize model performance using prompt engineering, fine-tuning, RAG, etc.
- Gave practical code examples on how to begin using these models
- Discussed pitfalls and limitations of these algorithms



Questions?

sambit.panda@utsa.edu; sampan.me

christian.cruz@utsa.edu