

生成式AI在企業永續發展的應用

Generative AI for Corporate ESG and Sustainable Development

Time: 15:00-17:00, Saturday, December 2, 2023

Place: NTPU USR Hub

<https://learningcollaboration.org/index.php/2023/09/21/collegeco1120103/>

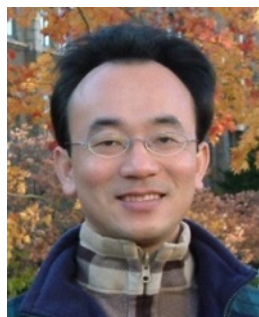
戴敏育 副教授

Min-Yuh Day, Ph.D.,
Associate Professor

Institute of Information Management, National Taipei University

<https://web.ntpu.edu.tw/~myday>

2023-12-02





戴敏育 博士

Min-Yuh Day, Ph.D.

Associate Professor, Information Management, NTPU

Visiting Scholar, IIS, Academia Sinica

Ph.D., Information Management, NTU

Director, Intelligent Financial Innovation Technology, IFIT Lab, IM, NTPU

Associate Director, Fintech and Green Finance Center, NTPU

**Publications Co-Chairs, IEEE/ACM International Conference on
Advances in Social Networks Analysis and Mining (ASONAM 2013-)**

**Program Co-Chair, IEEE International Workshop on
Empirical Methods for Recognizing Inference in Text (IEEE EM-RITE 2012-)**

**Publications Chair, The IEEE International Conference on
Information Reuse and Integration for Data Science (IEEE IRI 2007-)**



Outline

- 1. Generative AI**
- 2. ChatGPT and Large Language Models (LLMs)**
- 3. Popular Generative AI and Applications**
- 4. Generative AI for Corporate ESG and Sustainable Development**

衡量企業永續關鍵指標

臺北大學獨創ESG永續評鑑系統

社會(S)	經濟(E)	環境(E)	揭露(D)
1.人權 2.員工溝通與福利 3.人力資本發展 4.多元組成與包容性 5.供應鏈社會面向控管 6.客戶關係管理 7.產品安全 8.企業公民與慈善	1.股東權益 2.董事會結構與運作 3.行為準則與內控 4.風險及危機管理 5.永續金融 6.ESG創新	1.環境系統與治理 2.空氣管理 3.能源與氣候變遷 4.水管理 5.原物料與廢棄物管理/ 資源與廢棄物管理 6.生物多樣性 7.供應商及產品生命週期管理/ 供應鏈環境面向管理	1.ESG 揭露

Generative AI

(Gen AI)

AI Generated Content

(AIGC)

Generative AI (Gen AI)

AI Generated Content (AIGC)

Image Generation

Instruction 1:

An astronaut riding a horse in a photorealistic style.

Instruction 2:

Teddy bears working on new AI research on the moon in the 1980s.

 **OpenAI DALL·E 2**

Figure 1



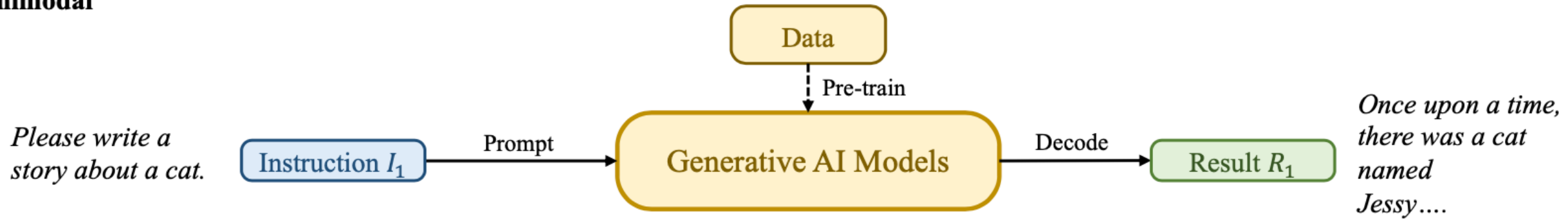
Figure 2



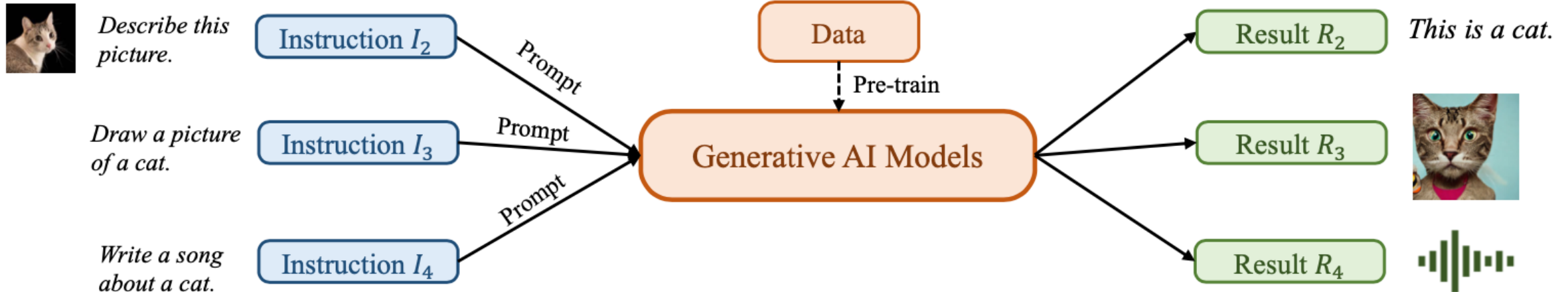
Generative AI (Gen AI)

AI Generated Content (AIGC)

Unimodal

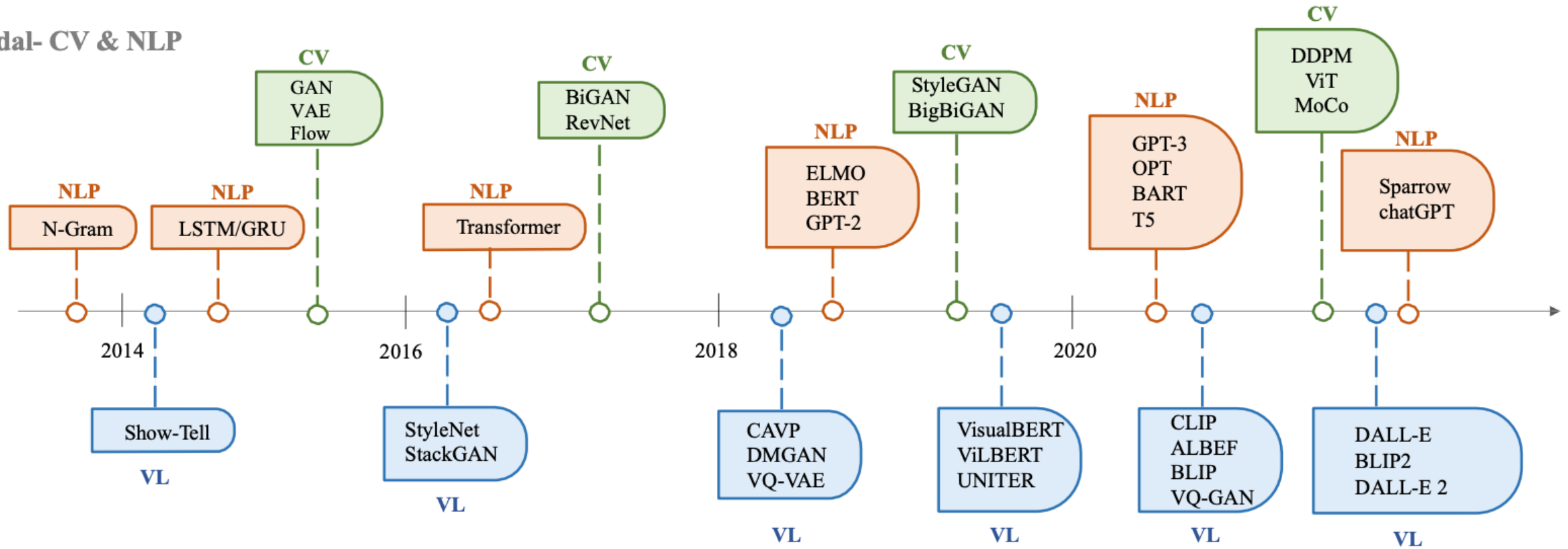


Multimodal



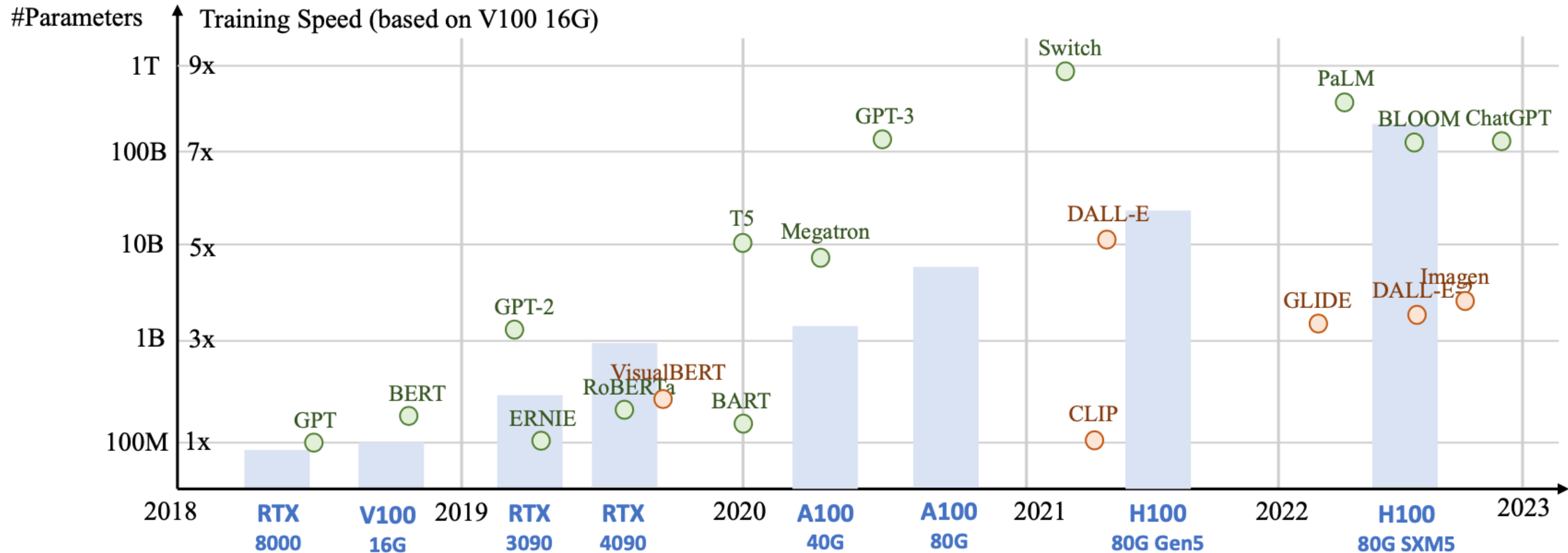
The history of Generative AI in CV, NLP and VL

Unimodal- CV & NLP

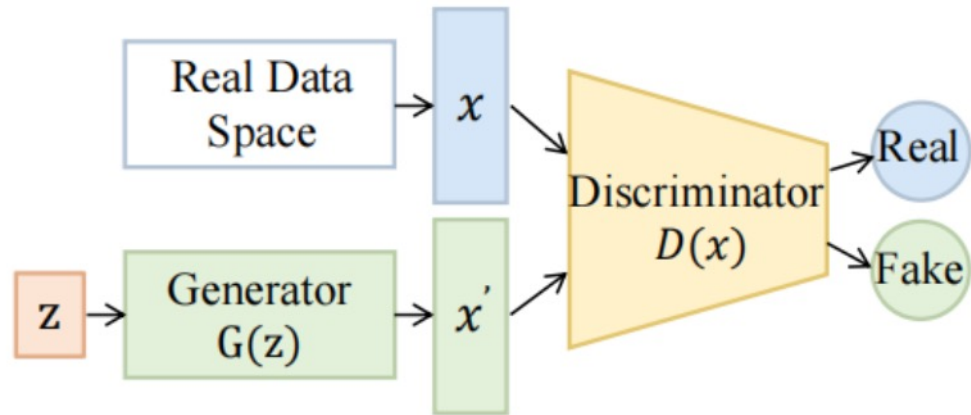


Multimodal – Vision Language

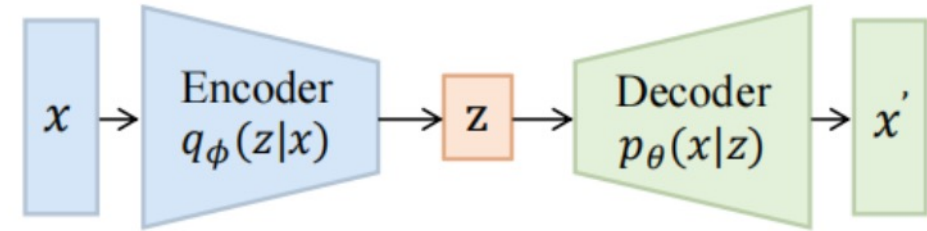
Generative AI Foundation Models



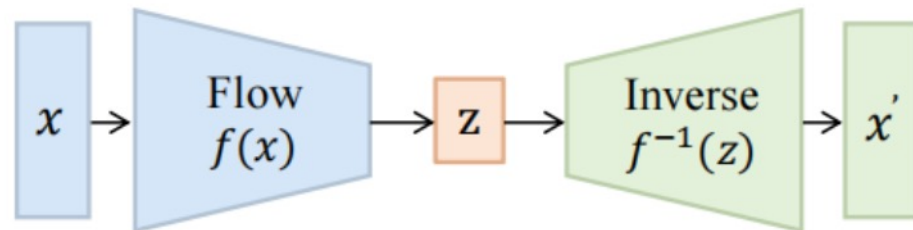
Categories of Vision Generative Models



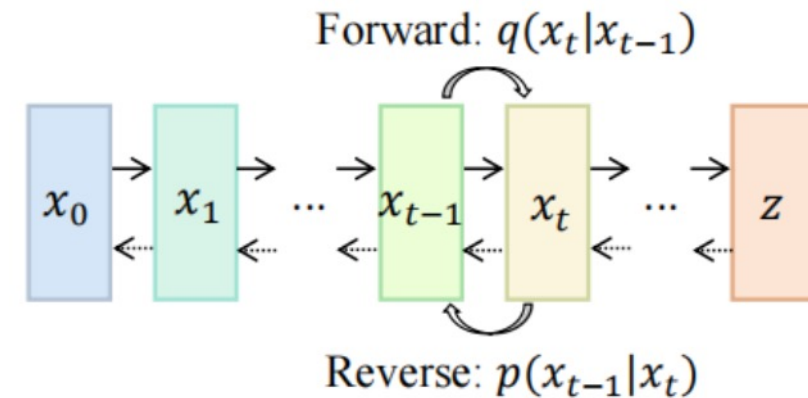
(1) Generative adversarial networks



(2) Variational autoencoders

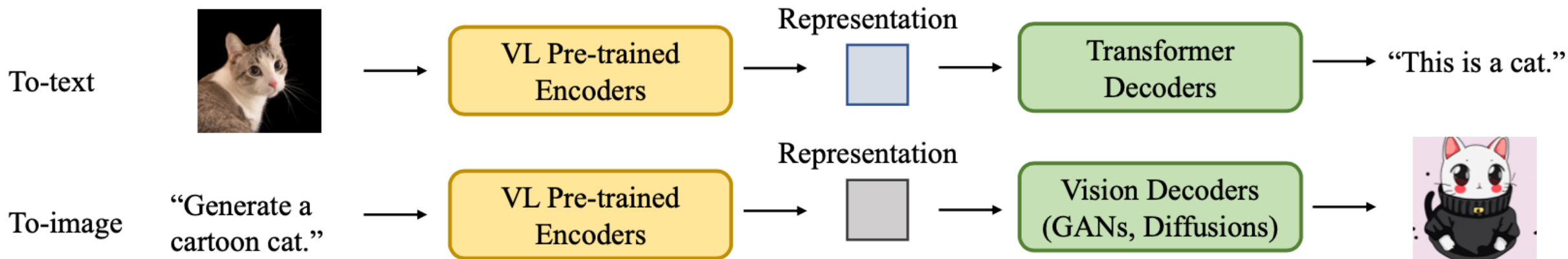
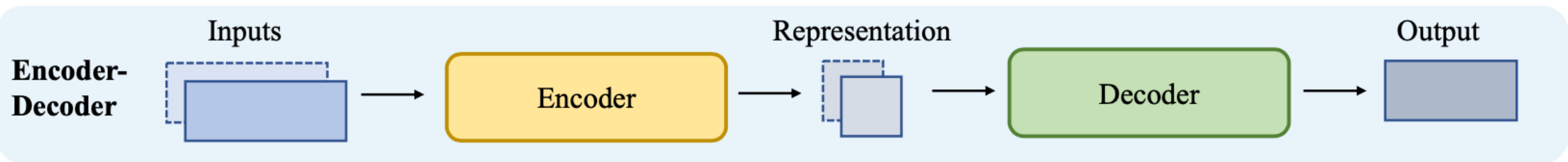


(3) Normalizing flows

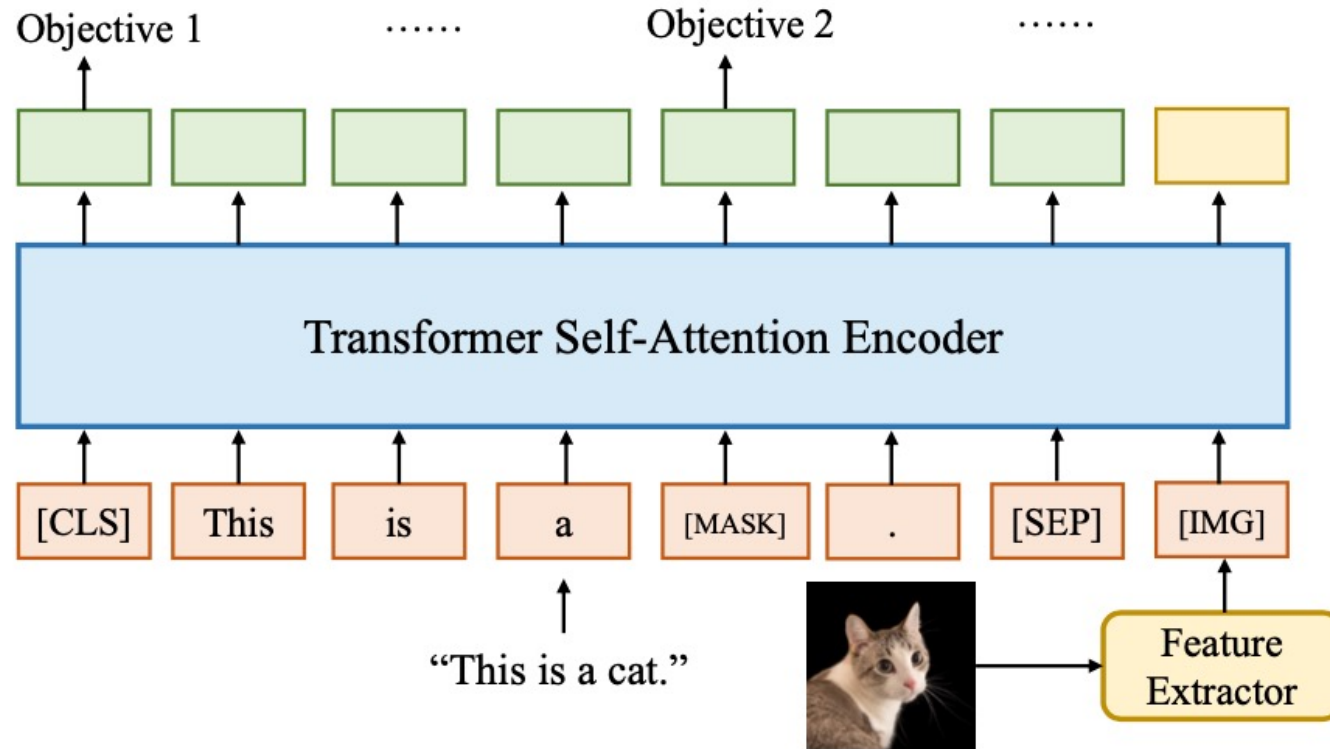


(4) Diffusion models

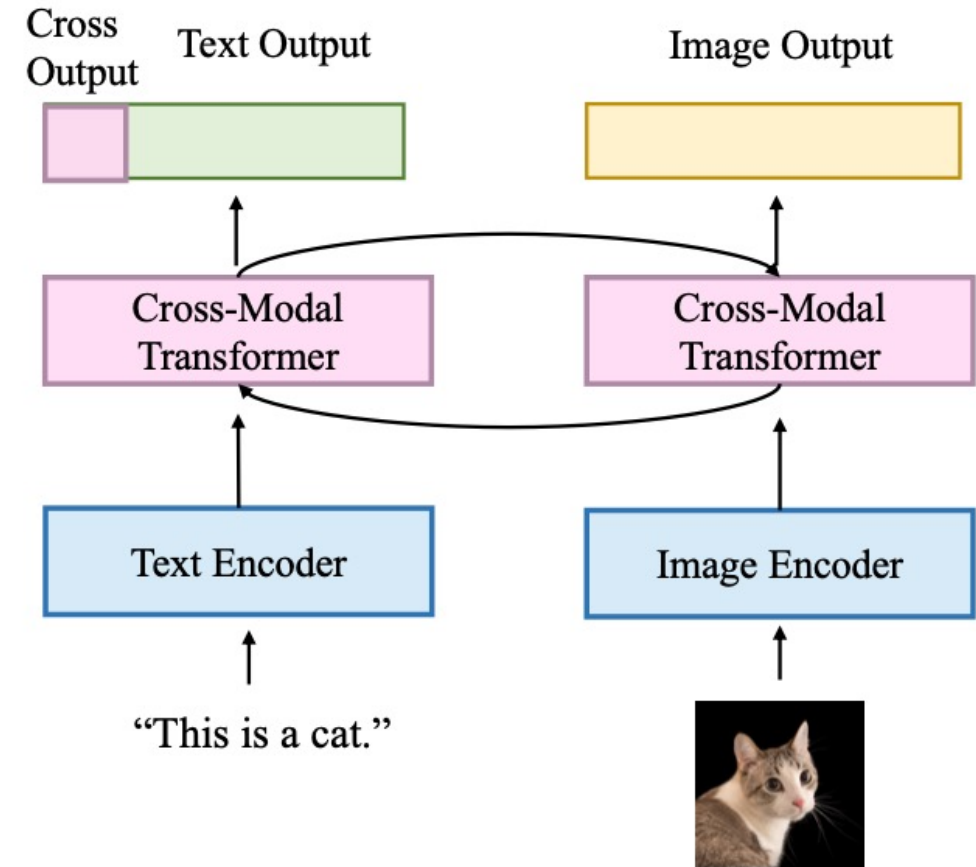
The General Structure of Generative Vision Language



Two Types of Vision Language Encoders: Concatenated Encoders and Cross-aligned Encoders

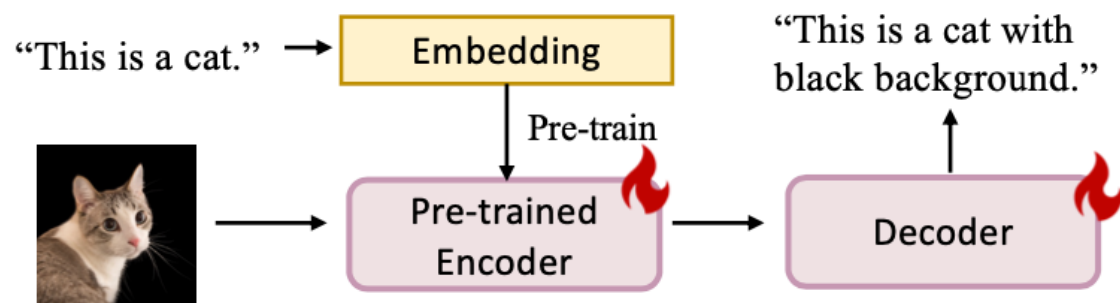


(a) Concatenated Encoder

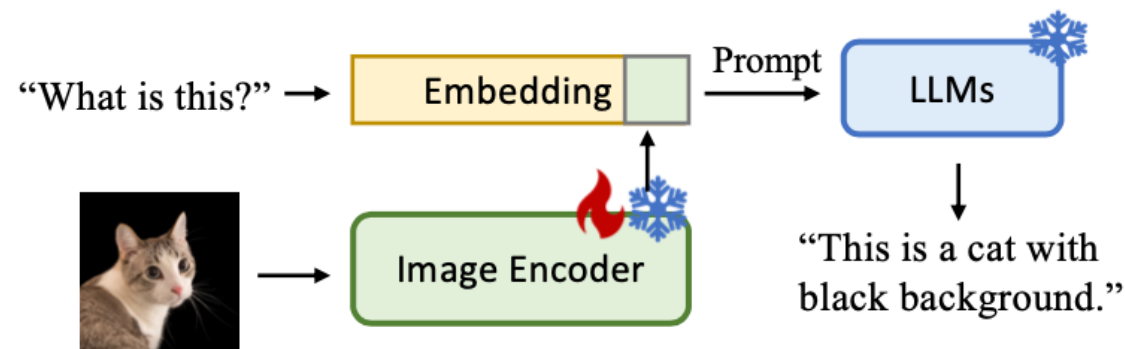


(b) Cross-aligned Encoder

Two Types of to-language Decoder Models: Jointly-trained Models and Frozen Models



(a) Jointly-trained Models



(b) Frozen Models

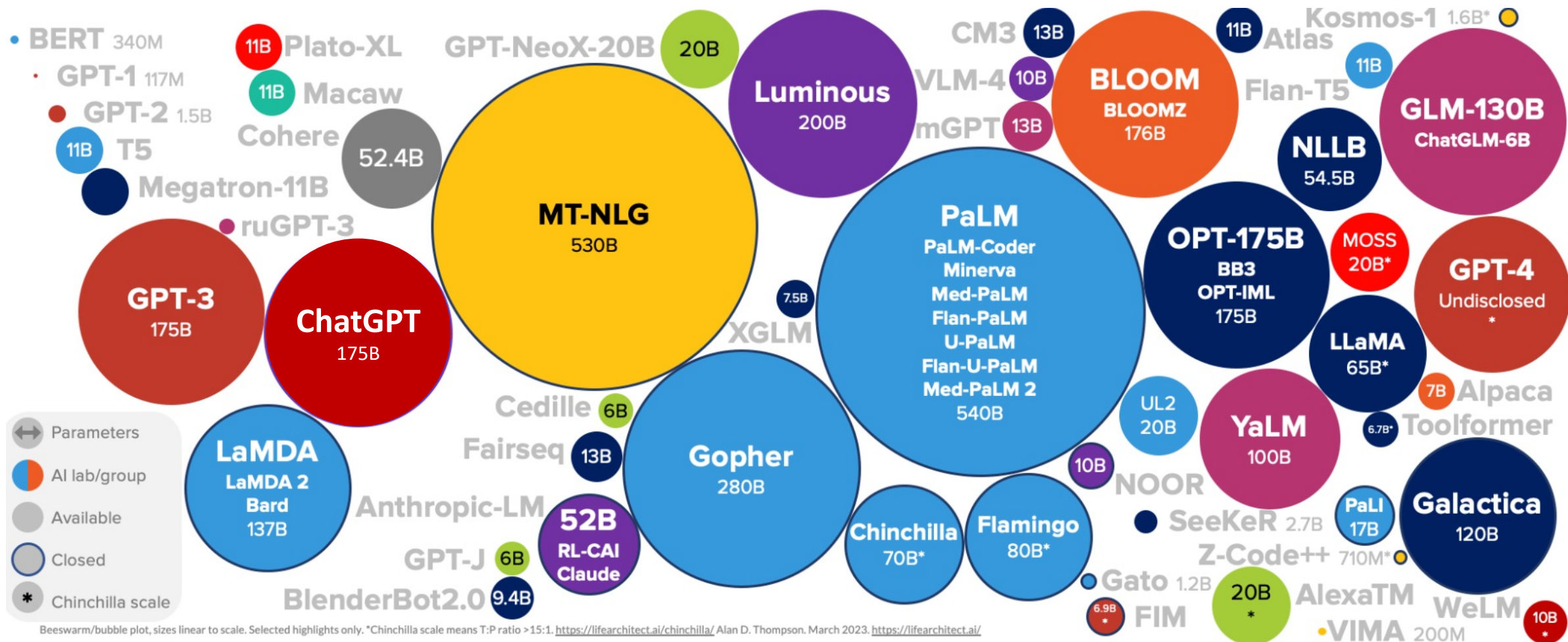
ChatGPT

Large Language Models (LLMs)

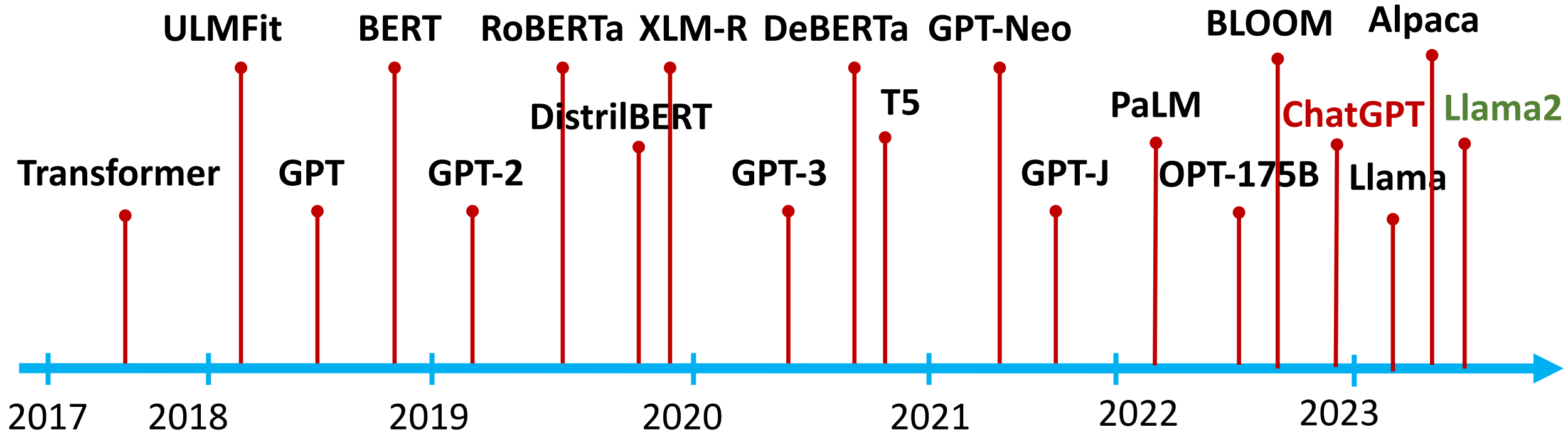
Foundation Models

Large Language Models (LLM)

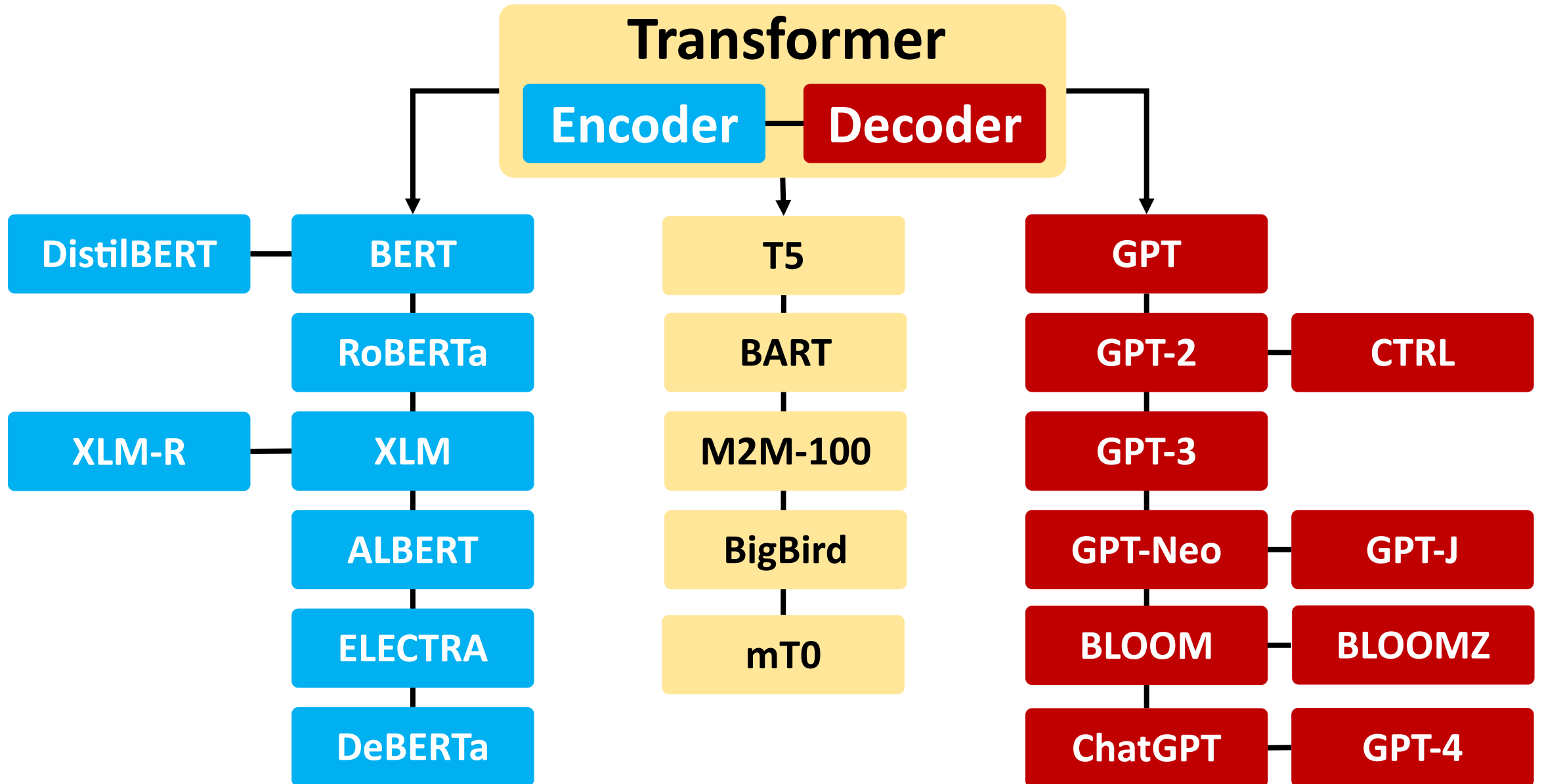
(GPT-3, ChatGPT, PaLM, BLOOM, OPT-175B, LLaMA)



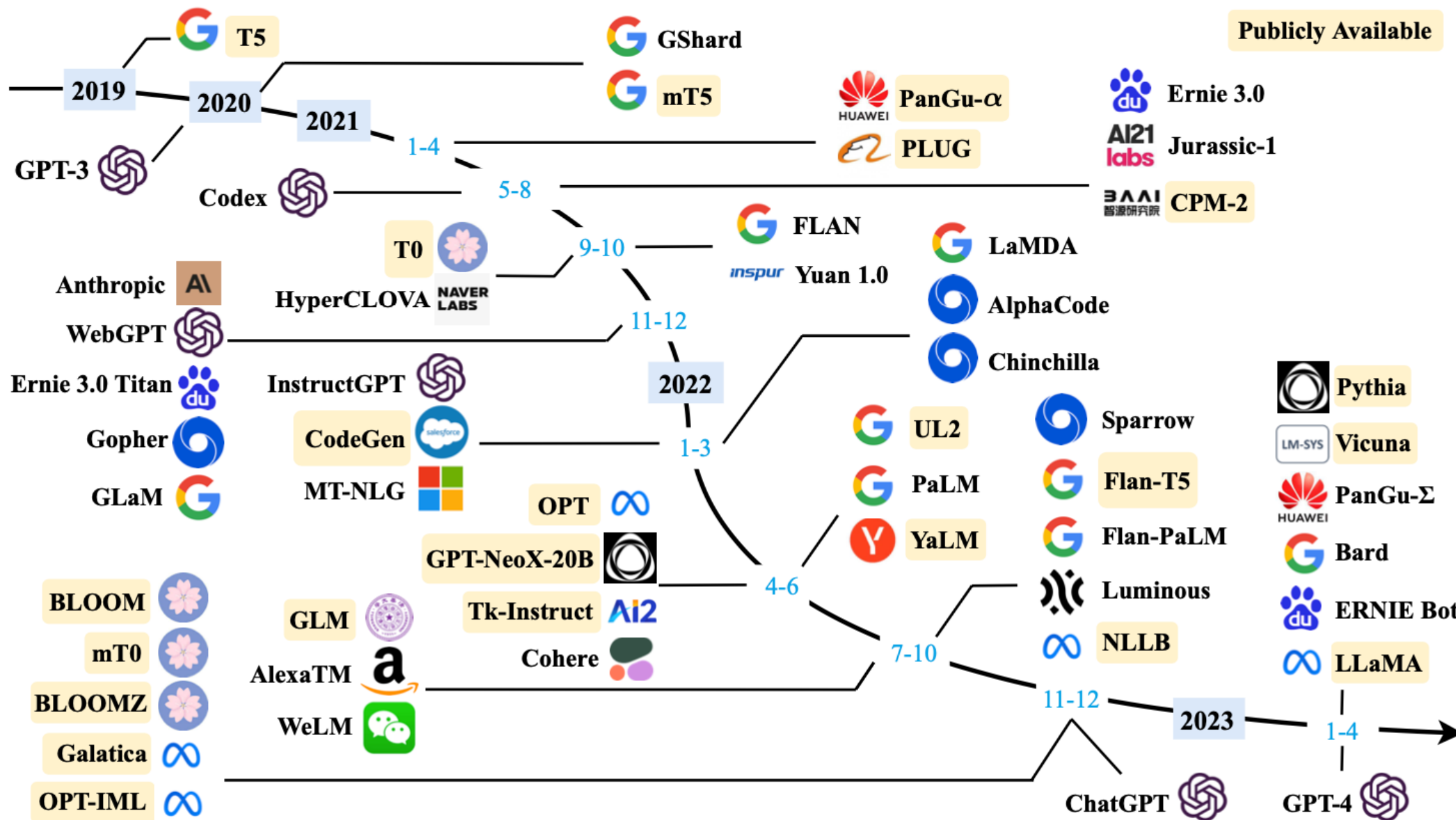
The Transformers Timeline



Transformer Models



Large Language Models (LLMs) (larger than 10B)



Large Language Models (LLMs) (larger than 10B)

	Model	Release Time	Size (B)	Base Model	Adaptation		Pre-train Data Scale	Latest Data Timestamp	Hardware (GPUs / TPUs)	Training Time	Evaluation	
					IT	RLHF					ICL	CoT
Publicly Available	T5 [72]	Oct-2019	11	-	-	-	1T tokens	Apr-2019	1024 TPU v3	-	✓	-
	mT5 [73]	Oct-2020	13	-	-	-	1T tokens	-	-	-	✓	-
	PanGu- α [74]	Apr-2021	13*	-	-	-	1.1TB	-	2048 Ascend 910	-	✓	-
	CPM-2 [75]	Jun-2021	198	-	-	-	2.6TB	-	-	-	-	-
	T0 [28]	Oct-2021	11	T5	✓	-	-	-	512 TPU v3	27 h	✓	-
	CodeGen [76]	Mar-2022	16	-	-	-	577B tokens	-	-	-	✓	-
	GPT-NeoX-20B [77]	Apr-2022	20	-	-	-	825GB	-	96 40G A100	-	✓	-
	Tk-Instruct [78]	Apr-2022	11	T5	✓	-	-	-	256 TPU v3	4 h	✓	-
	UL2 [79]	May-2022	20	-	-	-	1T tokens	Apr-2019	512 TPU v4	-	✓	✓
	OPT [80]	May-2022	175	-	-	-	180B tokens	-	992 80G A100	-	✓	-
	NLLB [81]	Jul-2022	54.5	-	-	-	-	-	-	-	✓	-
	GLM [82]	Oct-2022	130	-	-	-	400B tokens	-	768 40G A100	60 d	✓	-
	Flan-T5 [83]	Oct-2022	11	T5	✓	-	-	-	-	-	✓	✓
	BLOOM [68]	Nov-2022	176	-	-	-	366B tokens	-	384 80G A100	105 d	✓	-
	mT0 [84]	Nov-2022	13	mT5	✓	-	-	-	-	-	✓	-
	Galactica [35]	Nov-2022	120	-	-	-	106B tokens	-	-	-	✓	✓
	BLOOMZ [84]	Nov-2022	176	BLOOM	✓	-	-	-	-	-	✓	-
	OPT-IML [85]	Dec-2022	175	OPT	✓	-	-	-	128 40G A100	-	✓	✓
	LLaMA [57]	Feb-2023	65	-	-	-	1.4T tokens	-	2048 80G A100	21 d	✓	-
	Pythia [86]	Apr-2023	12	-	-	-	300B tokens	-	256 40G A100	-	✓	-

Large Language Models (LLMs) (larger than 10B)

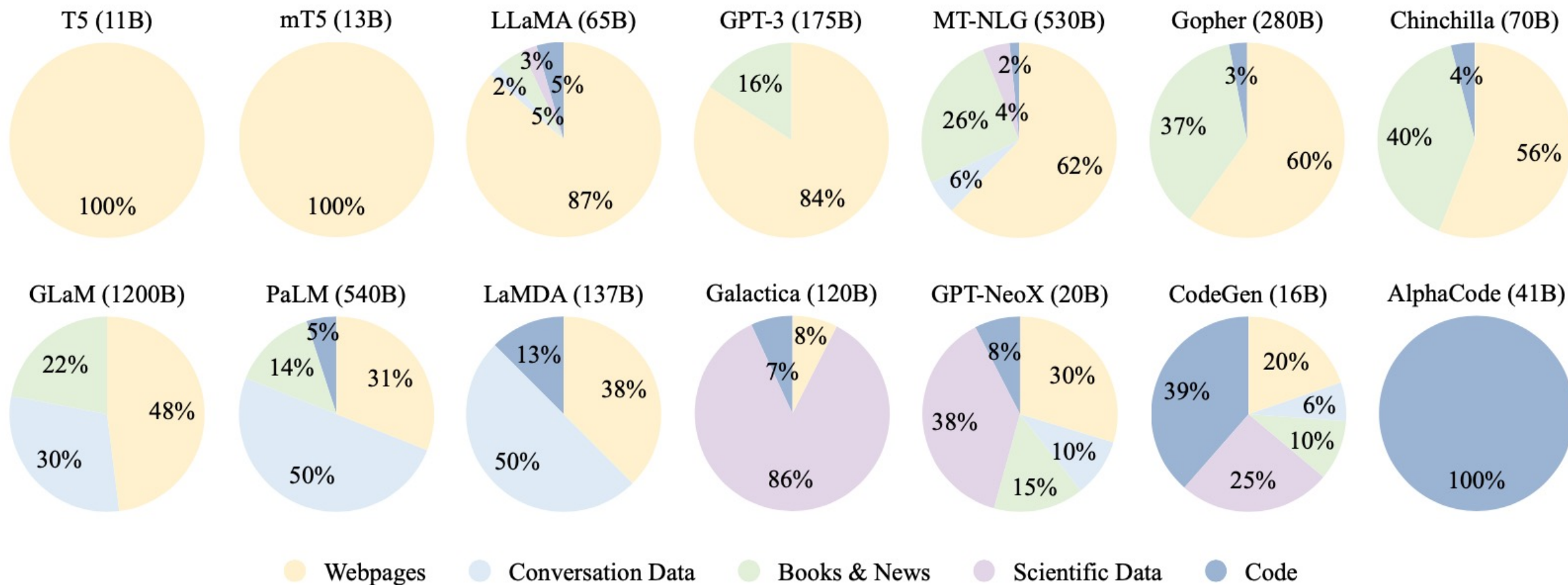
	Model	Release Time	Size (B)	Base Model	Adaptation		Pre-train Data Scale	Latest Data Timestamp	Hardware (GPUs / TPUs)	Training Time	Evaluation	
					IT	RLHF					ICL	CoT
Closed Source	GPT-3 [55]	May-2020	175	-	-	-	300B tokens	-	-	-	✓	-
	GShard [87]	Jun-2020	600	-	-	-	1T tokens	-	2048 TPU v3	4 d	-	-
	Codex [88]	Jul-2021	12	GPT-3	-	-	100B tokens	May-2020	-	-	✓	-
	ERNIE 3.0 [89]	Jul-2021	10	-	-	-	375B tokens	-	384 V100	-	✓	-
	Jurassic-1 [90]	Aug-2021	178	-	-	-	300B tokens	-	800 GPU	-	✓	-
	HyperCLOVA [91]	Sep-2021	82	-	-	-	300B tokens	-	1024 A100	13.4 d	✓	-
	FLAN [62]	Sep-2021	137	LaMDA	✓	-	-	-	128 TPU v3	60 h	✓	-
	Yuan 1.0 [92]	Oct-2021	245	-	-	-	180B tokens	-	2128 GPU	-	✓	-
	Anthropic [93]	Dec-2021	52	-	-	-	400B tokens	-	-	-	✓	-
	WebGPT [71]	Dec-2021	175	GPT-3	-	✓	-	-	-	-	✓	-
	Gopher [59]	Dec-2021	280	-	-	-	300B tokens	-	4096 TPU v3	920 h	✓	-
	ERNIE 3.0 Titan [94]	Dec-2021	260	-	-	-	300B tokens	-	2048 V100	28 d	✓	-
	GLaM [95]	Dec-2021	1200	-	-	-	280B tokens	-	1024 TPU v4	574 h	✓	-
	LaMDA [96]	Jan-2022	137	-	-	-	2.81T tokens	-	1024 TPU v3	57.7 d	-	-
	MT-NLG [97]	Jan-2022	530	-	-	-	270B tokens	-	4480 80G A100	-	✓	-
	AlphaCode [98]	Feb-2022	41	-	-	-	967B tokens	Jul-2021	-	-	-	-
	InstructGPT [61]	Mar-2022	175	GPT-3	✓	✓	-	-	-	-	✓	-
	Chinchilla [34]	Mar-2022	70	-	-	-	1.4T tokens	-	-	-	✓	-
	PaLM [56]	Apr-2022	540	-	-	-	780B tokens	-	6144 TPU v4	-	✓	✓
	AlexaTM [99]	Aug-2022	20	-	-	-	1.3T tokens	-	128 A100	120 d	✓	✓
	Sparrow [100]	Sep-2022	70	-	-	✓	-	-	64 TPU v3	-	✓	-
	WeLM [101]	Sep-2022	10	-	-	-	300B tokens	-	128 A100 40G	24 d	✓	-
	U-PaLM [102]	Oct-2022	540	PaLM	-	-	-	-	512 TPU v4	5 d	✓	✓
	Flan-PaLM [83]	Oct-2022	540	PaLM	✓	-	-	-	512 TPU v4	37 h	✓	✓
	Flan-U-PaLM [83]	Oct-2022	540	U-PaLM	✓	-	-	-	-	-	✓	✓
	GPT-4 [46]	Mar-2023	-	-	✓	✓	-	-	-	-	✓	✓
	PanGu- Σ [103]	Mar-2023	1085	PanGu- α	-	-	329B tokens	-	512 Ascend 910	100 d	✓	-

Statistics of Commonly-used Data Sources for LLMs

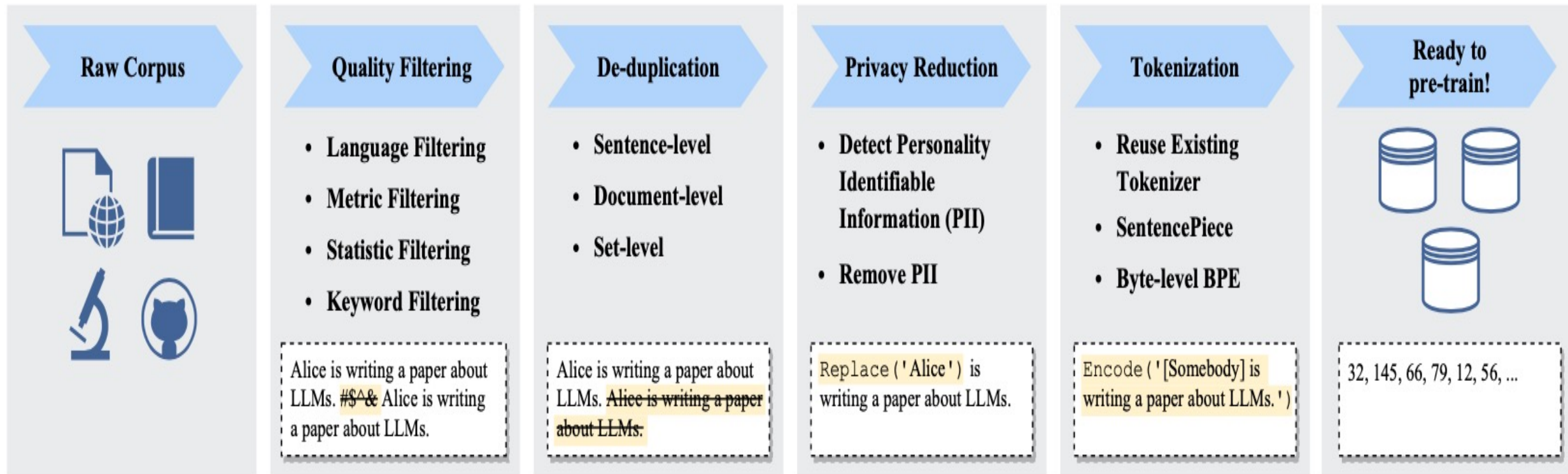
Corpora	Size	Source	Latest Update Time
BookCorpus [109]	5GB	Books	Dec-2015
Gutenberg [110]	-	Books	Dec-2021
C4 [72]	800GB	CommonCrawl	Apr-2019
CC-Stories-R [111]	31GB	CommonCrawl	Sep-2019
CC-NEWS [27]	78GB	CommonCrawl	Feb-2019
REALNEWS [112]	120GB	CommonCrawl	Apr-2019
OpenWebText [113]	38GB	Reddit links	Mar-2023
Pushift.io [114]	-	Reddit links	Mar-2023
Wikipedia [115]	-	Wikipedia	Mar-2023
BigQuery [116]	-	Codes	Mar-2023
the Pile [117]	800GB	Other	Dec-2020
ROOTS [118]	1.6TB	Other	Jun-2022

Source: Wanyin Liu Zhao, Kun Zhao, Junyi Li, Hanyu Tang, Xiaohu Wang, Lupeng Hou, Mingqian Wang et al. (2023). A Survey of Large Language Models. arXiv preprint arXiv:2305.10223.

Ratios of various data sources in the pre-training data for existing LLMs



Typical Data Preprocessing Pipeline for Pre-training Large Language Models (LLMs)



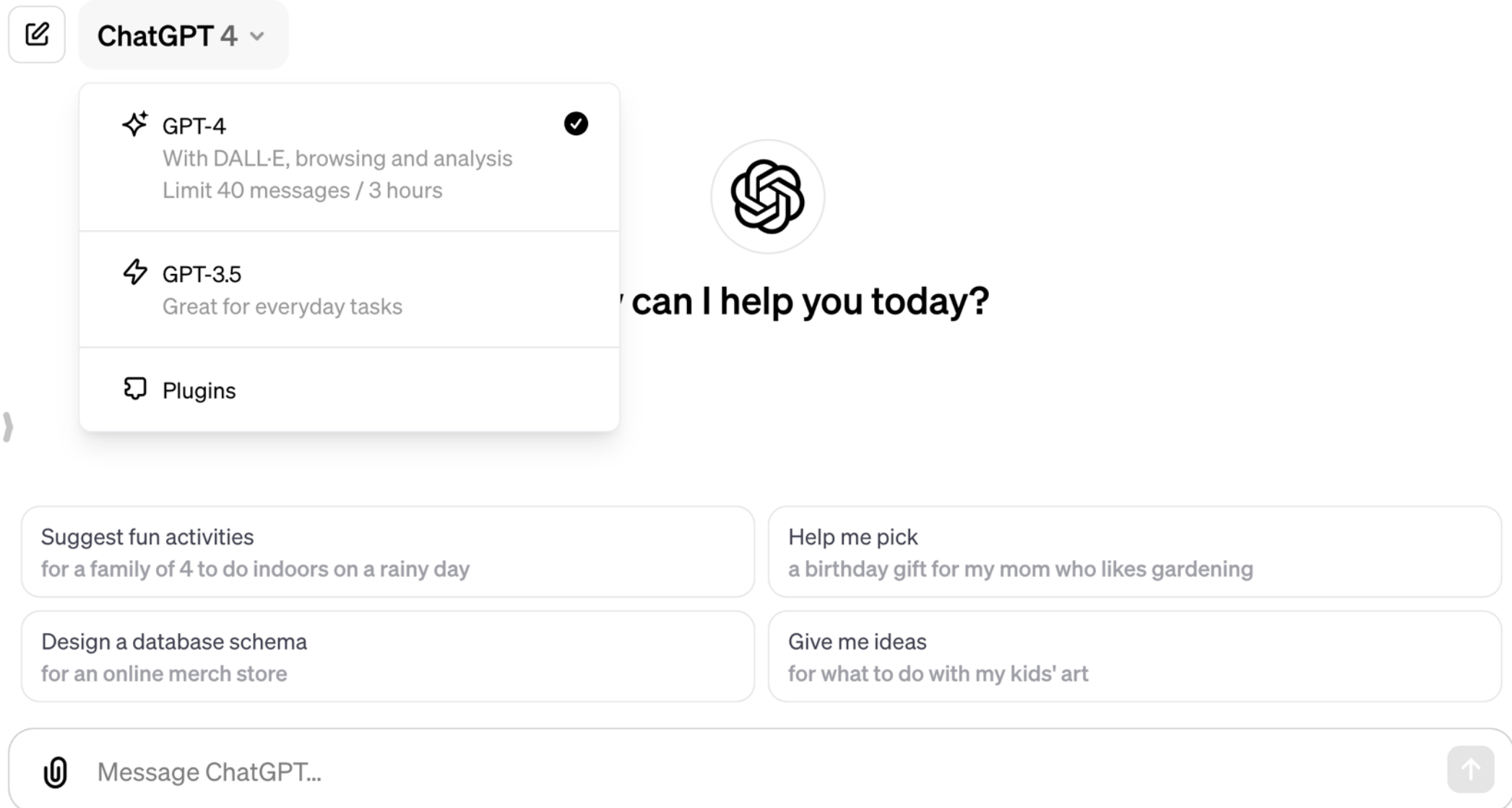
Generative AI

**Text, Image, Video, Audio
Applications**

Popular Generative AI

- **OpenAI ChatGPT (GPT-3.5, GPT-4)**
- **OpenAI DALL·E 3**
- **Perplexity.ai**
- **Chat.LMSys.org**
 - **Llama 2, Mistral, Zephyr, Vicuna, MPT, Falcon**
- **ChatPDF**
- **Stable Diffusion**
- **Video: D-ID, Synthesia**
- **Audio: Speechify**

OpenAI ChatGPT (GPT-4, GPT-3.5)



ChatGPT can make mistakes. Consider checking important information.


<https://chat.openai.com/>

OpenAI ChatGPT (GPT-4) DALL·E 3


Vector graphic of a flowchart depicting the integration of generative AI in the education process, from content creation to virtual experiments, personalized learning, and innovative learning.



Perplexity.ai

 perplexity

New Thread

 K

Home

Discover

Library

Login


Sign Up


Try Pro

Upgrade to Claude-2 or GPT-4, boost Copilot uses, and upload more files.

Learn More

Download






Anonymous Like

Share


Large Language Models for Data Science

Sources


What Is a Large Language Model and What Can...

 snowflake • 1

What are Large Language Models? | NVIDIA

 nvidia • 2

Catch Up On Large Language Models | by...

 towardsdatascien... • 3

View 3 more

Answer

Large Language Models (LLMs) are advanced AI systems designed to understand human language intricacies and generate intelligent, creative responses to queries ¹. They are deep learning algorithms that can recognize, summarize, translate, predict, and generate content using very large datasets ². LLMs are trained on enormous data sets typically measured in petabytes, sourced from books, articles, websites, and other text-based sources ¹. They are dynamic tools that are continuously updated and refined as new data becomes available ¹.

In the field of data science, LLMs can potentially transform how teams source, manage, and analyze data ¹. They are used in numerous data science applications, and their ability to process and interpret vast amounts of text data has made them an indispensable part of

Ask follow-up...

Copilot

↑

social media posts or customer reviews, to determine whether the overall sentiment is

<https://www.perplexity.ai/>

Chat with Open Large Language Models:

<https://chat.lmsys.org/>

- [Llama 2](#): open foundation and fine-tuned chat models by Meta
- [Mistral](#): a large language model by Mistral AI team
- [WizardLM](#): an instruction-following LLM using evol-instruct by Microsoft
- [Zephyr](#): a chatbot fine-tuned from Mistral by Hugging Face
- [Code Llama](#): open foundation models for code by Meta
- [ChatGLM](#): an open bilingual dialogue language model by Tsinghua University
- [Qwen](#): a large language model by Alibaba Cloud
- [Vicuna](#): a chat assistant fine-tuned on user-shared conversations by LMSYS
- [Falcon](#): TII's flagship series of large language models

Chat with Open Large Language Models

chat.lmsys.org

⚔ Chatbot Arena ⚔ : Benchmarking LLMs in the Wild

[GPT-3.5](#): GPT-3.5 by OpenAI

[GPT-3.5-Turbo-1106](#): GPT-3.5-Turbo-1106 by OpenAI

[GPT-4-Turbo](#): GPT-4-Turbo by OpenAI

[GPT-4](#): ChatGPT-4 by OpenAI

[Claude](#): Claude 2 by Anthropic

[Claude Instant](#): Claude Instant by Anthropic

[Vicuna](#): a chat assistant fine-tuned on user-shared conversations by LMSYS

[Llama 2](#): open foundation and fine-tuned chat models by Meta

[ChatGLM](#): an open bilingual dialogue language model by Tsinghua University

[OpenChat 3.5](#): OpenChat 3.5 is a versatile, open-source language model fine-tuned using C-RLFT

[Mistral](#): a large language model by Mistral AI team

[Zephyr](#): a chatbot fine-tuned from Mistral by Hugging Face

[Qwen](#): a large language model by Alibaba Cloud

[Code Llama](#): open foundation models for code by Meta

[WizardLM](#): an instruction-following LLM using evol-instruct by Microsoft

[PaLM 2 Chat](#): PaLM 2 for Chat (chat-bison@001) by Google

[Falcon](#): TII's flagship series of large language models

<https://chat.lmsys.org/>

Large Language Models for Data Science

Llama 2-70b-chat

llama-2-70b-chat

Model A

Large Language Models for Data Science

Large language models have gained significant attention in recent years for their ability to process and analyze natural language data. These models have been successfully applied in various fields such as natural language processing, text classification, sentiment analysis, and machine translation. In data science, large language models can be used for a variety of tasks, including data cleaning, data transformation, and data analysis.

In this article, we will explore the use of large language models for data science tasks and discuss their advantages and limitations. We will also provide examples of large language models that can be used for data science applications.

Advantages of Large Language Models for Data

Mistral-7b-instruct

mistral-7b-instruct

Model B

Large Language Models for Data Science

Large Language Models (LLMs) are becoming increasingly important in data science due to their ability to understand and generate natural language text. Here are some of the key ways that LLMs can be used in data science:

1. Text classification: LLMs can be trained on large datasets of labeled text and then used to classify new text into different categories. This can be useful in a variety of applications, such as sentiment analysis, spam filtering, and topic classification.
2. Named entity recognition: LLMs can also be used to identify and extract specific pieces of information from text, such as names of people, organizations, and locations. This can be useful in applications such as information extraction

<https://chat.lmsys.org/>

Chat with Open Large Language Models: Chatbot Arena

Chatbot Arena Leaderboard

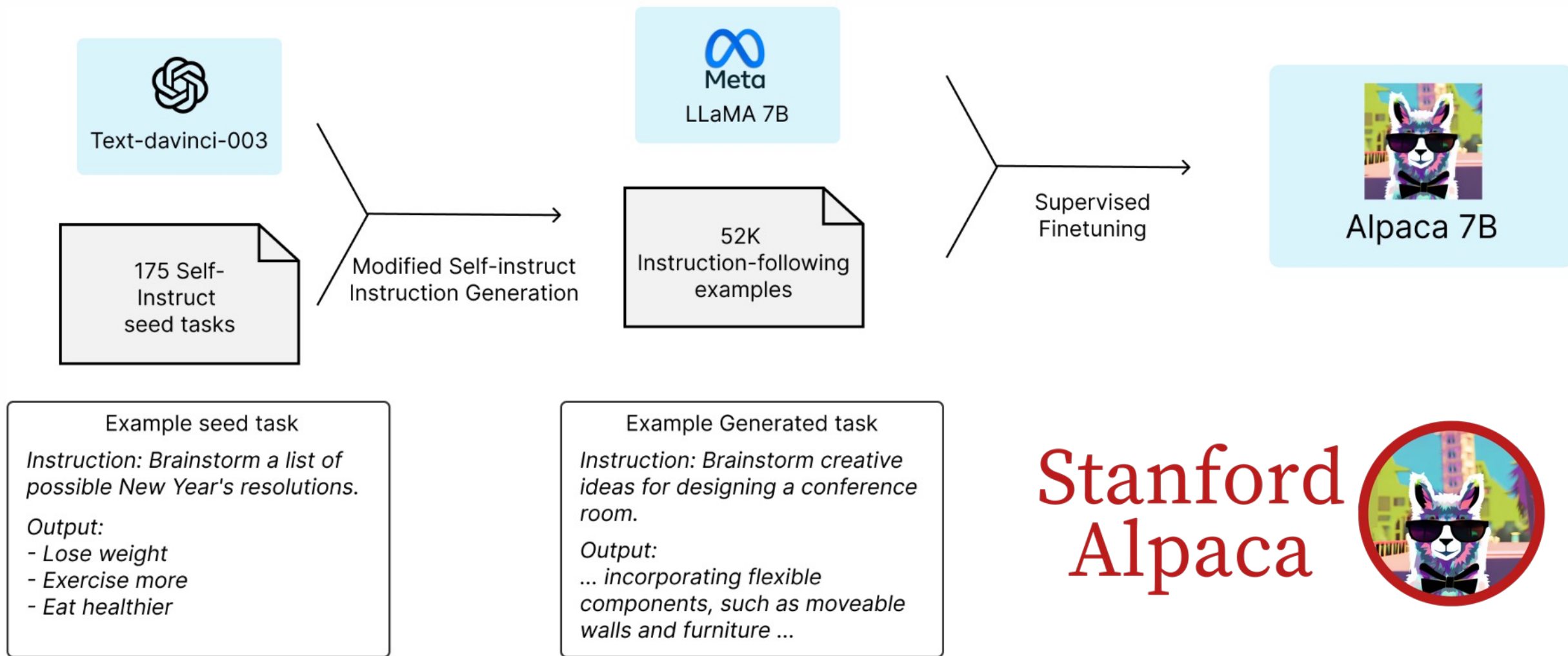
LLM Leaderboard

Model	★ Arena Elo rating	📈 MT-bench (score)	MMLU	License
GPT-4-Turbo	1210	9.32		Proprietary
GPT-4	1159	8.99	86.4	Proprietary
Claude-1	1146	7.9	77	Proprietary
Claude-2	1125	8.06	78.5	Proprietary
Claude-instant-1	1106	7.85	73.4	Proprietary
GPT-3.5-turbo	1103	7.94	70	Proprietary
WizardLM-70b-v1.0	1093	7.71	63.7	Llama 2 Community
Vicuna-33B	1090	7.12	59.2	Non-commercial
OpenChat-3.5	1070	7.81	64.3	Apache-2.0
Llama-2-70b-chat	1065	6.86	63	Llama 2 Community
WizardLM-13b-v1.2	1047	7.2	52.7	Llama 2 Community
zephyr-7b-beta	1042	7.34	61.4	MIT
MPT-30B-chat	1031	6.39	50.4	CC-BY-NC-SA-4.0
Vicuna-13B	1031	6.57	55.8	Llama 2 Community
QWen-Chat-14B	1030	6.96	66.5	Qianwen LICENSE

<https://chat.lmsys.org/>

Stanford Alpaca:

A Strong, Replicable Instruction-Following Model

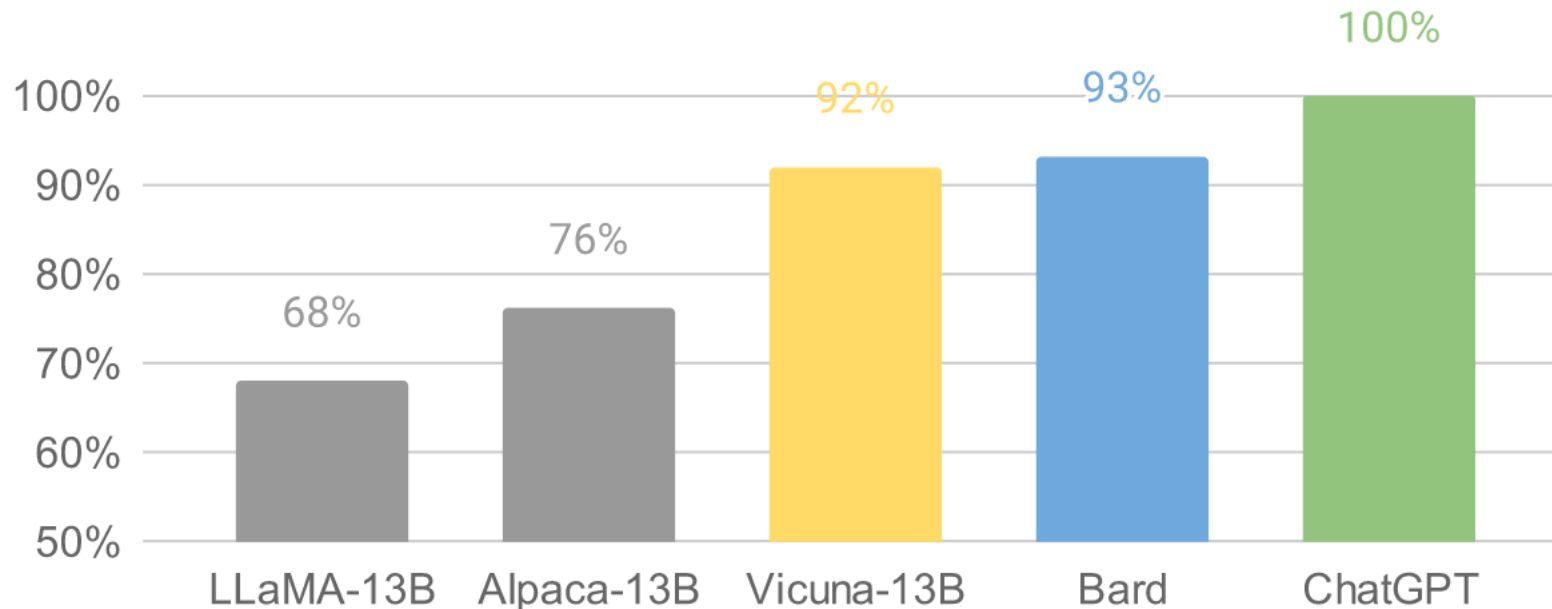


Vicuna: An Open-Source Chatbot

Impressing GPT-4 with 90%* ChatGPT Quality

by the Team with members from UC Berkeley, CMU, Stanford, and UC San Diego

- **Vicuna-13B: an open-source chatbot trained by fine-tuning LLaMA on user-shared conversations collected from ShareGPT.**
- **The cost of training Vicuna-13B is around \$300.**



Source: <https://vicuna.lmsys.org/>

Chinese-Vicuna:

A Chinese Instruction-following LLaMA-based Model

一個中文低資源的 llama+loram方案

Dataset Preview

Size: 409 MB

</> API

Go to dataset viewer

instruction (string)	input (string)	output (string)
"用一句话描述地球为什么是独一无二的。 \n "	" "	"地球上存在生命存在的条件和多样化的生命形式。 "
"给出一段对话，要求GPT模型使用合适的语气和回答方式继续对话。 对话： A：你今...	" "	"A： 这太好了！你的新工作听起来很令人兴奋。你对接接下来的日子有什么期望吗？ B： 是啊，我非...
"基于以下提示填写以下句子的空格。 提示： - 提供多种现实世界的场景 - 空格...	" "	"多种形容词可填，以下是其中一些例子： - 愉快的 - 惬意的 - 轻松的 - 安静的 - 美妙的"
"请生成一个新闻标题，描述一场正在发生的大型自然灾害。 \n "	" "	"强烈飓风肆虐，数百万人疏散！ "
"为指定的词汇创建一个关于该词汇的简短解释。 "人工智能"	" "	"人工智能是一种科技，通过使用计算机程序和算法实现类似于人类智能的能力，例如学习、推理和自...
"在给定的文本中查找并替换特定的单词。 原始文本： "这个句子中有一个错误的单词..."	" "	"替换后的文本： "这个句子中有一个修改的单词"
"根据给定的年份，计算该年是否为闰年。 \n \n1996\n"	" "	"是闰年。 "



Bactrian camel



Wild Bactrian camel



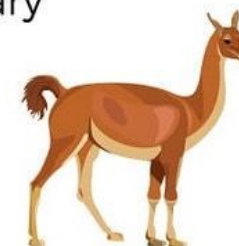
Dromedary



Llama



Alpaca



Guanaco



Vicuña

Chinese-Vicuna based on Guanaco Dataset and Belle Dataset

Source: https://huggingface.co/datasets/Chinese-Vicuna/guanaco_belle_merge_v1.0

Source: <https://github.com/Facico/Chinese-Vicuna>

RedPajama

a project to create leading open-source models,
starts by reproducing LLaMA training dataset of over 1.2 trillion tokens



Dataset	RedPajama	LLaMA*
CommonCrawl	878 billion	852 billion
C4	175 billion	190 billion
Github	59 billion	100 billion
Books	26 billion	25 billion
ArXiv	28 billion	33 billion
Wikipedia	24 billion	25 billion
StackExchange	20 billion	27 billion
Total Tokens	1.2 trillion	1.25 trillion

ChatPDF

www.chatpdf.com



ChatPDF & Jenni AI: Write your next paper with AI - Unlimited access to the [Jenni AI Writer](#) with ChatPDF Plus! [i](#)

Chat with any PDF

Join millions of students, researchers and professionals to instantly answer questions and understand research with AI



Drop PDF here

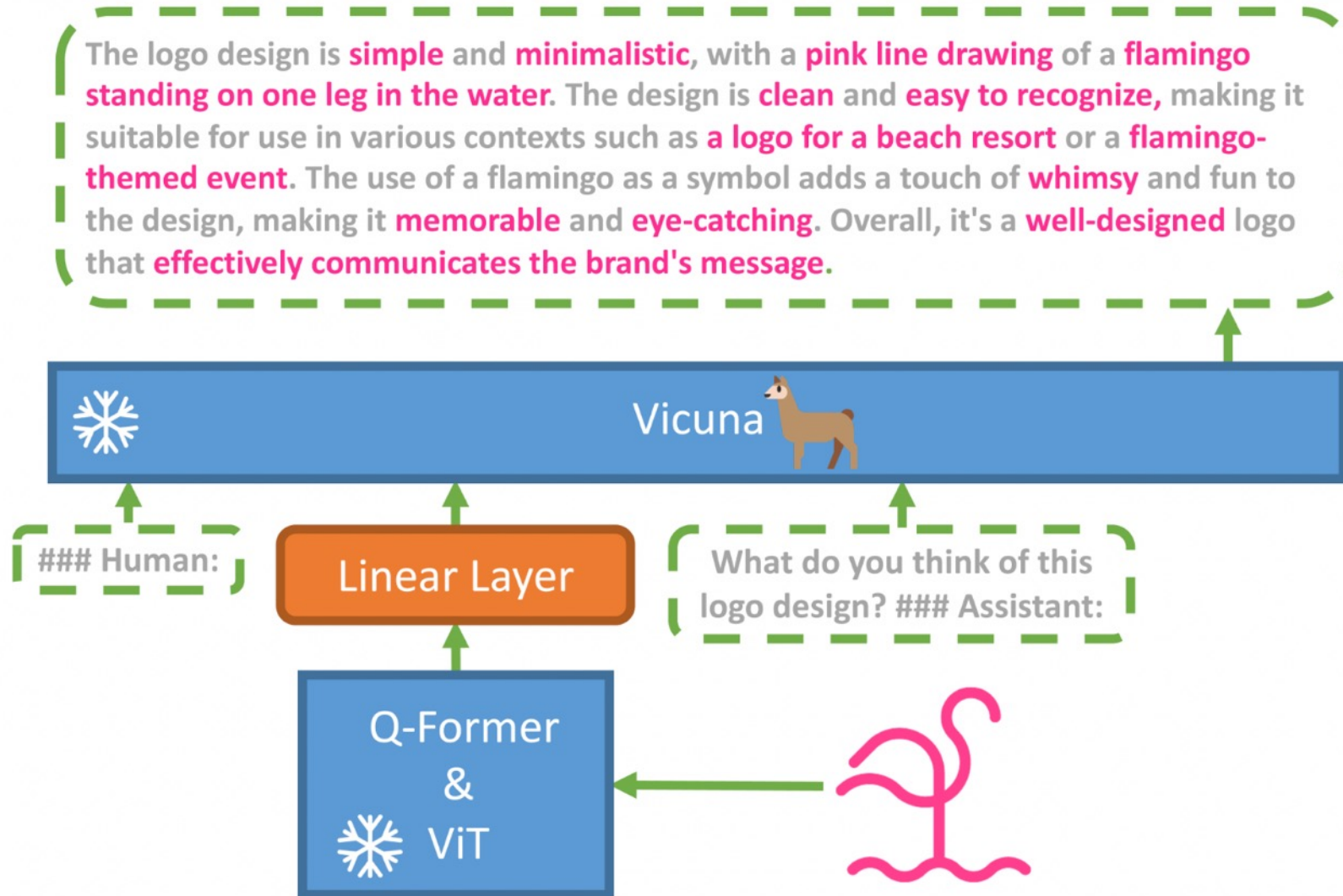
[Browse my Computer](#)

[From URL](#)

<https://www.chatpdf.com/>

MiniGPT-4:

Enhancing Vision-language Understanding with Advanced Large Language Models



LLaVA: Large Language and Vision Assistant

llava-vl.github.io



LLaVA: Large Language and Vision Assistant

[\[Project Page\]](#) [\[Paper\]](#) [\[Code\]](#) [\[Model\]](#)

LLaVA-13B-v0

Image

Drop Image Here

- or -

Click to Upload

Preprocess

When the image is not square, you may try adjusting this .

☒ Crop

☐ Resize

☐ Pad

Parameters

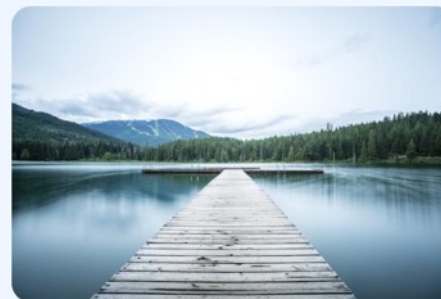
Terms of use

By using this service, users are required to agree to the following terms:

The service is a research preview intended for non-commercial use only. It only provides limited safety measures and may generate offensive content. It must not be used for any illegal, harmful, violent, racist, or sexual purposes. The service may collect user dialogue data for future research.

LLaVA Chatbot

What are the things I should be cautious about when I visit here?



When visiting this location, one should be cautious about the weather conditions and the state of the wooden pier. In the image, the wooden pier is extended over a large, still lake under cloudy skies, which could suggest the possibility of sudden changes in weather, such as wind, rain, or even high waves. Visitors should pay attention to weather forecasts and remain vigilant while on the pier to ensure their safety. Additionally, they should be mindful of the pier's condition and the surrounding water, as wood can be slippery, especially when wet or in windy conditions. It is crucial to maintain proper footing and follow any posted safety guidelines, as the quiet waters of the lake can be deceptively

Enter text and press ENTER

Submit

👍 Upvote

👎 Downvote

⚠️ Flag

🔄 Regenerate

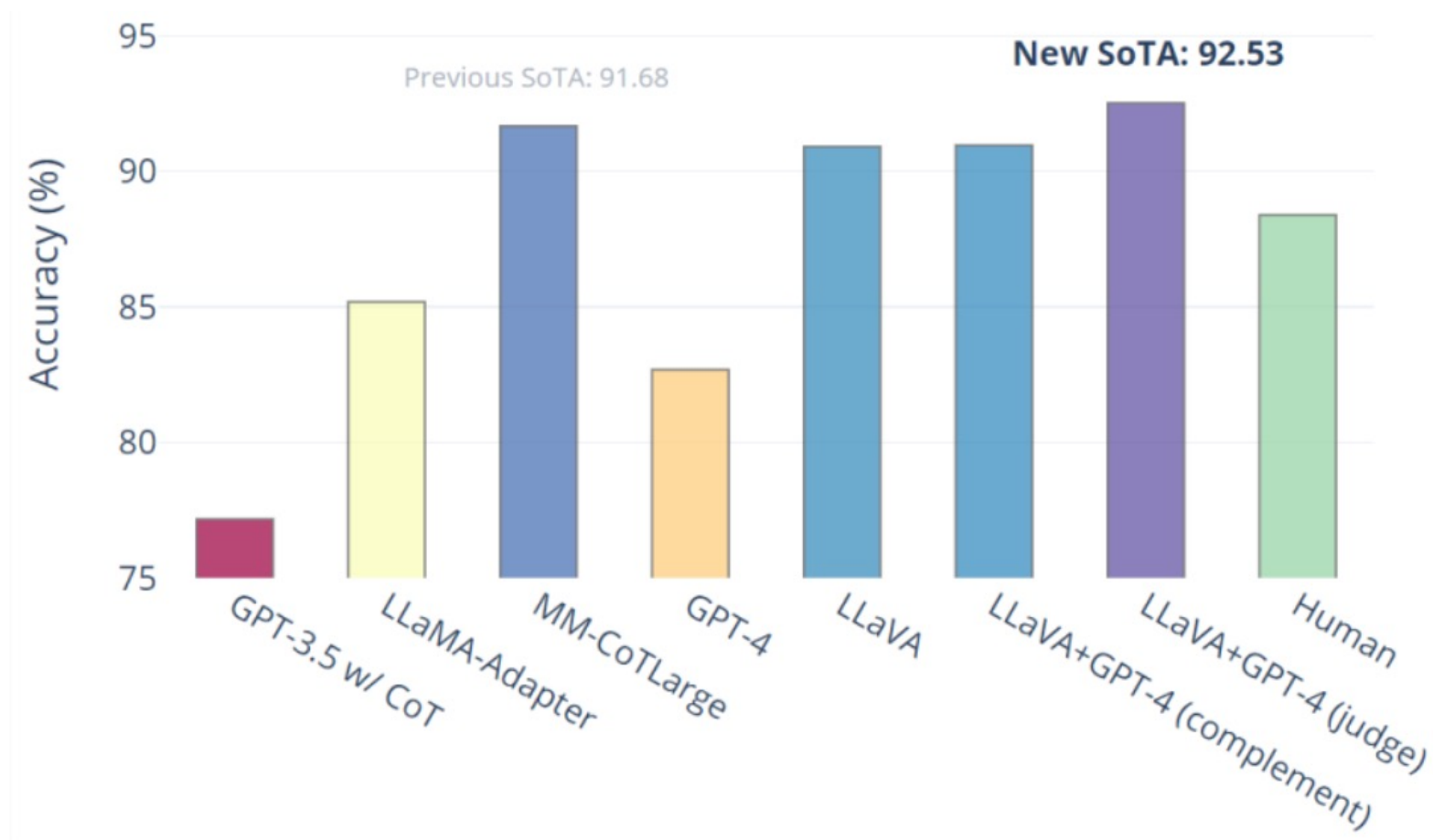
Source: <https://llava-vl.github.io/>

Visual Instruction Tuning

LLaVA: Large Language and Vision Assistant

University of Wisconsin-Madison, Microsoft Research, Columbia University

Science QA:
New SoTA with
the synergy of
LLaVA with
GPT-4

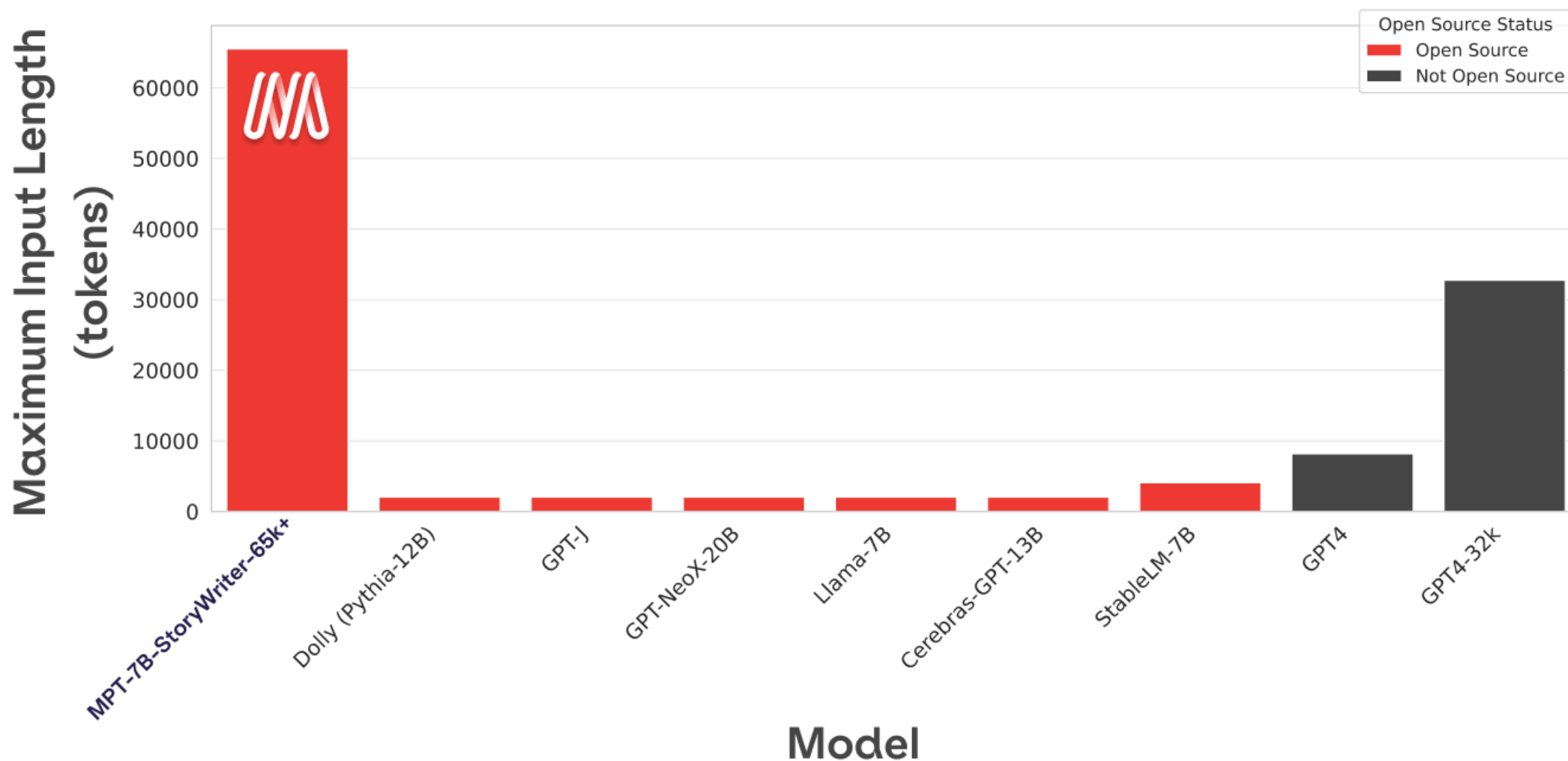


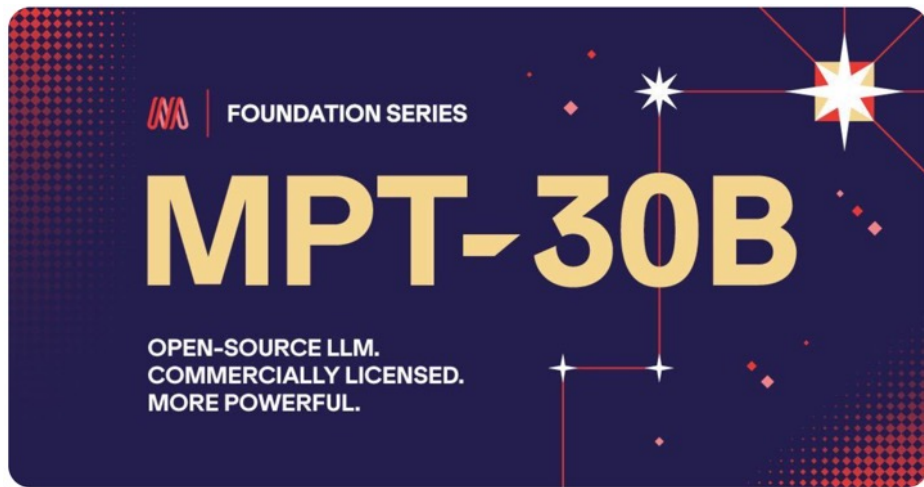
LLaVA represents a novel end-to-end trained large multimodal model that combines a vision encoder and Vicuna for general-purpose visual and language understanding, achieving impressive chat capabilities mimicking spirits of the multimodal GPT-4 and setting a new state-of-the-art accuracy on Science QA.

Source: <https://llava-vl.github.io/>

MPT-7B-StoryWriter-65k+

Maximum Input Lengths of Different LLMs





MPT-30B, MPT-7B LLaMa-30B, LLaMa-7B

Model Purpose	Model Series	Model	Sequence Length	Accuracy (Pass@1)	Externally Reported Pass@1 & [Source]
General Purpose	MPT	MPT-30B	1024	25.00%	N/A
		MPT-30B Chat	1024	37.20%	N/A
		MPT-30B Instruct	1024	26.20%	N/A
		MPT-7B	1024	15.90%	N/A
		MPT-7B Instruct	1024	16.50%	N/A
	LLaMa	LLaMa-7B	1024	10.10%	10.5% [1]
		LLaMa-13B	1024	16.50%	15.8% [1]
		LLaMa-30B	1024	20.10%	21.7% [1]
	Falcon	Falcon-40B	1024	1.2%* (did not generate code)	N/A
		Falcon-40B Instruct	1024	0.6%* (did not generate code)	18.9% [2]

Meta Llama-2 70B: Best Open Source and Commercial LLM (Llama-2, Falcon, MPT)

MODEL SIZE (PARAMETERS)	PRETRAINED	FINE-TUNED FOR CHAT USE CASES
7B	Model architecture:	Data collection for helpfulness and safety:
13B	Pretraining Tokens: 2 Trillion	Supervised fine-tuning: Over 100,000
70B	Context Length: 4096	Human Preferences: Over 1,000,000

Llama 2 pretrained models are trained on 2 trillion tokens, and have double the context length than Llama 1. Its fine-tuned models have been trained on over 1 million human annotations.

**Meta
Llama-2 70B:
Best
Open Source
and
Commercial
LLM
(Llama-2,
Falcon, MPT)**

Benchmark (Higher is better)	MPT (7B)	Falcon (7B)	Llama-2 (7B)	Llama-2 (13B)	MPT (30B)	Falcon (40B)	Llama-1 (65B)	Llama-2 (70B)
MMLU	26.8	26.2	45.3	54.8	46.9	55.4	63.4	68.9
TriviaQA	59.6	56.8	68.9	77.2	71.3	78.6	84.5	85.0
Natural Questions	17.8	18.1	22.7	28.0	23.0	29.5	31.0	33.0
GSM8K	6.8	6.8	14.6	28.7	15.2	19.6	50.9	56.8
HumanEval	18.3	N/A	12.8	18.3	25.0	N/A	23.7	29.9
AGIEval (English tasks only)	23.5	21.2	29.3	39.1	33.8	37.0	47.6	54.2
BoolQ	75.0	67.5	77.4	81.7	79.0	83.1	85.3	85.0

Llama 2 outperforms other open source language models on many external benchmarks, including reasoning, coding, proficiency, and knowledge tests.

Source: <https://ai.meta.com/llama/>

Llama-2: Comparison to closed-source models (GPT-3.5, GPT-4, PaLM) on academic benchmarks

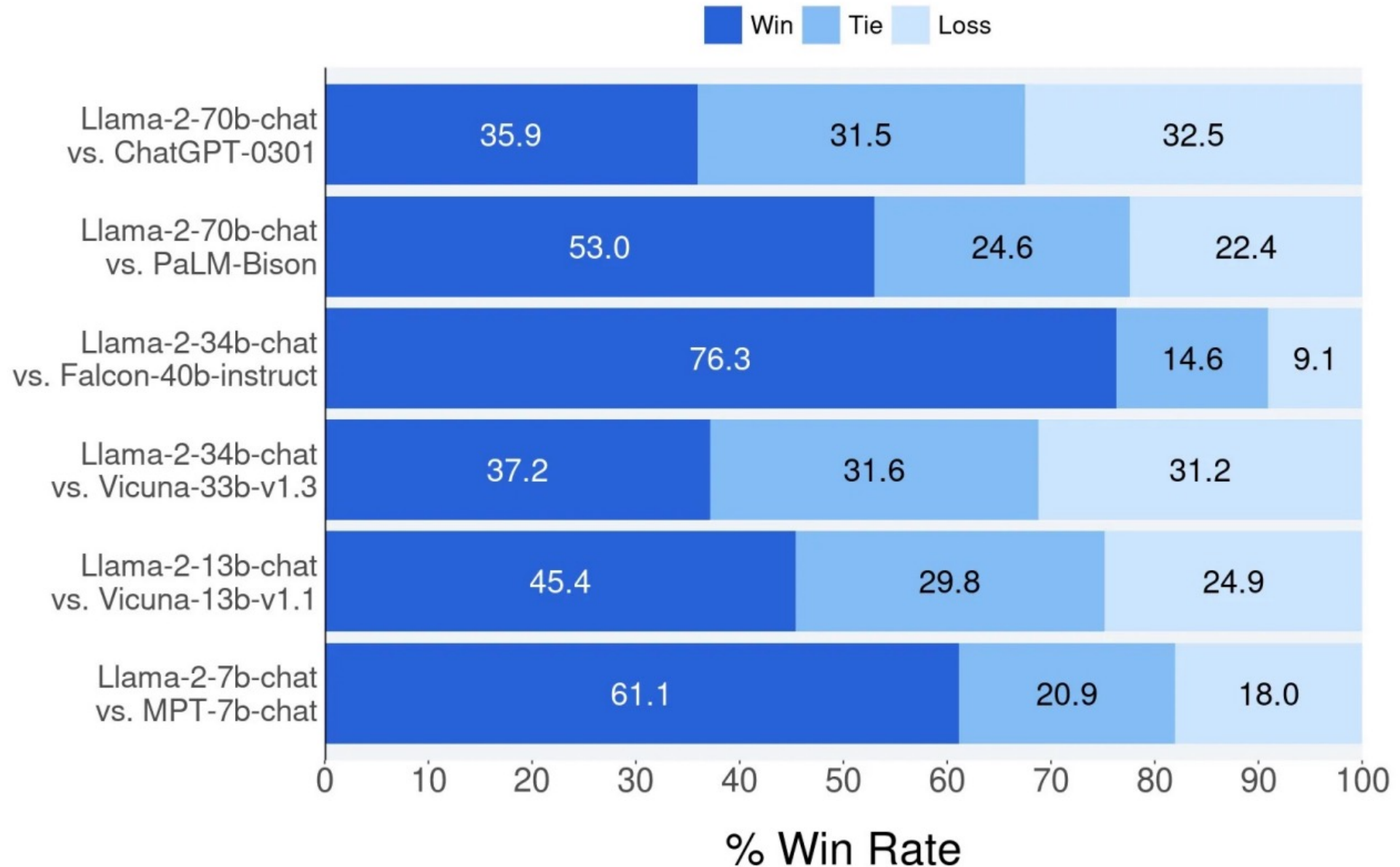
Benchmark (shots)	GPT-3.5	GPT-4	PaLM	PaLM-2-L	LLAMA 2
MMLU (5-shot)	70.0	86.4	69.3	78.3	68.9
TriviaQA (1-shot)	–	–	81.4	86.1	85.0
Natural Questions (1-shot)	–	–	29.3	37.5	33.0
GSM8K (8-shot)	57.1	92.0	56.5	80.7	56.8
HumanEval (0-shot)	48.1	67.0	26.2	–	29.9
BIG-Bench Hard (3-shot)	–	–	52.3	65.7	51.2

Results for GPT-3.5 and GPT-4 are from OpenAI (2023).

Results for the PaLM model are from Chowdhery et al. (2022).

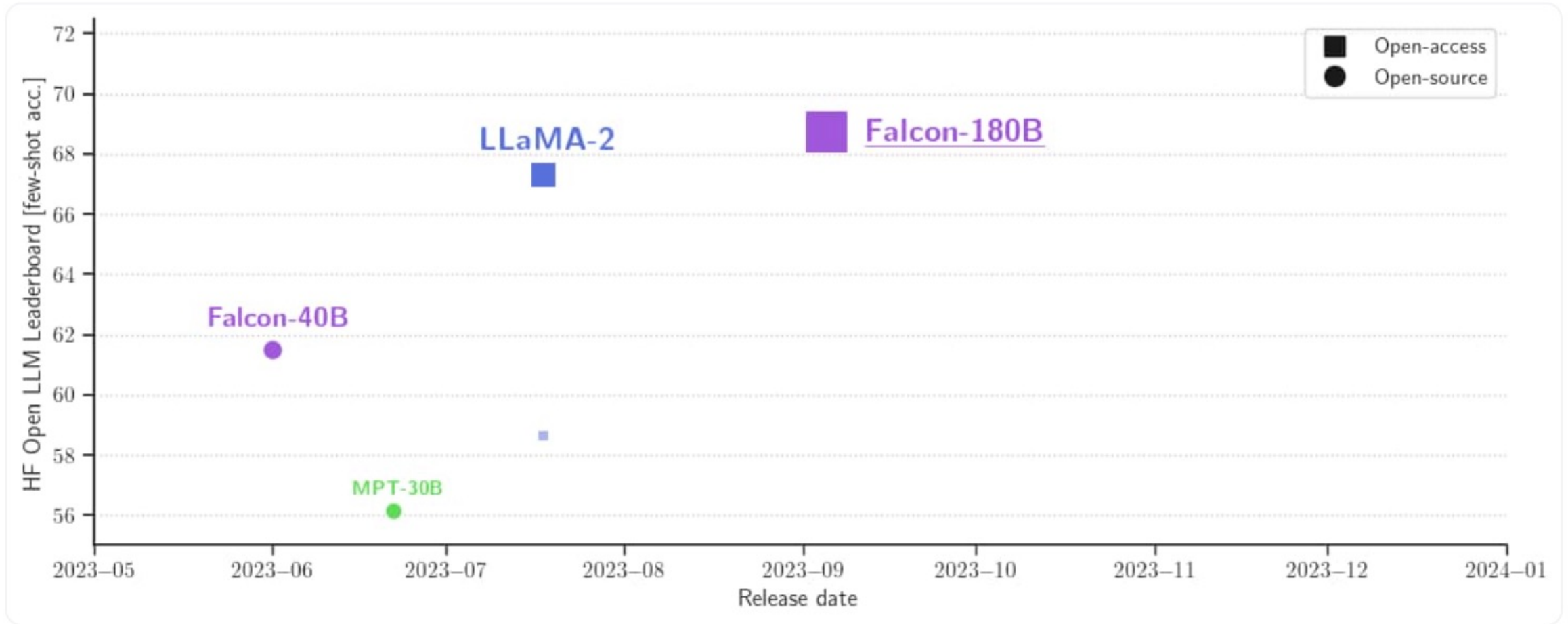
Results for the PaLM-2-L are from Anil et al. (2023).

Llama-2 Chat: Helpfulness Human Evaluation





Falcon 180B





Falcon 180B, LLaMA 65B, MPT 30B

Model	Size	Leaderboard score	Commercial use or license	Pretraining length
Falcon	180B	68.74	🟡	3,500B
Llama 2	70B	67.35	🟡	2,000B
LLaMA	65B	64.23	🔴	1,400B
Falcon	40B	61.48	🟢	1,000B
MPT	30B	56.15	🟢	1,000B



Falcon 180B

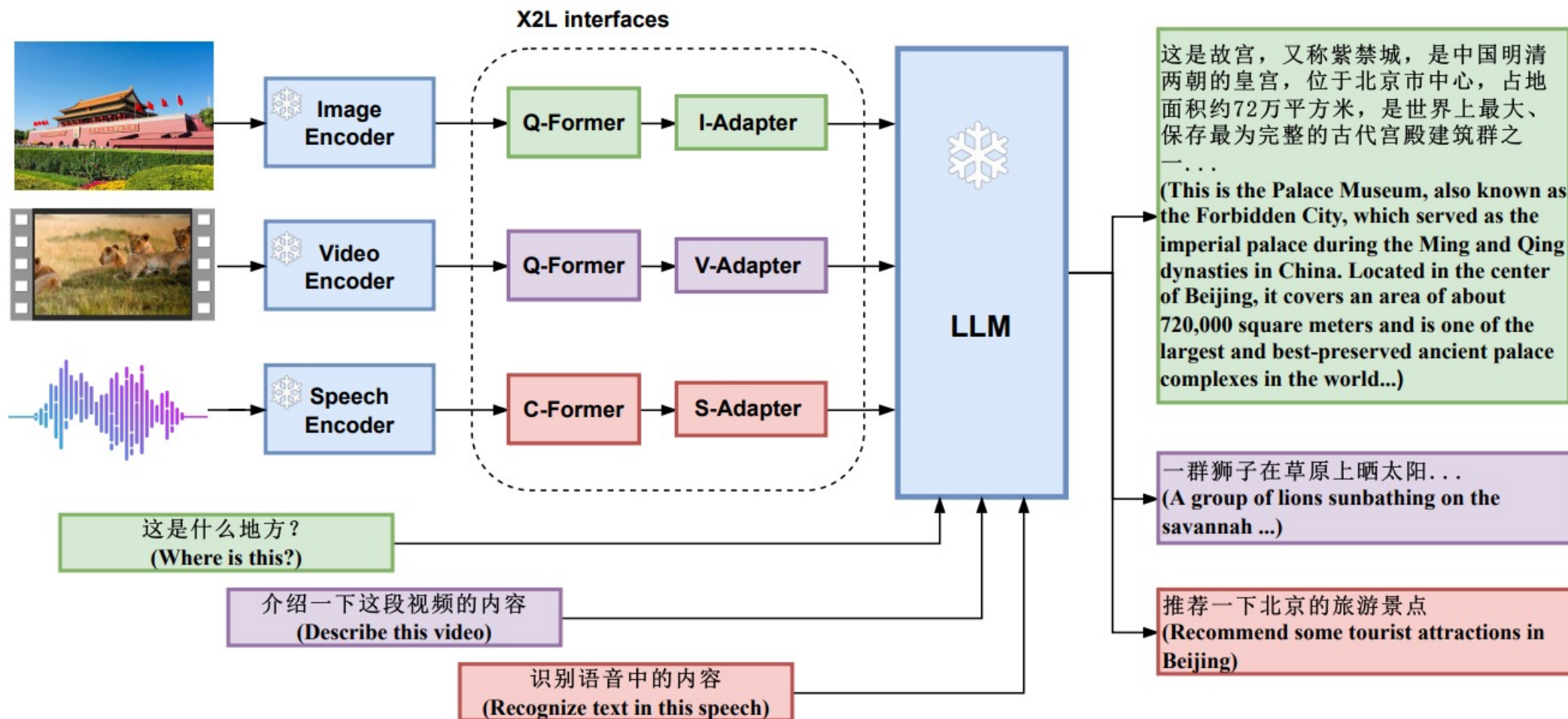
Hardware requirements

NVIDIA A100 80 GB:
\$16,135

	Type	Kind	Memory	Example
Falcon 180B	Training	Full fine-tuning	5120GB	8x 8x A100 80GB
Falcon 180B	Training	LoRA with ZeRO-3	1280GB	2x 8x A100 80GB
Falcon 180B	Training	QLoRA	160GB	2x A100 80GB
Falcon 180B	Inference	BF16/FP16	640GB	8x A100 80GB
Falcon 180B	Inference	GPTQ/int4	320GB	8x A100 40GB

X-LLM:

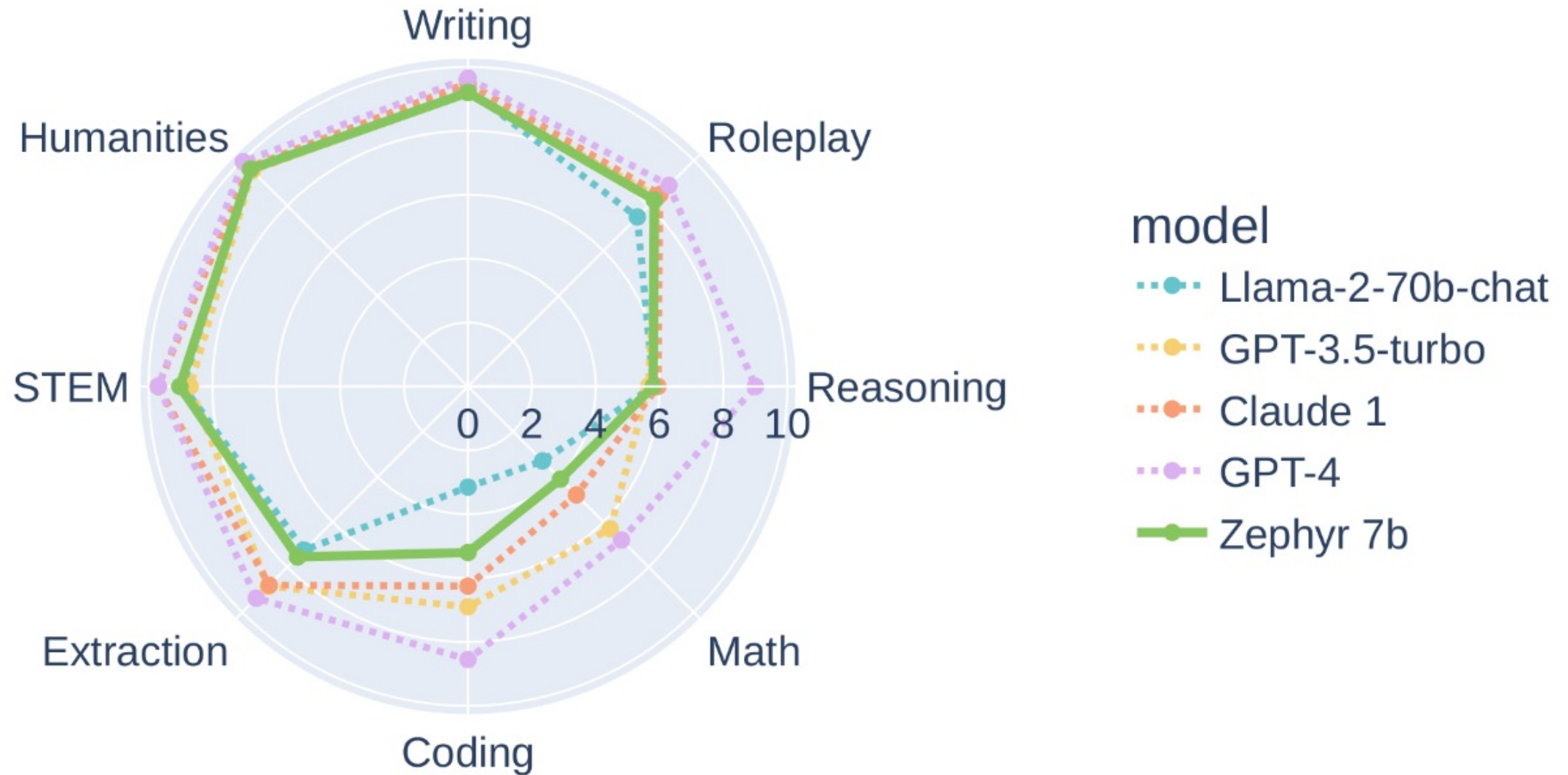
Bootstrapping Advanced Large Language Models by Treating Multi-Modalities as Foreign Languages



Zephyr-7B- β , Llama2-Chat-70B, GPT-4

Model	Size	Alignment	MT-Bench (score)	AlpacaEval (win rate %)
StableLM-Tuned- α	7B	dSFT	2.75	-
MPT-Chat	7B	dSFT	5.42	-
Xwin-LMv0.1	7B	dPPO	6.19	87.83
Mistral-Instructv0.1	7B	-	6.84	-
Zephyr-7b- α	7B	dDPO	6.88	-
Zephyr-7b-β 🌂	7B	dDPO	7.34	90.60
Falcon-Instruct	40B	dSFT	5.17	45.71
Guanaco	65B	SFT	6.41	71.80
Llama2-Chat	70B	RLHF	6.86	92.66
Vicuna v1.3	33B	dSFT	7.12	88.99
WizardLM v1.0	70B	dSFT	7.71	-
Xwin-LM v0.1	70B	dPPO	-	95.57
GPT-3.5-turbo	-	RLHF	7.94	89.37
Claude 2	-	RLHF	8.06	91.36
GPT-4	-	RLHF	8.99	95.28

Zephyr-7B- β , Llama2-Chat-70B, GPT-4



Zephyr: Direct Distillation of LM Alignment

distilled supervised fine-tuning (dSFT)

AI Feedback (AIF)

distilled Direct Preference Optimization (dDPO)

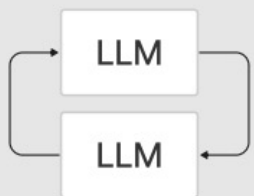
Step 1 - dSFT

Generate multi-turn AI dialogues

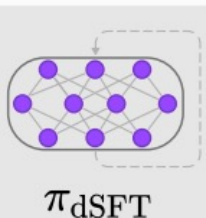
Prompt sampled from dataset of prompts.

Create a scenario for a game about space exploration

LLM simulates multi-turn user-assistant interactions.



Dialogues are used for supervised fine-tuning.



Step 2 - AIF

Response generation and AI ranking

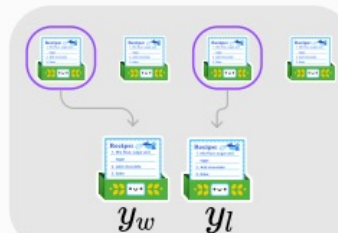
Prompt sampled from dataset of prompts.

Describe how to make chocolate brownies

4 different language models generate responses.



GPT-4 ranks the responses.



Step 3 - dDPO

Distillation of AI preferences

Prompt sampled from dataset of prompts.

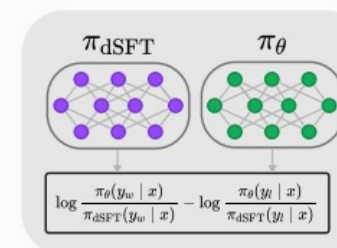
Describe how to make chocolate brownies

Best and another random response are selected.



y_w y_l

Direct Preference Optimization



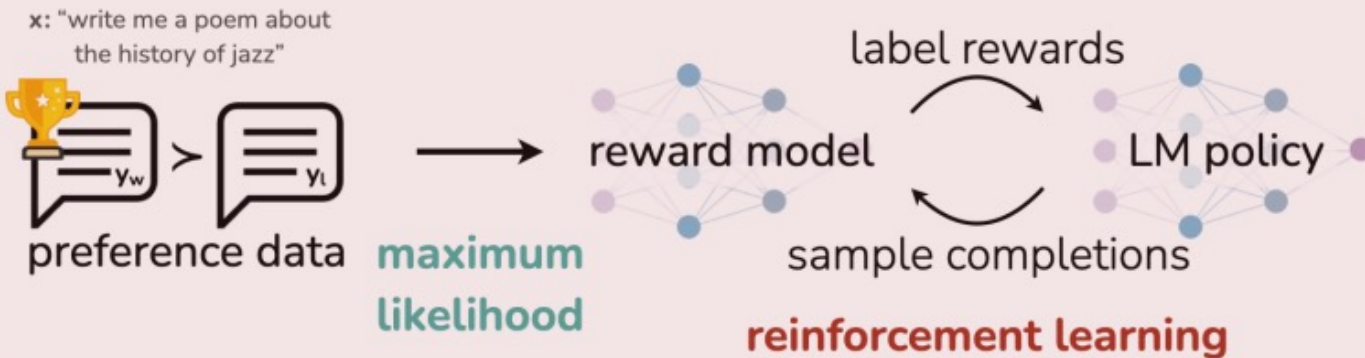
Zephyr: Direct Distillation of LM Alignment

The three steps of **Zephyr**:

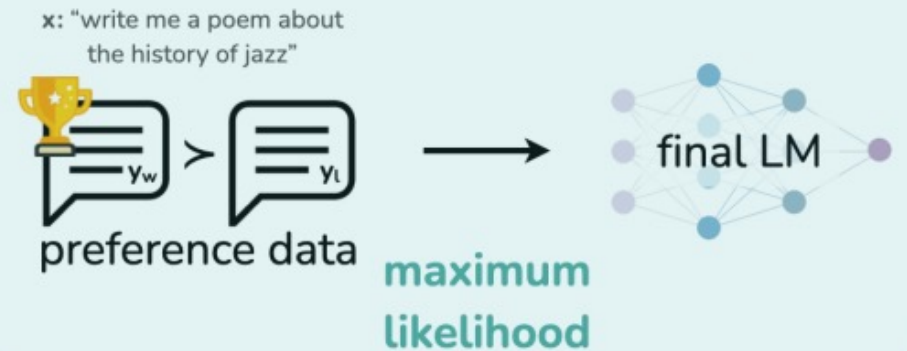
- (1) large scale, self-instruct-style dataset construction (UltraChat), followed by **distilled supervised fine-tuning (dSFT)**,
- (2) **AI Feedback (AIF)** collection via an ensemble of chat model completions, followed by scoring by GPT-4 (UltraFeedback) and binarization into preferences, and
- (3) **distilled direct preference optimization (dDPO)** of the dSFT model utilizing the feedback data.

DPO optimizes for human preferences while avoiding reinforcement learning

Reinforcement Learning from Human Feedback (RLHF)



Direct Preference Optimization (DPO)



Direct Preference Optimization (DPO)

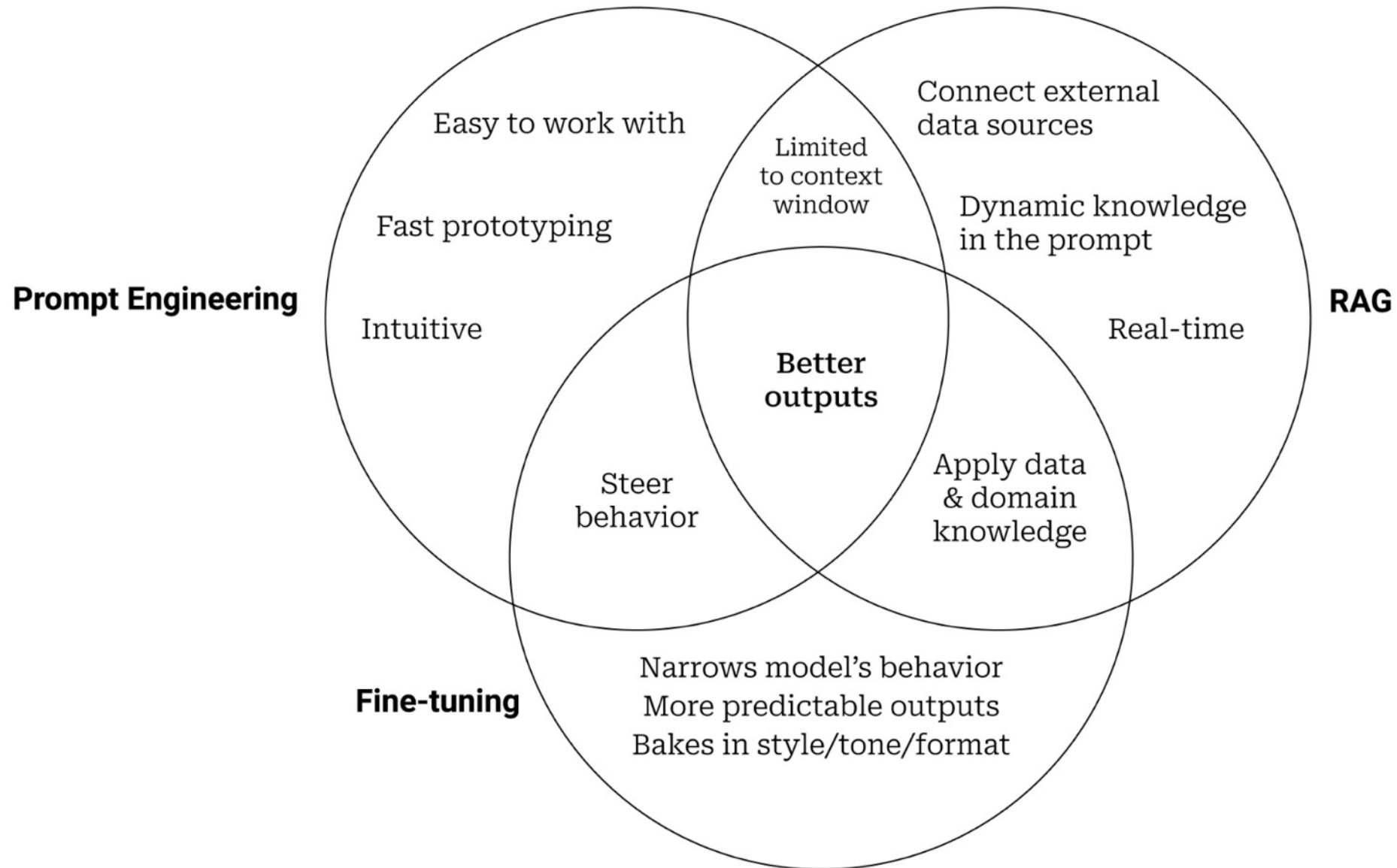
Direct Preference Optimization (DPO)

x : "write me a poem about
the history of jazz"

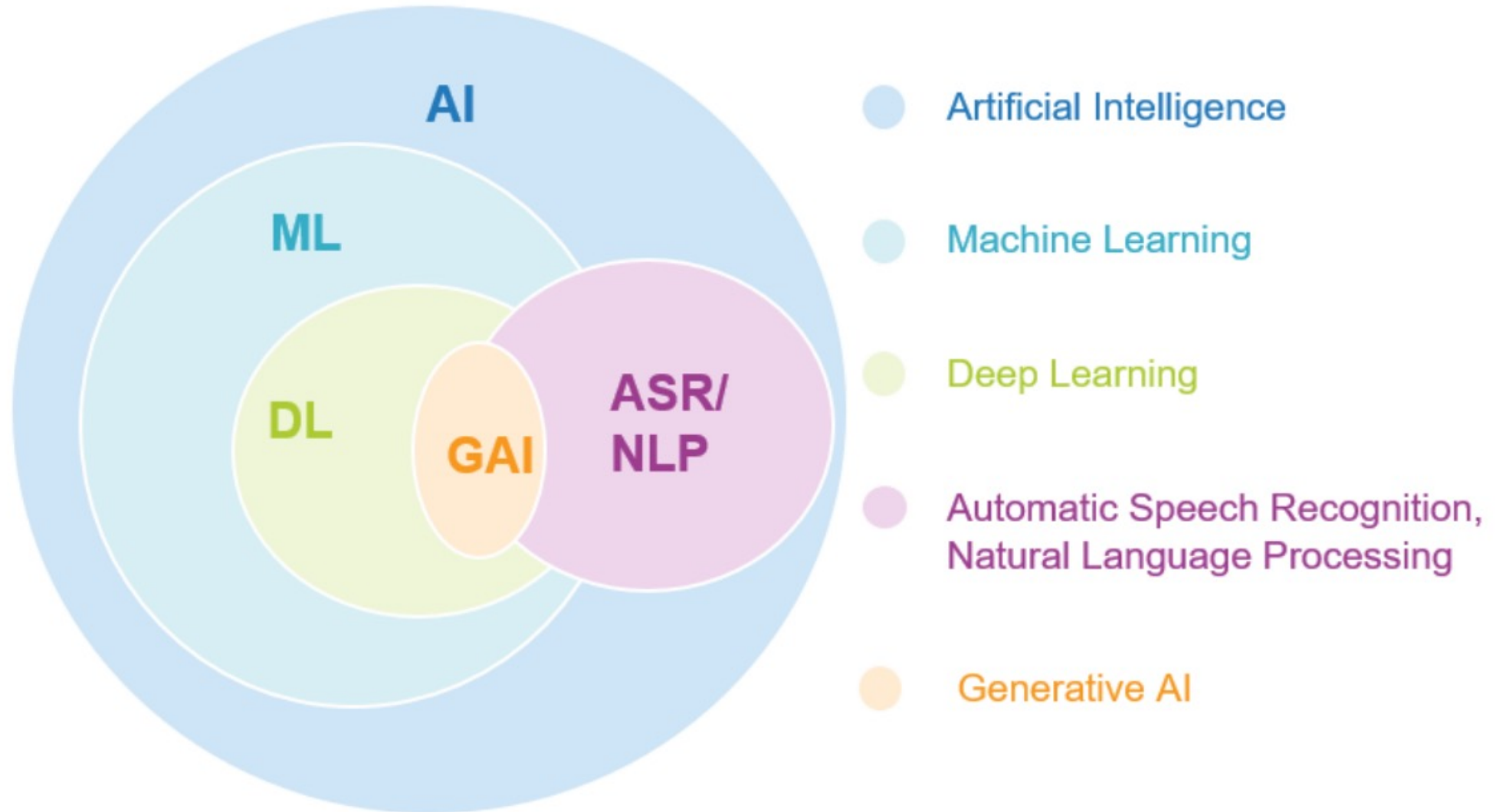


**maximum
likelihood**

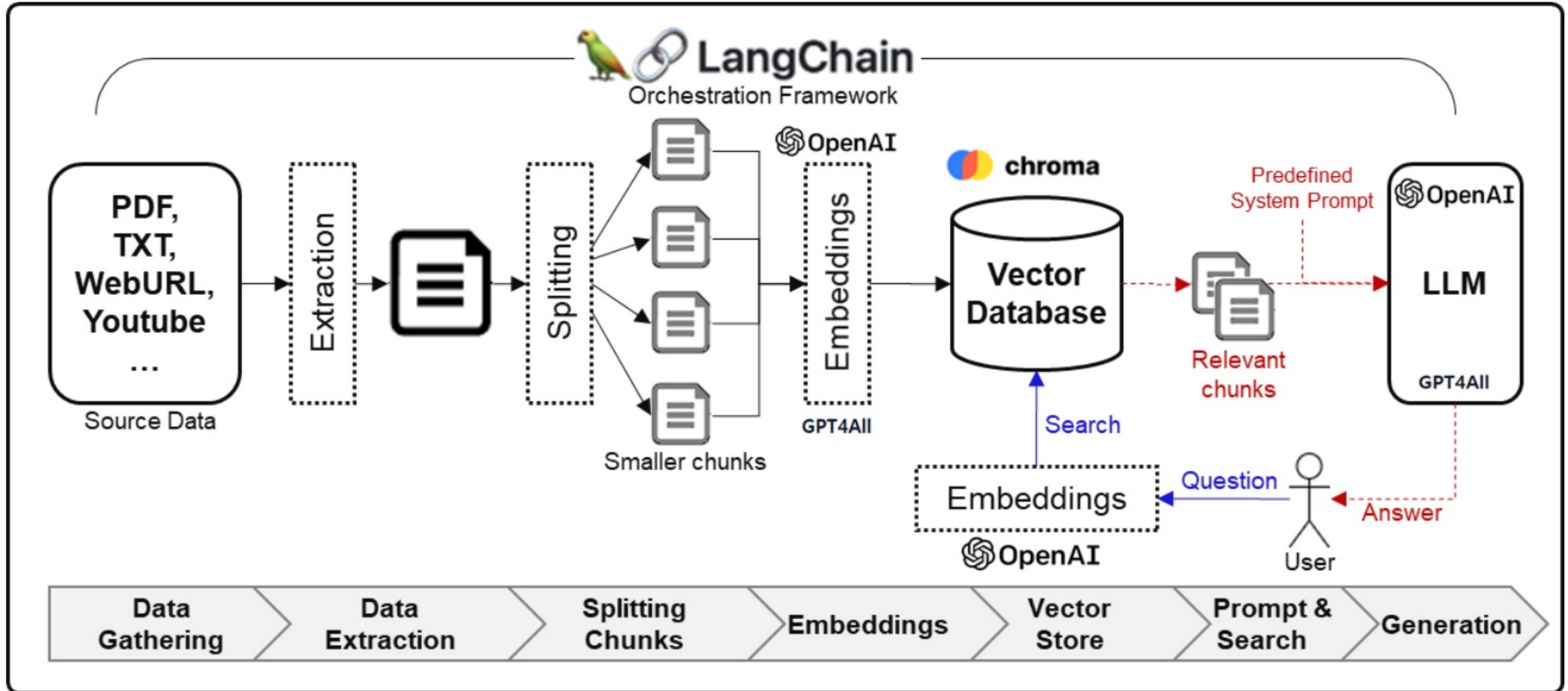
Prompt Engineering, Fine-tuning, and RAG



Generative AI



Framework for Implementing Generative AI Services using RAG Model



Factuality Enhancement of Large Language Models (LLMs)

Factuality Enhancement of Large Language Models

Factuality Enhancement of Large Language Models						
Standalone LLMs		Retrieval Augmented Generation			Domain Factuality Enhancement	
Supervised Finetuning		Normal RAG Setting			Domain enhancement techniques	
Continual SFT	Model Editing	Post-editing			Continue-SFT	Continue Pretraining
Pretraining-based		Interactive Retrieval			Train From Scratch	External Knowledge
Initial Pretraining	Continual Pretraining	CoT-based Retrieval	Agent-based Retrieval		Domains	
Prompt Engineering		Retrieval Adaption			Healthcare and medicine	Finance and Ecommerce
Multi-Agent		Prompt-based	SFT-based	RLHF-based	Legal/Law	Geoscience and Environment
Inference and Decoding		Retrieval on External Memory			Education	Food Industry
		Retrieval on KGs/Databases				Home Renovation

Stable Diffusion



Hugging Face

Search models, datasets, users...



Models



Datasets



Spaces



Docs



Solutions

Pricing



Spaces: stabilityai/

stable-diffusion



like 1.89k



Running



App



Files



Community 241



Linked Models

Stable Diffusion Demo

Stable Diffusion is a state of the art text-to-image model that generates images from text.

For faster generation and forthcoming API access you can try [DreamStudio Beta](#)

an insect robot preparing a delicious meal

Generate image



<https://huggingface.co/spaces/stabilityai/stable-diffusion>

Stable Diffusion Colab

woctezuma / **stable-diffusion-colab** Public

Notifications

Fork 7

Star 31

<> Code Issues Pull requests Actions Projects Wiki Security Insights

main

1 branch 0 tags

Go to file

Code



woctezuma README: add a reference for sampler schedules

37bc02d 24 days ago 18 commits



LICENSE

Initial commit

27 days ago



README.md

README: add a reference for sampler schedules

24 days ago



stable_diffusion.ipynb

Allow to choose the scheduler

25 days ago

README.md

Stable-Diffusion-Colab

The goal of this repository is to provide a Colab notebook to run the text-to-image "Stable Diffusion" model [1].

Usage

- Run `stable_diffusion.ipynb` . [Open in Colab](#)

About

Colab notebook to run Stable Diffusion.

github.com/CompVis/stable-diffusion

deep-learning colab image-generation

text-to-image diffusion text2image

colaboratory google-colab

colab-notebook google-colaboratory

google-colab-notebook

text-to-image-synthesis huggingface

diffusion-models

text-to-image-generation latent-diffusion

stable-diffusion huggingface-diffusers

diffusers stable-diffusion-diffusers

Readme

MIT license

31 stars

2 watching

<https://github.com/woctezuma/stable-diffusion-colab>

Stable Diffusion Reimagine



Clipdrop
by stability.ai

► Stable diffusion Reimagine

Apps ▾

API

Blog

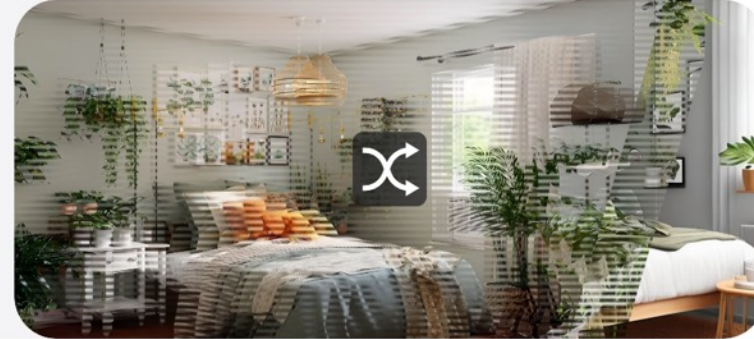
Pricing

Sign-in / Sign-up



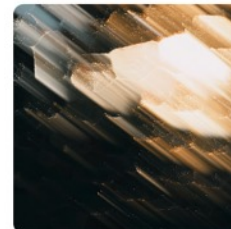
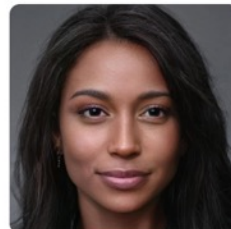
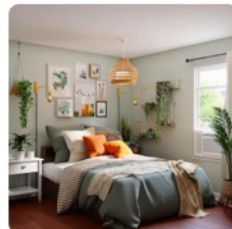
Stable diffusion reimagine

Create multiple variations from a single image.



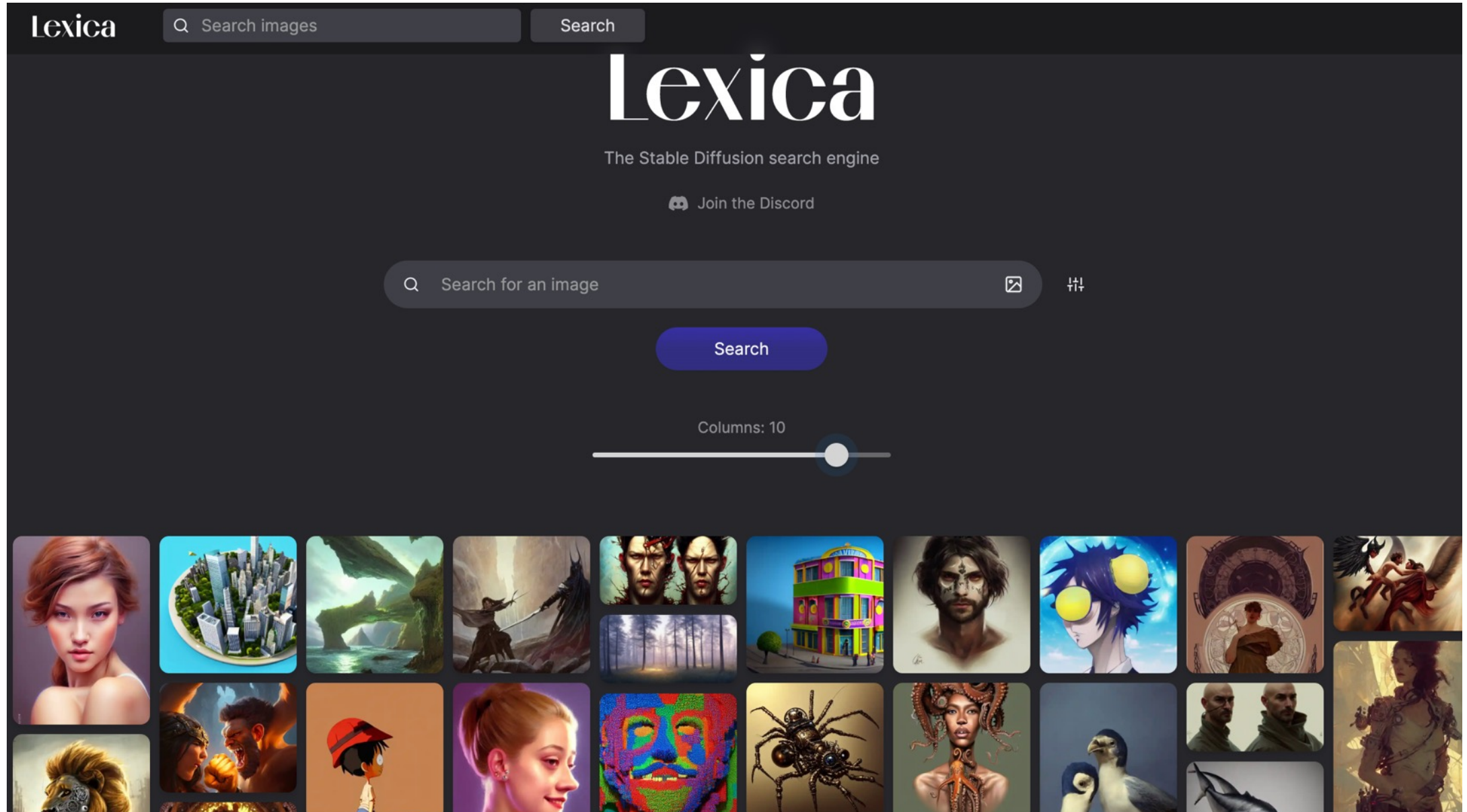
Click, paste, or drop a file here to start.

↓ Or click on an example below



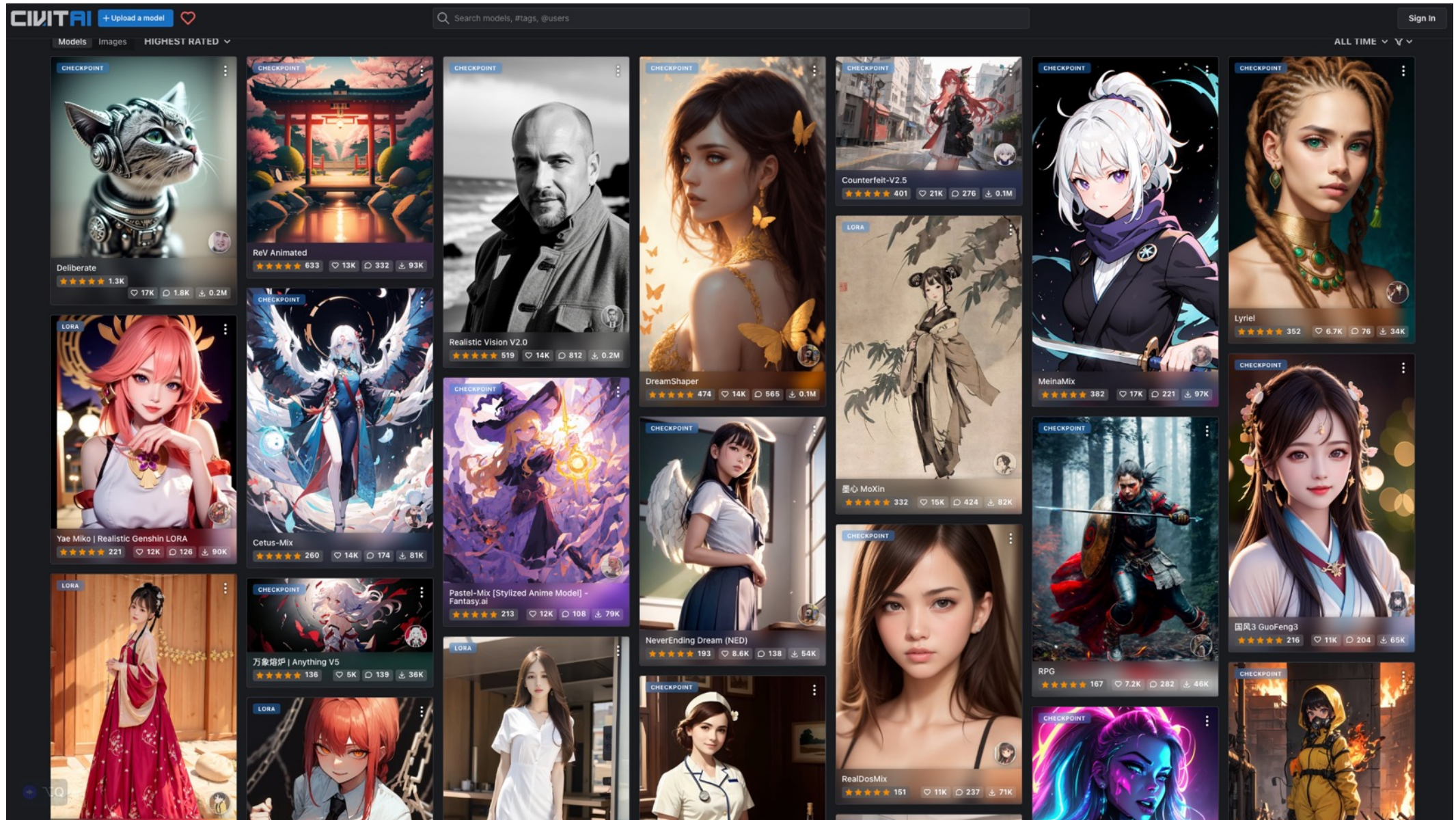
<https://clipdrop.co/stable-diffusion-reimagine>

Lexica Art: Search Stable Diffusion images and prompts



<https://lexica.art/>

Civitai: Stable Diffusion AI Art Models



<https://civitai.com/>

D-ID Text to Video

[Products](#)[Technology](#)[Ethics](#)[Pricing](#)[Company](#)[Start Free Trial](#)[Log in](#)

Turn Text To Video In 30 Seconds

Save time and money and enrich your content with engaging videos. Try it today!

[Start Free Trial](#)

<https://www.d-id.com/text-to-video/>

Synthesia: #1 AI Video Generation



Features ▾

Use cases ▾

Pricing

Resources ▾

Company ▾

Log in

Create account

#1 AI VIDEO GENERATION PLATFORM ⓘ

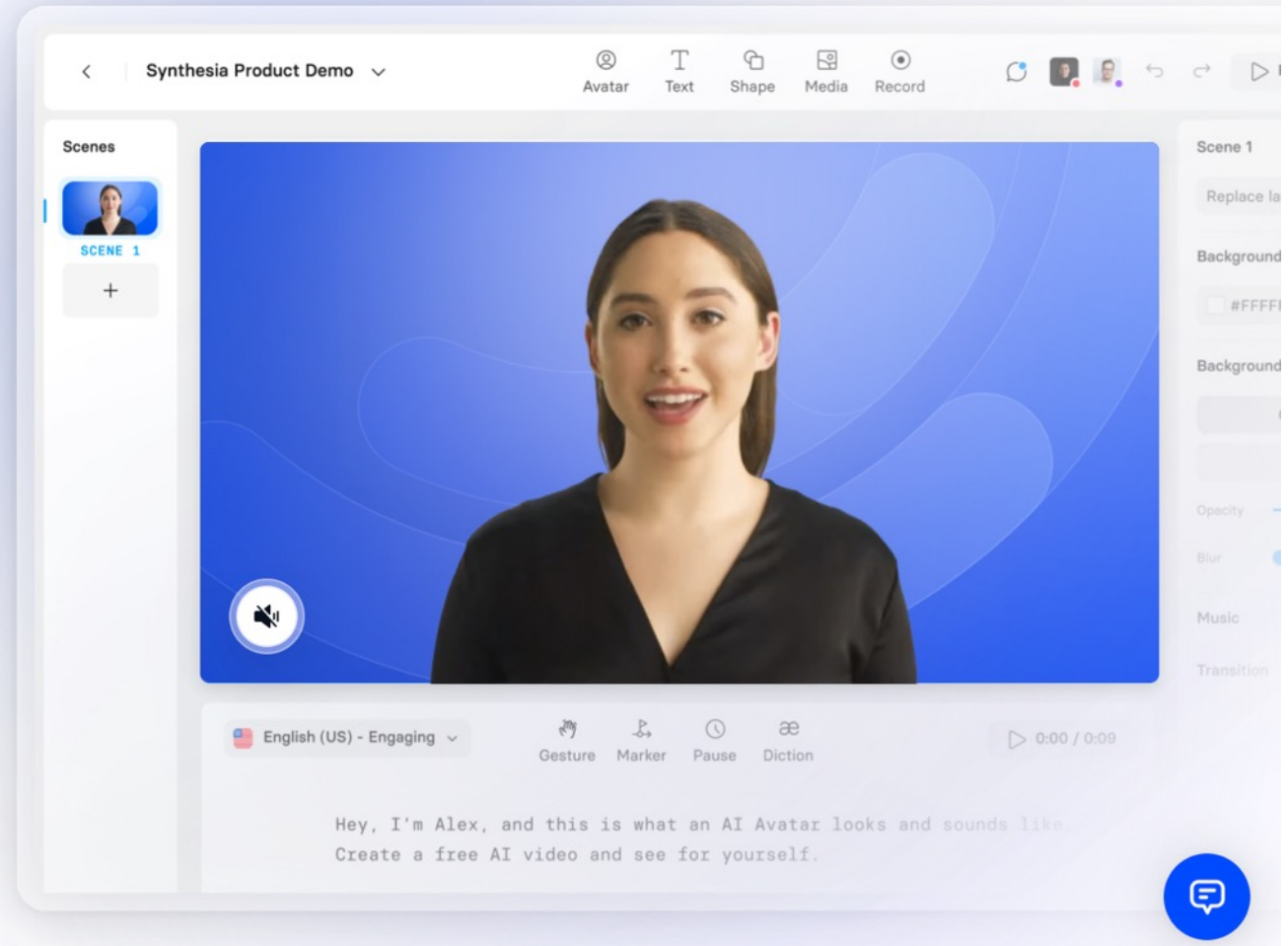
Turn your text into videos in minutes

- Get natural sounding AI voices in 120+ languages
- Make your videos more engaging with 140+ AI Avatars
- Edit as simply as a slide-deck, no experience required

Create a free AI video

▶ Watch 2 min demo

No credit card required.



<https://www.synthesia.io/>

Speechify: #1 AI Voice Over Generator

[Text to Speech](#)[AI Voice Studio](#)[Products](#)[Teams](#)[Edu](#)[About](#)[Log in](#)[Talk to Sales](#)[Try for free](#)

The #1 AI Voice Over Generator

Natural sounding, human-quality voice generator for all your content. Try our AI voice today, for free!



Type text here



Introducing the ultimate voiceover tool for professionals and amateurs – a powerful and easy-to-use software that lets you easily create high-quality voiceovers.



Narrate text, videos, explainers, slides, books – anything – in any style.



Our voiceover product is perfect for businesses, content creators, podcasters, video editors, and anyone else who needs to add professional-quality voiceovers to their projects.



Select Voice

[More voices >](#)

Davis

General



Aria

Chat



Guy

Friendly



Clone My
Voice



Add Pause

Pause for a few seconds to add emphasis.



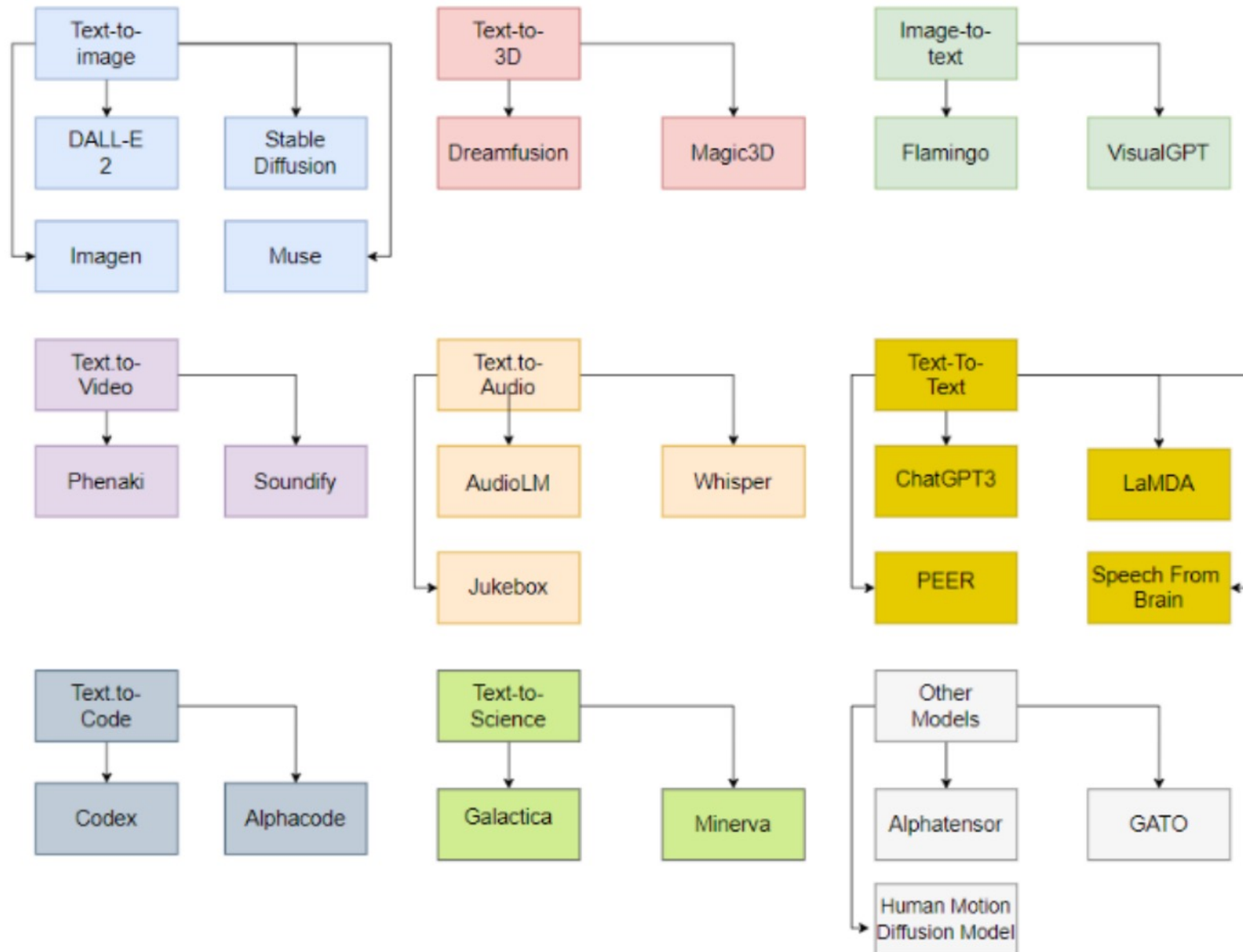
Listen With Music

Change music v



Create an account to access 200+ high-quality voices and Granular controls on the pitch, tone and speed.

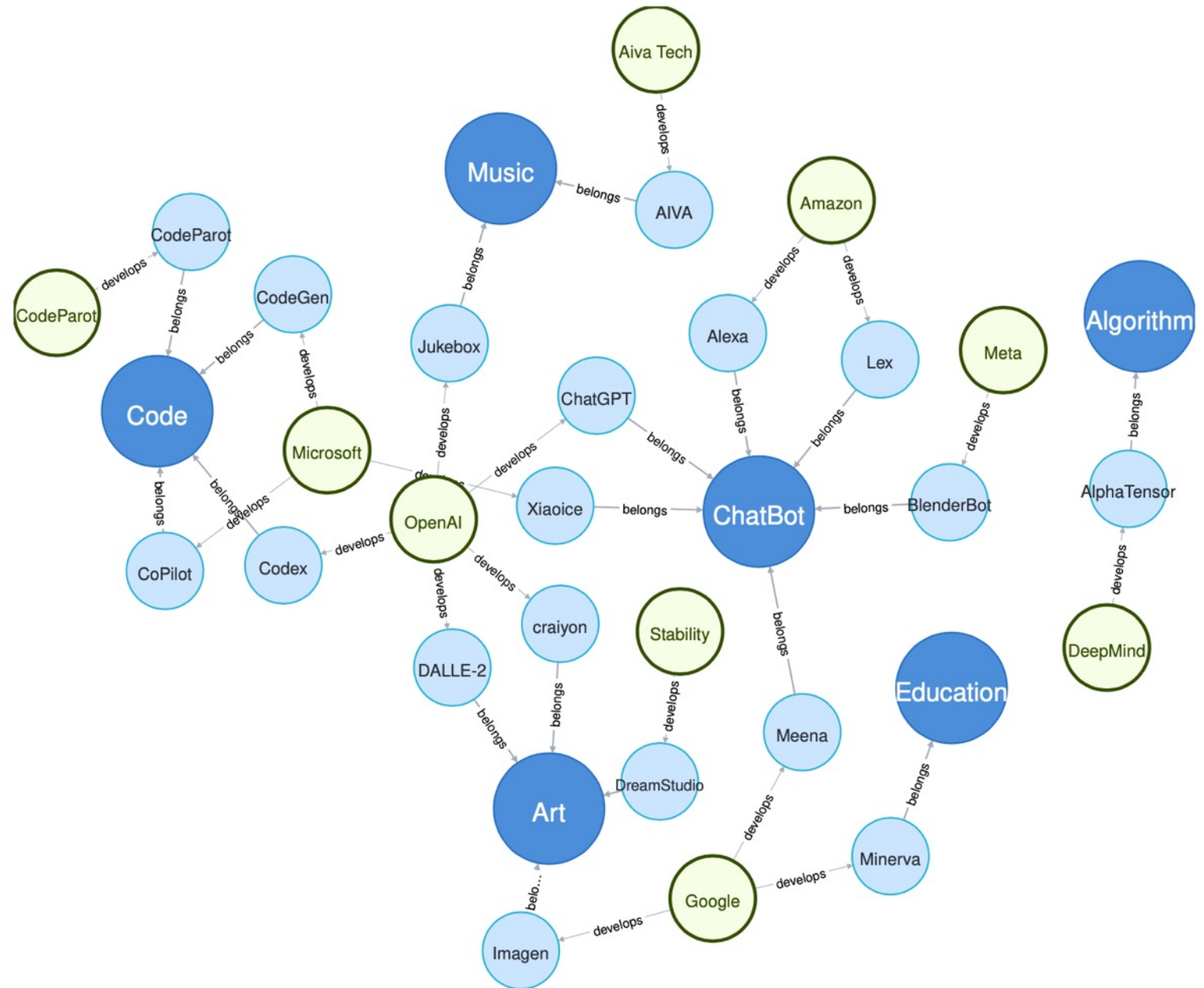
Generative AI Models



**ChatGPT
is not
all you need**

**Attention
is
all you need**

Generative AI Research Areas, Applications and Companies



Applications of Generative AI Models

Application	Platform/Software	Company	Year	Papaer	Link
ChatBot	Xiaoice	Microsoft	2018	[200]	Xiaoice
ChatBot	Meena	Google	2020	[201]	Meena Blog
ChatBot	BlenderBot	Meta	2022	[202]	Blenderbot
ChatBot	ChatGPT	OpenAI	2022	[10]	ChatGPT
ChatBot	Alexa	Amazon	2014	-	Amazon Alexa
ChatBot	Lex	Amazon	2017	-	Amazon Lex
Music	AIVA	Aiva Tech	2016	-	AIVA
Music	Jukebox	OpenAI	2020	[203]	Jukebox
Code	CodeGPT	Microsoft	2021	[204]	CodeGPT
Code	CodeParrot	CodeParrot	2022	[205]	CodeParrot
Code	Codex	OpenAI	2021	[206]	Codex blog
Code	CoPilot	Microsoft	2021	[206]	CoPilot
Art	DALL-E-2	OpenAI	2022	[5]	DALL-E-2 Blog
Art	DreamStudio	Stability	2022	[13]	Dreamstudio
Art	craiyon	OpenAI	2021	[1]	Craiyon
Art	Imagen	Google	2022	[152]	Imagen
Education	Minerva	Google	2022	[207]	Minerva Blog
Algorithm	AlphaTensor	DeepMind	2022	[208]	AlphaTensor

Generative AI for Corporate ESG and Sustainable Development

ESG:

Environmental

Social

Governance

CSR: Corporate Social Responsibility

Sustainable Development Goals (SDGs)



Source: <https://sdgs.un.org/goals>

Sustainable Development Goals (SDGs) and 5P

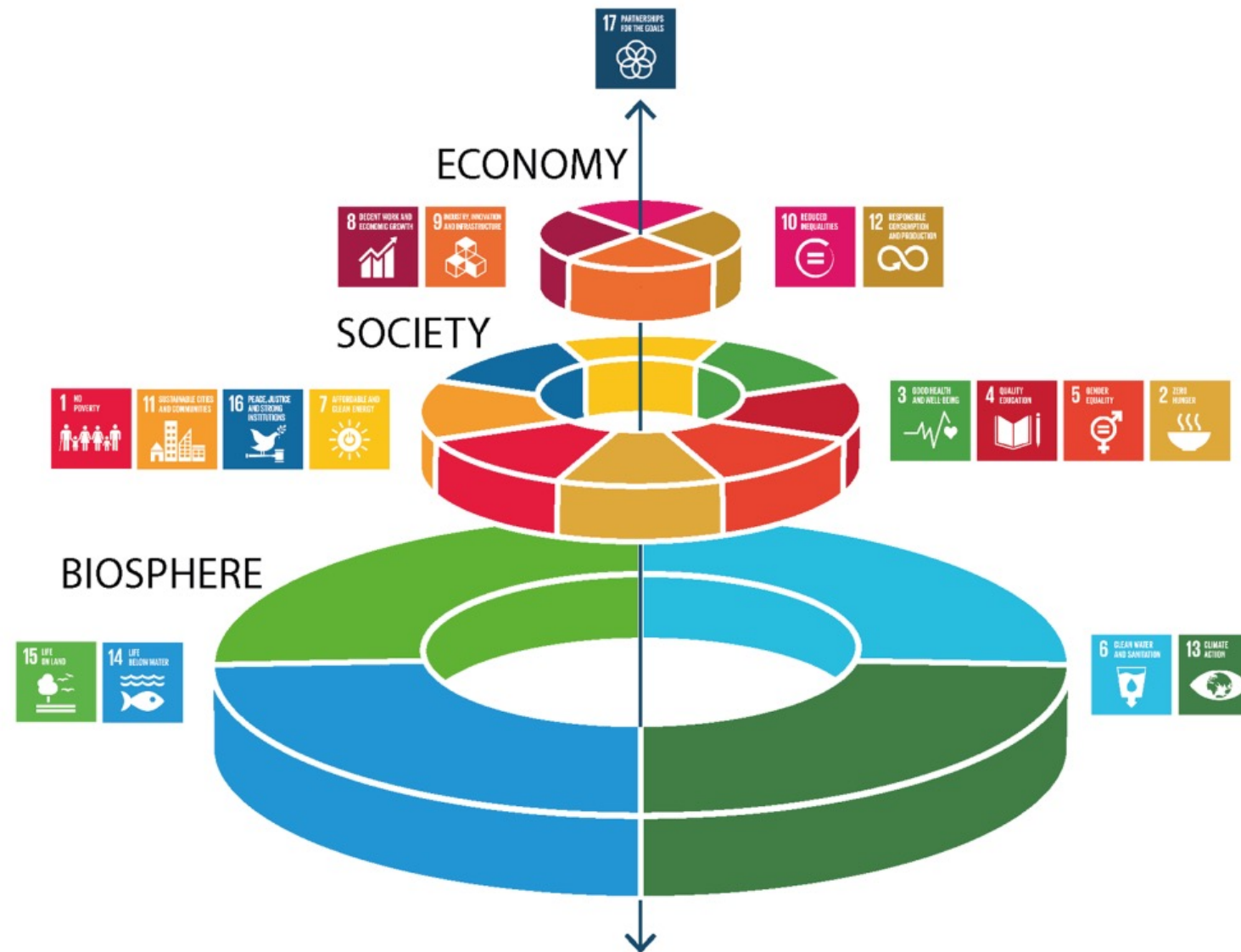
Partnership

Peace

Prosperity

People

Planet



Evolution of Sustainable Finance Research

SDGs:

Sustainable Development Goals

SDGs

Innovative Financial Instrument

Impact Investing

ESG: Environmental, Social, and Governance

CSR: Corporate Social Responsibility

Conscious Capitalism

Climate Financing

Carbon Financing

Green Financing

Ethical Investing

Socially Responsible Investing

Topic

1986

1995

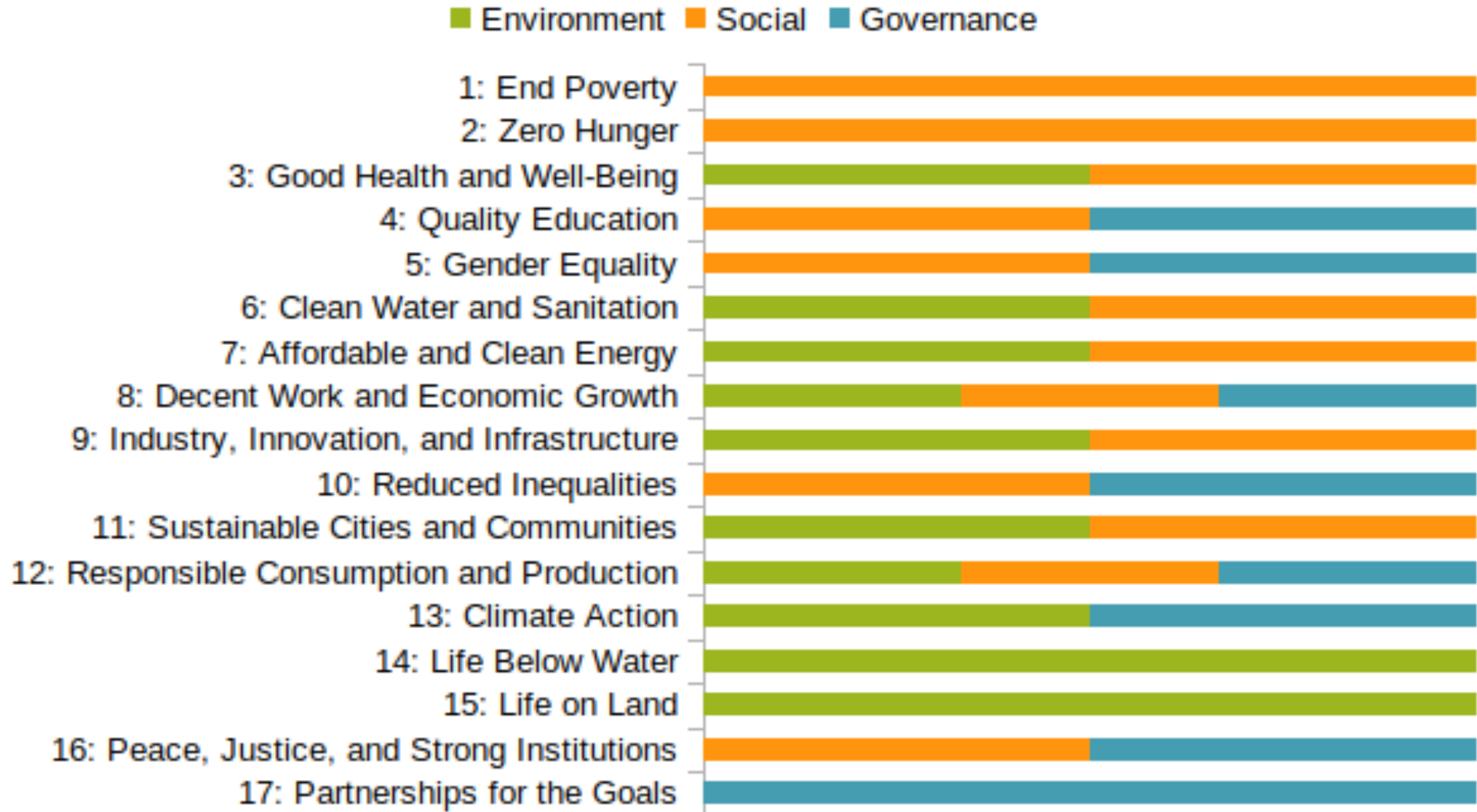
2005

2015

2020

Source: Kumar, S., Sharma, D., Rao, S., Lim, W. M., & Mangla, S. K. (2022). Past, present, and future of sustainable finance: Insights from big data analytics through machine learning of scholarly research. *Annals of Operations Research*, 1-44.

ESG to 17 SDGs



Generative AI and LLMs for Sustainability and ESG Data Analytics



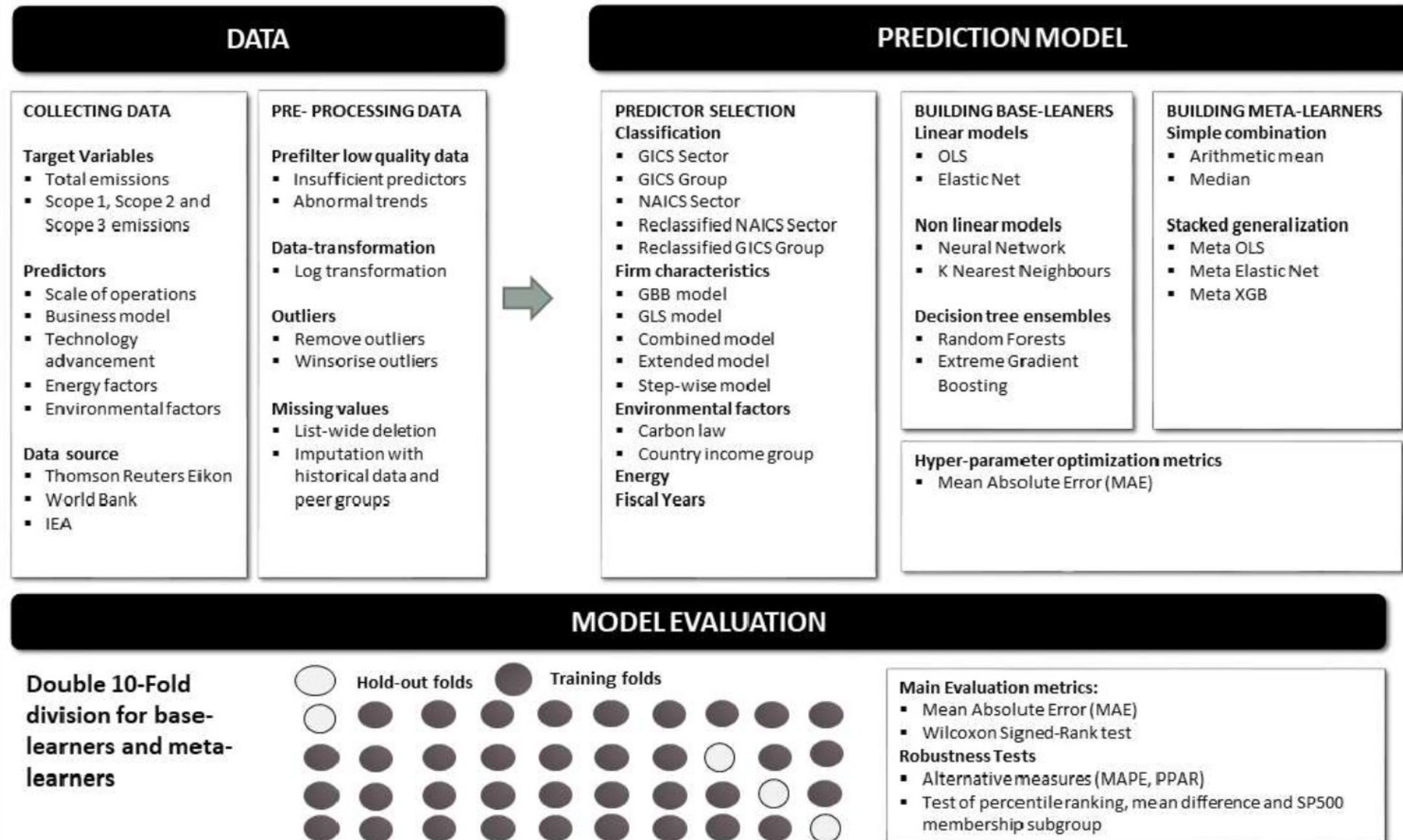
Sustainability and ESG Data Analytics



Generative AI for ESG Rating and Reporting Generation

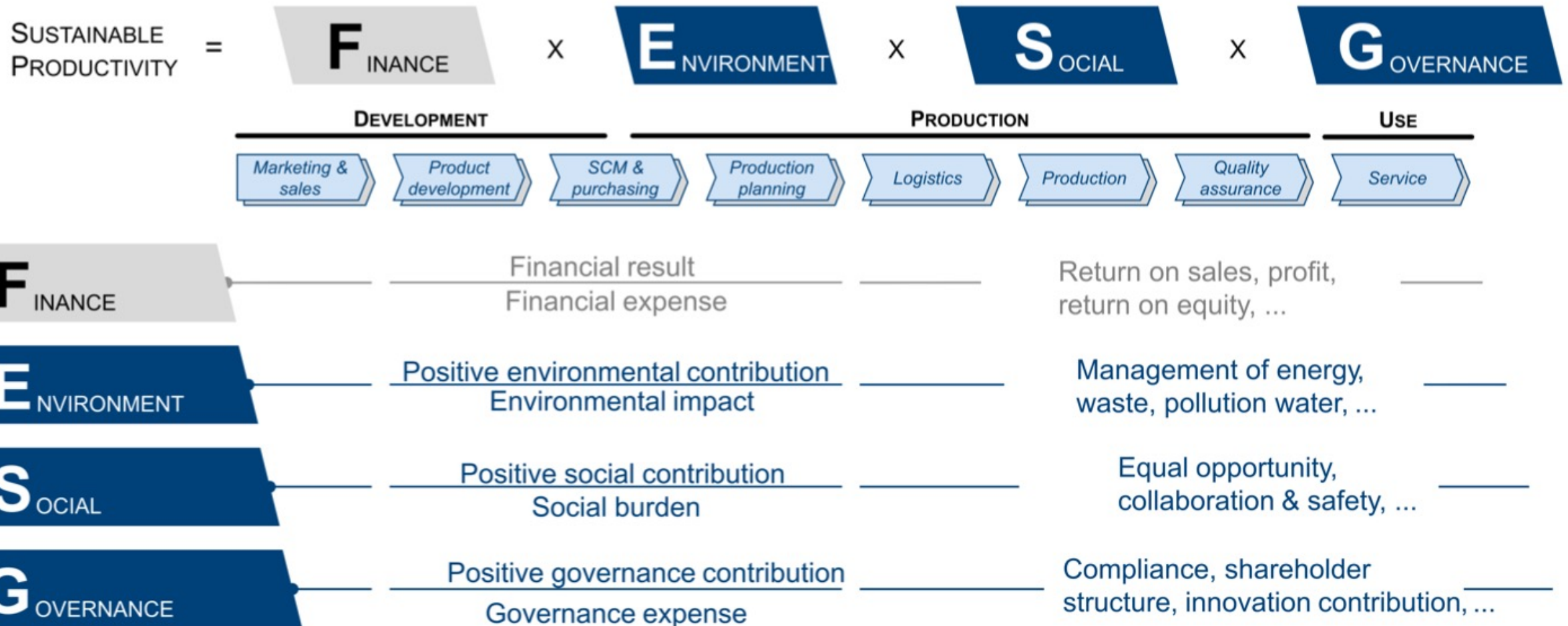


Modelling Strategy to Forecast Carbon Emissions with AI



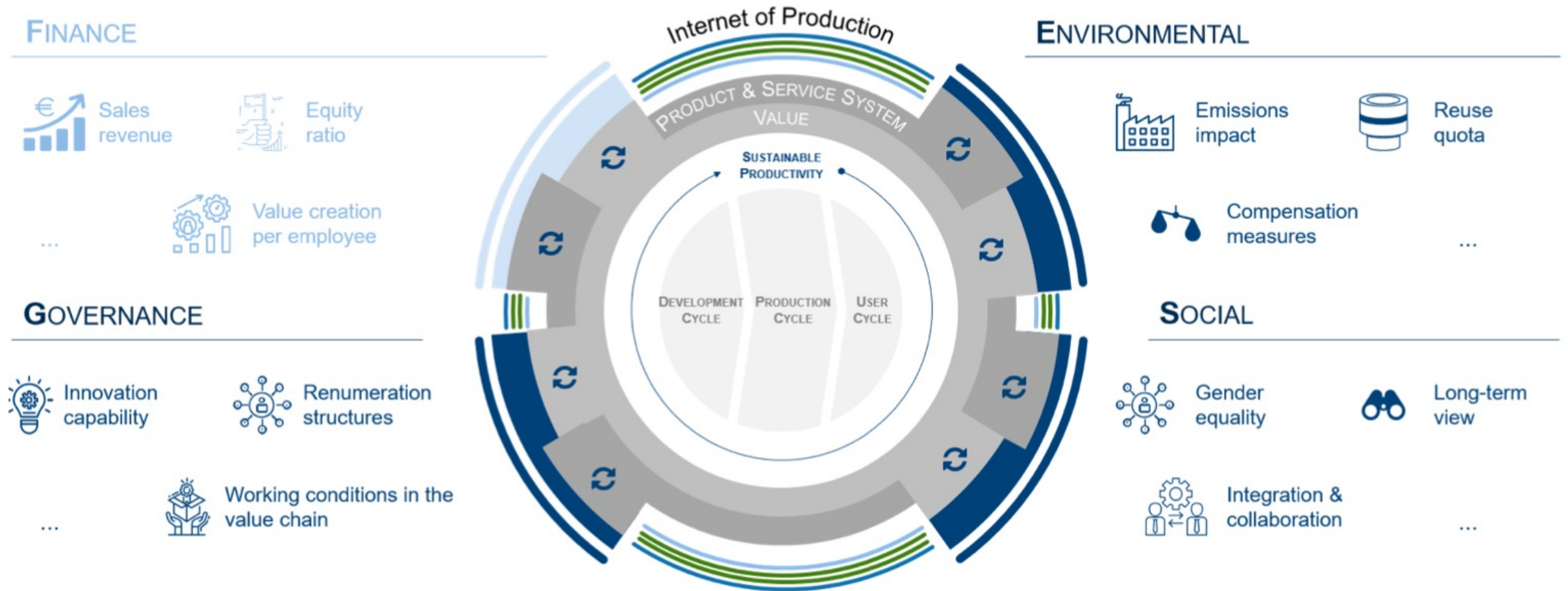
Sustainable Productivity:

Finance ESG



Sustainable Resilient Manufacturing

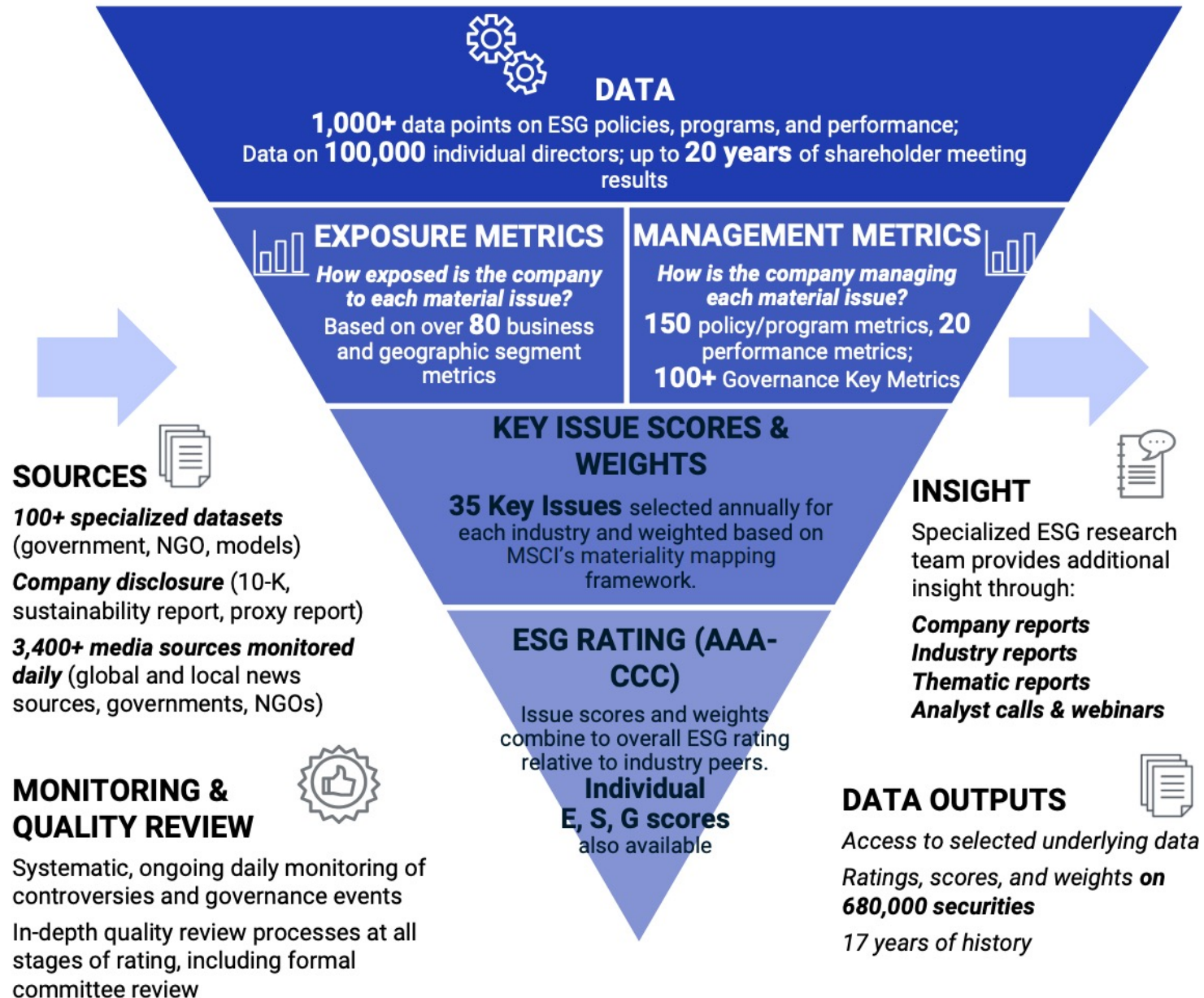
ESG



ESG Indexes

- **MSCI ESG Index**
- **Dow Johns Sustainability Indices (DJSI)**
- **FTSE ESG Index**

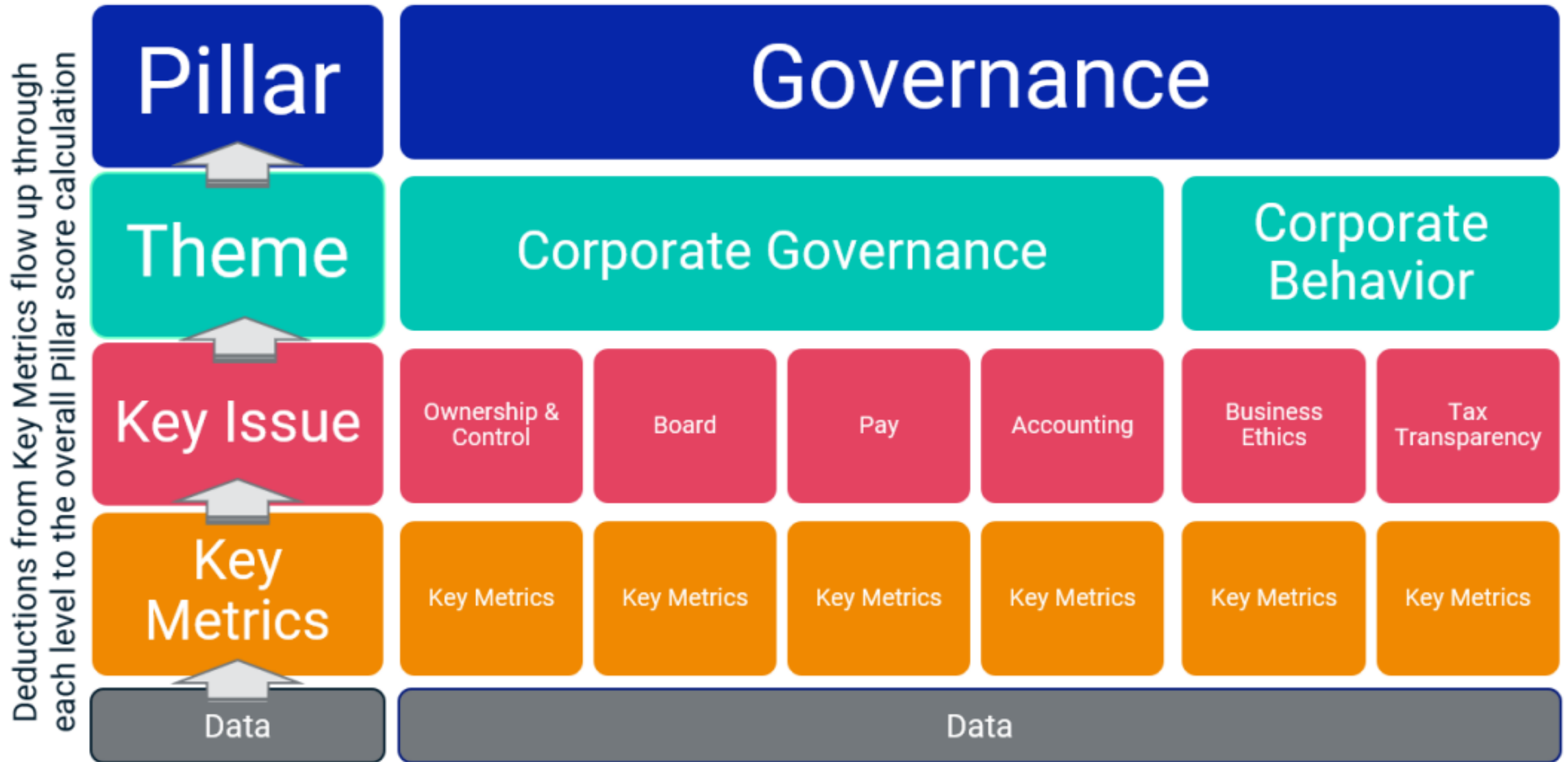
MSCI ESG Rating Framework



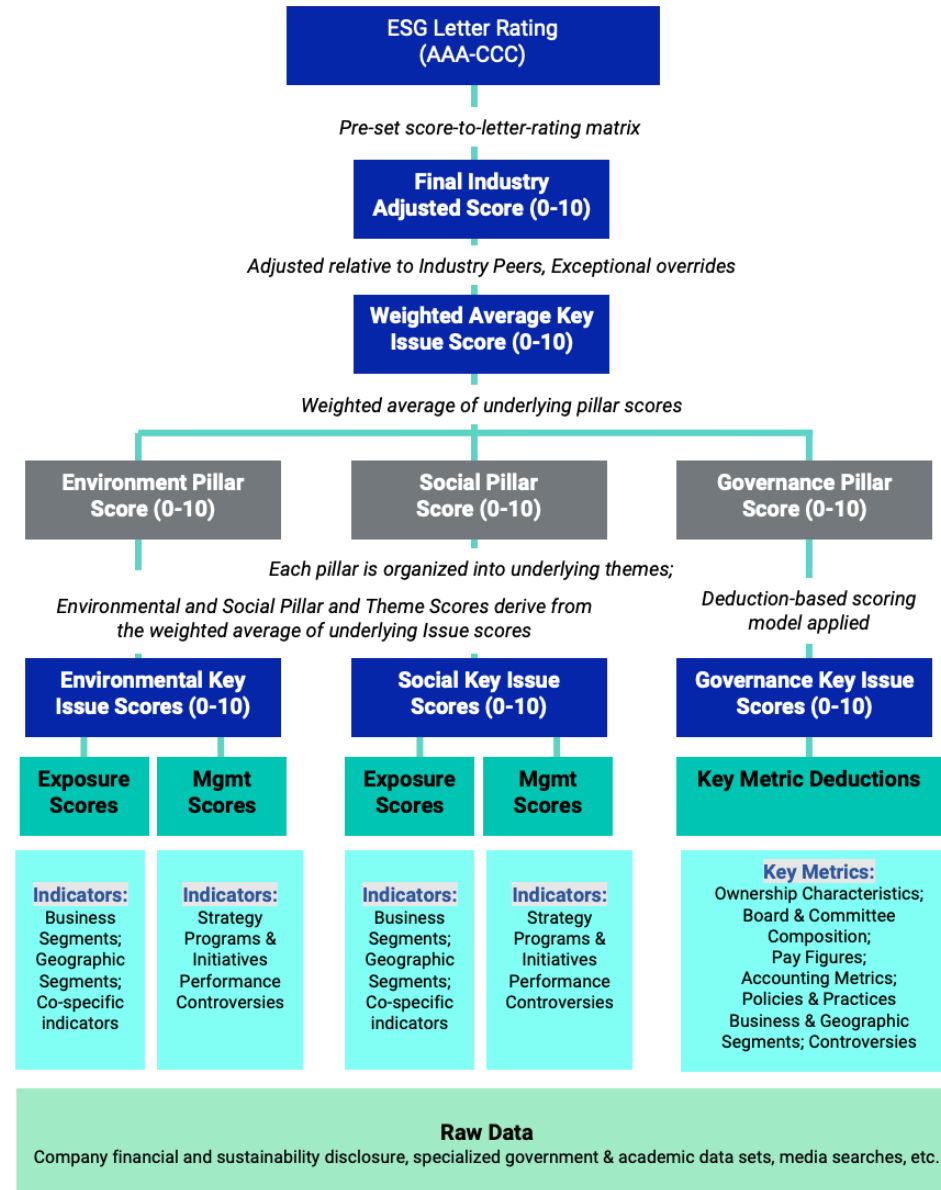
MSCI ESG Key Issue Hierarchy

3 Pillars	10 Themes	35 ESG Key Issues	
Environment	Climate Change	Carbon Emissions Product Carbon Footprint	Financing Environmental Impact Climate Change Vulnerability
	Natural Capital	Water Stress Biodiversity & Land Use	Raw Material Sourcing
	Pollution & Waste	Toxic Emissions & Waste Packaging Material & Waste	Electronic Waste
	Environmental Opportunities	Opportunities in Clean Tech Opportunities in Green Building	Opportunities in Renewable Energy
Social	Human Capital	Labor Management Health & Safety	Human Capital Development Supply Chain Labor Standards
	Product Liability	Product Safety & Quality Chemical Safety Consumer Financial Protection	Privacy & Data Security Responsible Investment Health & Demographic Risk
	Stakeholder Opposition	Controversial Sourcing Community Relations	
	Social Opportunities	Access to Communications Access to Finance	Access to Health Care Opportunities in Nutrition & Health
Governance	Corporate Governance	Ownership & Control Board	Pay Accounting
	Corporate Behavior	Business Ethics Tax Transparency	

MSCI Governance Model Structure



MSCI Hierarchy of ESG Scores



FTSE Russell ESG Ratings



Sustainalytics

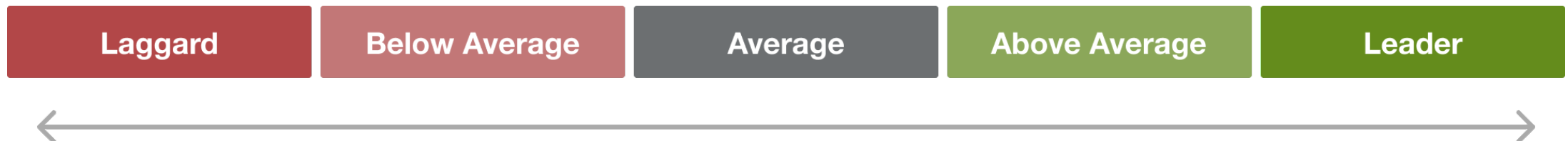
ESG Risk Ratings

Sustainalytics' ESG Risk Ratings measure a company's exposure to industry-specific material ESG risks and how well a company is managing those risks.

Negligible	Low	Medium	High	Severe
0 - 10	10 - 20	20 - 30	30 - 40	40+

Truvalue ESG Ranks

- **Truvalue Labs** applies **AI** to analyze over **100,000 sources** and uncover **ESG risks** and opportunities hidden in **unstructured text**.
- The ESG Ranks data service produces an overall company rank based on industry percentile leveraging the **26 ESG categories** defined by the **Sustainability Accounting Standards Board (SASB)**.
- The data feed covers 20,000+ companies with more than 13 years of history.



Analyst-driven vs. AI-driven ESG

Analyst-driven ESG research

Derives ratings in a structured data model



Sustainalytics

Analyst role at the end of the process allows subjectivity to color results

AI-driven ESG research

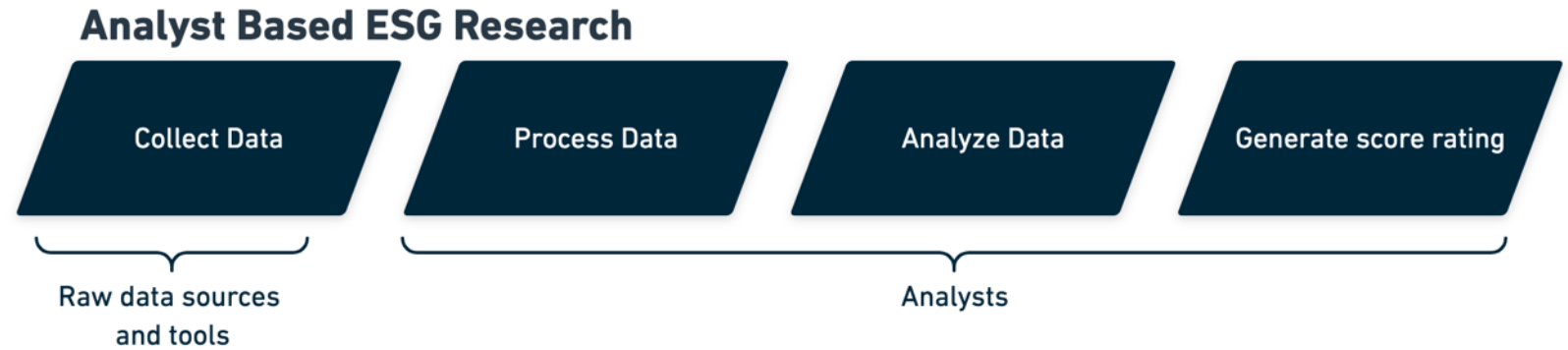
Derives signals from unstructured data



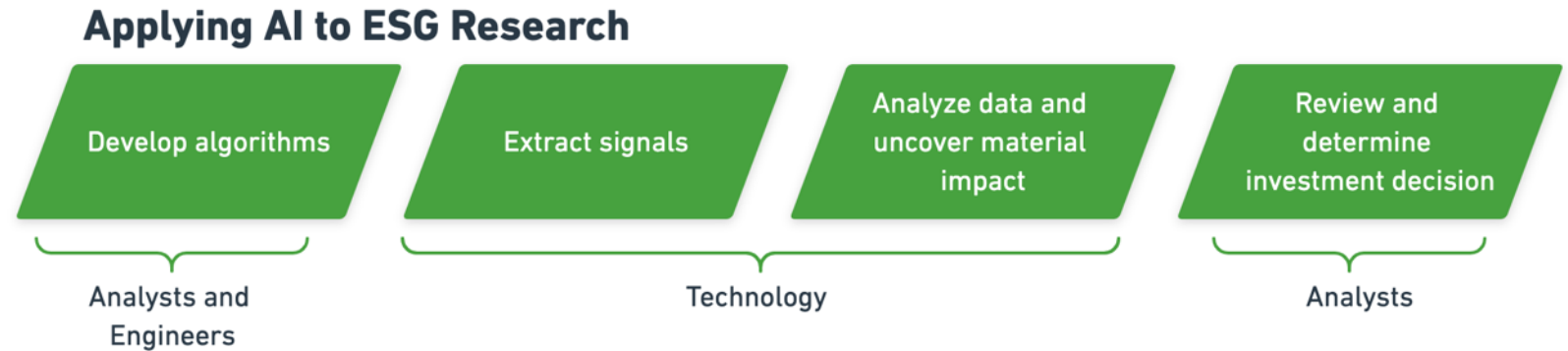
Truvalue Labs

Analyst expertise at the beginning of the process produces consistent results

Analyst based ESG Research

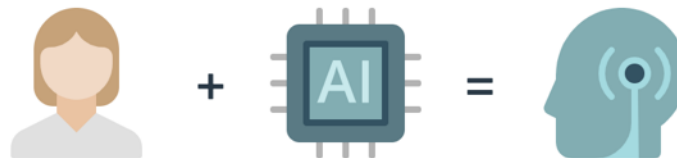


AI based ESG Research



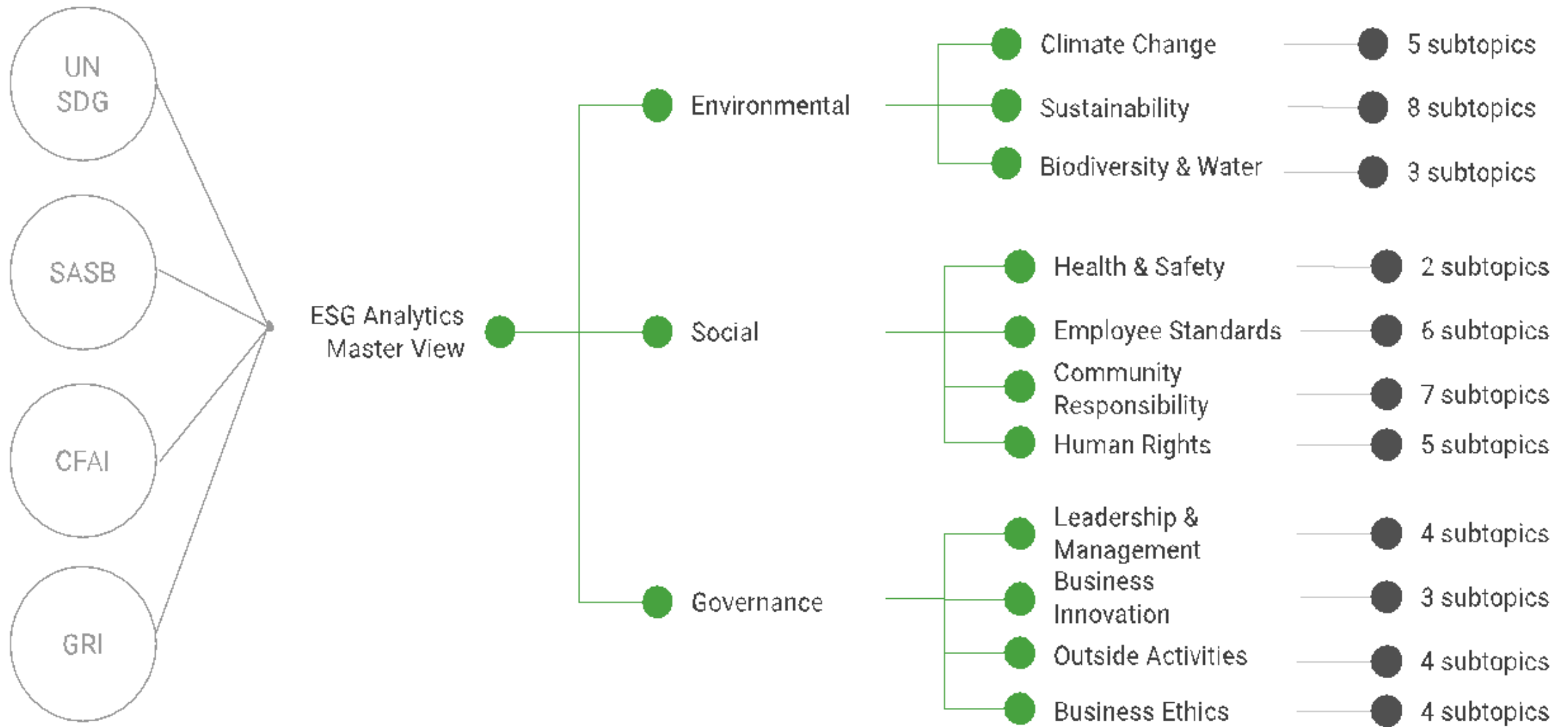
It would take an analyst over 5 years to do what our AI can in 1 week

Combining analysts with AI creates gives you the full picture



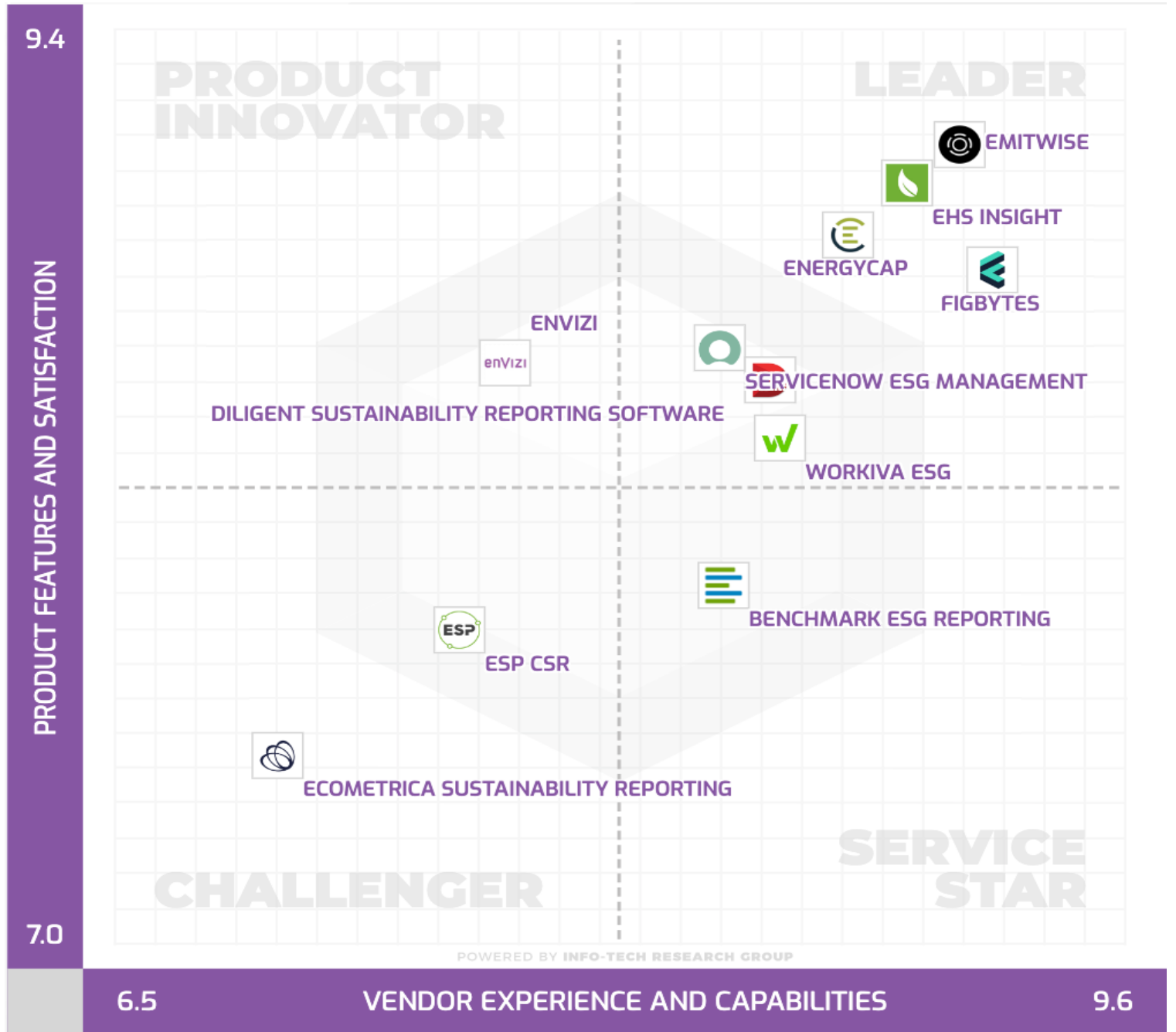
ESG ANALYTICS
Invest where it matters.

ESG Analytics: NLP Taxonomy



Top ESG Reporting Software

Environmental, Social and Governance (ESG) Reporting software or **Sustainability software** helps organizations **manage their operational data, evaluate their impact on the environment and provide reporting to perform audits.**



Summary

- 1. Generative AI**
- 2. ChatGPT and Large Language Models (LLMs)**
- 3. Popular Generative AI and Applications**
- 4. Generative AI for Corporate ESG and Sustainable Development**

Acknowledgments: Research Projects

1. **Applying AI technology to construct knowledge graphs of cryptocurrency anti-money laundering: a few-shot learning model**
 - MOST, 110-2410-H-305-013-MY2, 2021/08/01~2023/07/31
2. **Fintech Green Finance for Carbon Market Index, Corporate Finance, and Environmental Policies. Carbon Emission Sentiment Index with AI Text Analytics**
 - NTPU, 112-NTPU_ORDA-F-003 , 2023/01/01~2024/12/31
3. **Digital Support, Unimpeded Communication: The Development, Support and Promotion of AI-assisted Communication Assistive Devices for Speech Impairment. Multimodal Cross-lingual Task-Oriented Dialogue System for Inclusive Communication Support**
 - NSTC 112-2425-H-305-002-, 2023/05/01-2026/04/30
4. **Establishment and Implement of Smart Assistive Technology for Dementia Care and Its Socio-Economic Impacts. Intelligent, individualized and precise care with smart AT and system integration**
 - NSTC, NSTC, 112-2627-M-038-001-, 2023/08/01~2024/07/31
5. **Use deep learning to identify commercially dental implant systems - observational study**
 - USTP-NTPU-TMU, USTP-NTPU-TMU-112-01, 2023/01/01~2023/12/31
6. **Metaverse AI Multimodal Cross-Language Task-Oriented Dialogue System**
 - ATEC Group x NTPU, NTPU-112A413E01, 2023/05/01~2026/04/30
7. **Metaverse Avatar Automatic Metadata Generation Module**
 - FormosaVerse x NTPU, NTPU-111A413E01, 2022/12/01~2023/11/30
8. **Pilot Study on Universal Data Processing for Code Generation Engine**
 - III x NTPU, NTPU-112A513E01, 2023/08/01~2023/12/22

References

- Yihan Cao, Siyu Li, Yixin Liu, Zhiling Yan, Yutong Dai, Philip S. Yu, and Lichao Sun (2023). "A Comprehensive Survey of AI-Generated Content (AIGC): A History of Generative AI from GAN to ChatGPT." arXiv preprint arXiv:2303.04226.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. (2023) "Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing." ACM Computing Surveys 55, no. 9 (2023): 1-35.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min et al. (2023) "A Survey of Large Language Models." arXiv preprint arXiv:2303.18223.
- Touvron, Hugo, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov et al. (2023) "Llama 2: Open Foundation and Fine-Tuned Chat Models." arXiv preprint arXiv:2307.09288 (2023).
- Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., Manning, C. D., & Finn, C. (2023). Direct preference optimization: Your language model is secretly a reward model. arXiv preprint arXiv:2305.18290.
- Tunstall, Lewis, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang et al. "Zephyr: Direct Distillation of LM Alignment." arXiv preprint arXiv:2310.16944 (2023).
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. arXiv preprint arXiv:2203.02155.
- Gozalo-Brizuela, Roberto, and Eduardo C. Garrido-Merchan (2023). "ChatGPT is not all you need. A State of the Art Review of large Generative AI models." arXiv preprint arXiv:2301.04655 (2023).
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. (2023) "InstructBLIP: Towards General-purpose Vision-Language Models with Instruction Tuning." arXiv preprint arXiv:2305.06500 (2023).
- Shahab Saquib Sohail, Faiza Farhat, Yassine Himeur, Mohammad Nadeem, Dag Øivind Madsen, Yashbir Singh, Shadi Atalla, and Wathiq Mansoor (2023). "The Future of GPT: A Taxonomy of Existing ChatGPT Research, Current Challenges, and Possible Future Directions." Current Challenges, and Possible Future Directions (April 8, 2023) (2023).
- Longbing Cao (2022). "Decentralized ai: Edge intelligence and smart blockchain, metaverse, web3, and desc." IEEE Intelligent Systems 37, no. 3: 6-19.
- Qinglin Yang, Yetong Zhao, Huawei Huang, Zehui Xiong, Jiawen Kang, and Zibin Zheng (2022). "Fusing blockchain and AI with metaverse: A survey." IEEE Open Journal of the Computer Society 3 : 122-136.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. arXiv preprint arXiv:2203.02155.

Q & A

生成式AI在企業永續發展的應用

Generative AI for Corporate ESG and Sustainable Development

Time: 15:00-17:00, Saturday, December 2, 2023

Place: NTPU USR Hub

<https://learningcollaboration.org/index.php/2023/09/21/collegeco1120103/>

戴敏育 副教授

Min-Yuh Day, Ph.D,
Associate Professor

Institute of Information Management, National Taipei University

<https://web.ntpu.edu.tw/~myday>

2023-12-02

