

# The Data Flywheel for FMware

Gopi Krishnan Rajbahadur



# How to cite this session?

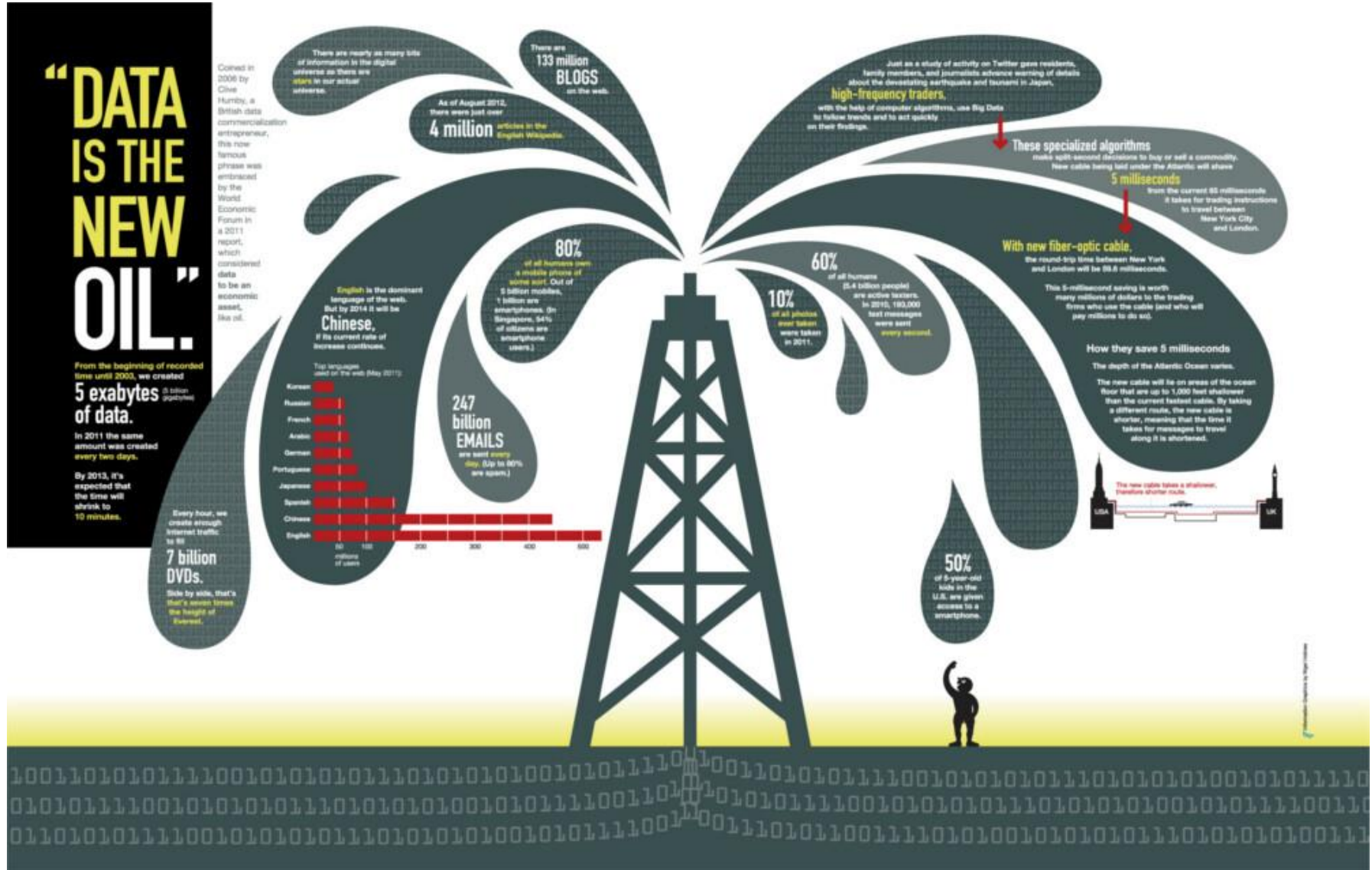
```
@misc{Rajbahadur2024dataflywheel,  
author = {Gopi Krishnan Rajbahadur},  
title = {The Data Flywheel for FMware},  
howpublished = {Tutorial presented at the AIware Leadership Bootcamp 2024},  
month = {November},  
year = {2024},  
address = {Toronto, Canada},  
note = {Part of the AIware Leadership Bootcamp series.},  
url = {https://aiwarebootcamp.io/slides/2024_aiwarebootcamp_rajbahadur_dataflywheelforfmware.pdf } }
```



# NETFLIX



# Data is the new oil





# Data is the new oil

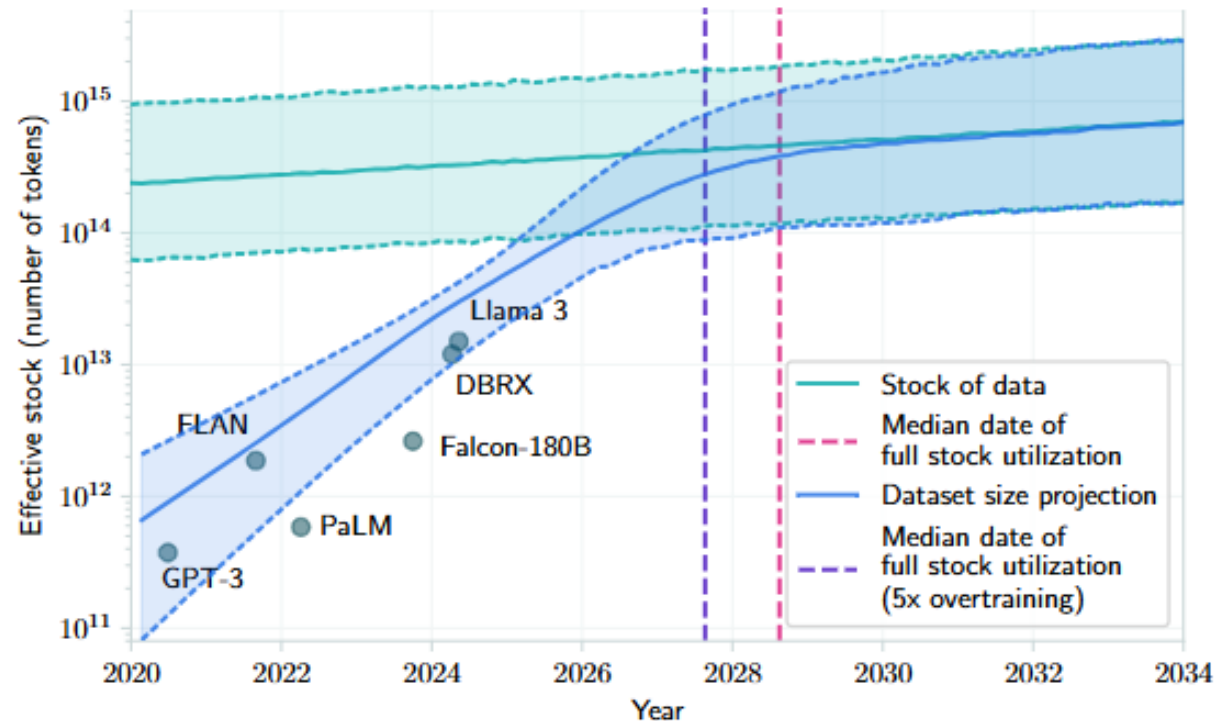




# What happens when we run out of data?

Will we run out of data? Limits of LLM scaling based on human-generated data

Pablo Villalobos<sup>1</sup> Anson Ho<sup>1</sup> Jaime Sevilla<sup>1,2</sup> Tamay Besiroglu<sup>1,3</sup> Lennart Heim<sup>1,4</sup> Marius Hobbhahn<sup>1,5</sup>



“ Our findings indicate that if current LLM development trends continue, models will be trained on datasets roughly equal in size to the available stock of public human text data **between 2026 and 2032, or slightly earlier if models are overtrained.** ”



# What happens when we run out of data?



**STORE  
CLOSING**

**EVERYTHING  
MUST GO!**





# We turn to renewable resources



But what does it mean in terms of data for  
FMware?





# Say Hello to our savior!



# Can anyone tell me what this is?





**Can anyone tell me what this is?**

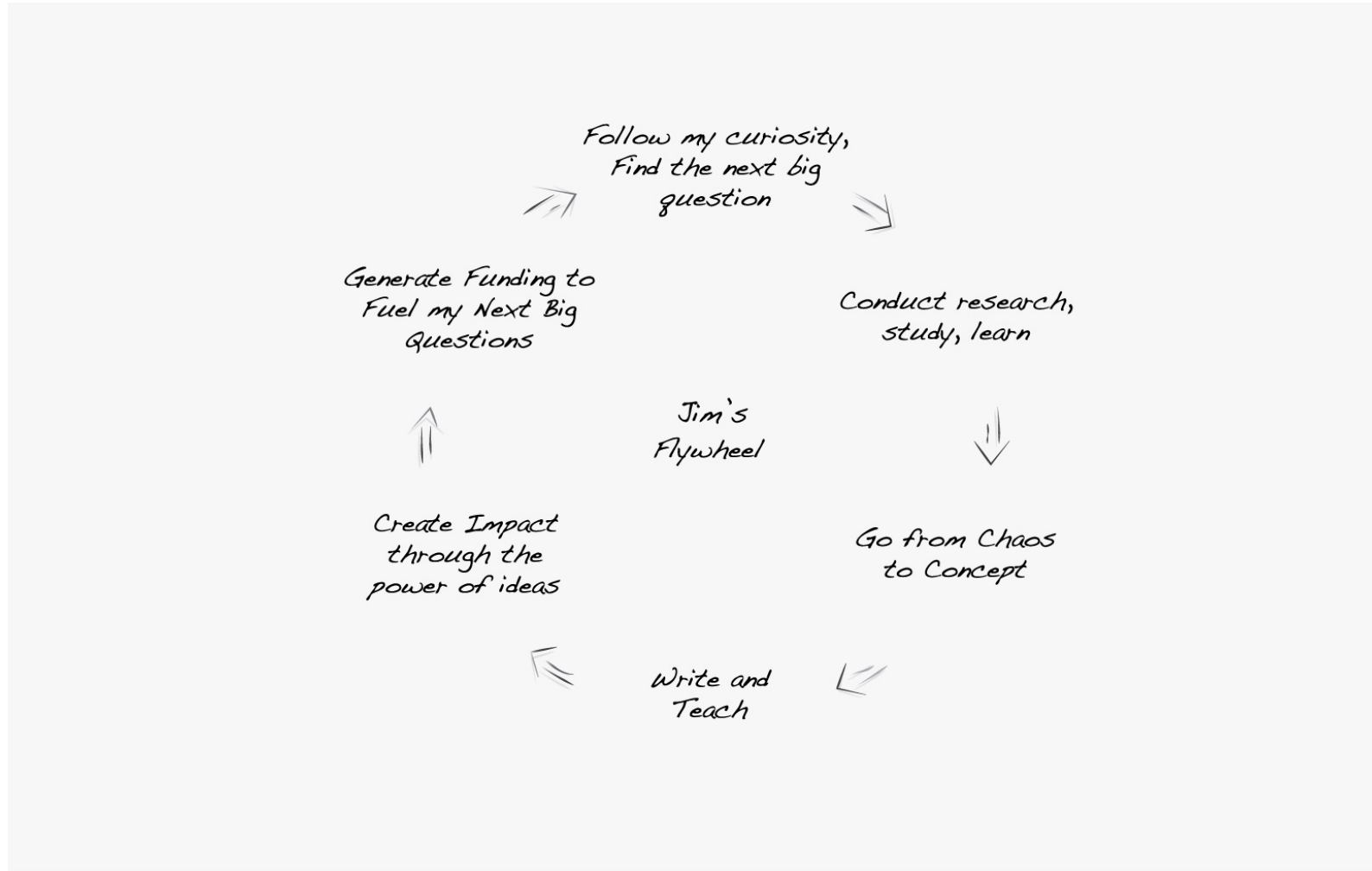


**Flywheel!**





# Jim Collin's Flywheel effect from the book "Good to Great"

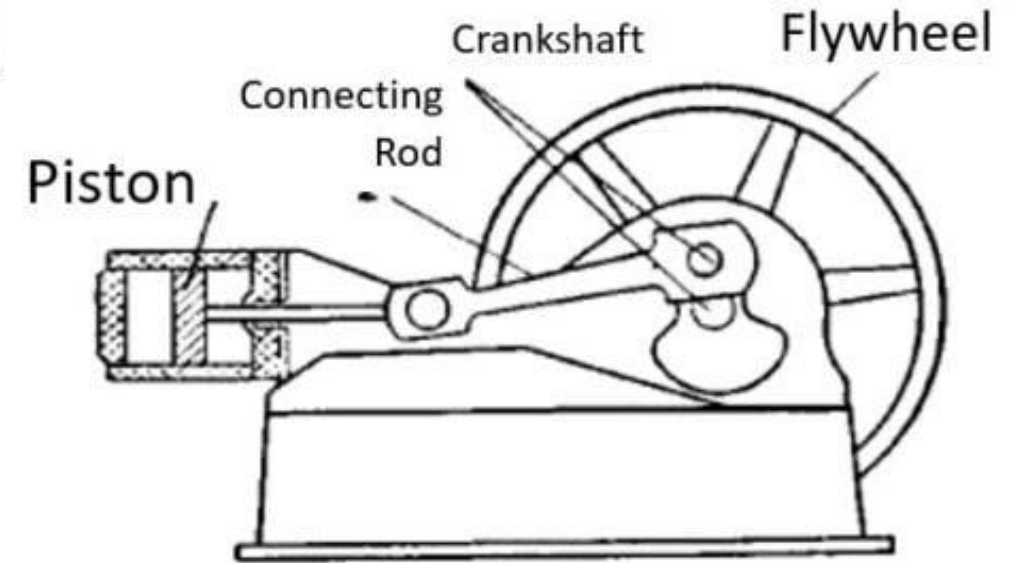
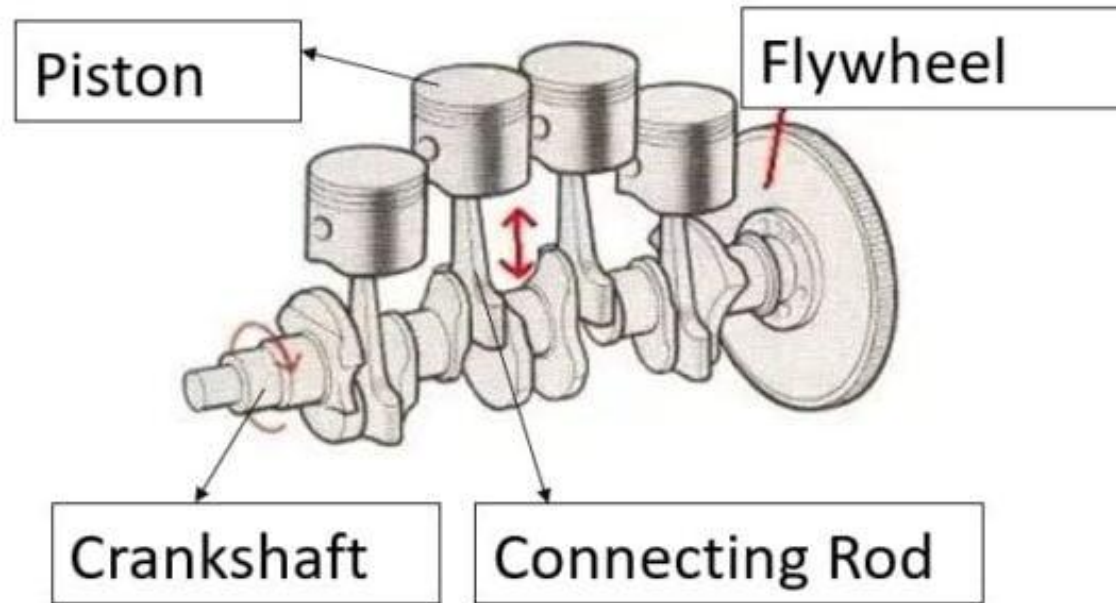


## The Data Flywheel: Building momentum by putting your data to work

Take a comprehensive approach to  
getting the most value from your data,  
one project at a time



# How does a Flywheel work?



An engine flywheel stores rotational energy, maintaining a steady speed by **smoothing out fluctuations from power strokes**. It conserves momentum, helping the engine run smoothly and **providing energy to keep it turning between cycles**.





# What is the first thing we need to make an engine work?



# What is the first thing we need to make an engine work?



# What is the first thing we need to make an engine work?



## Crude Oil

*['krüd 'oi(-ə)l]*

A raw natural resource that is extracted from the earth and used to propel vehicles, heat buildings, produce electricity, and make everyday products.



# Let's see if we can identify all the different types of data involved in the FMware Lifecycle



# Let's see if we can identify all the different types of data involved in the FMware Lifecycle

Pretraining data

Instruction data

Finetuning data

Preference data

Grounding data

Guarding data

Legal data

Examples data

Implicit and  
Explicit Feedback  
data

Performance data

Benchmark data

Log data

Telemetry data

Output data



# Evolution of synthetic data generation methods

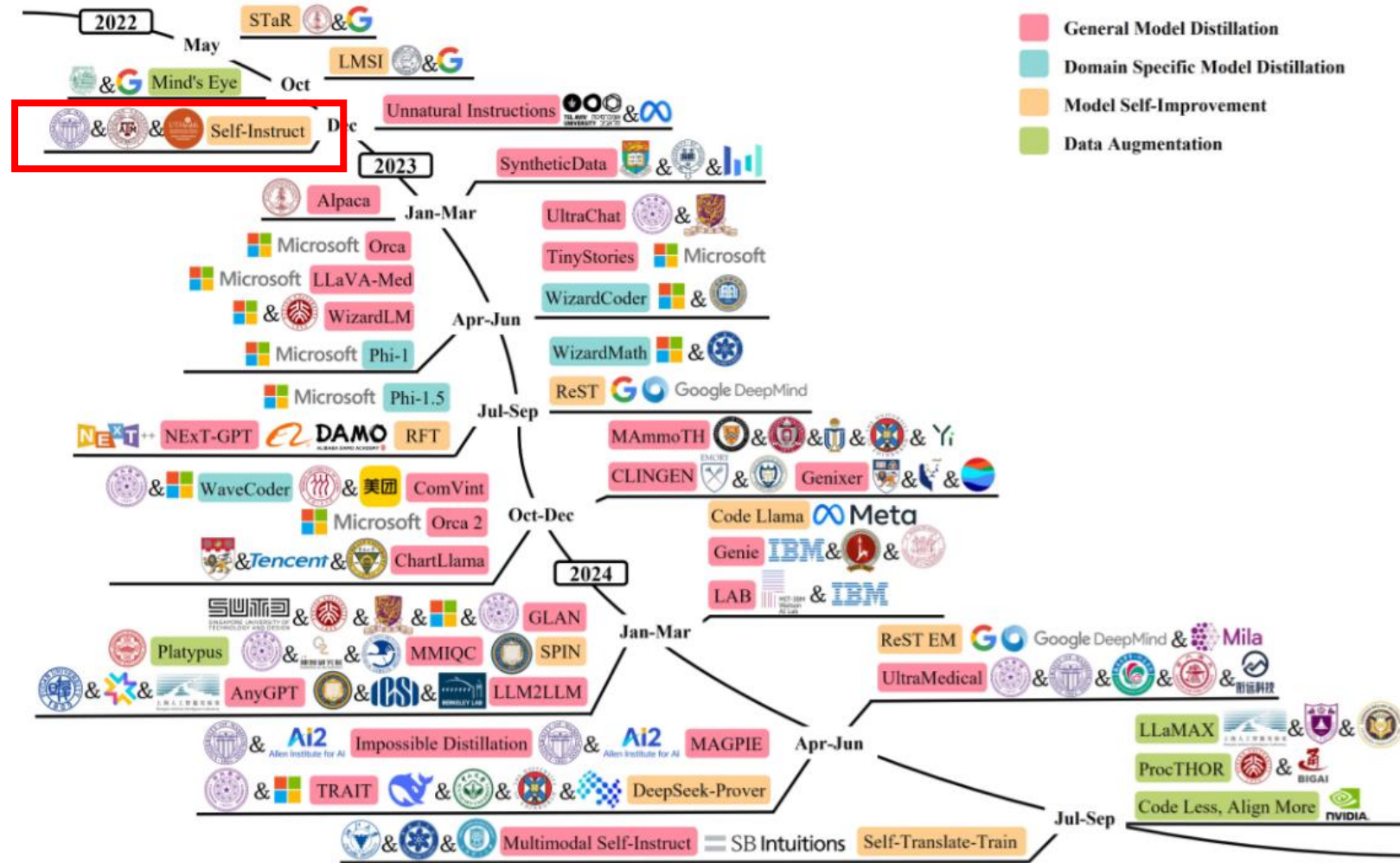


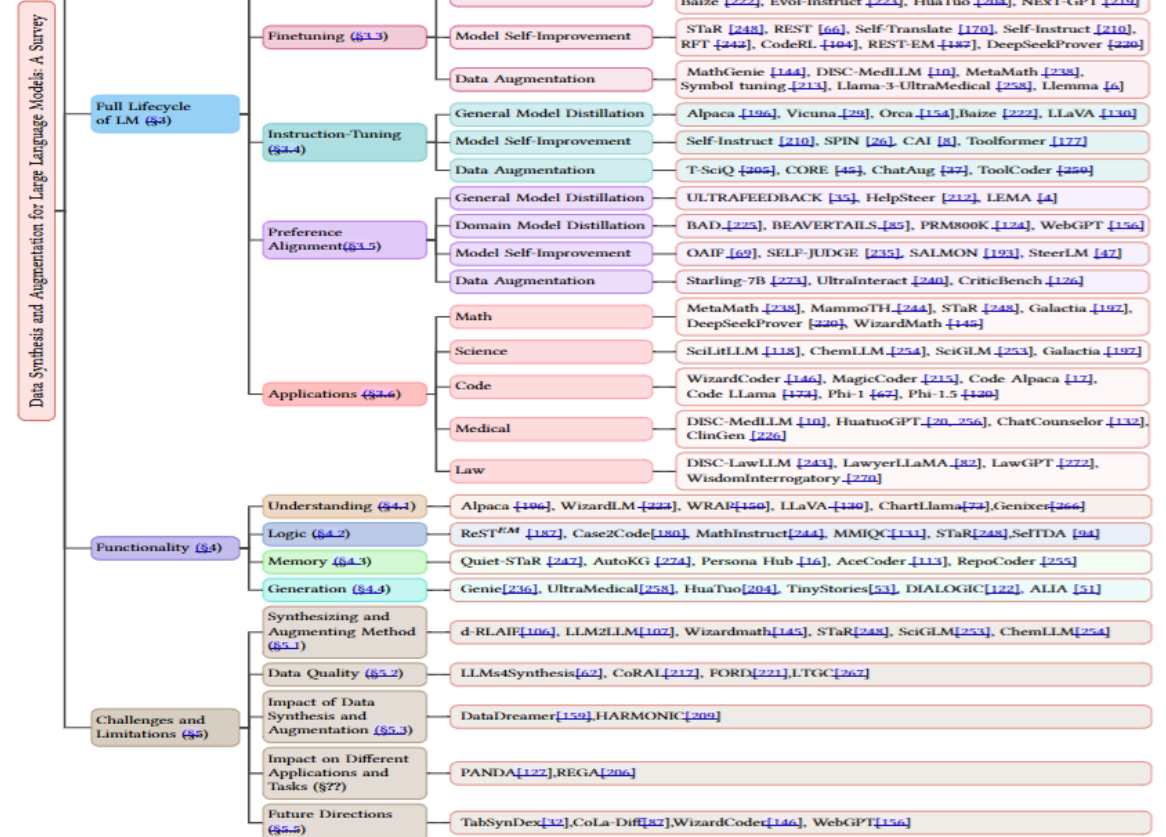
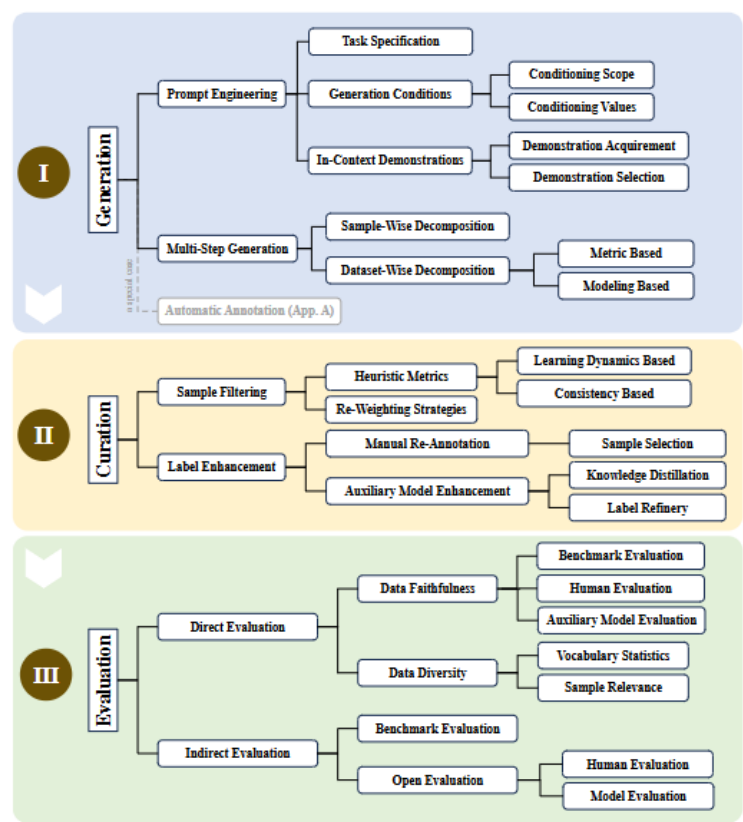
Figure 4: Illustration of the evolutionary steps in the development of data synthesis and augmentation techniques for large models.





# Much of this data can be generated synthetically

I can't cover all of these techniques in 35 mins!



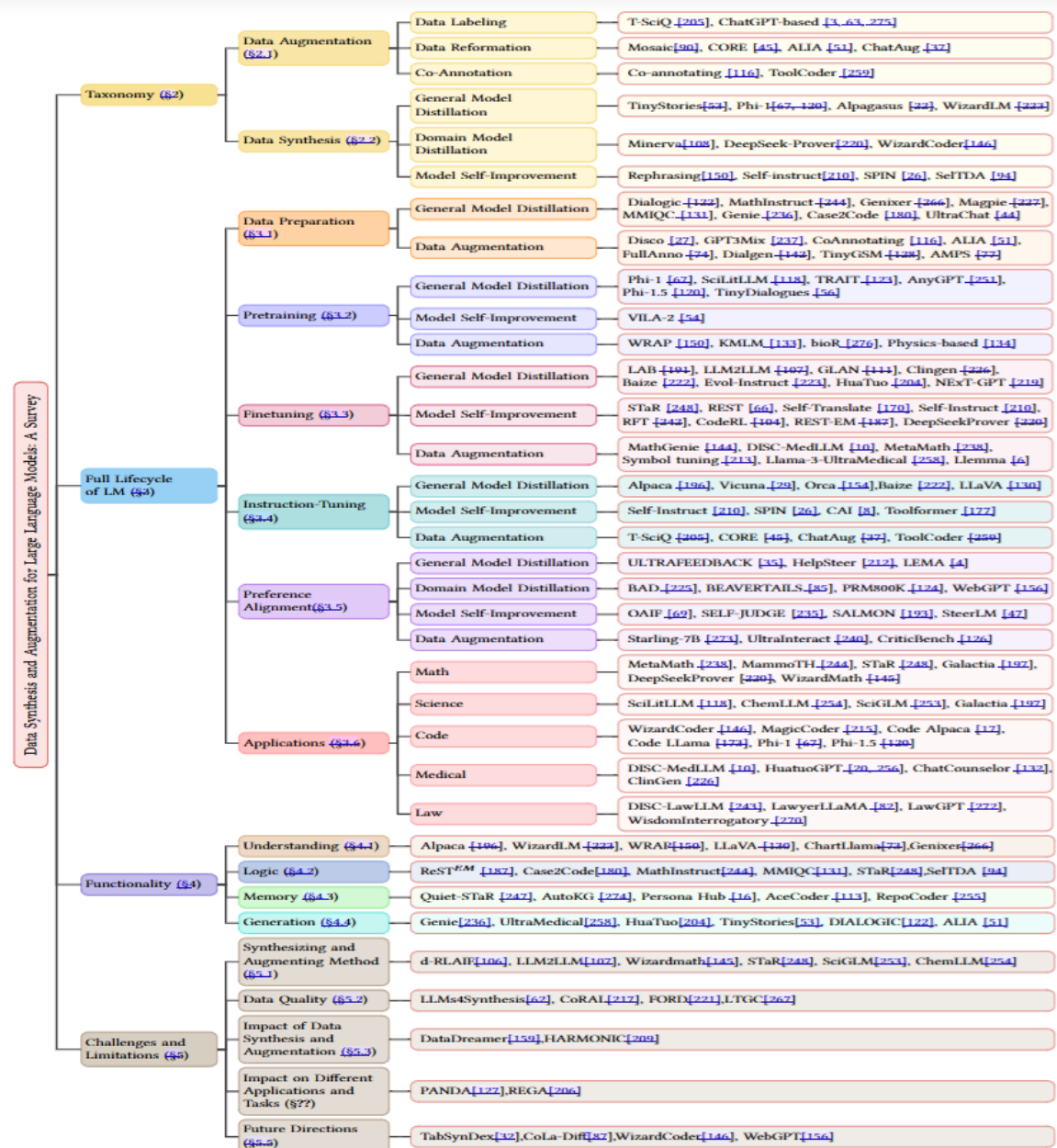
# Much of this data can be generated synthetically

I can't cover all of these techniques in 30 mins! So I am going to give you guys home work!

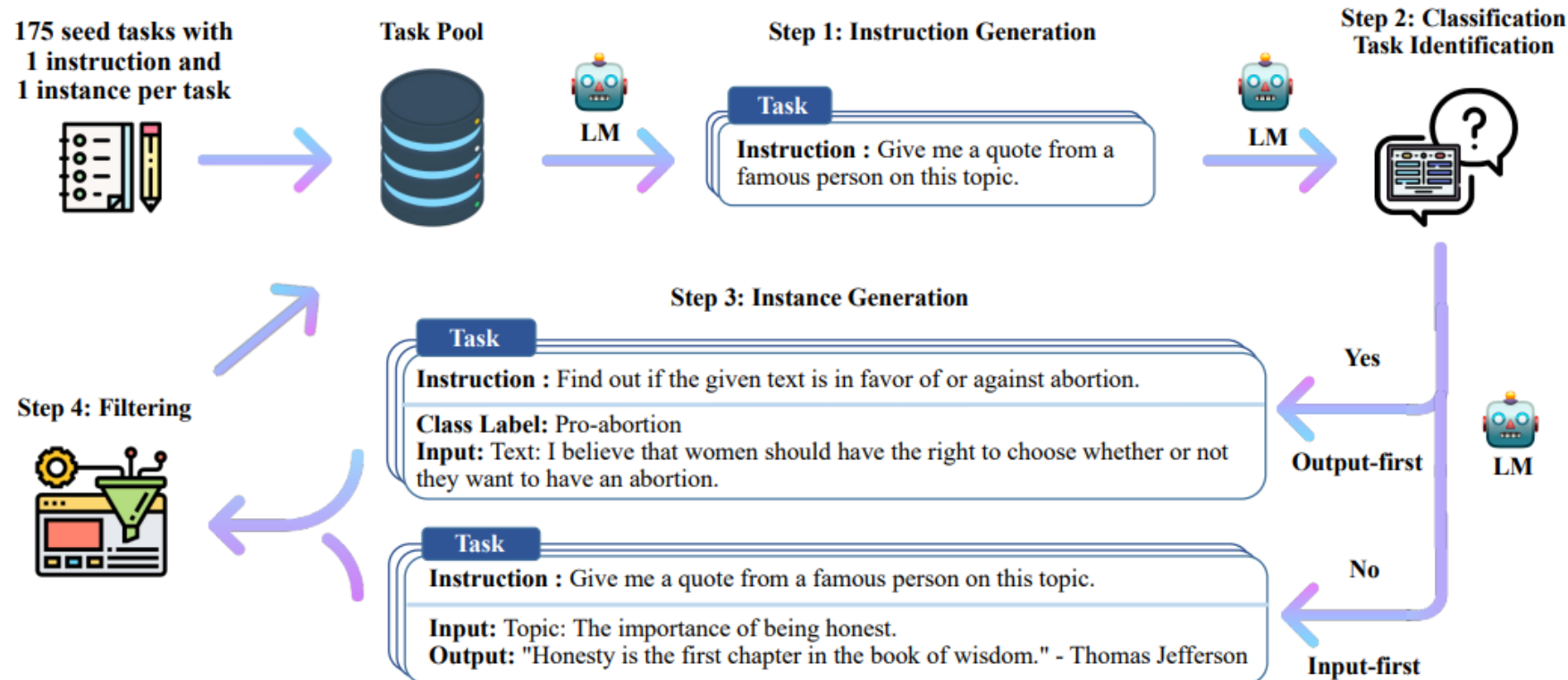
Home work

## A Survey on Data Synthesis and Augmentation for Large Language Models

## On LLMs-Driven Synthetic Data Generation, Curation, and Evaluation: A Survey



# Self-Instruct – Leveraging FMs to synthetically generate data

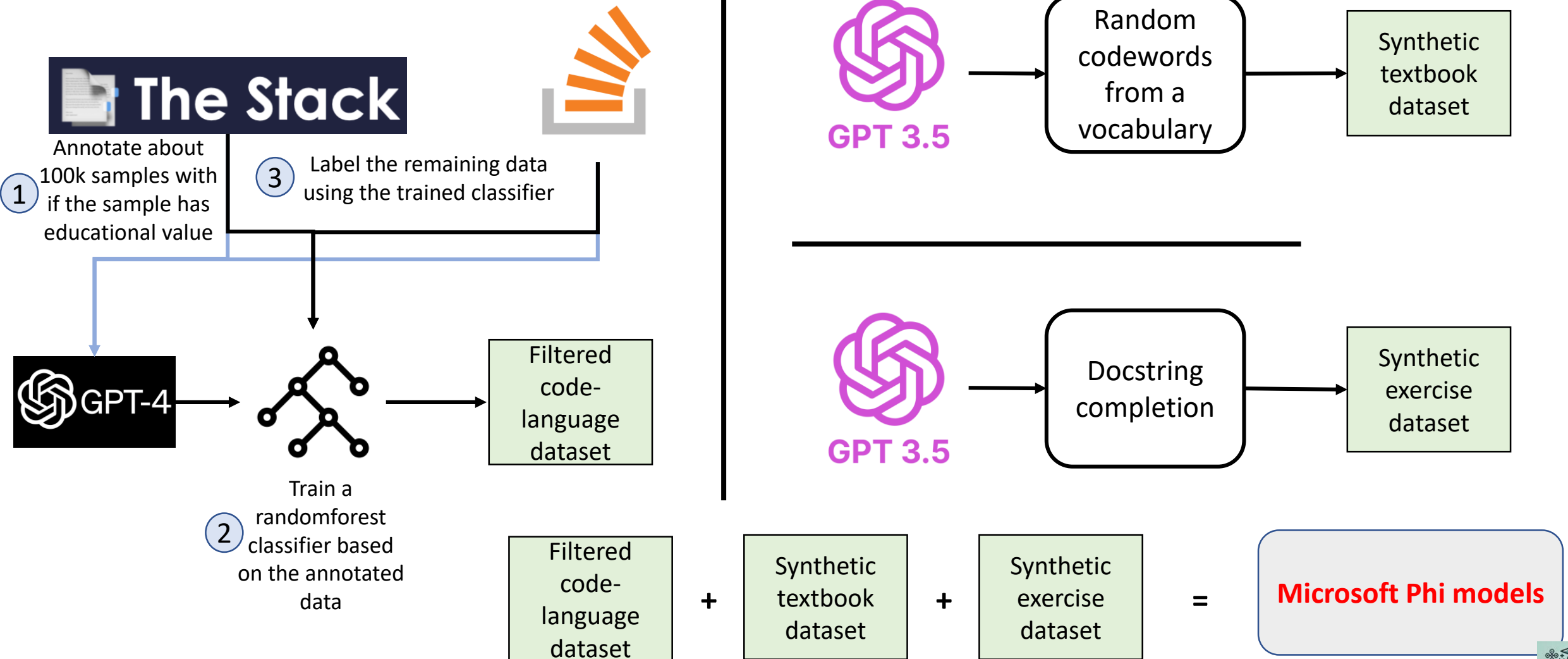


- Use FMs to generate data
- Start with a seed pool of prompts and then iteratively generate and refine task-specific data

Ensuring information diversity and overcoming FM-bias is challenging



# Microsoft Phi Models



**We now have all the data we need; Crisis averted – can we go home now?**

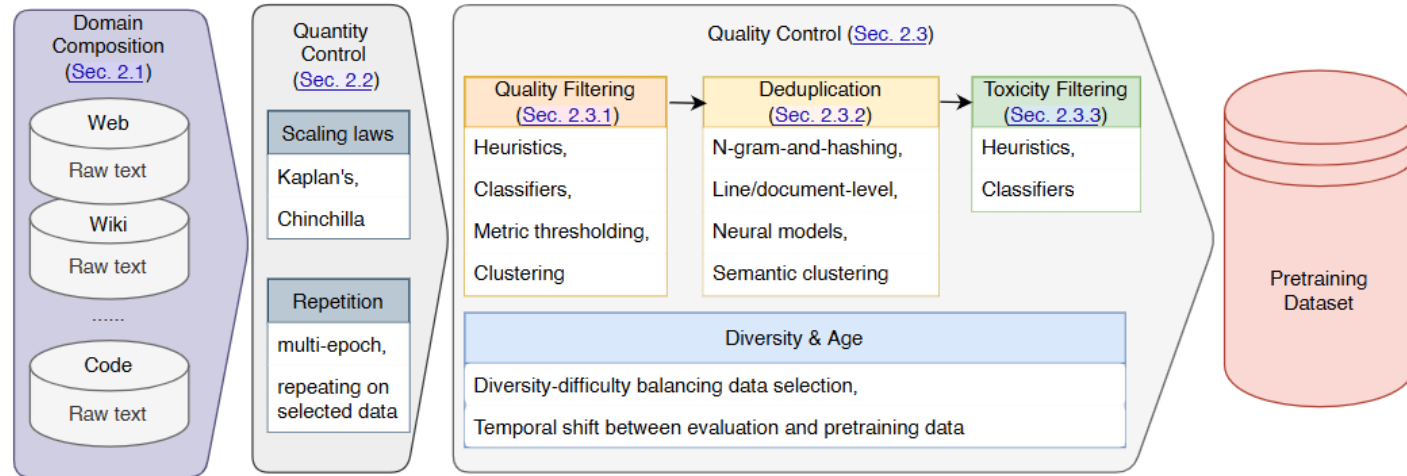


**Not all data are created equal – For the flywheel to work, we need the right kind of data**

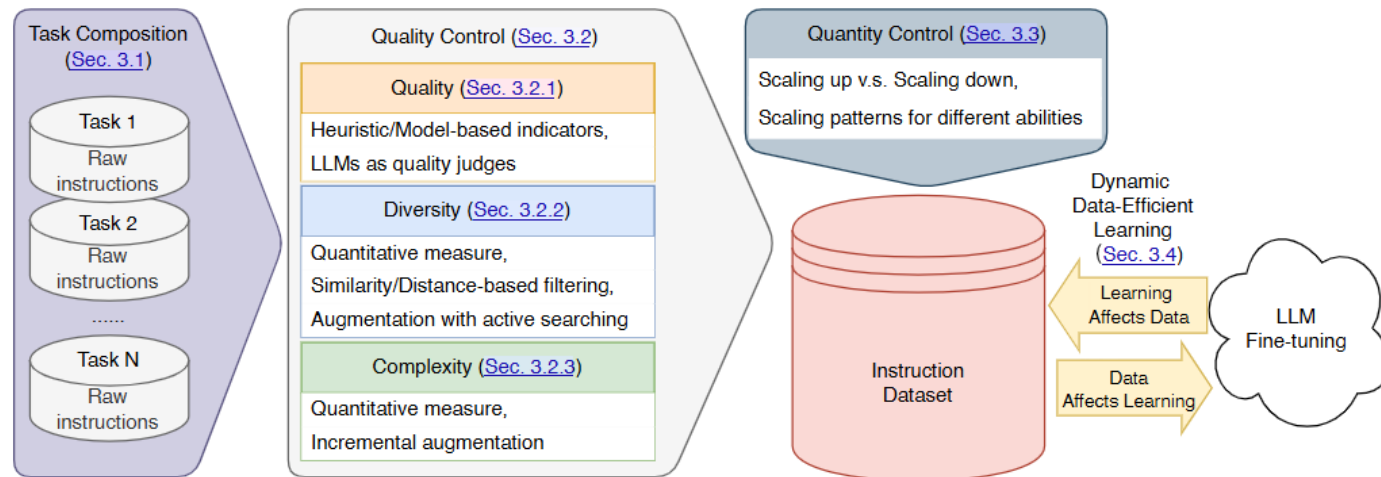




# Step 1: Ensuring the quality of data crucial for Data Flywheel's success



(a) Data management pipeline in the pretraining stage of LLMs

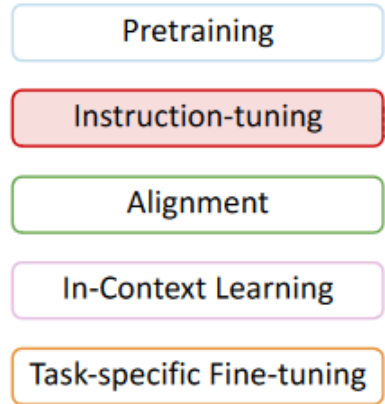


(b) Data management pipeline in the supervised fine-tuning stage of LLMs

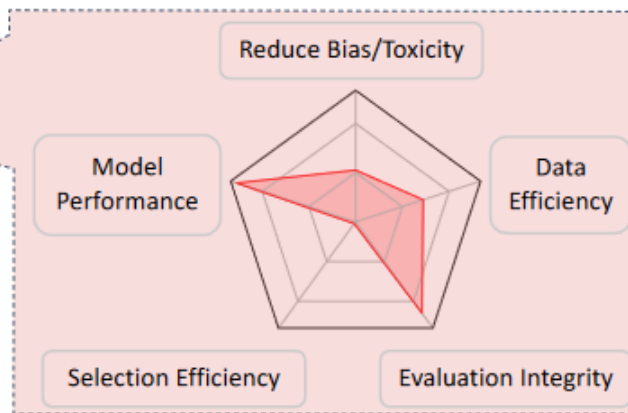


# Managing the quality, quantity, informativeness and diversity of the data

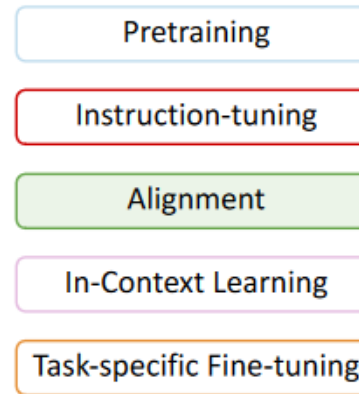
## Learning Stage



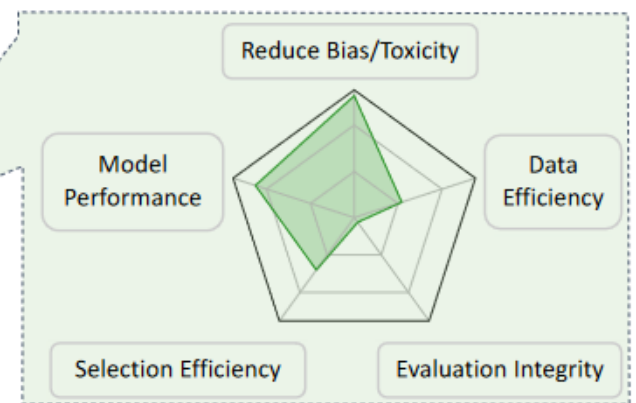
## Selection Objective



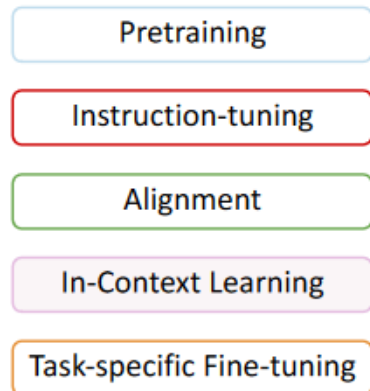
## Learning Stage



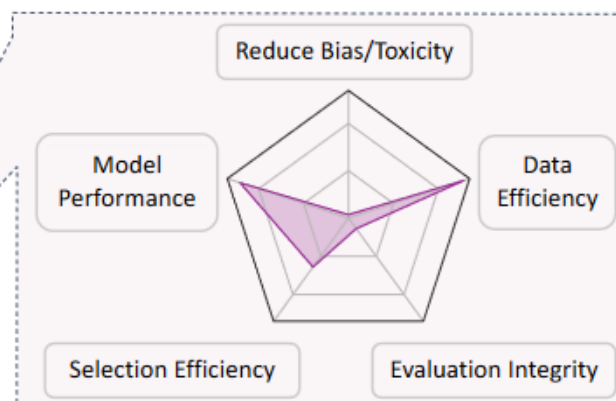
## Selection Objective



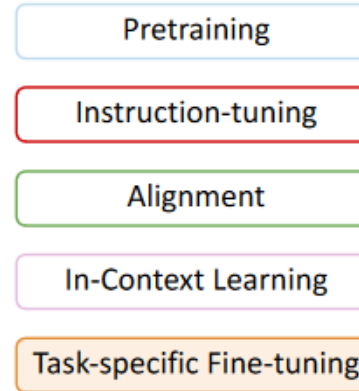
## Learning Stage



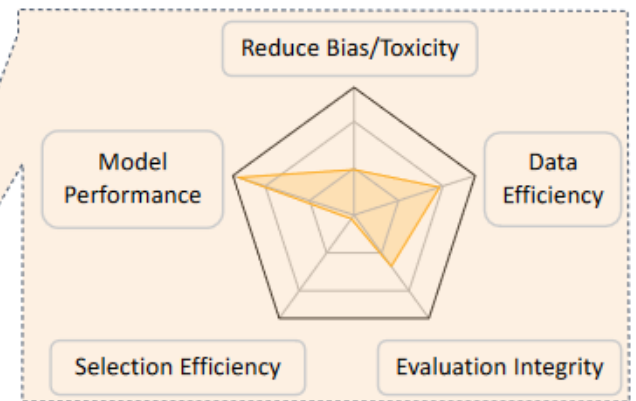
## Selection Objective



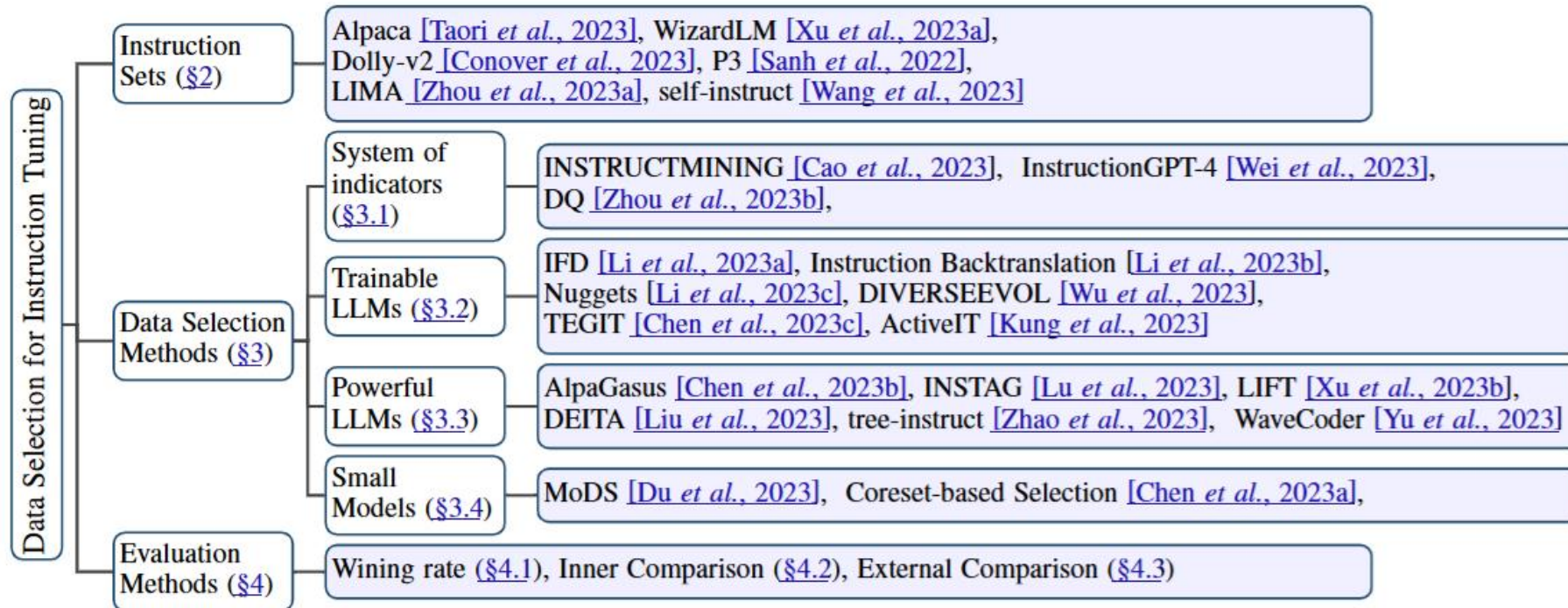
## Learning Stage



## Selection Objective



# Step 2: Select the right type and quantity of data



---

## LIMA: Less Is More for Alignment

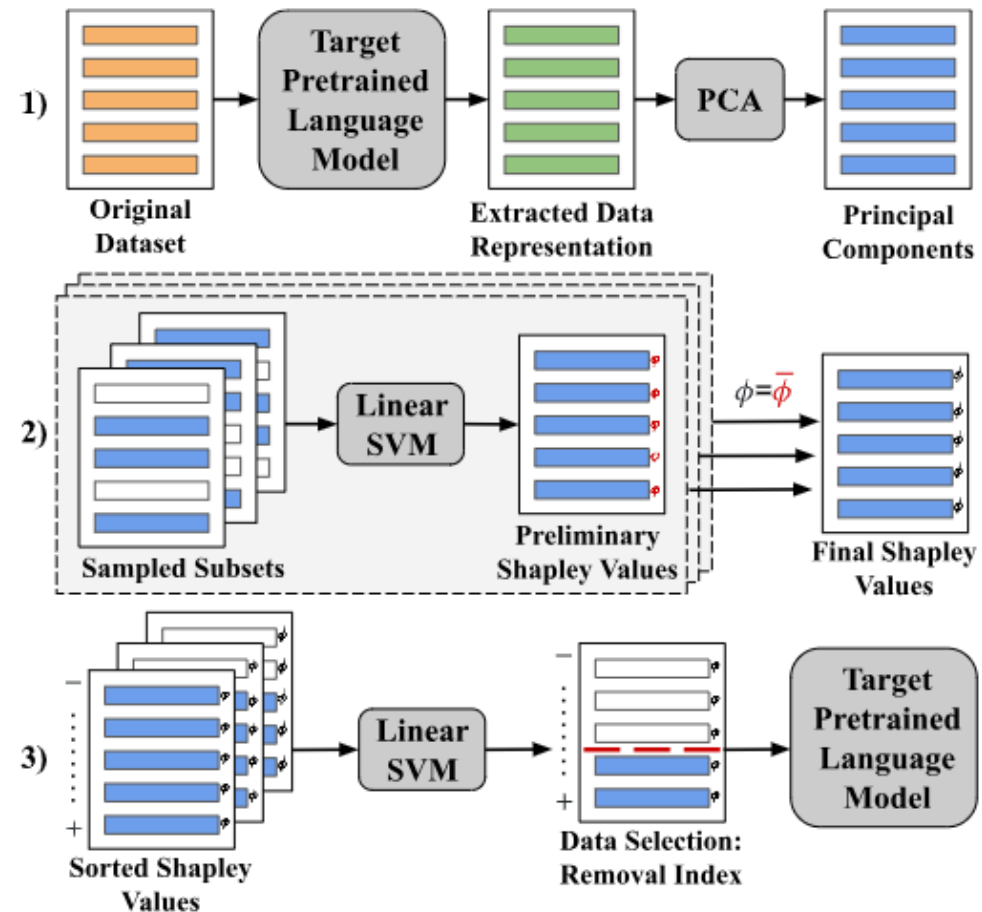
---





# Step 3: Ensure that the data that we use to build the FMware is useful

Data valuation methods estimate the contributive value of each data point so that the data we include is useful



# Do we now have all the ingredients?

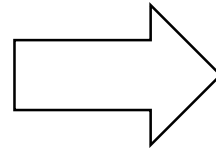


## Crude Oil

*['krüd 'oi(-ə)]*

A raw natural resource that is extracted from the earth and used to propel vehicles, heat buildings, produce electricity, and make everyday products.

 Investopedia



# Not Yet!



# **We can't always rely on just Synthetic data**

FMs are limited by the data that it was trained on ...  
If we run out of data, all the variety of data that we  
can synthetically generate **might not have  
enough information diversity!**

**We need to get the human generated  
data!**






# But can you trust data from anyone!




# Lots of feedback data is available

Which response do you prefer?  
Your choice will help make ChatGPT better.




 Response 1


Your approach of looking at a subject holistically and then breaking it down into manageable chunks is an effective way to tackle complex topics or problems. It allows for a comprehensive understanding of the overall context before diving into the details. This method can help with problem-solving, decision-making, and learning. If you ever have specific topics or questions you'd like to explore in this manner, feel free to bring them up, and I'll be here to assist!

 Response 2

That's a wise approach, especially when dealing with complex subjects. Starting with a holistic view allows you to understand the big picture and context, while breaking it down into smaller, manageable chunks helps in addressing each aspect systematically and thoroughly. It's a valuable method for learning, problem-solving, and making informed decisions.

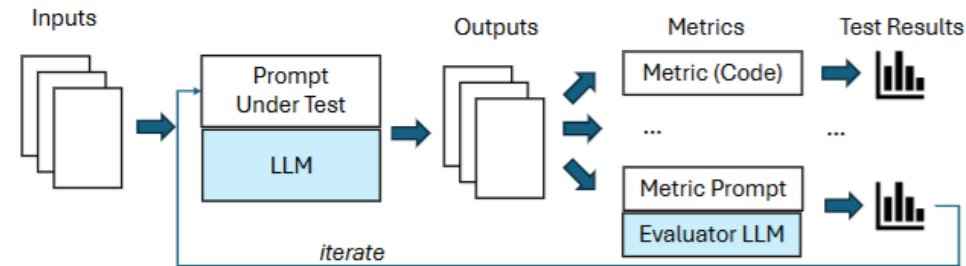
revealed data of the

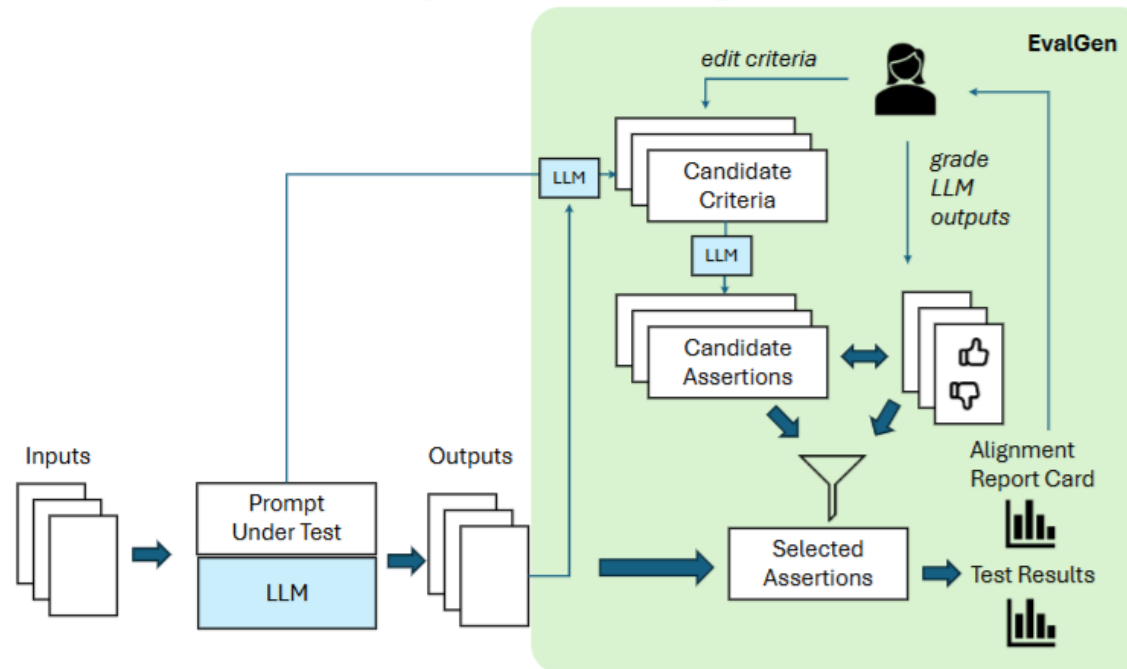
 Message ChatGPT...



# Flywheel requires robust ways of integrating user generated data



(a) Typical Evaluation Pipeline





# Flywheel requires robust ways of integrating user generated data

**a** Prompt Node: You will be doing named entity recognition (NER). Extract up to 3 well-known entities from the following tweet: {tweet\_full\_text}. For each entity, write one sentence describing the person or entity. All the entities you extract should be found in a knowledge base like Wikipedia, so don't make up.

**b** Multi-Evaluator: Let an AI help you generate criteria and implement evaluation functions.

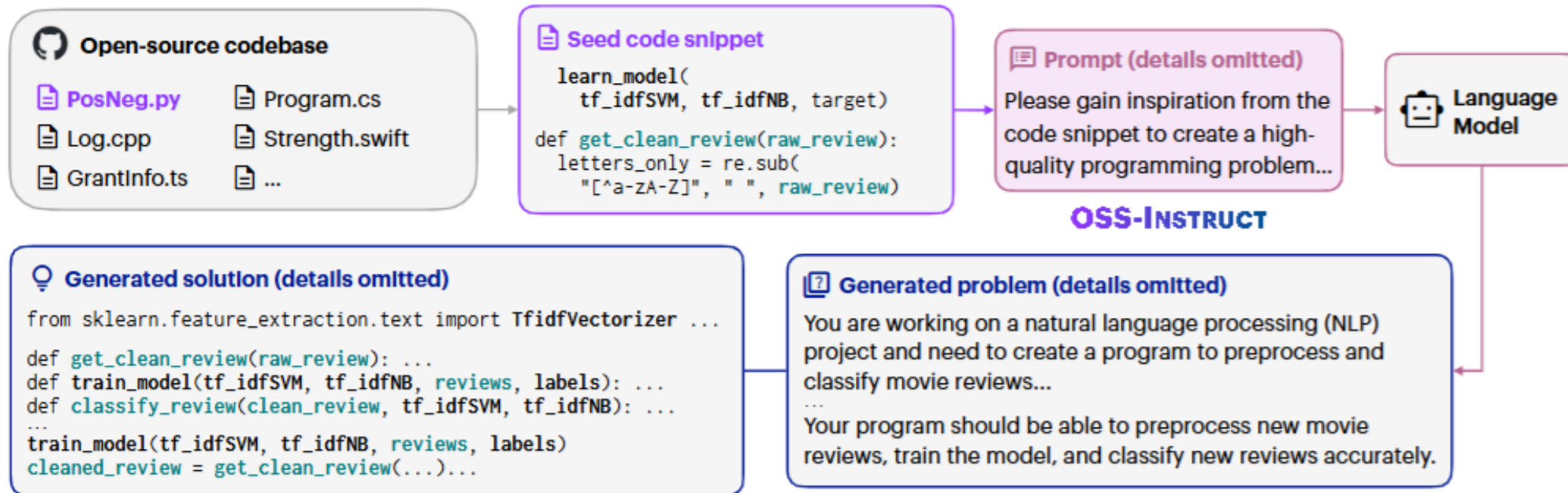
**c** Type a new criteria to add, then press Enter: the response is valid JSON. Suggest more.

**d** Is this response **Bad!** or **Good!**?  
- Bravotv: A television network that focuses on reality TV shows, including popular franchises like The Real Housewives.  
- BravoWHL: Stands for Bravo's Watch What Happens Live, a late-night talk show hosted by Andy Cohen that features celebrity interviews, games, and discussions about Bravo's reality TV shows.  
- Paris Hilton: A well-known American socialite, businesswoman, and media personality, known for her appearance on the reality TV show The Simple Life and her work as a singer, actress, and entrepreneur.

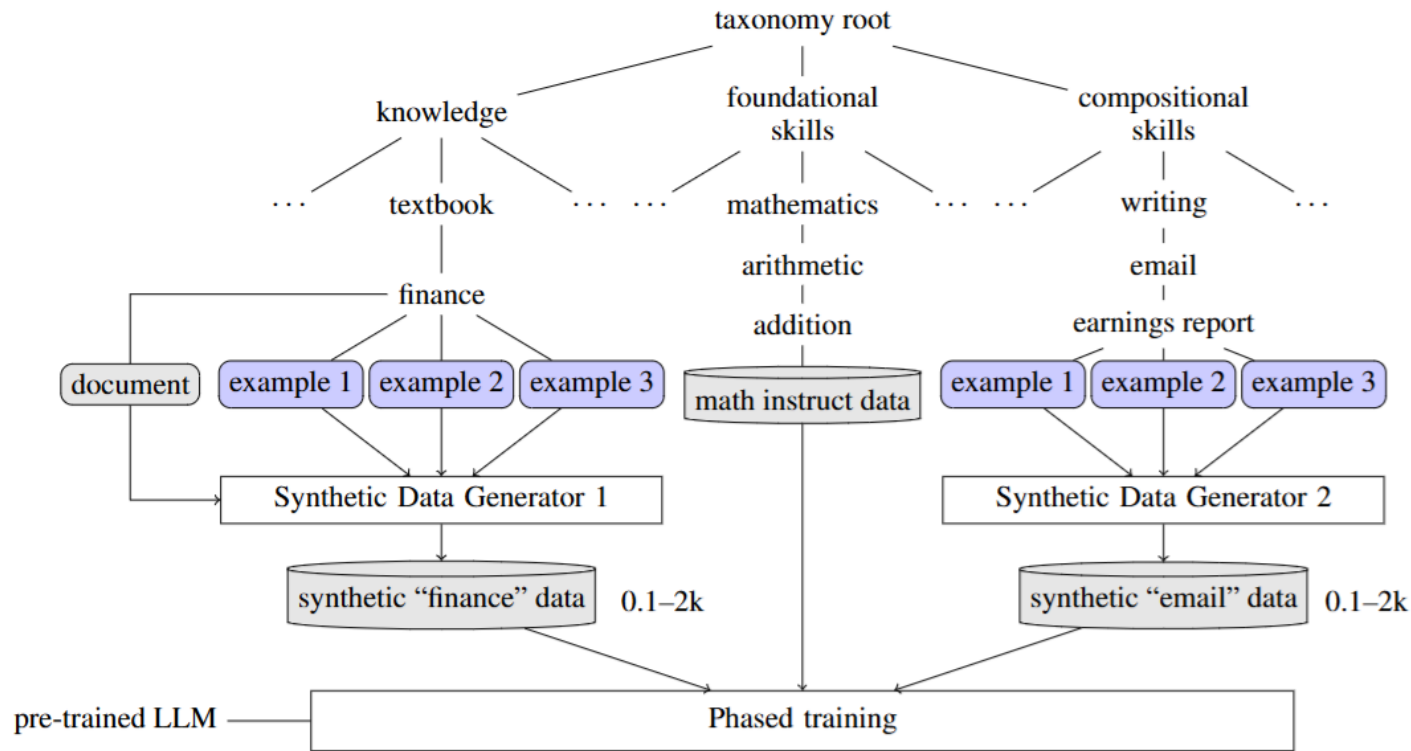
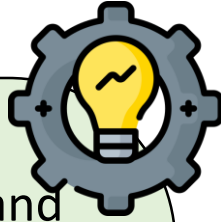
**e** Chosen Functions and Alignment: Coverage of Bad Responses 77.78%, False Failure Rate 28.57%, Entity Notability 56%.



# Flywheel in Action – OSS Instruct and Magicoder



# IBM InstructLab



- Decompose knowledge into Knowledge, foundational skills and compositional skills
- Leverage crowdsourcing to collect such Knowledge, skills and compositional skills for multiple domains in the form of question and answer pairs
- These skills and knowledge together acts as the curriculum for synthetic data generation
- The synthetically trained data is then used to fine-tune the model in a two-phase process





