

Microsoft Learn
Student Ambassadors

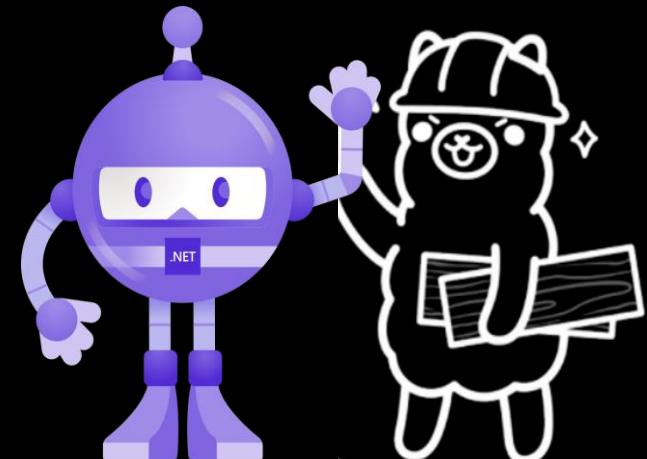
Microsoft
Research

Developing Local AI Models with .NET and AutoGen



Charunthon Limseelo (@boatchrnthn)
Microsoft Learn Student Ambassador | Thai Student Tech Community Lead

.NET Conf Thailand at Seven Peaks Software (November 23rd)



Charunthon Limseelo (Boat)

Beta Microsoft Learn Student Ambassadors (Thai Student Tech Lead) and Google Developer Group (GDG Cloud Bangkok) Member

- + Microsoft Office Specialist (Excel)
- + Open-source small AI and NLP Interest, with BDE Applications for business aspect
- + Applied Skills Challenger (Azure AI Document Intelligence and NLP)



Charunthon Limseelo



@boatchrnthn



Charunthon Limseelo



Boat Charunthon (boatchrnthn)



Acknowledgement To All AI Leads and Specialists

Engineering, Research-Based, and Practical Leads



Business, Financial, Commercial-based, and Daily Users



Proprietary Gemini

(adj.) used, made, or marketed by one having the exclusive legal right, privately owned and managed and run as a **profit-making** organization, Close-Source

Open-Source

(adj.) having the source code **freely available** for possible modification and redistribution, along with publicly available for use by the community at large



Bad gateway

The web server reported a bad gateway error.

[SUBSCRIBE TO UPDATES](#)

ChatGPT is unavailable for some users.

[Subscribe](#)

Investigating - We are currently investigating this issue.

Jun 04, 2024 - 07:33 PDT

Uptime over the past 90 days. [View historical uptime.](#)

API [?](#)



Operational

ChatGPT [?](#)



Major Outage

Labs [?](#)



Operational

Tasks Libraries Datasets Languages Licenses
Other

Filter Tasks by name

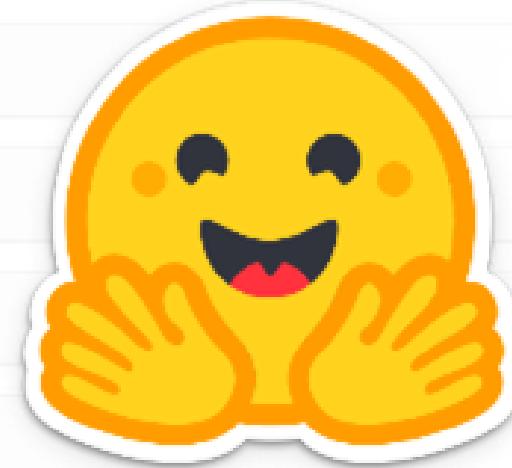
Multimodal

- Image-Text-to-Text
- Visual Question Answering
- Document Question Answering
- Video-Text-to-Text
- Any-to-Any

Computer Vision

- Depth Estimation
- Image Classification
- Object Detection
- Image Segmentation
- Text-to-Image
- Image-to-Text
- Image-to-Image
- Image-to-Video
- Unconditional Image Generation
- Video Classification
- Text-to-Video
- Zero-Shot Image Classification
- Mask Generation
- Zero-Shot Object Detection
- Text-to-3D
- Image-to-3D
- Image Feature Extraction
- Keypoint Detection

- Models 1,045,019 Filter by name
- Full-text search Sort: Trending
- openai/whisper-large-v3-turbo
Automatic Speech Recognition • Updated 7 days ago • ↓ 102k • ⚡ • ❤ 859
 - nvidia/NVLM-D-72B
Image-Text-to-Text • Updated 2 days ago • ↓ 18.8k • ❤ 590
 - rain1011/pyramid-flow-sd3
Text-to-Video • Updated about 2 hours ago • ❤ 287
 - black-forest-labs/FLUX.1-dev
Text-to-Image • Updated Aug 16 • ↓ 1.13M • ⚡ • ❤ 5.33k
 - ostris/OpenFLUX.1
Text-to-Image • Updated 7 days ago • ↓ 11.1k • ❤ 453
 - rhymes-ai/Aria
Text Generation • Updated 1 day ago • ↓ 172 • ❤ 208
 - apple/DepthPro
Depth Estimation • Updated 1 day ago • ❤ 193
 - jxm/cde-small-v1
Feature Extraction • Updated 1 day ago • ↓ 1.68k • ❤ 185



Data Leak

Data Centers' Power Consumption

Server's Down

Prevent from Paying Premium Services

No Internet Connection

Tasks Libraries Datasets Languages Licenses

Other

Filter Tasks by name

Multimodal

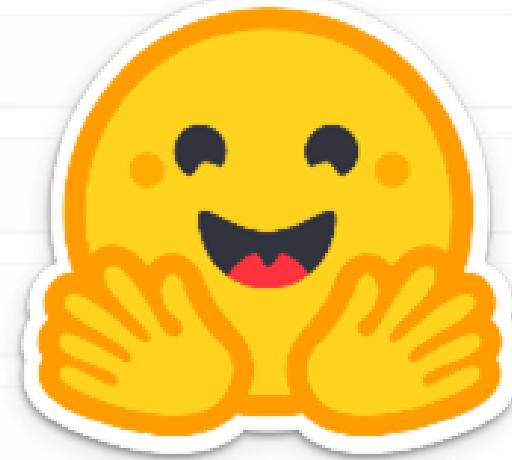
- Image-Text-to-Text
- Visual Question Answering
- Document Question Answering
- Video-Text-to-Text
- Any-to-Any

Computer Vision

- Depth Estimation
- Image Classification
- Object Detection
- Image Segmentation
- Text-to-Image
- Image-to-Text
- Image-to-Image
- Image-to-Video
- Unconditional Image Generation
- Video Classification
- Text-to-Video
- Zero-Shot Image Classification
- Mask Generation
- Zero-Shot Object Detection
- Text-to-3D
- Image-to-3D
- Image Feature Extraction
- Keypoint Detection

Models 1,045,019 Filter by name Full-text search Sort: Trending

- openai/whisper-large-v3-turbo
Automatic Speech Recognition • Updated 7 days ago • 102k • 859
- nvidia/NVLM-D-72B
Image-Text-to-Text • Updated 2 days ago • 18.8k • 590
- rain1011/pyramid-flow-sd3
Text-to-Video • Updated about 2 hours ago • 287
- black-forest-labs/FLUX.1-dev
Text-to-Image • Updated Aug 16 • 1.13M • 5.33k
- ostris/OpenFLUX.1
Text-to-Image • Updated 7 days ago • 11.1k • 453
- rhymes-ai/Aria
Text Generation • Updated 1 day ago • 172 • 208
- apple/DepthPro
Depth Estimation • Updated 1 day ago • 193
- jxm/cde-small-v1
Feature Extraction • Updated 1 day ago • 1.68k • 185

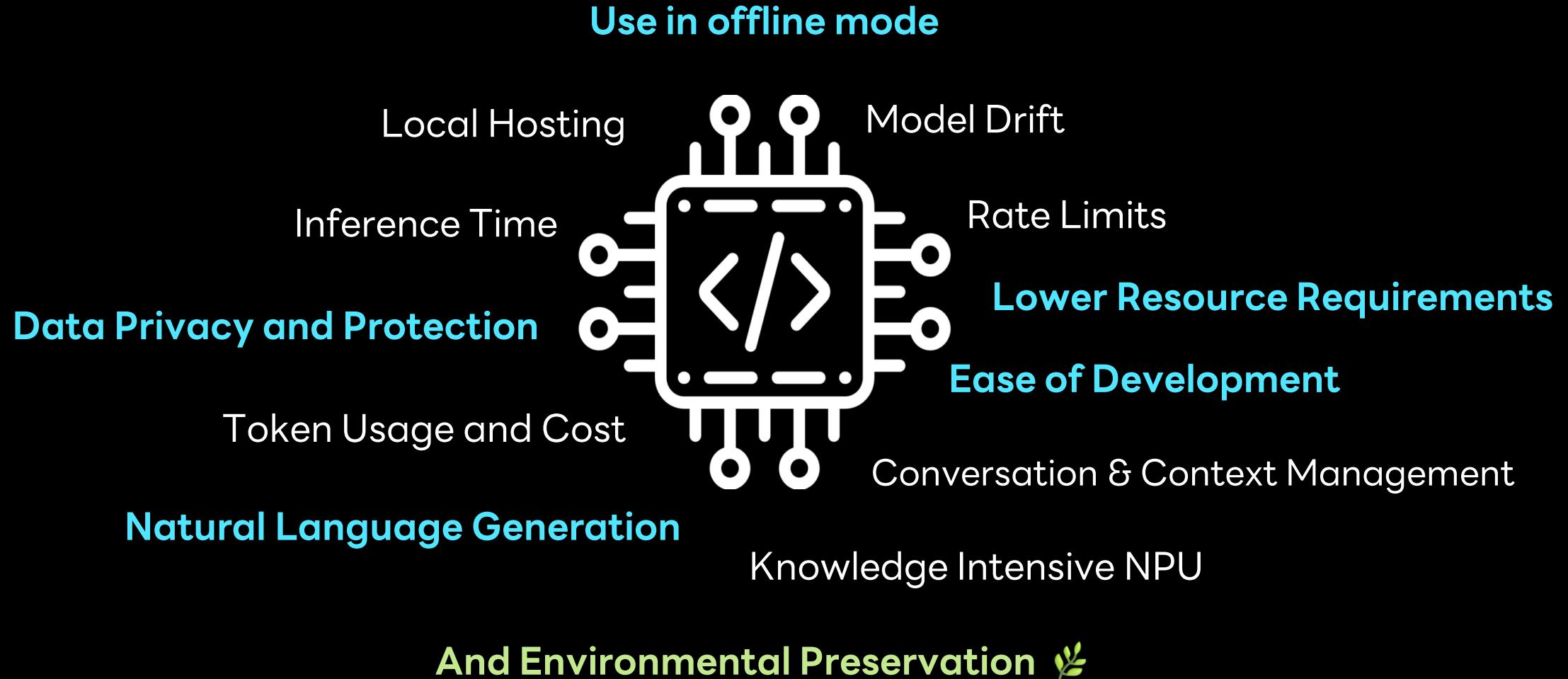


Found out some LLMs are too big to download, need to use API to connect (Require internet)

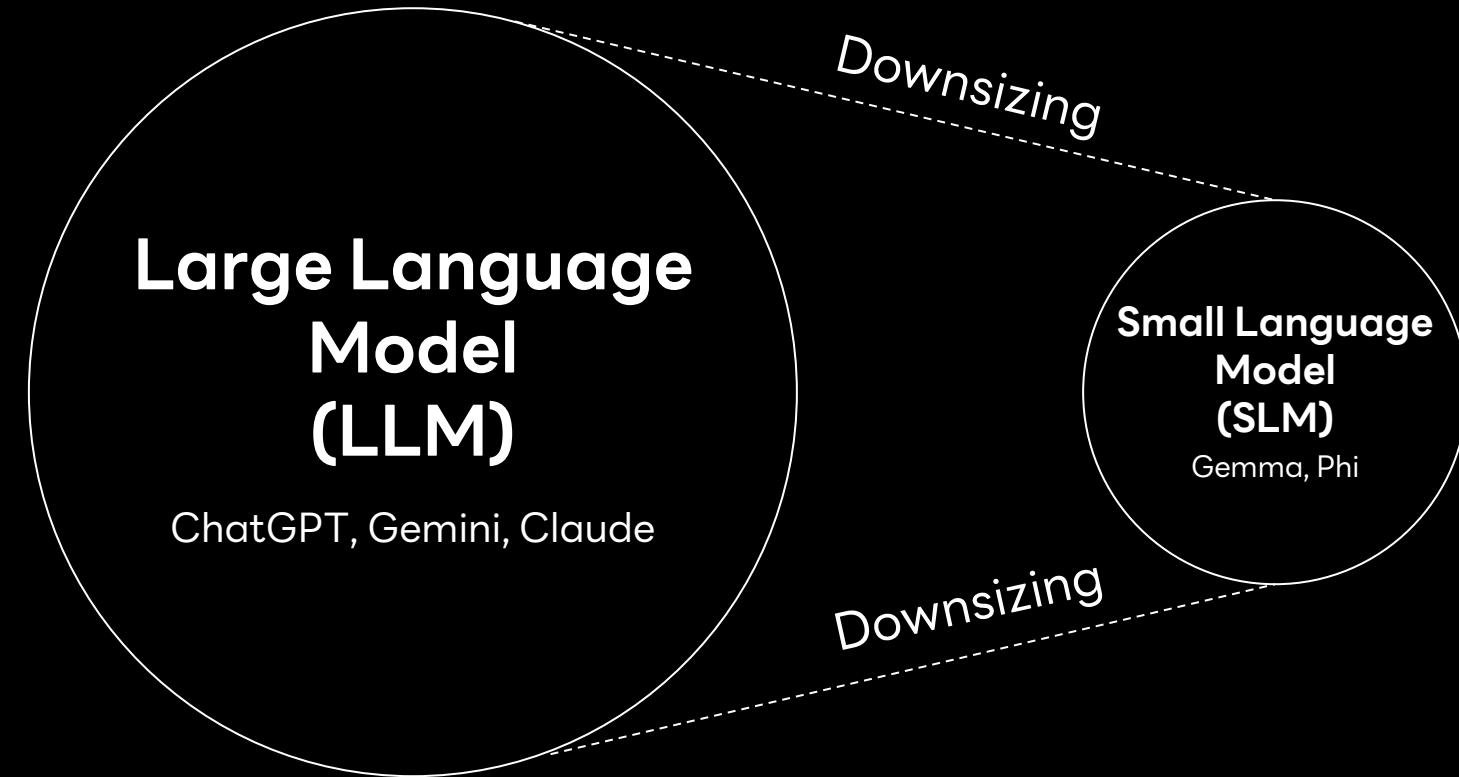
I have an interesting way to prevent all these problems for using each device offline!

Small Language Models (SLMs)

What is Small Language Model?

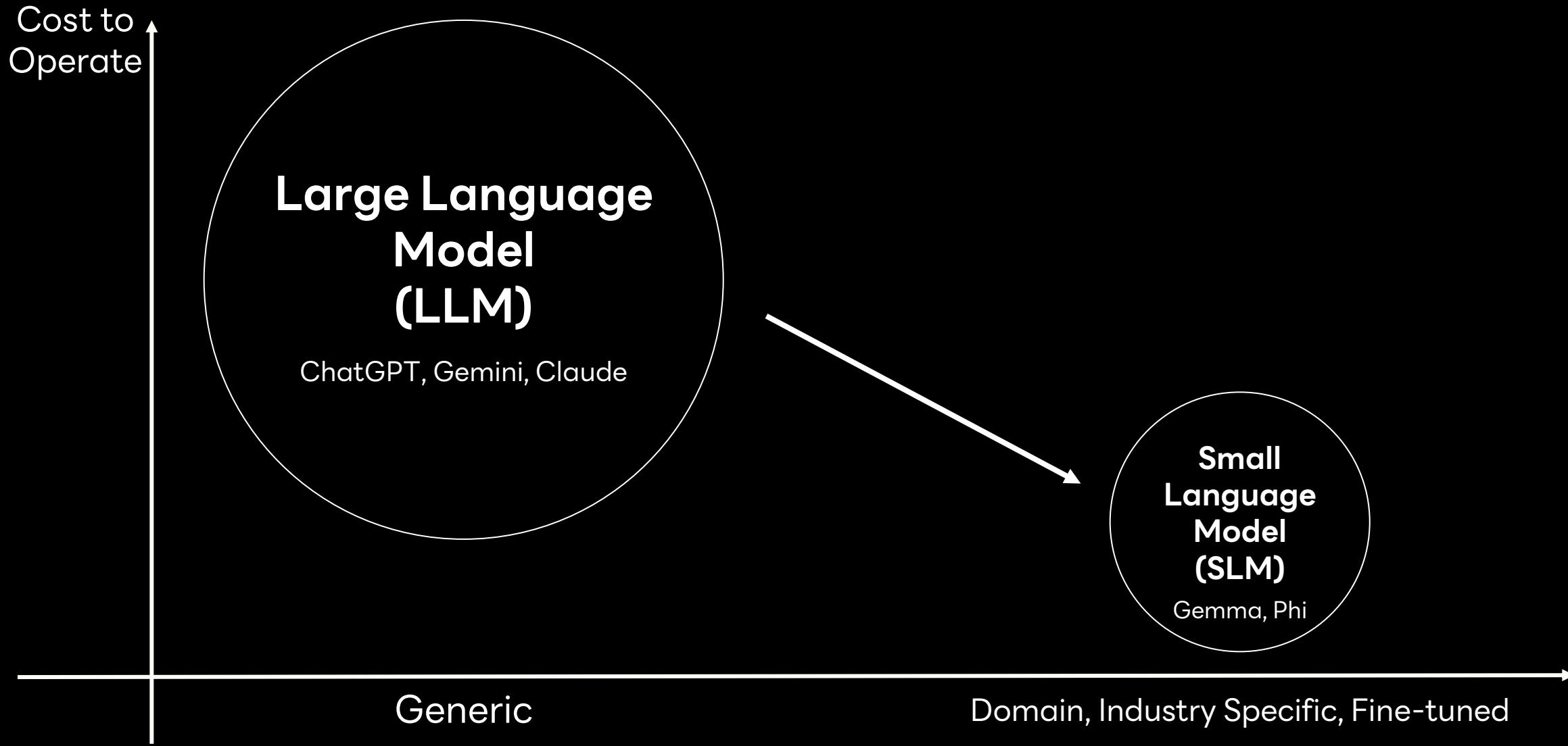


Language Model Comparison



- Ten of billions in parameters
 - Requires substantial computational power for training and development
 - Higher performance in broader and more complex tasks
-
- Millions to few billions in parameters
(should be lower than 7B)
 - Capable of being trained with consumer GPUs and lower budgets
 - Effective for specific and narrow tasks

Language Model Comparison



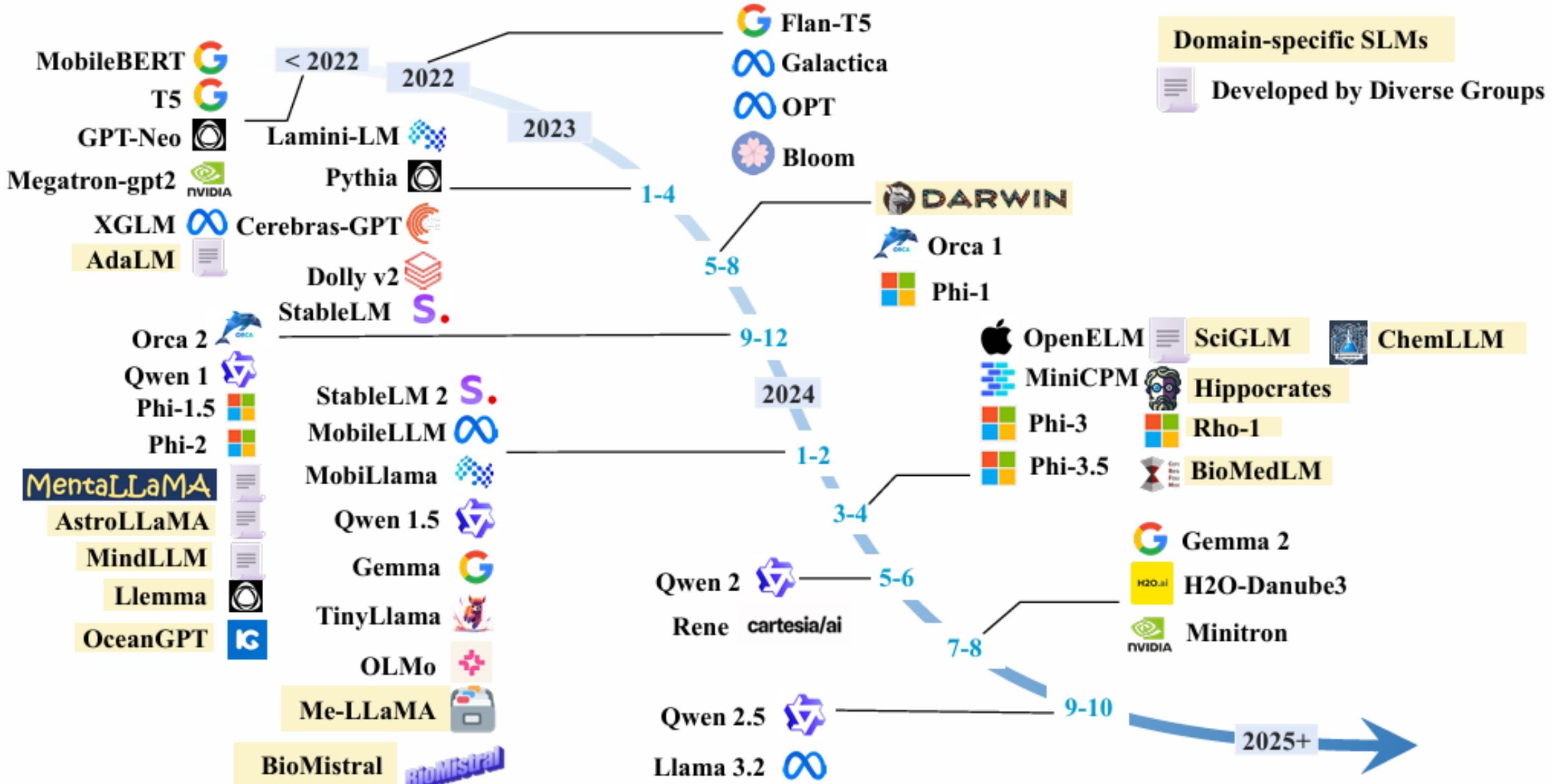
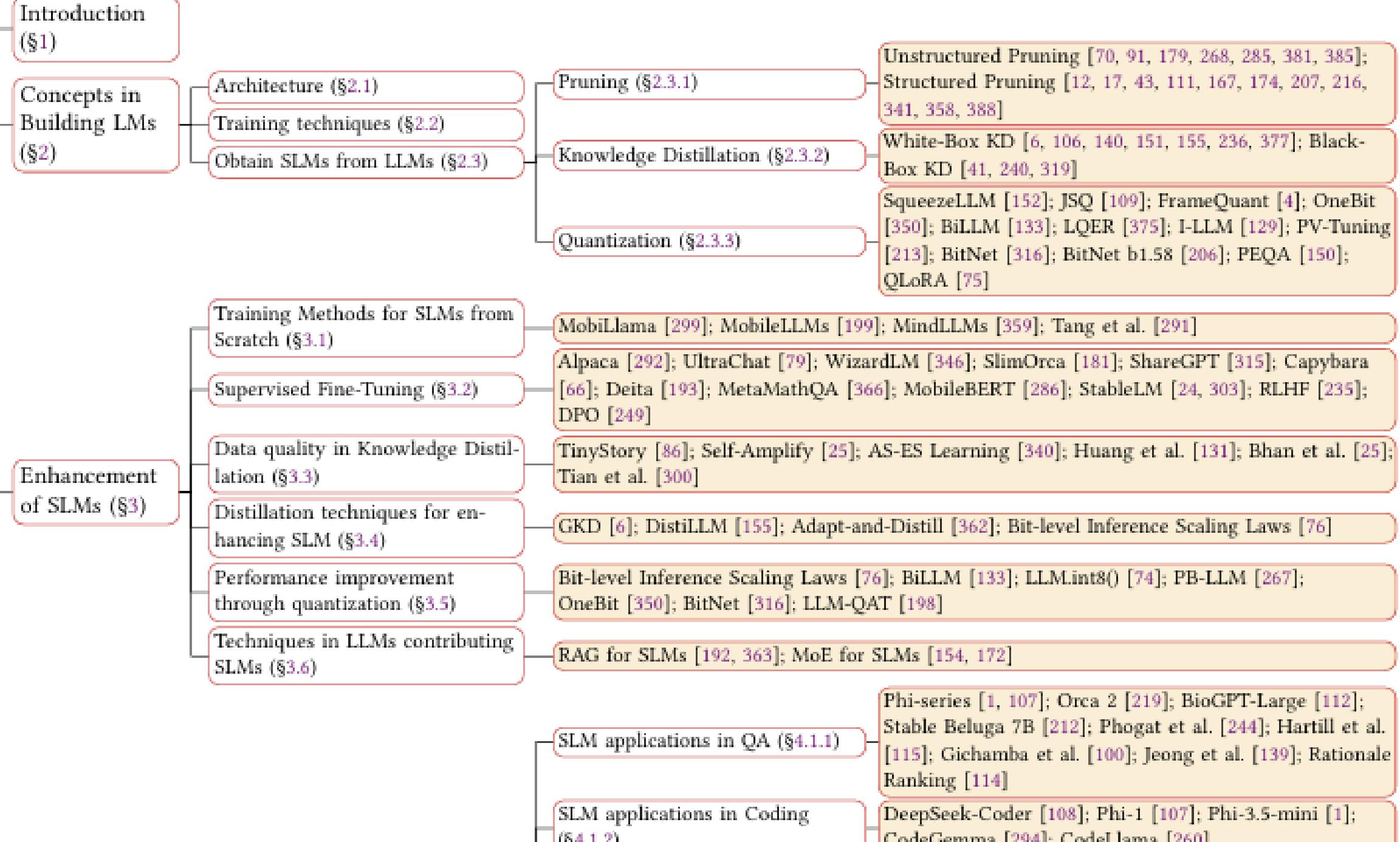
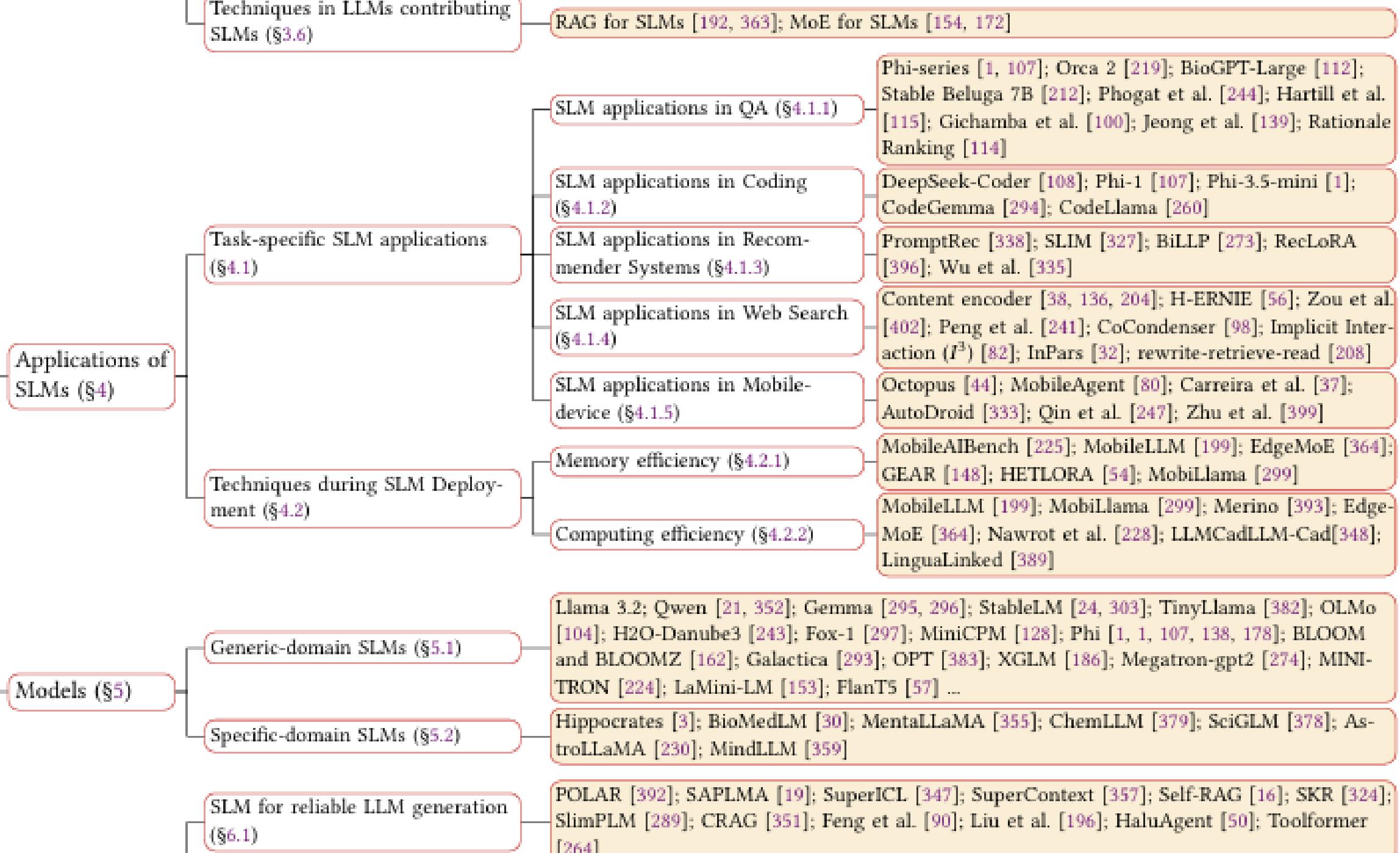


Fig. 3. A timeline of existing small language models.



Small Language Models



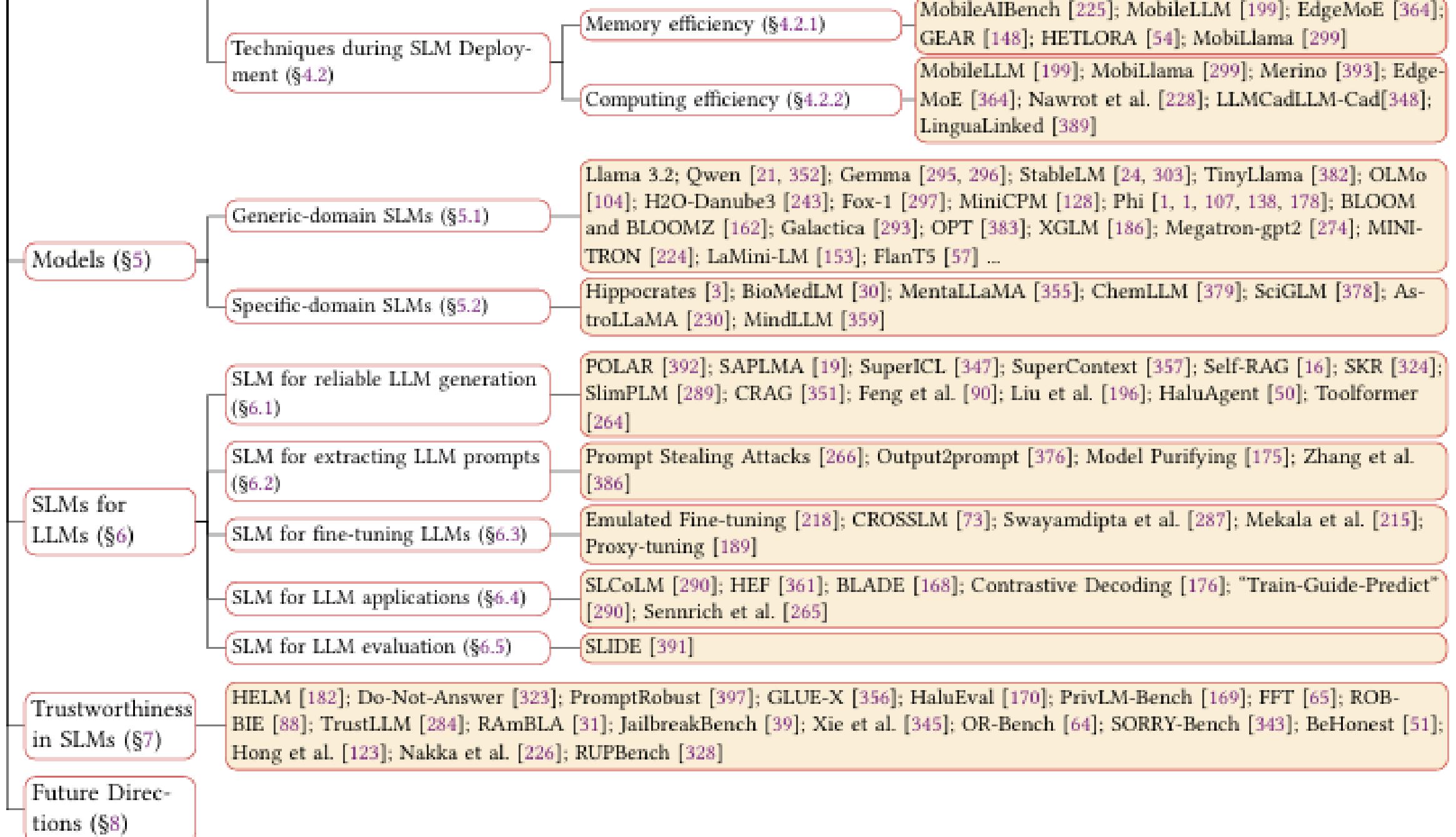


Fig. 1. Overview of Small Language Models.

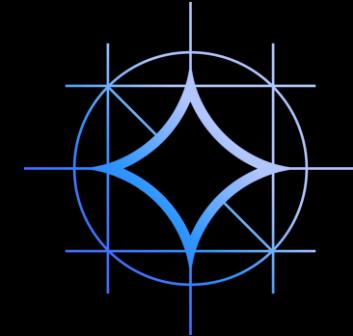
Less-Parameters Language Models Arena



Phi-1 : 1.3B
Phi-2 : 2.7B
Phi-3.5: 3.8B (\approx 2.2 GB)



LLaMA 3.2 : 3B (\approx 2.02 GB)
Thai-Lang Supported



Deepmind's Gemma 2
2B: \approx 1.6 GB
9B: \approx 5.4 GB (LLM)



SCB^X

KBTG

Qwen2.5 3B (\approx 1.93 GB)
Thai-Lang Supported

Typhoon 8B (\approx 4.7-4.9 GB)
(Consider as LLM)
Thai-Lang Supported

THaLLE 7B (\approx 4.4 GB-5.4 GB)
(Consider as LLM)
Thai-Lang Supported

Case Examples as Using Small Language Model



Anuchit Sapanpong (O) -

Actor

Biography Trivia with Small Language Model (using RAG Technique)



Jirachai Chansivanon (P' Job)

- Developer

TinyLM for Internet of Things, OfflineGPT on Airplane



Wichayada Chamnansilp (Namin) - Programmer

Offline AI Interaction Research



Wit Sitthivaekin (The Standard) - Influencer

History Trivia on Local Device



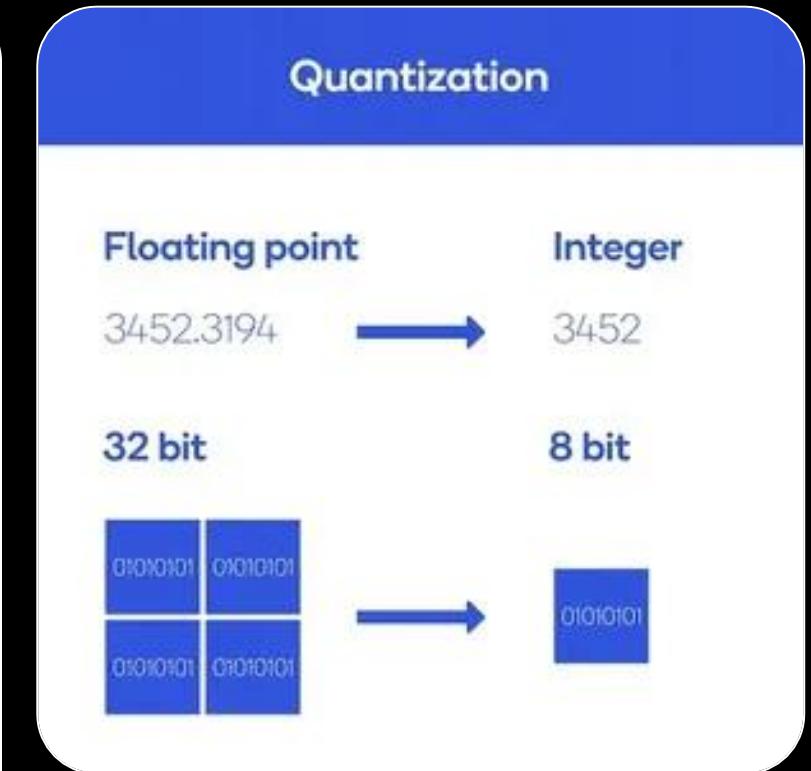
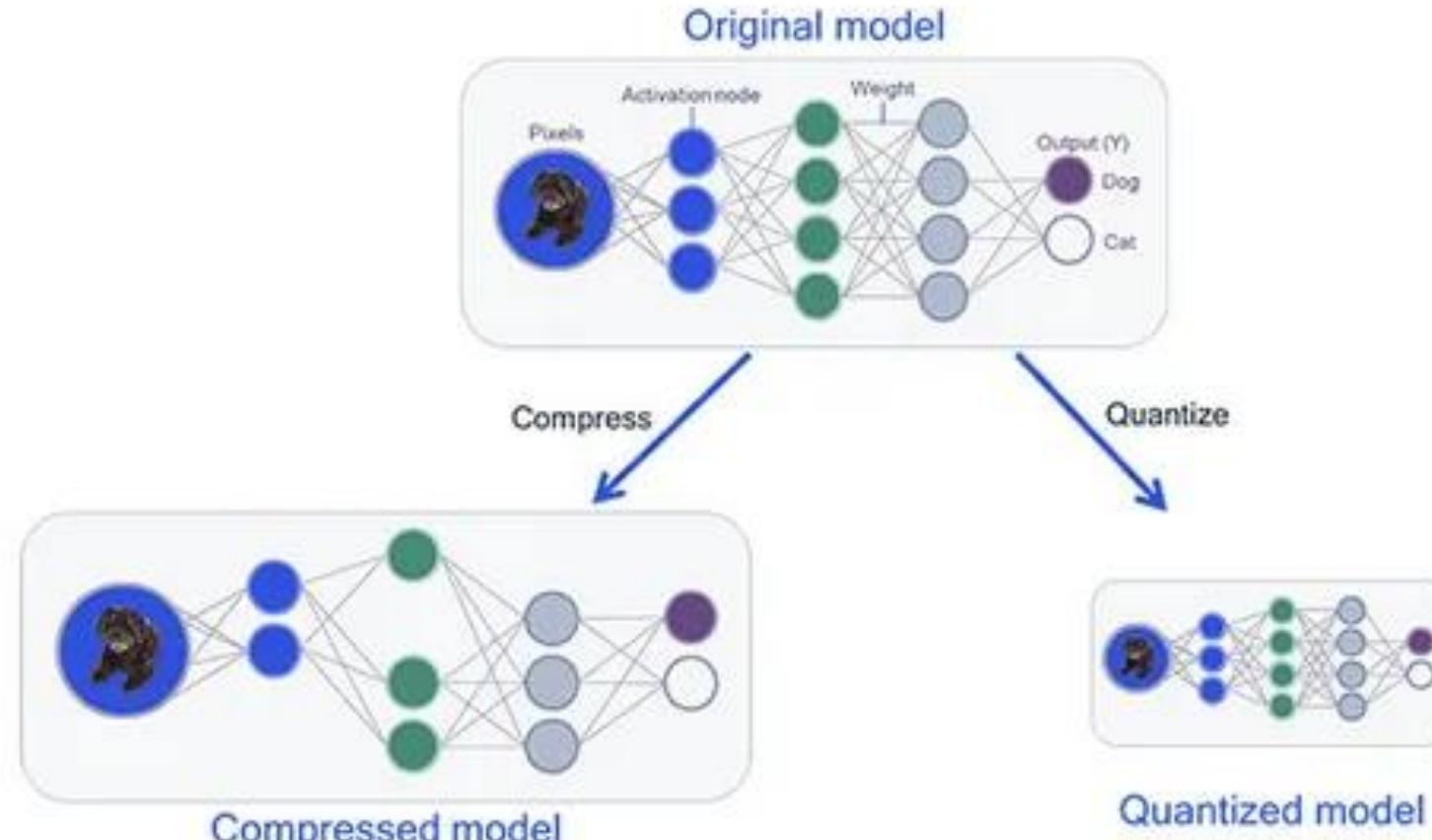
CK Cheong (Fastwork) -

Influencer

Research and Local ESL Practice for his hobby

How are you going to use
Small Language Models
on your device at the first step?

SLMs Quantization



Reduce size and computation cost

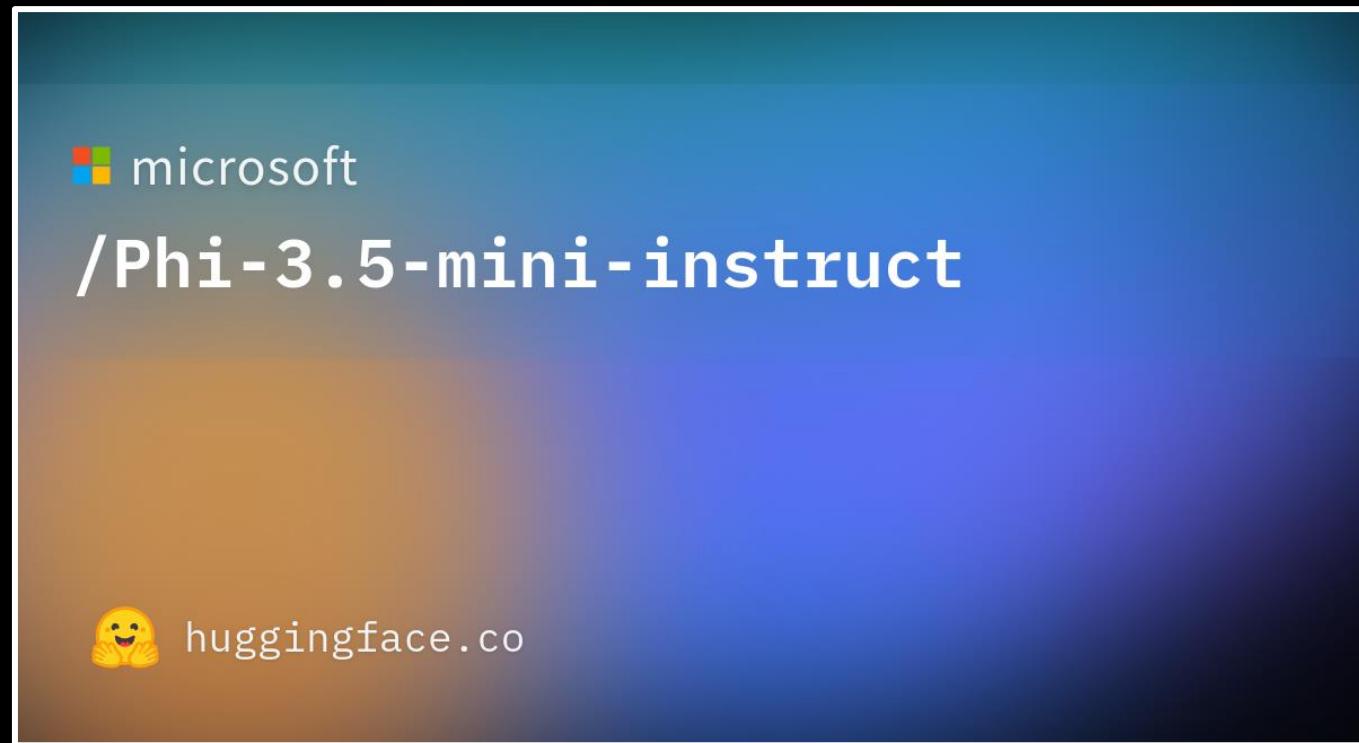
Deploy model with limited device resources

Reduce Memory Footprint + Accelerate inference time

May loss some accuracy (but could improve by fine tuning)

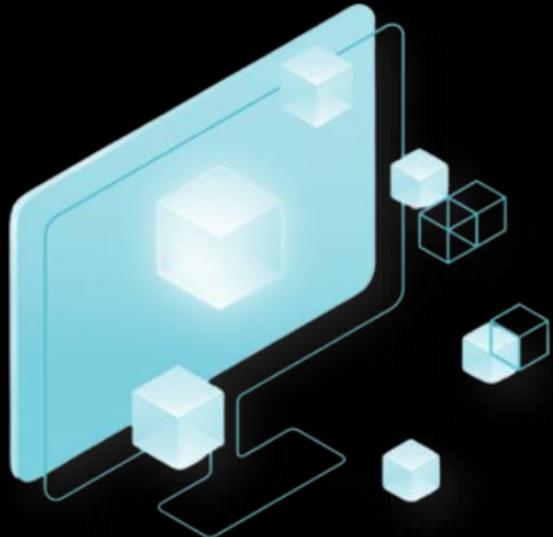
Generally Available

Microsoft's Phi-3.5



[microsoft/Phi-3.5-mini-instruct · Hugging Face](https://huggingface.co/microsoft/Phi-3.5-mini-instruct)

Microsoft's Phi-3.5 Family

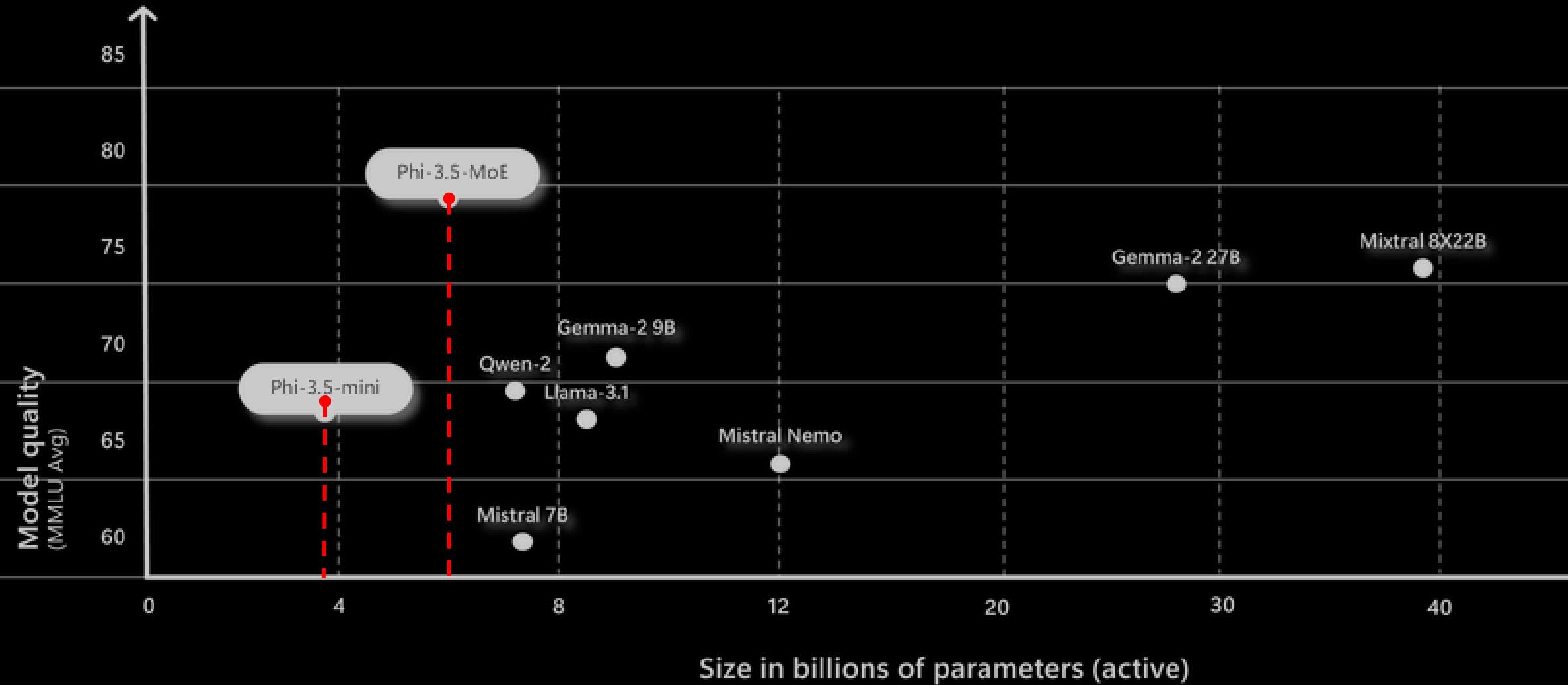


Model For Today's Experiment

Phi-3.5 mini-4k instruct (SLM)

- **Architecture:** 3.8B parameters and is a dense decoder-only Transformer model. The model is fine-tuned with Supervised fine-tuning (SFT) and Direct Preference Optimization (DPO) to ensure alignment with human preferences and safety guidelines.
- **Inputs:** Text. It is best suited for prompts using chat format.
- **Context Window Length:** 4K tokens
- **GPUs:** 512 H100-80G
- **Training time:** 7 days
- **Training data:** 3.3T tokens
- **Outputs:** Generated text in response to the input
- **Dates:** Our models were trained between February and April 2024
- **Status:** This is a static model trained on an offline dataset with cutoff date October 2023. Future versions of the tuned models may be released as we improve models.

Phi-3.5 Performance



Textbooks Are All You Need

Evolution of Phi family

Textbooks Are All You Need

Suriva Gunasekar
Allie Del Gismo
Gustavo de Rosa
Xia Wang
Soband

Textbooks Are All You Need II: phi-1.5 technical report

Yuanzhi Li
Sébastien Bubeck
Ronen Eldan
Sunoya Gunasekar
Yin Tat Lee
Allie Del Gismo

Microsoft Research

Abstract

We continue the investigation into the power of smaller Transformer-based language models as initiated by [“Textbooks Are All You Need”](#), a 10 million parameter model that can perform coherent reasoning tasks and the follow-up work on phi-1, a 1.3 billion parameter model with Python coding performance close to the state-of-the-art. The latter work proposed to use existing Large Language Models (LLMs) to process mathematical reasoning tasks, such as solving equations, performing arithmetic operations, and analyzing web data. We follow the “Textbooks Are All You Need” approach, focusing this time on common sense reasoning in natural language, and create a new 1.3 billion parameter model named phi-1.5, with approximately 10% of the parameters of phi-1. The new model is trained on a mix of common sense and few-shot LLMs on more complex reasoning tasks such as grade-school mathematics and basic science. The results show that phi-1.5 can solve a wide range of tasks from the largest few-shot benchmarks, including halving sums and the potential for text-based generations—encouragingly though, we are seeing improvements in that front thanks to the absence of web data. We open-source phi-1.5 to invite further research on these exciting topics.

1 Introduction

The art of training large open-domain language models remains limited. And not just at the same time as Transformers have become the dominant paradigm. In fact, the term “large language model” was first coined in 2016, and since then, the quality of data leads to better RASH²⁰, and it can yield quality dataset verifications of high quality data sets, the scaling laws, potentialities, and limitations of the data can even improve dataset size and training times, and so on.

We focus our attention from their doings as it has been widely adopted

arXiv:2306.11644v2 [cs.CL] 2 Oct 2023

arXiv:2309.05638v1 [cs.LG] 11 Sep 2023

Figure 1: Benchmark results comparing phi-1.5 vs phi-1. The figure consists of three bar charts side-by-side. The left chart is ‘Common Sense Reasoning’, the middle is ‘Language Understanding and Knowledge’, and the right is ‘Multi-Step Reasoning’. Each chart compares phi-1.5 (blue bars) with phi-1 (red bars) across four benchmarks: LLaMA-13B (green), LLaMA-7B (orange), Qwen-1.5B (purple), and Qwen-7B (yellow). In all categories, phi-1.5 consistently outperforms phi-1 across all benchmarks.

1

Phi-1 and Phi-1.5 (First Development)

Phi-3 Technical Report:
A Highly Capable Language Model Locally on Your Phone

Microsoft

Abstract

We introduce **phi-3-mini**, a 3.8 billion parameter language model trained on 3.3 trillion tokens, whose overall performance, as measured by both academic benchmarks and internal testing, rivals that of models such as Mixtral 8x7B and GPT-3.5 (e.g., phi-3-mini achieves 99% on MMU-U and 8.38 on MT-Bench, while setting a new state-of-the-art on both). The model can be fine-tuned directly on our dataset for training, a scaled-up version of the one used for **phi-2**. The composed of heavily filtered publicly available with data and synthetic data. The model is also further aligned for robustness, safety, and privacy. We also introduce **phi-3-video**, a 4.2 billion parameter video model based on 1.1B models trained for 1.8T tokens, called **phi-3-small** and **phi-3-medium**, both significantly more capable than **phi-3-mini** (e.g., respectively 75% and 78% on MMU-U, and 8.7 and 8.9 on MT-Bench). Moreover, we also introduce **phi-3-video**, a 4.2 billion parameter model based on **phi-3-mini** with strong reasoning capabilities for image and text prompts.

1 Introduction

The striking progress of AI in the last few years can be largely attributed to major efforts throughout the world towards scaling-up to ever-larger models and datasets. Large Language Models (LLMs) have steadily increased in size from a mere billion parameters just five years ago (GPT-2 had 1.5 billion parameters) [RPM+19] to more than 1000x today. The expectation for the effect of scaling in the seemingly predictable direction one obtains by training larger models on the available corpora [KMF+20, HBM+22, MRB+23]. However, these laws assume a “fixed” data source. This assumption is now significantly disrupted by the existence of frontier LLMs themselves, which allow us to interact with data in novel ways. In this paper, we introduce **phi-3** (Figure 12.2), which on May 23, 2024, was shown that it can outperform a wide range of LLM-based filters on a variety of web, data, and LLM-critical tasks, enabling performance in smaller language models that were typically seen only in much larger models. For example, our previous model trained on this data recipe, **phi-2** (2.7B parameters), matched the performance of Mixtral 25 trillion parameter model on similar benchmarks. This report also presents new and improved variants of **phi-3** (**phi-3-mini** (3.8B parameters), trained for 3.3B tokens on larger and more diverse datasets used in **phi-2**). With its small size, **phi-3-mini** can easily be inference locally on a modern phone (see Figure 2), yet it achieves a quality that seems on-par with models such as Mixtral 8x7B [JSR+24] and GPT-3.5.

1

Phi-3 (Latest & Stable) – Before #MSBuild 2024

Source: Microsoft Research (ArXiv)

Phi-1 and Phi-1.5 (First Development)

- 1.3 Billion Parameters
- Focusing on Python Coding and Some Commonsense reasoning
- First state-of-the-art performance on Python coding tasks (test with HumanEval and MBPP)
- Comparable performance to model 5x larger in language understanding

Phi-2 (First Public Released) – Ignite 2023

- 2.7 Billion Parameters
- More efficient on reasoning and language
- Matches or outperforms models up to 25x larger on complex benchmarks. Ideal for research and exploration

Phi-3 (Stable) – Before #MSBuild 2024

- At least 3.8B parameters
- More efficient on scientific reasoning and outperform on mathematical problems
- Providing multimodal features on Phi-3-vision-128k-instruct

Phi-3.5 - Fine Tuning from Phi-3 (Latest - Released on August 2024)

- Supporting Multilingual Context (including Thai Language) - MoE model
- More Logical Thinking for specific tasks

Microsoft Cloud

AI Infrastructure

Foundation models

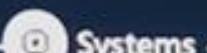
Your data

AI toolchain

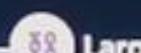
AI for science



Small molecules



Systems



Large molecules

Chemistry

New compounds and materials

Physics

Semiconductors and chip development

Life sciences

AI-driven drug discovery

AI for science: Scientific reasoning • Simulation agents • Scientific data

Cloud supercomputing: HPC • AI •



When you cannot purely use Multilingual Model

The **phi-3.5-MoE** adopts an Mixture-of-Experts (MoE) architecture to selectively activate parts of modules on specific inputs to improve the model efficiency. It incorporates MoE layer as its feedforward models, employing the top2 routing among 16 expert networks. Particularly, each expert network is a separate GLU network and the routing module will selectively activate 2 expert networks out of the 16 expert networks for each token, leaving $16 \times 3.8B$ model to have 6.6B activated parameters with 42B

model-00001-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00002-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00003-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00004-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00005-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00006-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00007-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00008-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00009-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00010-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00011-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00012-of-00017.safetensors	Safe	LFS	4.96 GB	upload initial files	3 months ago
model-00013-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00014-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00015-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00016-of-00017.safetensors	Safe	LFS	4.99 GB	upload initial files	3 months ago
model-00017-of-00017.safetensors	Safe	LFS	3.91 GB	upload initial files	3 months ago

We don't know
which one is
Thai Language

Phi-3.5 family of open models

Now also in MaaS
for Azure AI Studio

New

New

New

Phi-3.5-mini

Phi-3.5-small

Phi-3.5-medium

Phi-3.5-vision

Generally Available At

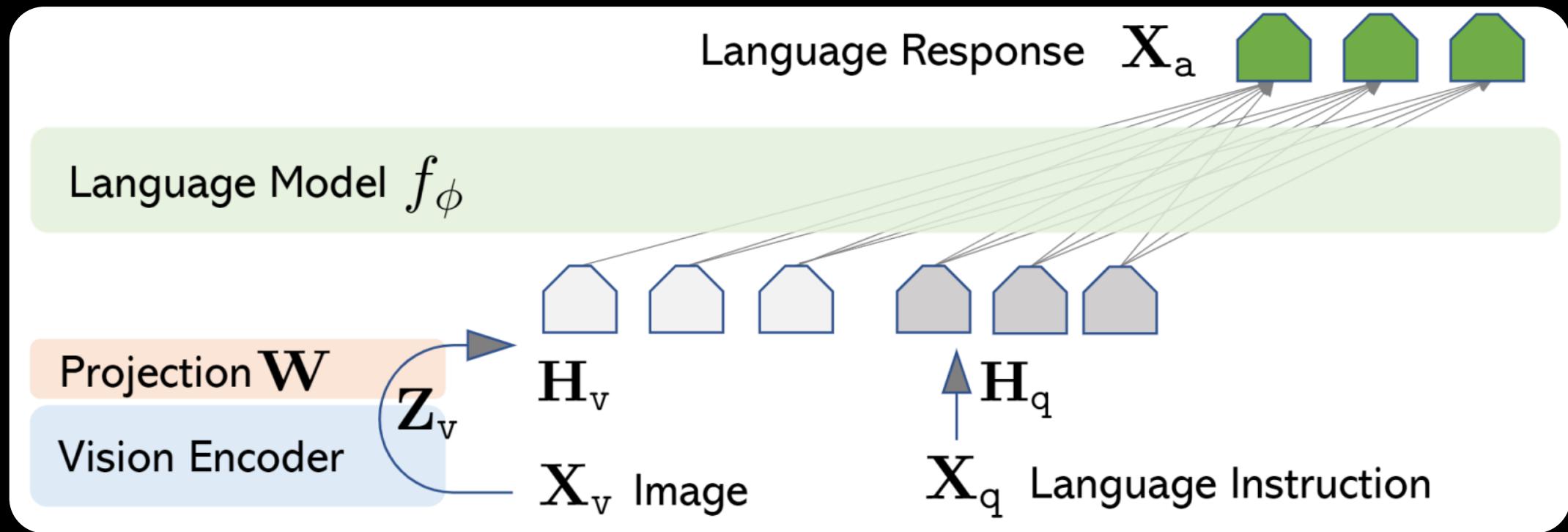


Hugging Face



LlaVA - Large Language-and-Vision Assistant

- An end-to-end trained large multimodal model that connects a vision encoder and LLM for general-purpose visual and language understanding.
- **Open-source.** Making GPT-4 generated visual instruction tuning data, our model and code base publicly available.



Microsoft's Responsible AI principles



Fairness



Reliability
& Safety



Privacy &
Security



Inclusiveness



Transparency



Accountability

P.Y.D. - Simple Responsible AI principles for SLMs



P stands for
Privacy:

- **Data Minimization:**
Collect only necessary data.
- **Consent:** Obtain informed consent for data use.
- **Anonymization:**
Protect personal information.

Y stands for **Your Responsibility:**

- **Bias Mitigation:**
Address biases in data and algorithms.
- **Transparency:**
Explain AI decisions and processes.
- **Accountability:** Take responsibility for AI outcomes.

D stands for
Diversity:

- **Inclusivity:** Ensure AI development and use benefits diverse populations.
- **Fairness:** Avoid discrimination in AI systems.
- **Equity:** Address disparities in AI access and outcomes.

Introducing

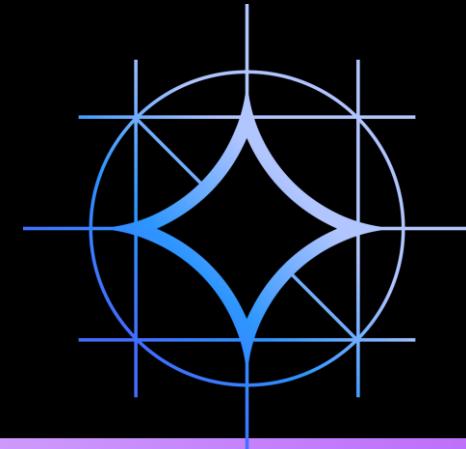
Ollama

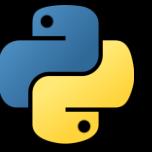
Running AI on Local Devices



What is Ollama?

- Ollama is a platform that makes local development with open-source language models a breeze.
- With Ollama, everything you need to run an LLM—model weights and all of the config—is packaged into a single Modelfile. Think it is Docker for LLMs. Act like GitHub at some points
- Run Llama 3.1, Phi 3.5, Mistral, Gemma 2, and other models from Hugging Face. Customize and create your own.





Introducing

OpenAI API Capability

Ollama now has built-in compatibility with the OpenAI Chat Completions API, making it possible to use more tooling and applications with Ollama locally.

```
from openai import OpenAI
```

```
client = OpenAI(  
    base_url = 'http://localhost:11434/v1',  
    api_key='ollama', # required, but unused  
)
```

```
response = client.chat.completions.create(  
    model="llama2",  
    messages=[  
        {"role": "system", "content": "You are a helpful  
assistant."}  
    ]  
)  
print(response.choices[0].message.content)
```

JS

```
import OpenAI from 'openai'
```

```
const openai = new OpenAI({  
    baseURL: 'http://localhost:11434/v1',  
    apiKey: 'ollama', // required but unused  
})
```

```
const completion = await  
openai.chat.completions.create({  
    model: 'llama2',  
    messages: [{ role: 'user', content: 'Why is the sky  
blue?' }],  
})
```

```
console.log(completion.choices[0].message.content)
```

Introducing

Download from Hugging Face

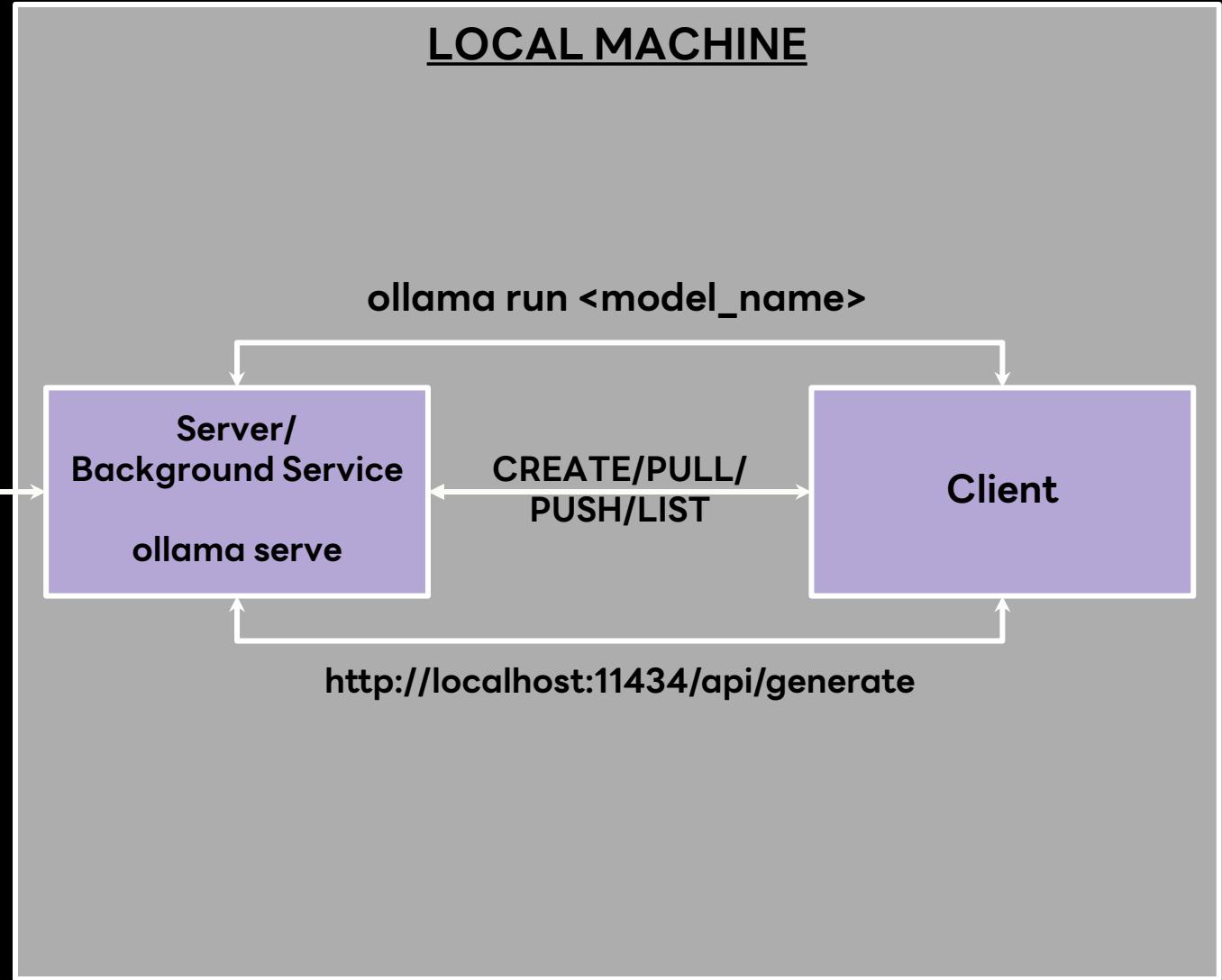
You can use **any GGUF quants** created by the community on Hugging Face directly with Ollama, without creating a new Modelfile. You can run any of them with a single **ollama run** command. We also provide customizations like choosing quantization type to improve your overall experience.



Hugging Face

```
ollama run hf.co/username/repository
```

**OLLM Model Registry
(Remote Server)
or
Hugging Face
(GGUF Model Only)**





First Demo

Ollama Implementation on .NET
with Semantic Kernel



Old
Logo

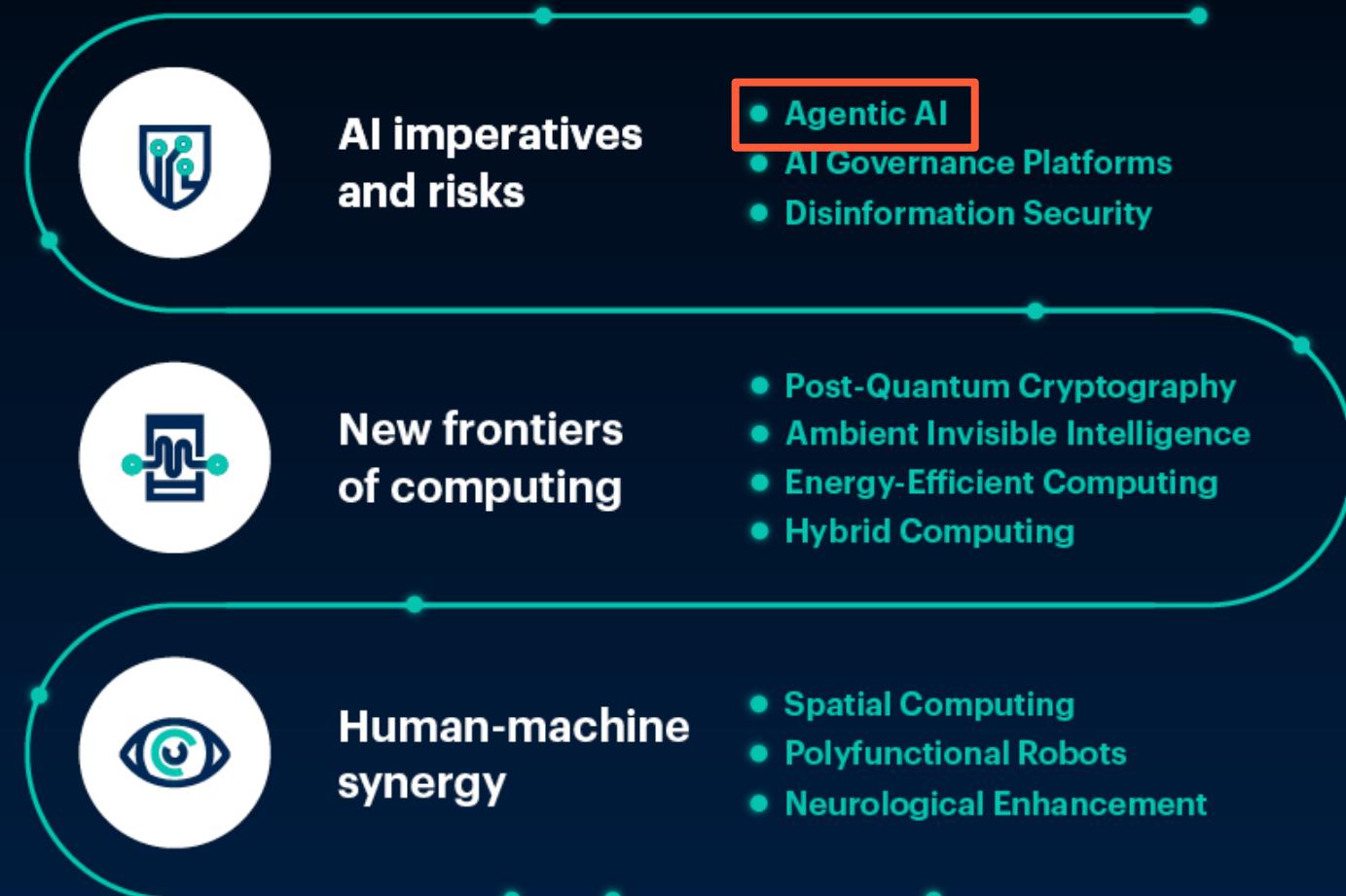


New
Logo



Why
Semantic Kernel important
for .NET developers along to the year 2025?

2025 Top 10 Strategic Technology Trends



Source: Gartner
© 2024 Gartner, Inc. and/or its affiliates.
All rights reserved. 3185862

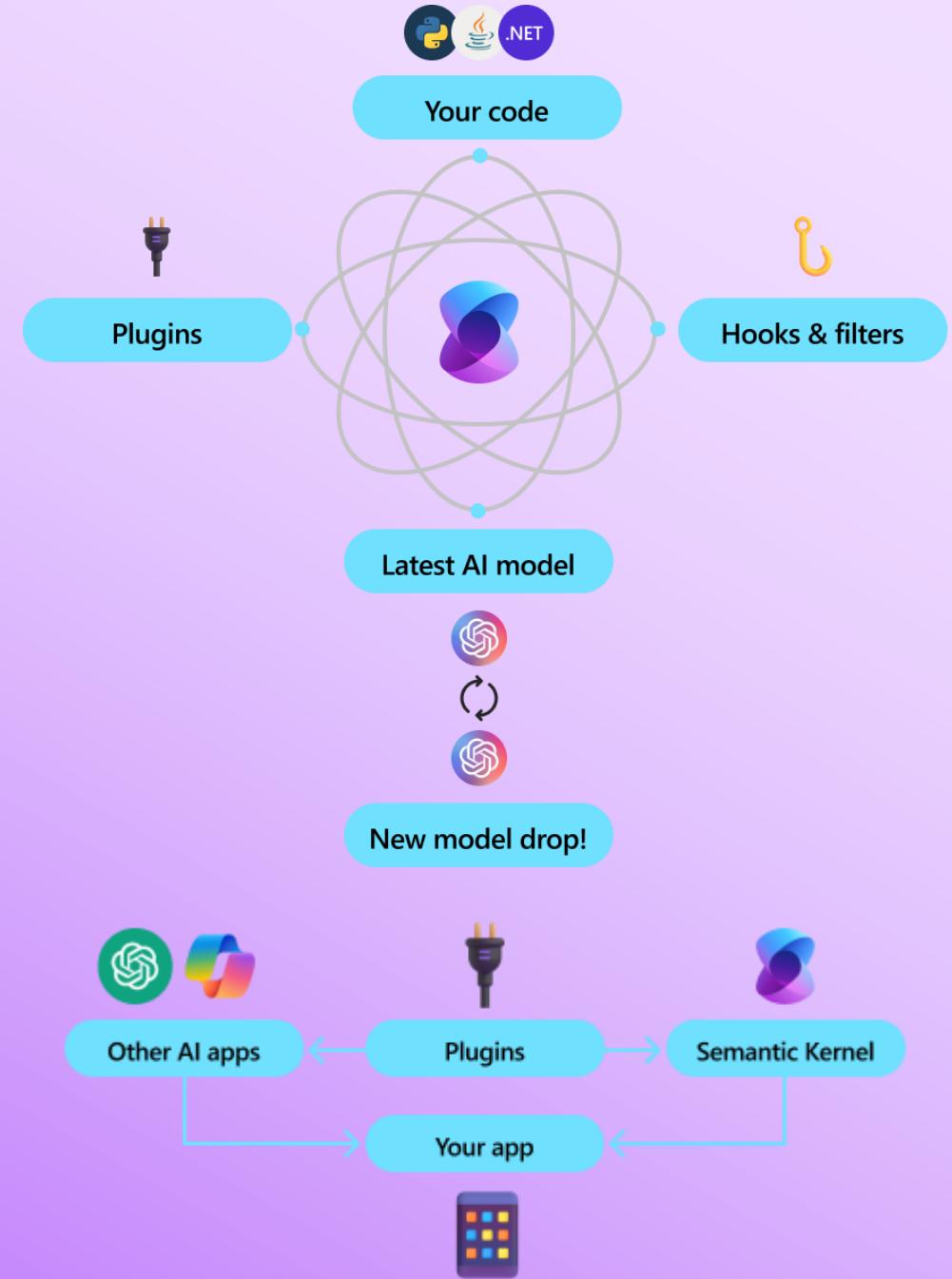
Gartner®

What is AI Agent and Multi-Agent AI System?

- An **AI agent** is an autonomous entity that can perceive its environment, reason about it, and take actions to achieve its goals. It functions through perception, reasoning, and action components.
- A **Multi-Agent AI System** is a system composed of multiple AI agents that interact and collaborate to achieve shared or individual goals. It exhibits decentralization, emergent behavior, collaboration, and competition. Applications include simulation, robotics, game AI, and e-commerce.

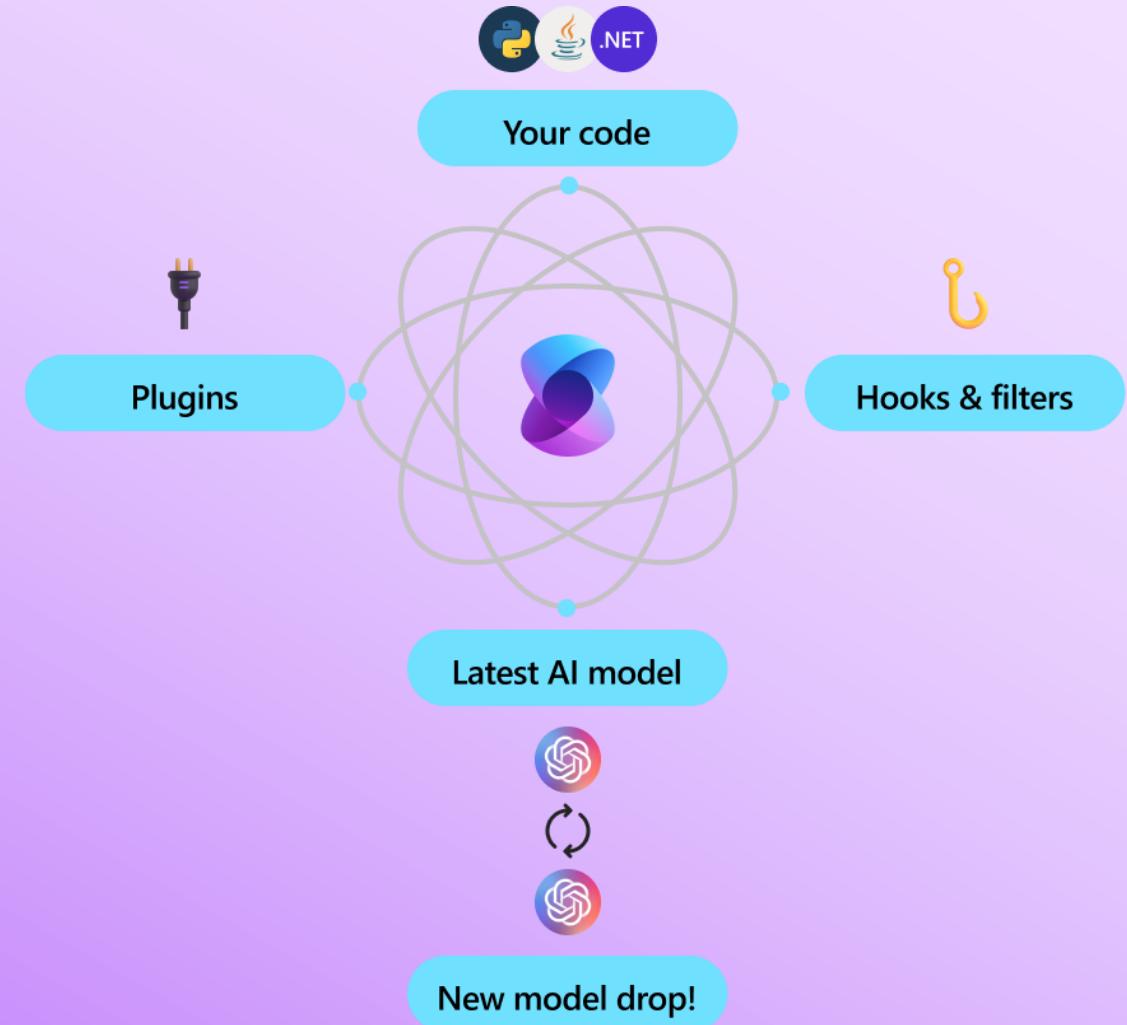
Semantic Kernel

Semantic Kernel is a lightweight, open-source development kit that lets you easily build **AI agents** and integrate the latest AI models into your **C# (.NET)**, Python, or Java codebase. It serves as an efficient middleware that enables rapid delivery of enterprise-grade solutions.



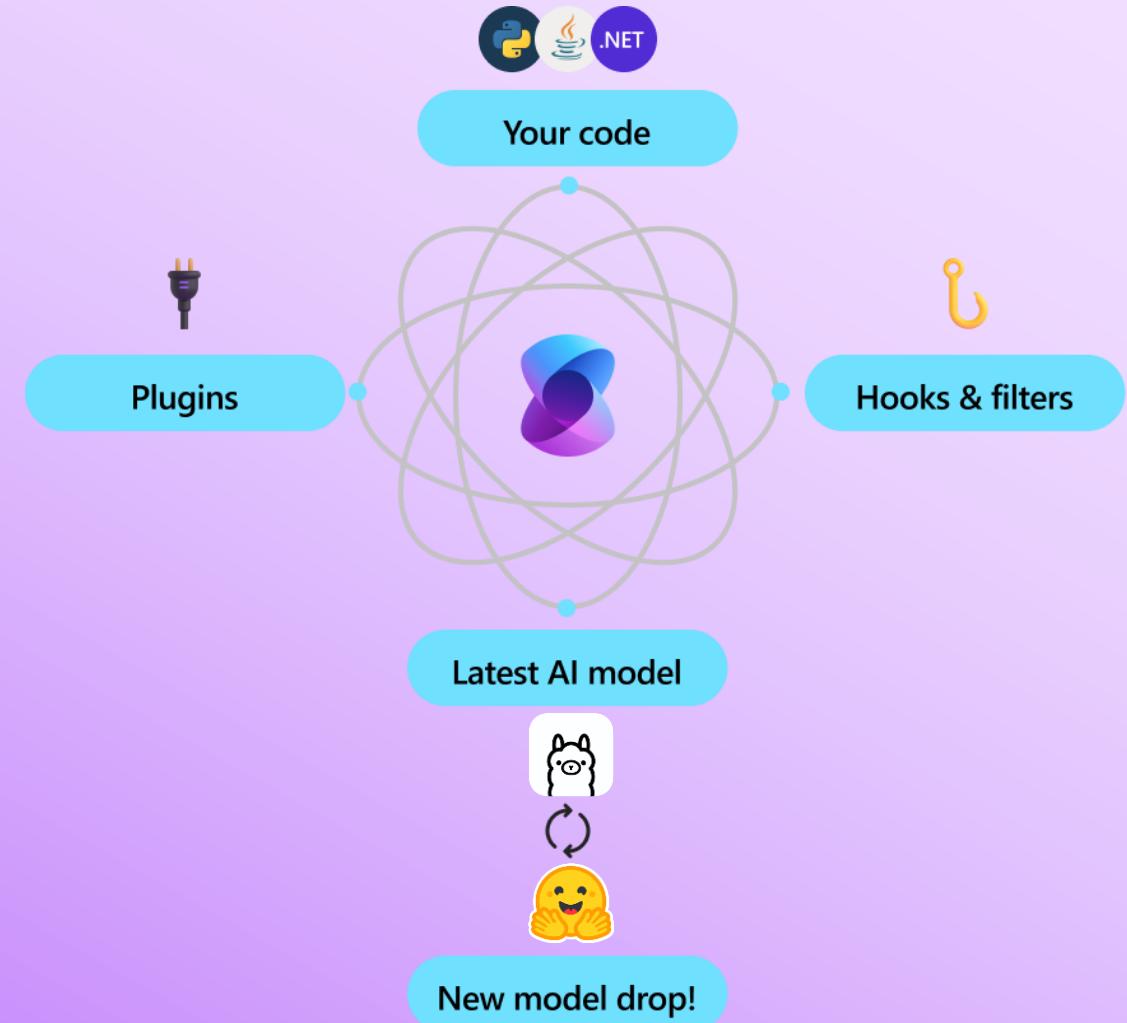
Semantic Kernel

Semantic Kernel is a lightweight, open-source development kit that lets you easily build **AI agents** and integrate the latest AI models into your **C# (.NET)**, Python, or Java codebase. It serves as an efficient middleware that enables rapid delivery of enterprise-grade solutions.



Semantic Kernel

Semantic Kernel is a lightweight, open-source development kit that lets you easily build **AI agents** and integrate the latest AI models into your **C# (.NET)**, Python, or Java codebase. It serves as an efficient middleware that enables rapid delivery of enterprise-grade solutions.



Installing all dependencies in Terminal

```
> dotnet new console  
  
> dotnet add package Microsoft.SemanticKernel  
> dotnet add package Microsoft.SemanticKernel.Connectors.Ollama  
  
> code .      # entering inside Visual Studio Code
```

Declaring Libraries, Models from Ollama (Plug-in)

```
using Microsoft.SemanticKernel;
using Microsoft.SemanticKernel.ChatCompletion;

// Create a kernel with OpenAI chat completion API
// Warning due to the experimental state of some Semantic Kernel SDK
features.

#pragma warning disable SKEXP0070
Kernel kernel = Kernel.CreateBuilder()
    .AddOllamaChatCompletion(
        modelId: "your ollama model name",
        endpoint: new Uri("http://localhost:11434"))
    .Build();
```

Chat Completion + Storing Chat History

```
var aiChatService = kernel.GetRequiredService<IChatCompletionService>();  
var chatHistory = new ChatHistory();  
  
while (true)  
{  
    // Get user prompt and add to chat history  
    Console.WriteLine("Your prompt:");  
    var userPrompt = Console.ReadLine();  
    chatHistory.Add(new ChatMessageContent(AuthorRole.User, userPrompt));
```

Chat Completion + Storing Chat History

```
// Stream the AI response and add to chat history
Console.WriteLine("AI Response:");
var response = "";
await foreach(var item in
    aiChatService.GetStreamingChatMessageContentsAsync(chatHistory))
{
    Console.Write(item.Content);
    response += item.Content;
}
chatHistory.Add(new ChatMessageContent(AuthorRole.Assistant, response));
Console.WriteLine();
}
```

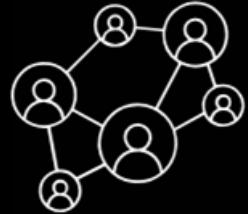


Ollama Implementation with AutoGen for .NET for making the AI agents

AG[★]
★

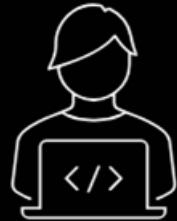
What is AutoGen?

AutoGen An Open-Source Programming Framework for Agentic AI



Multi-Agent Conversation Framework

AutoGen provides multi-agent conversation framework as a high-level abstraction. With this framework, one can conveniently build LLM workflows.



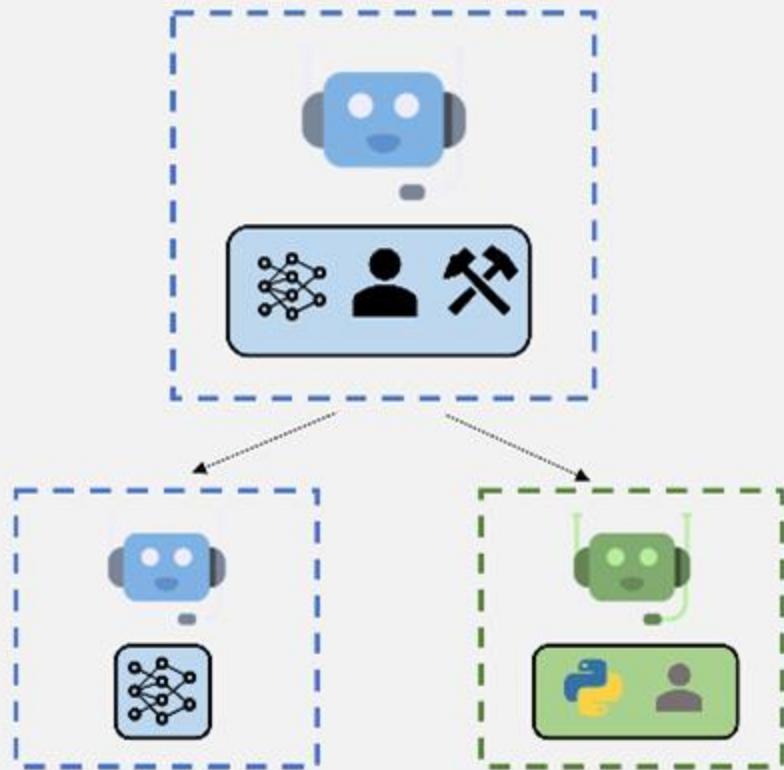
Easily Build Diverse Applications

AutoGen offers a collection of working systems spanning a wide range of applications from various domains and complexities.

Enhanced LLM Inference & Optimization

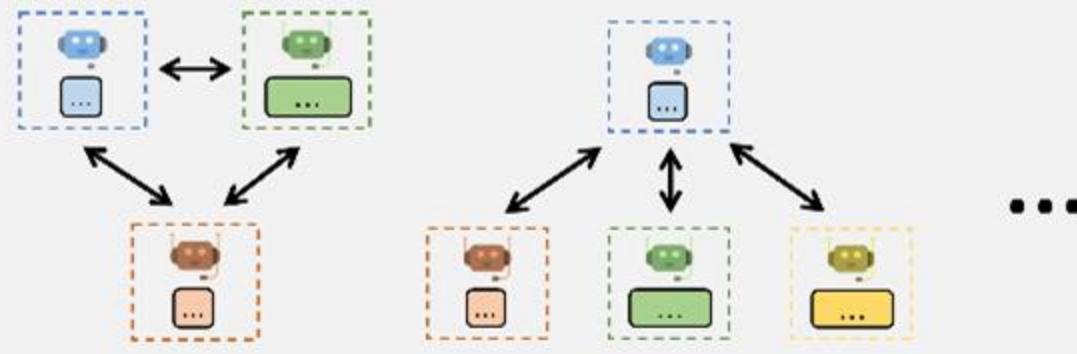
AutoGen supports enhanced LLM inference APIs, which can be used to improve inference performance and reduce cost.

Conversable agent



Agent Customization

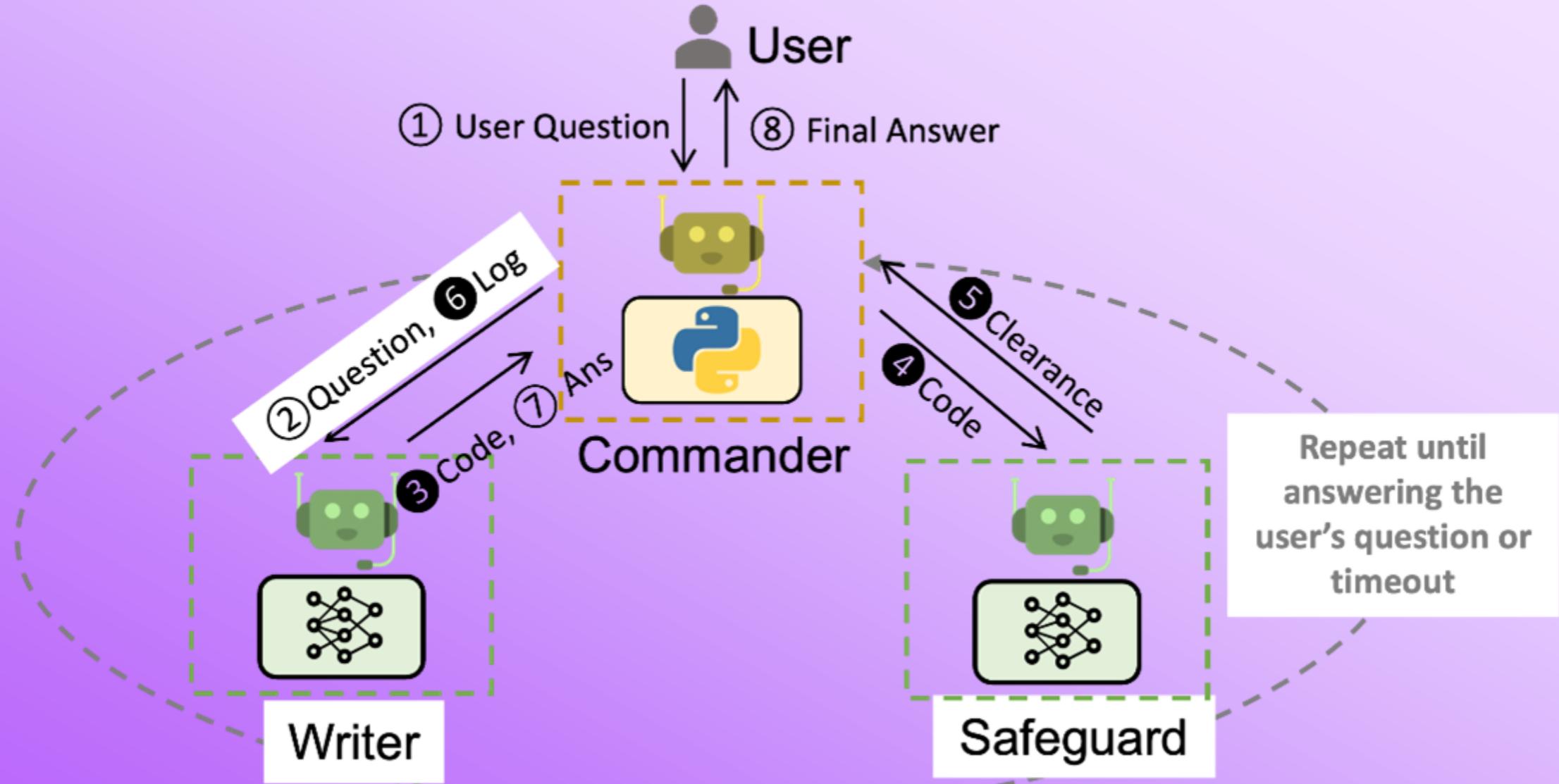
Multi-Agent Conversations



Joint chat

Hierarchical chat

Flexible Conversation Patterns



Some Syntaxes

Running Ollama Model with AutoGen for .NET
and Semantic Kernel

Pre-requisite

First, install the AutoGen.Ollama package using the following command:

```
dotnet add package AutoGen.Ollama
```

Step 2: Add using statement

```
using AutoGen.Core; using AutoGen.Ollama.Extension;
```

Create OllamaAgent: Text Based

```
using var httpClient = new HttpClient()
{
    BaseAddress = new Uri("http://localhost:11434"),
};

var ollamaAgent = new OllamaAgent(
    httpClient: httpClient,
    name: "ollama",
    modelName: "llama3:latest",
    systemMessage: "You are a helpful AI assistant")
    .RegisterMessageConnector()
    .RegisterPrintMessage();

var reply = await ollamaAgent.SendAsync("Can you write a piece of C# code to
calculate 100th of fibonacci?");
```

Create OllamaAgent: Multimodal Based (VLM)

```
using var httpClient = new HttpClient()
{
    BaseAddress = new Uri("http://localhost:11434"),
};

var ollamaAgent = new OllamaAgent(
    httpClient: httpClient,
    name: "ollama",
    modelName: "llava-phi3",
    systemMessage: "You are a helpful AI assistant")
    .RegisterMessageConnector()
    .RegisterPrintMessage();

var image = Path.Combine("resource", "images", "background.png");
var binaryData = BinaryData.FromBytes(File.ReadAllBytes(image), "image/png");
var imageMessage = new ImageMessage(Role.User, binaryData);
var textMessage = new TextMessage(Role.User, "what's in this image?");
var reply = await ollamaAgent.SendAsync(chatHistory: [textMessage, imageMessage]);
```

Code and Slides are available on my GitHub

chrnthonkmutt/ **BoatSLM_cs_experiment**

This repository is used to make the experiment of using ollama, along with custom small language model, in .NET framework

1

Contributor

0

Issues

0

Stars

0

Forks

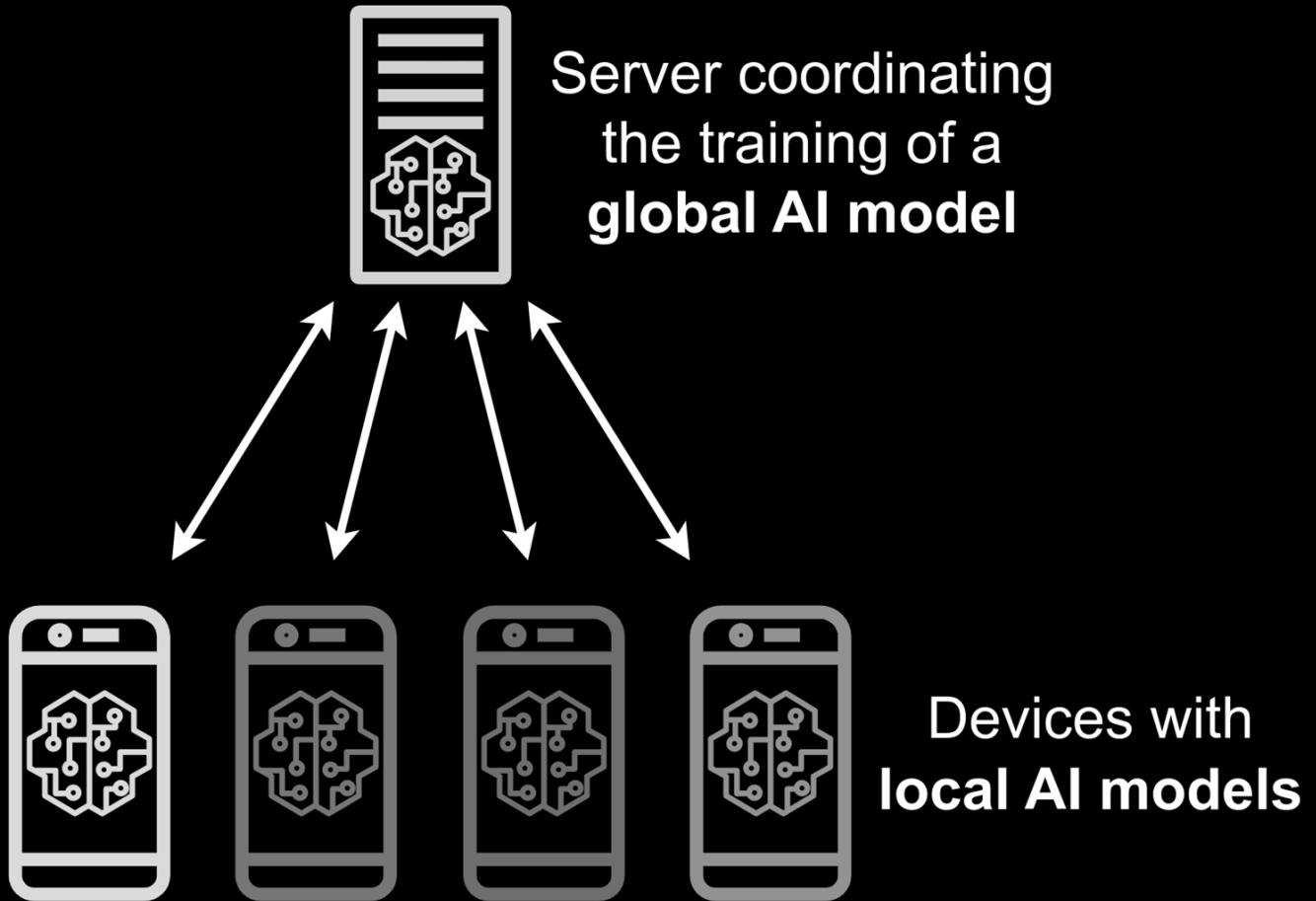


Key Takeaways

Upcoming Trends of Generative AI and NLPs Industry

- Shift to **multimodality**: Text → text, images, pictures, speech
- **Explainability** (Reason/Evidence for why the output is produced)
Interpretability (understanding how the model works)
Trustworthiness (safety and transparency of model usage)
- **Model Miniaturization**: separation between **language fluency** and **knowledge base** for AI on small devices (Laptop, Smartphones) with Multi Agents Compatibility

Federated Learning



The question is...

Why Thai People Still Doesn't Use Small Language Models?

Hallucination
& Performance

Limited Device
Resources

Languages Barrier

Low
Engagement?

Let me be there for you ❤️
to learn more on new innovations on the world

Our New Finetuned Small Language Model For Thai People



will be launched and available on **National
Coding Day next week in Hugging Face**





Hundreds of Thai Data & AI
Community Leaders are in 2024

We hope to see you more than
a million leaders and trainers in
2025 as a part of us :)

Escaping Yourself from
“Know-it-all” to **“Learn-it-all”**

Thank You!



Charunthon Limseelo



@boatchrnthn



Charunthon Limseelo



Boat Charunthon (boatchrnthn)

