

Microsoft Learn
Student Ambassadors



Small Language Models, Developing AI for Sustainable Future and Their Applications

Charunthon Limseelo (@boatchrnthn)

Microsoft Learn Student Ambassador | Thai Student Tech Community Lead
Computer Engineering, King Mongkut's University of Technology Thonburi

#MSFTStudentambassadors #SDGs

Retrieved from A Comprehensive Survey of Small Language Models in the Era of Large Language Models: Techniques, Enhancements, Applications, Collaboration with LLMs, and Trustworthiness



Charunthon Limseelo (Boat)

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

- + Microsoft Certified Office Specialist (Microsoft 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 of AI



Business, Financial, Commercial-based, and Daily Users of AI



3 Tracks will be discussed in 30 Minutes

1

**Pain Points of Using
Large Language
Models**

2

**Small Language
Models**
Capabilities + Use Cases

3

**Microsoft Learn
Student Ambassadors:**
Pathway to New Models

1



Pain points from using **Large Language Models**

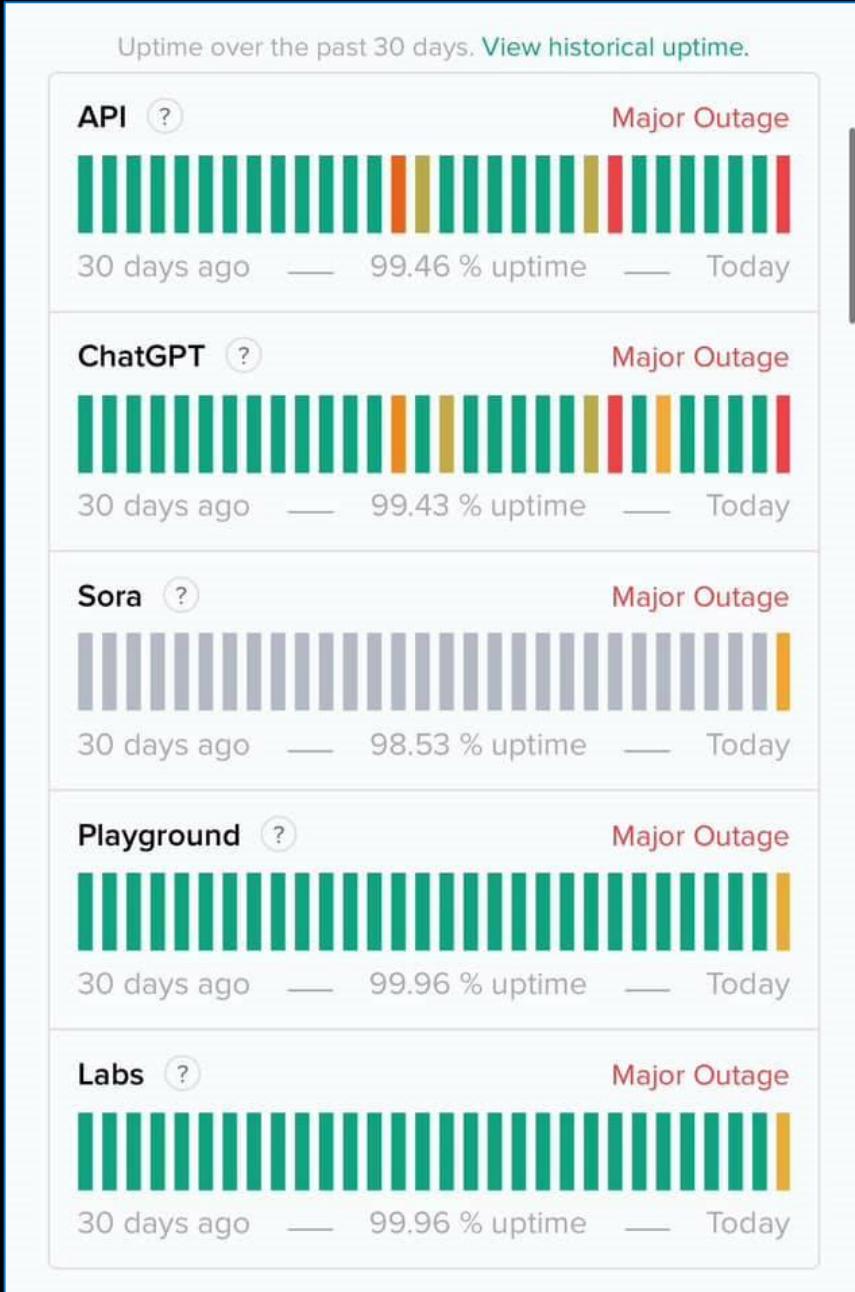
Proprietary



(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



ChatGPT is currently unavailable.

Status: Identified - We have identified the issue and are working to roll out a fix.

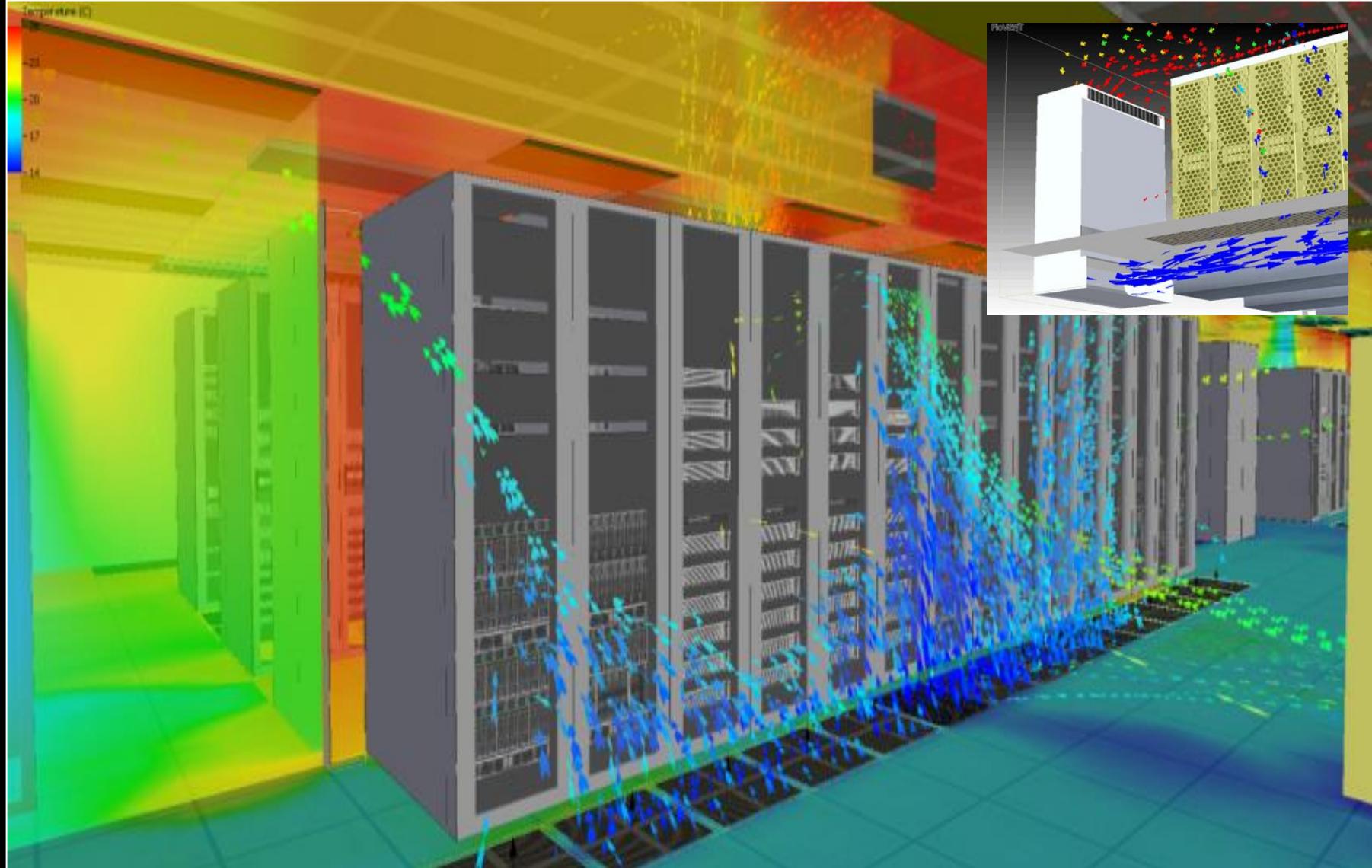
Major Outage Incident of OpenAI ecosystem on December 12th, 2024. (Dec 11 by US Time Zone)

11 Dec 2024

Major outage 4 hrs 10 mins

RELATED

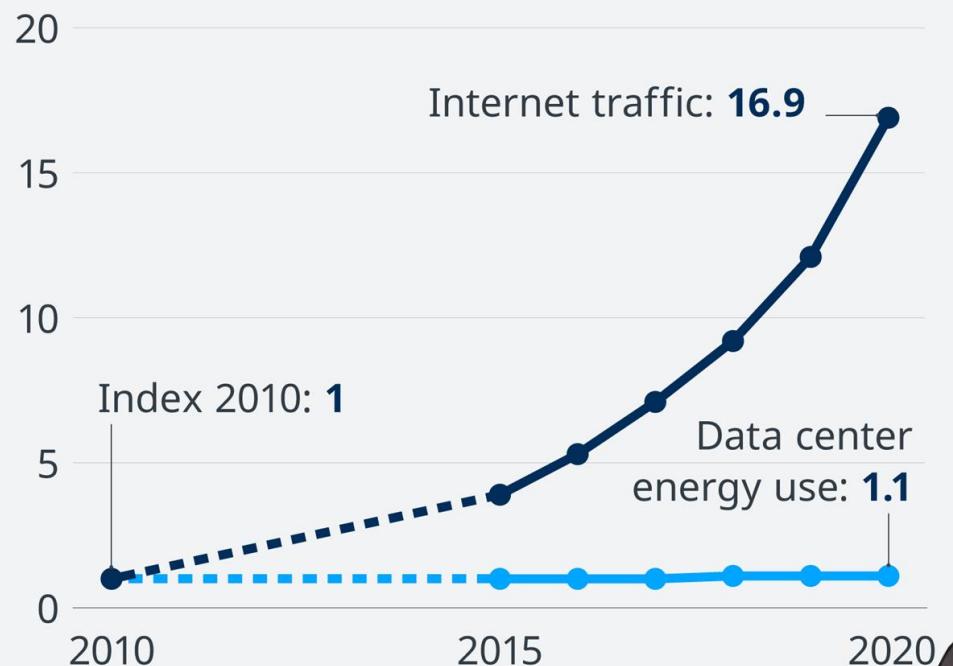
[API, ChatGPT & Sora Facing Issues](#)



Thermal Simulation Software on Visualizing Heat Flow of Datacenter

Global trends in internet traffic and data center energy use

Internet traffic and data center energy use compared to 2010



Source: IEA



Source: Enerdata, IEA

Data centers use more electricity than entire countries

Domestic electricity consumption of selected countries vs. data centers in 2020 in TWh



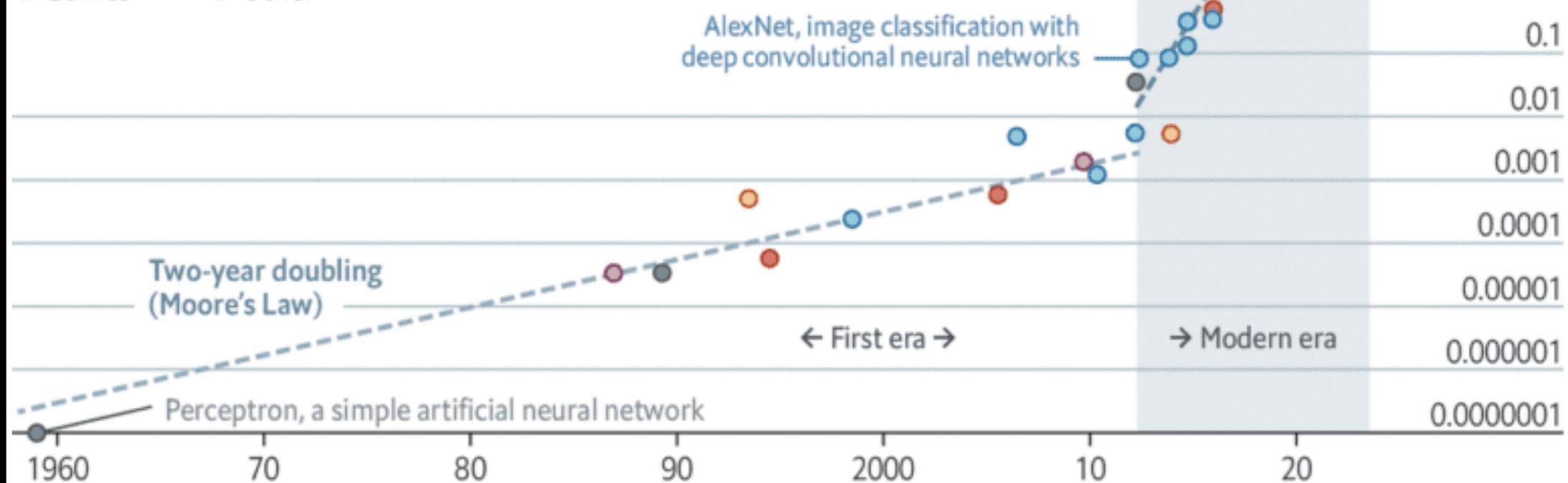
Deep and steep

Computing power used in training AI systems

Days spent calculating at one petaflop per second*, log scale

By fundamentals

- Language
- Games
- Speech
- Other
- Vision



Source: OpenAI

*1 petaflop=10¹⁵ calculations

The Economist



Climate change is the defining issue of our generation, and addressing it requires swift, collective action and technological innovation. We are committed to meeting our own goals while enabling others to do the same. That means taking responsibility for our operational footprint and accelerating progress through technology.

- **Satya Nadella**, CEO of Microsoft

Tasks Libraries Datasets Languages Licenses
Other

Filter Tasks by name

Multimodal

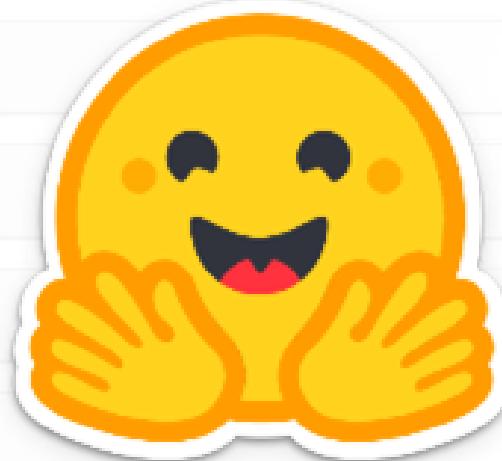
- Image-Text-to-Text
- Visual Question Answering
- Document Question Answering
- Video-Text-to-Text
- Any-to-Any

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

Big Tech Server's Down

Prevent from Paying Premium Services



No internet

Try:

- Checking the network cables, modem, and router
- Reconnecting to Wi-Fi

ERR_INTERNET_DISCONNECTED



Tasks Libraries Datasets Languages Licenses

Other

Filter Tasks by name

Multimodal

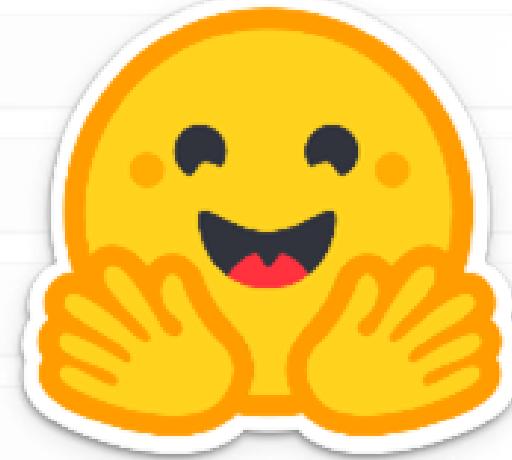
- Image-Text-to-Text
- Visual Question Answering
- Document Question Answering
- Video-Text-to-Text
- Any-to-Any

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

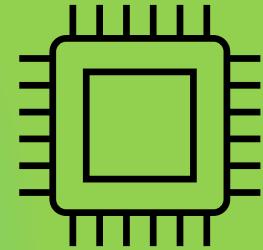


Found out some LLMs are too big to download, need to use API to connect (Require internet)

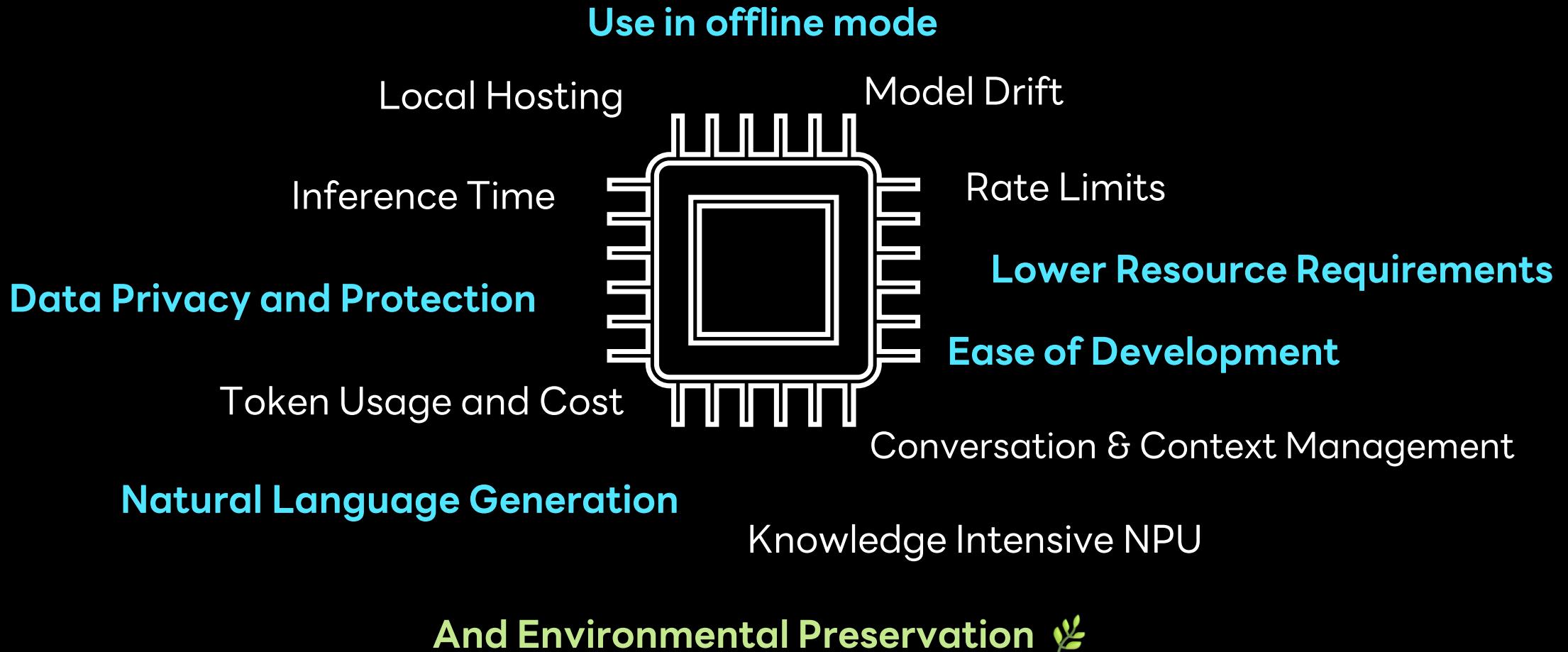
I have an interesting way to prevent all these problems for using each device offline!

2

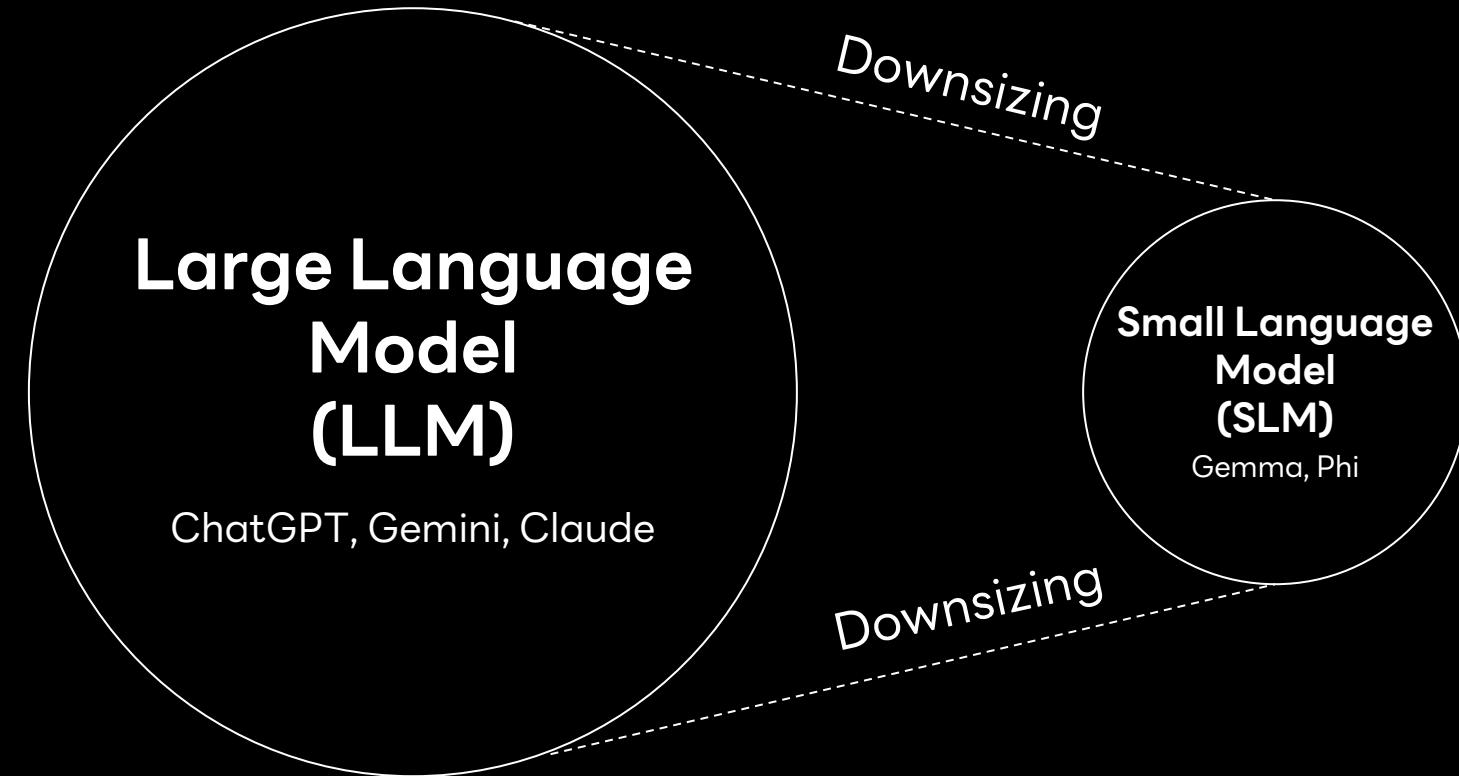
ทำความรู้จักกับ Small Language Models



What is Small Language Model?

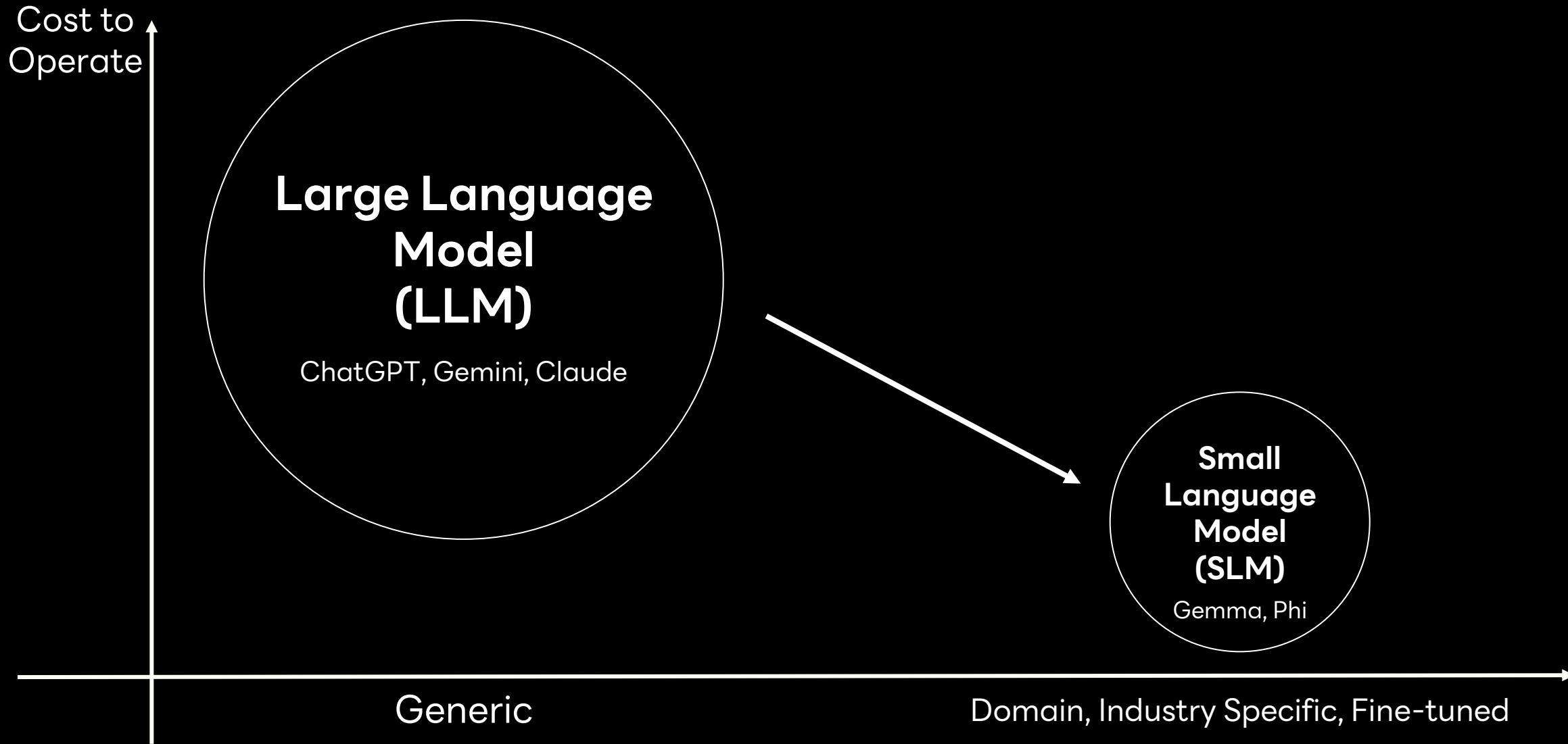


Language Model Comparison



- 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

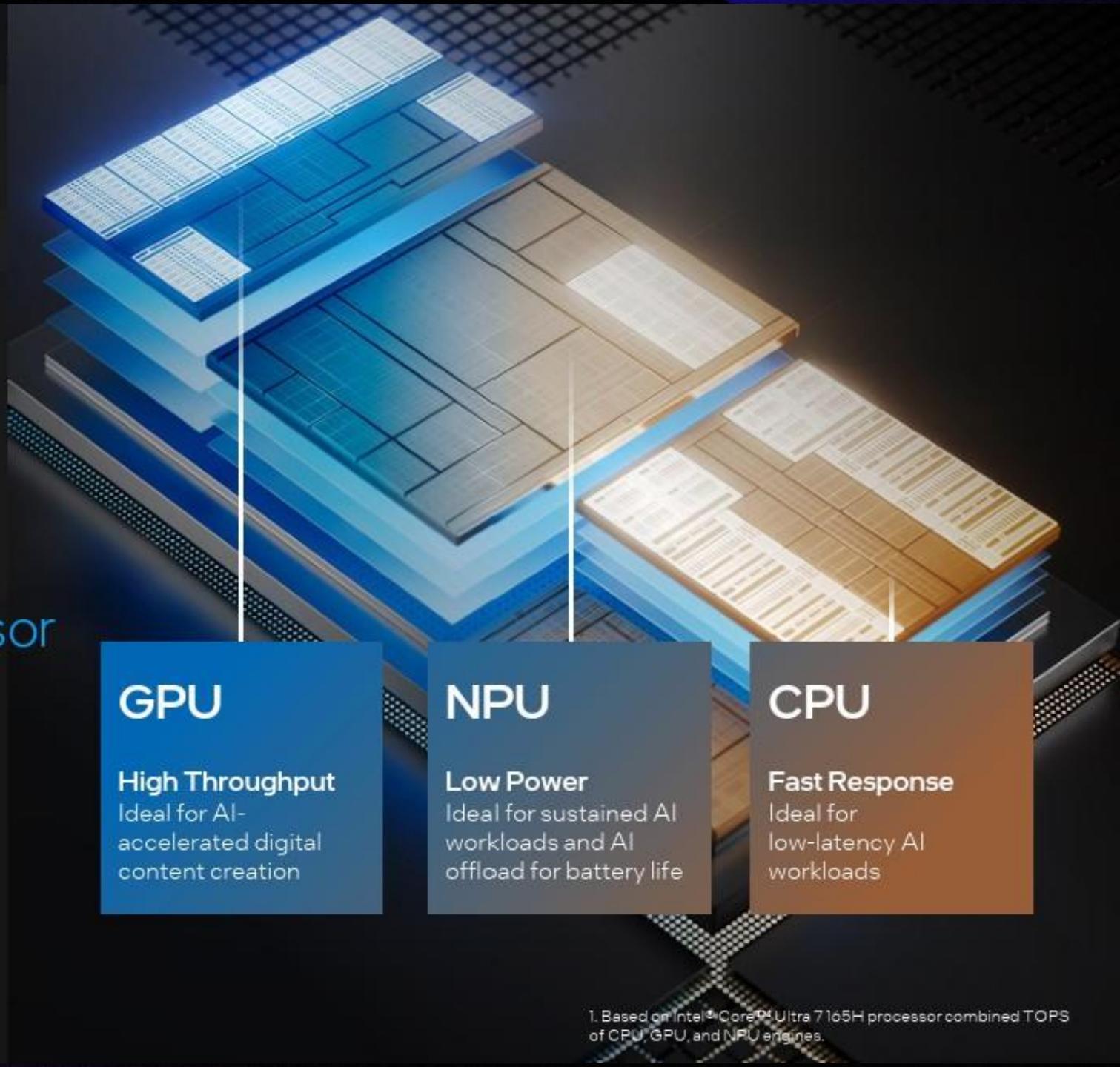


Three AI Engines

with Intel® Core™ Ultra Processor

Heterogenous execution of AI workloads embraces the best practices in AI software design

Deliver up to **34 TeraOPS¹**



GPU

High Throughput
Ideal for AI-accelerated digital content creation

NPU

Low Power
Ideal for sustained AI workloads and AI offload for battery life

CPU

Fast Response
Ideal for low-latency AI workloads

1. Based on Intel® Core™ Ultra 7 165H processor combined TOPS of CPU, GPU, and NPU engines.

Potential Applications of SLMs in Nowadays



Phone Smart
Device



Smart Home
Devices



Wearable
Technologies



Automative
System



Educational
Tools



Entertainment
Systems

Khanmigo from Khan Academy

The screenshot displays the Khanmigo Tools interface. At the top, a banner reads "Khanmigo Tools" and "Free AI powered tools designed to save you time and improve instruction!". Below the banner is a navigation bar with buttons for "All", "Plan", "Create", "Support", "Learn", and "Students". A search bar is also present. The main area contains ten tool cards arranged in two columns:

- Lesson Plan**: Create structured, detailed lesson plans tailored to your curriculum and students' needs.
- Leveler**: Adjust the complexity of a given text.
- Chunk Text**: Break complex texts into manageable sections for easier student comprehension.
- Clear Directions**: Generate concise, easy-to-follow instructions for assignments and activities.
- Exit Ticket**: Create quick end-of-lesson assessments to check student understanding.
- Learning Objective(s)**: Develop clear, measurable learning objectives to guide instruction.
- Real World Context Generator**: Connect lesson topics to engaging real-world examples and applications.
- Multiple Choice Assessment**: Create multiple-choice assessments on a variety of topics.
- Questions Generator**: Create questions for a specific piece of content.
- Discussion Prompts**: Craft engaging prompts to stimulate meaningful classroom discussions.
- Lesson Hook**: Plan compelling lesson starters to engage students.
- Make It Relevant**: Link lesson content to students' lives and interests to boost engagement.

Khanmigo: Free, AI-powered teacher assistant by Khan Academy



Copilot in NX X, by Siemens



Vinit Shukla
What is the significance of styling data in the BIW Design process?

Copilot

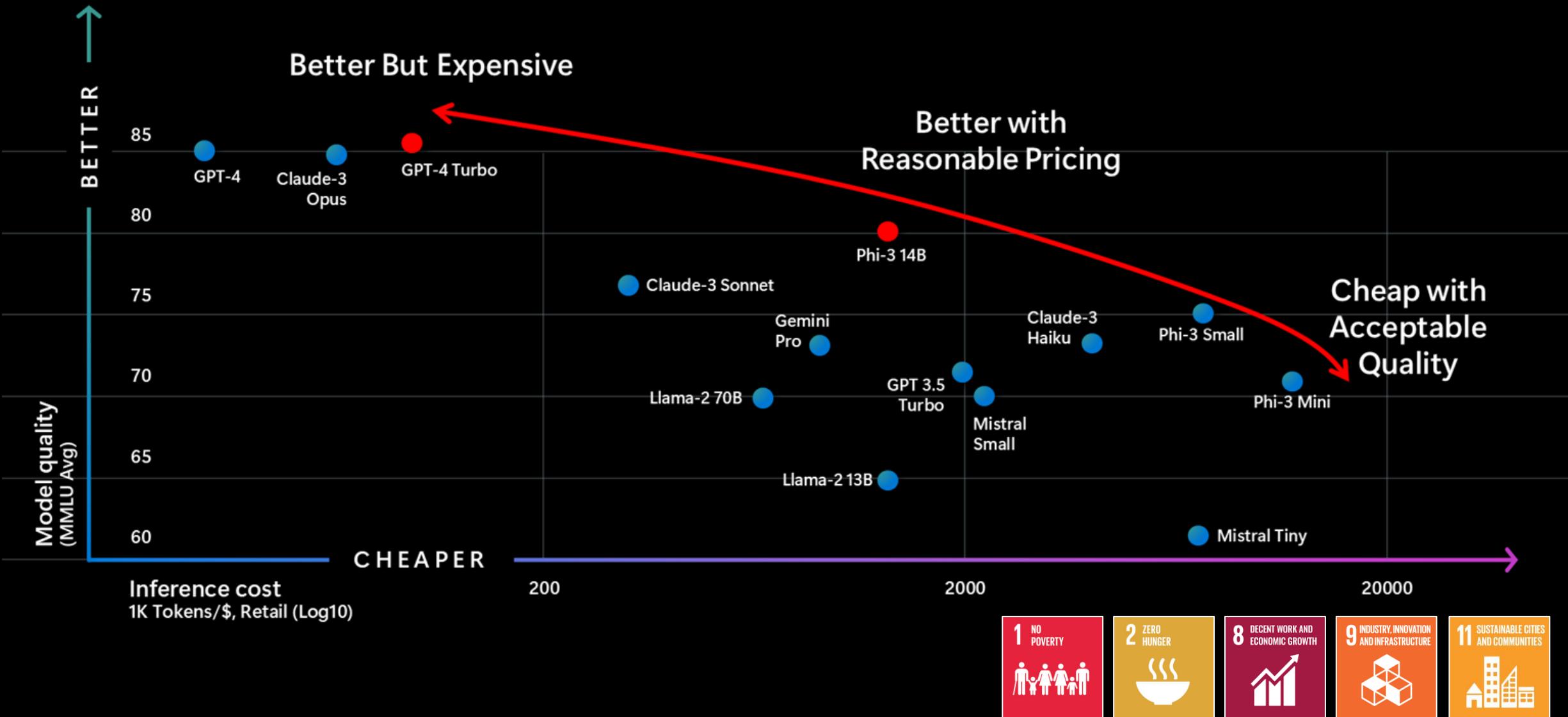
The significance of styling data in the Body-in-White (BIW) design process lies in its role as the foundation for the vehicle's outer body shell. Styling data, which includes the outer surface geometry, is crucial for defining the aesthetic and functional aspects of the vehicle. It guides the development of interface data, which connects the styling surfaces to the structural components. Key points include:

- Initial Concept:** Styling data represents the vehicle's intended appearance, influencing design decisions.
- Engineering Criteria:** It provides essential information for engineers to develop structural components while adhering to material and manufacturing constraints.
- Data Transformation:** Scanned clay models or rough styling data are refined into high-quality CAD surfaces, enabling accurate structural design.

In summary, styling data is integral to ensuring that the BIW design aligns with both aesthetic and engineering requirements.

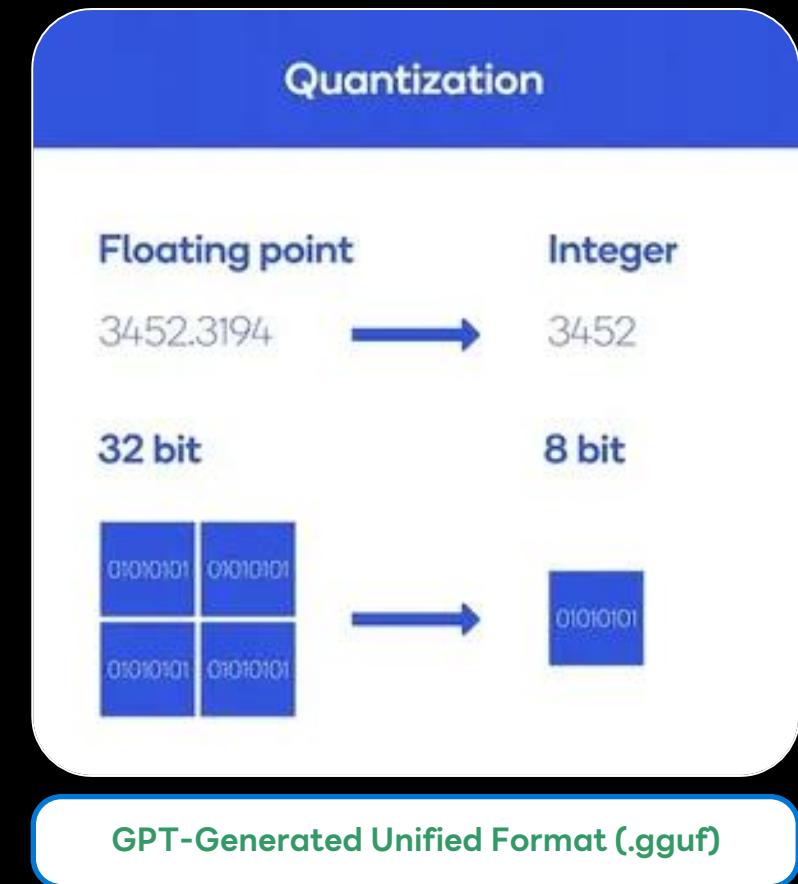
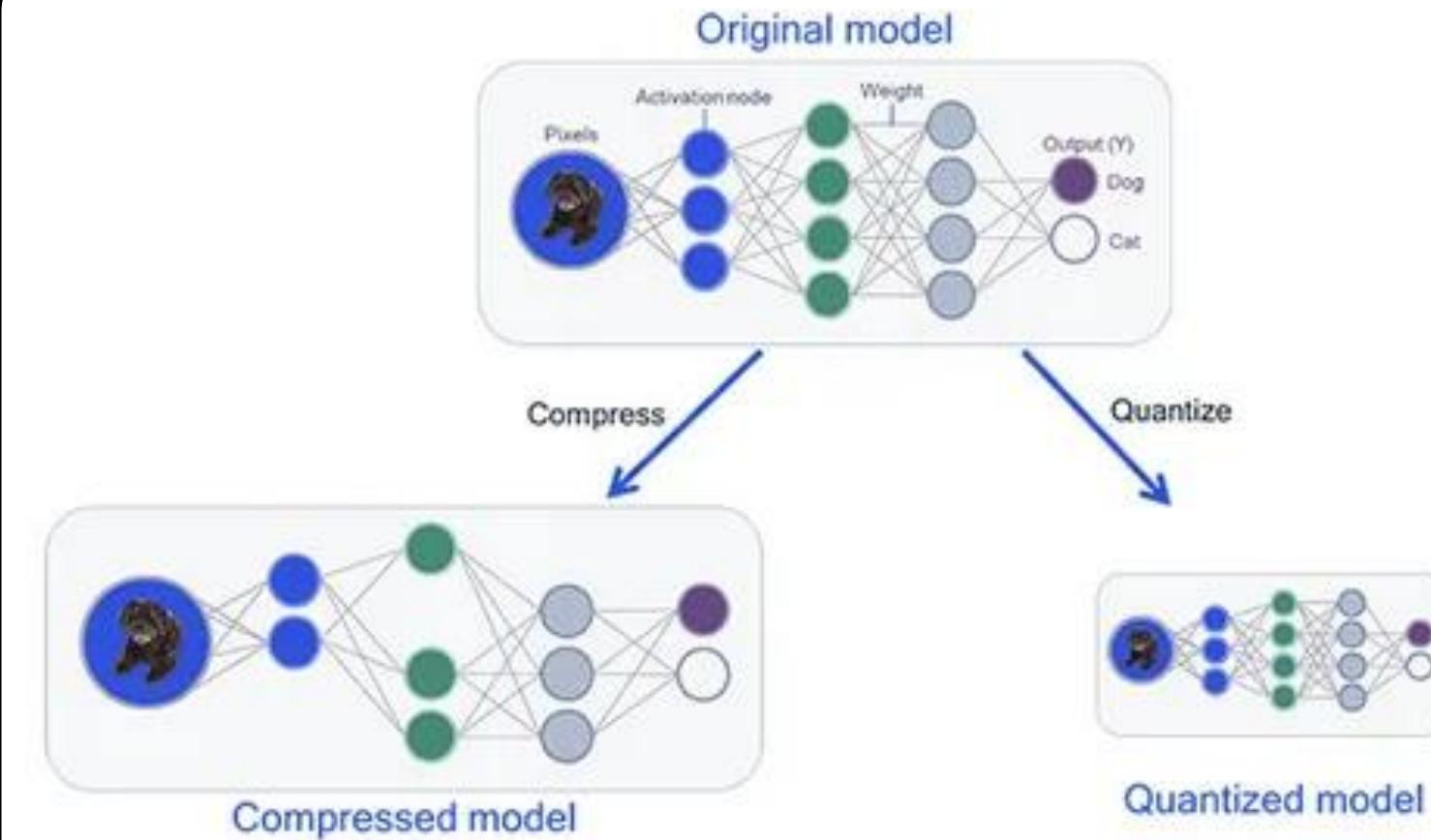
Enter your prompt here

Cost Estimation of Tokenization + Quality



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)

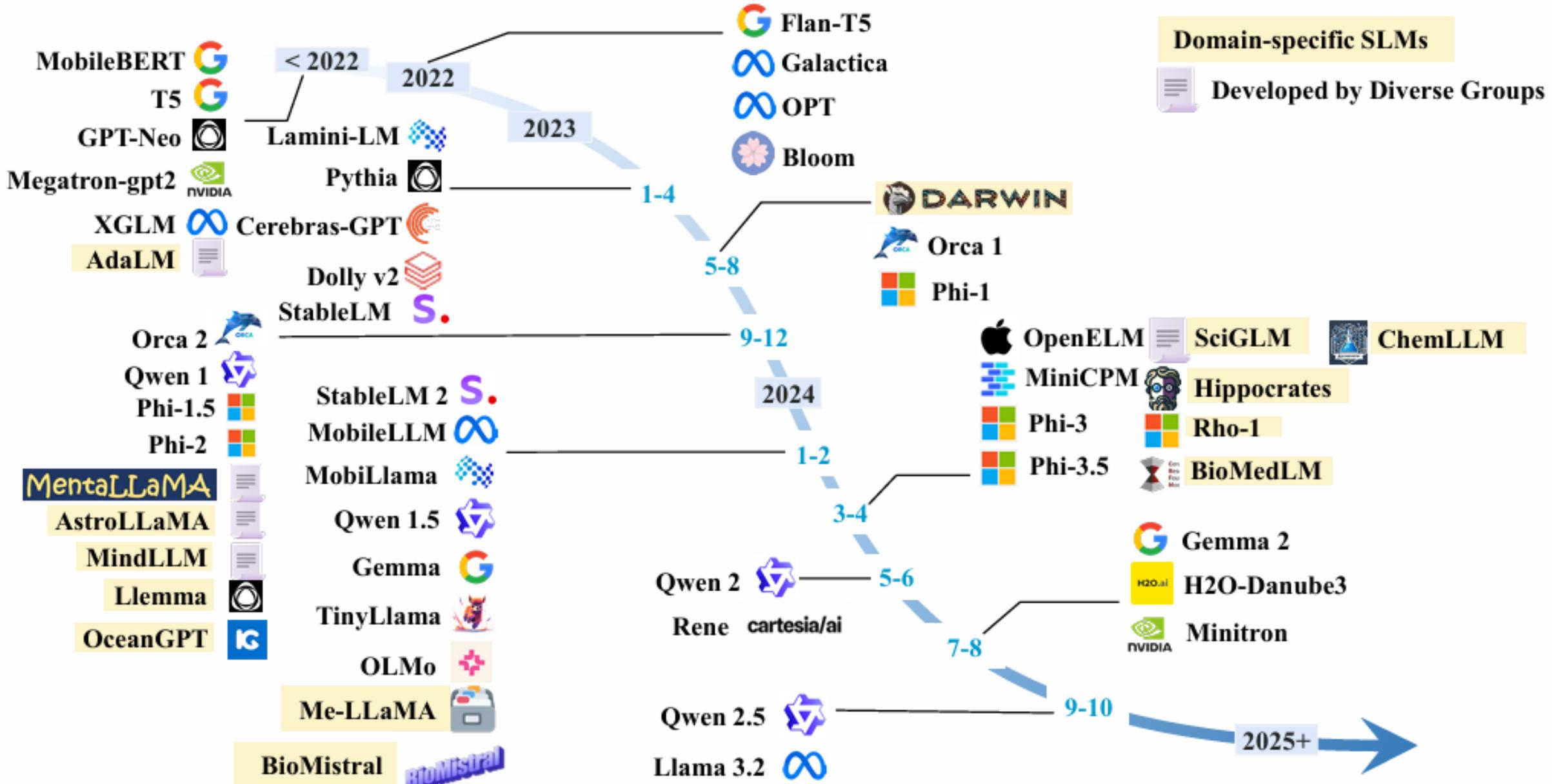
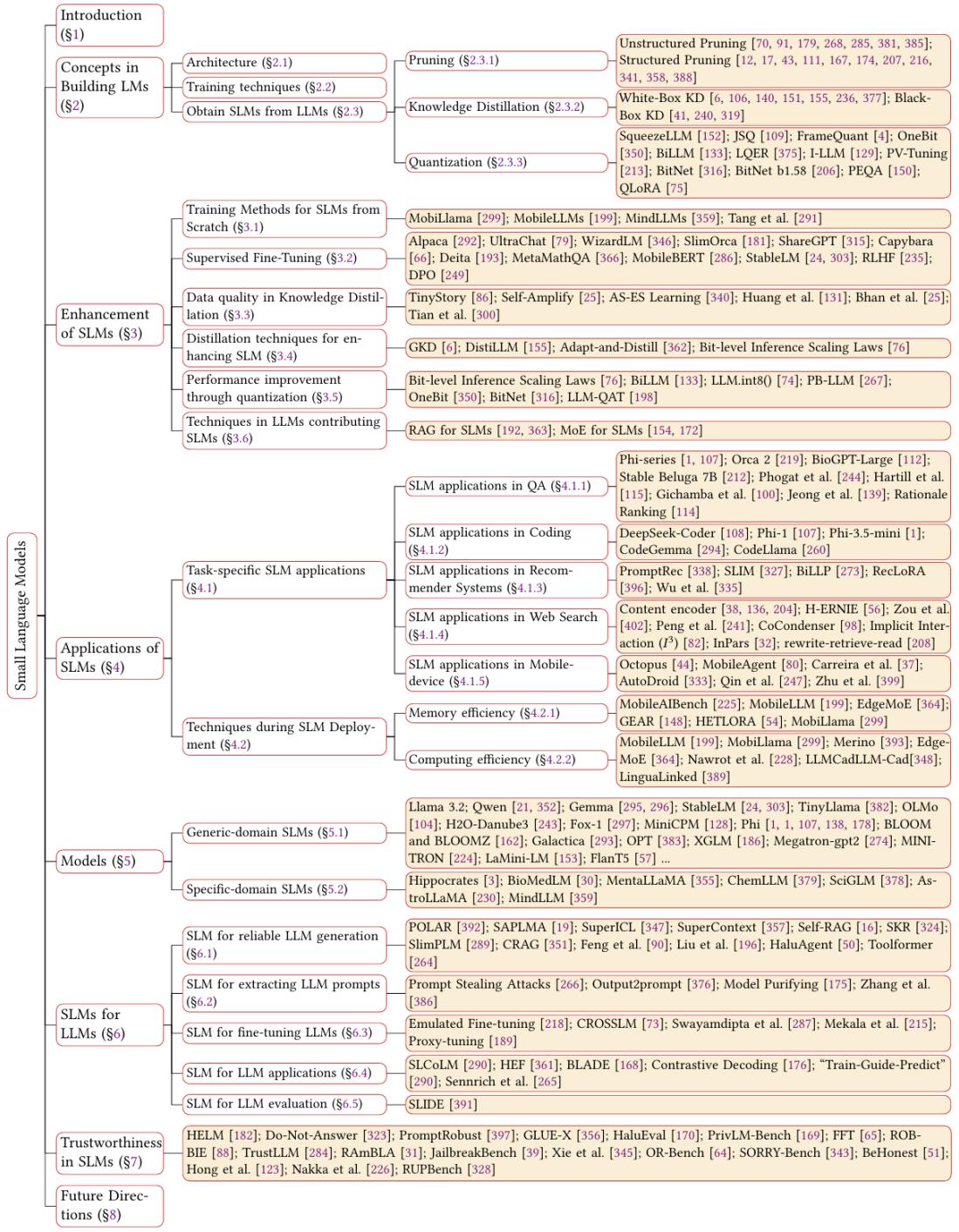
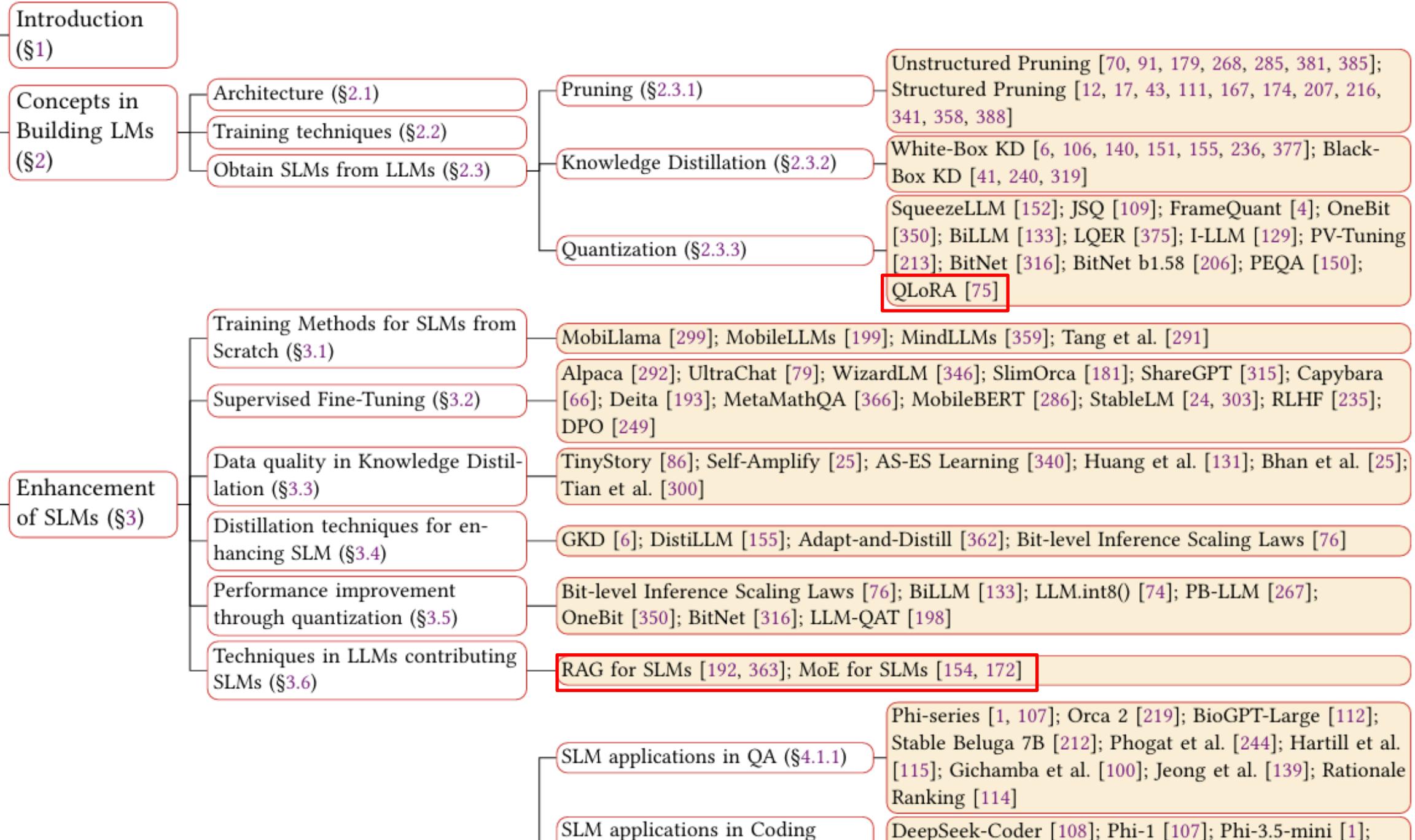


Fig. 3. A timeline of existing small language models.





SLMs (§3.6)

RAG for SLMs [192, 363], MoL for SLMs [154, 172]

Phi-series [1, 107]; Orca 2 [219]; BioGPT-Large [112]; Stable Beluga 7B [212]; Phogat et al. [244]; Hartill et al. [115]; Gichamba et al. [100]; Jeong et al. [139]; Rationale Ranking [114]

Task-specific SLM applications
(&4.1)

SLM applications in QA (&4.1.1)

DeepSeek-Coder [108]; Phi-1 [107]; Phi-3.5-mini [1]; CodeGemma [294]; CodeLlama [260]

SLM applications in Coding
(&4.1.2)

PromptRec [338]; SLIM [327]; BiLLP [273]; RecLoRA [396]; Wu et al. [335]

SLM applications in Recommender Systems (&4.1.3)

Content encoder [38, 136, 204]; H-ERNIE [56]; Zou et al. [402]; Peng et al. [241]; CoCondenser [98]; Implicit Interaction (I^3) [82]; InPars [32]; rewrite-retrieve-read [208]

SLM applications in Web Search
(&4.1.4)

Octopus [44]; MobileAgent [80]; Carreira et al. [37]; AutoDroid [333]; Qin et al. [247]; Zhu et al. [399]

SLM applications in Mobile-device (&4.1.5)

MobileAIBench [225]; MobileLLM [199]; EdgeMoE [364]; GEAR [148]; HETLORA [54]; MobiLlama [299]

MobileLLM [199]; MobiLlama [299]; Merino [393]; EdgeMoE [364]; Nawrot et al. [228]; LLMcadLLM-Cad [348]; LinguaLinked [389]

Applications of SLMs (&4)

Techniques during SLM Deployment (&4.2)

Memory efficiency (&4.2.1)

Computing efficiency (&4.2.2)

Models (&5)

Generic-domain SLMs (&5.1)

Llama 3.2; Qwen [21, 352]; Gemma [295, 296]; StableLM [24, 303]; TinyLlama [382]; OLMo [104]; H2O-Danube3 [243]; Fox-1 [297]; MiniCPM [128]; Phi [1, 1, 107, 138, 178]; BLOOM and BLOOMZ [162]; Galactica [293]; OPT [383]; XGLM [186]; Megatron-gpt2 [274]; MINI-TRON [224]; LaMini-LM [153]; FlanT5 [57] ...

Specific-domain SLMs (&5.2)

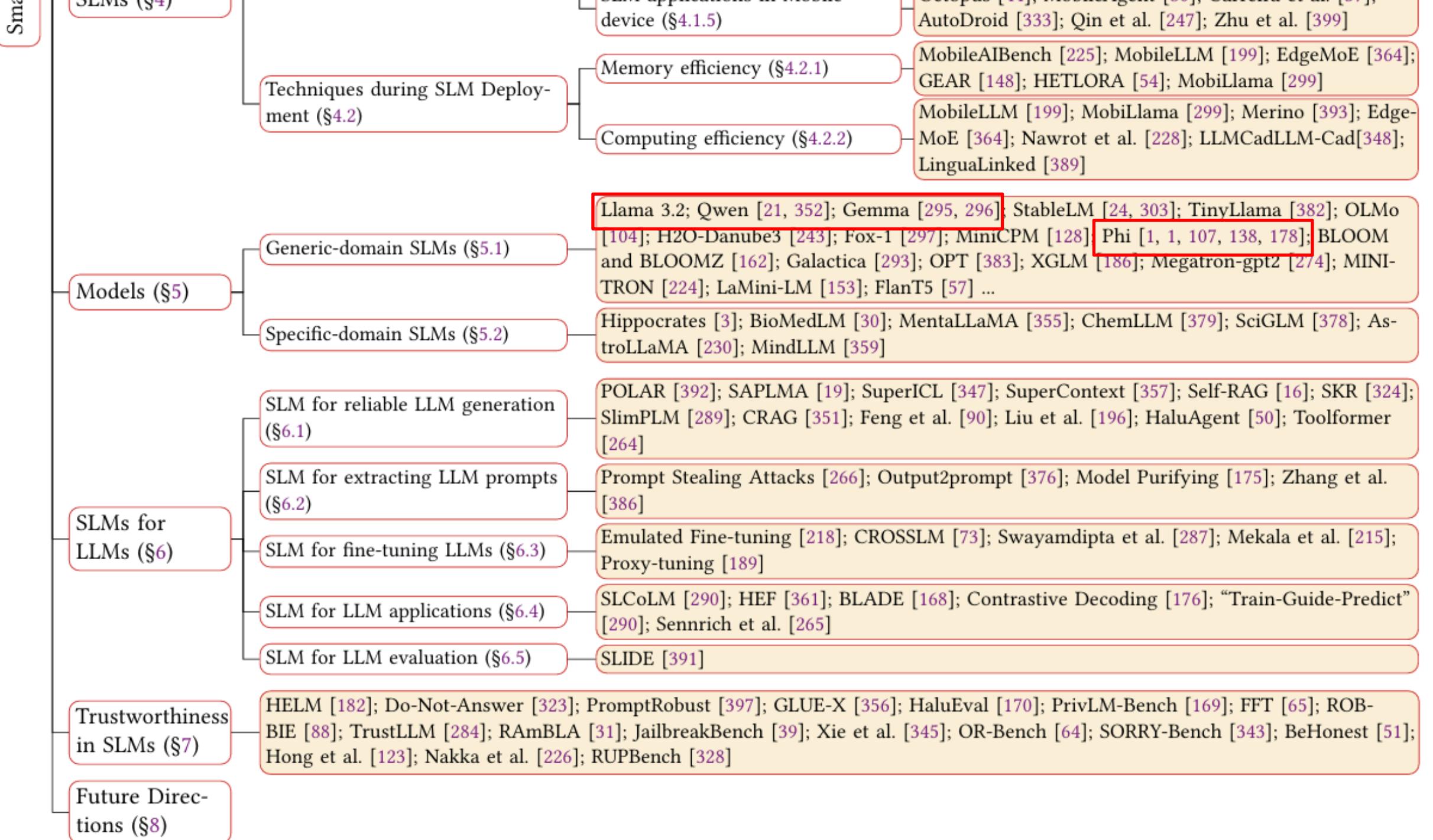
Hippocrates [3]; BioMedLM [30]; MentaLLaMA [355]; ChemLLM [379]; SciGLM [378]; AstroLLaMA [230]; MindLLM [359]

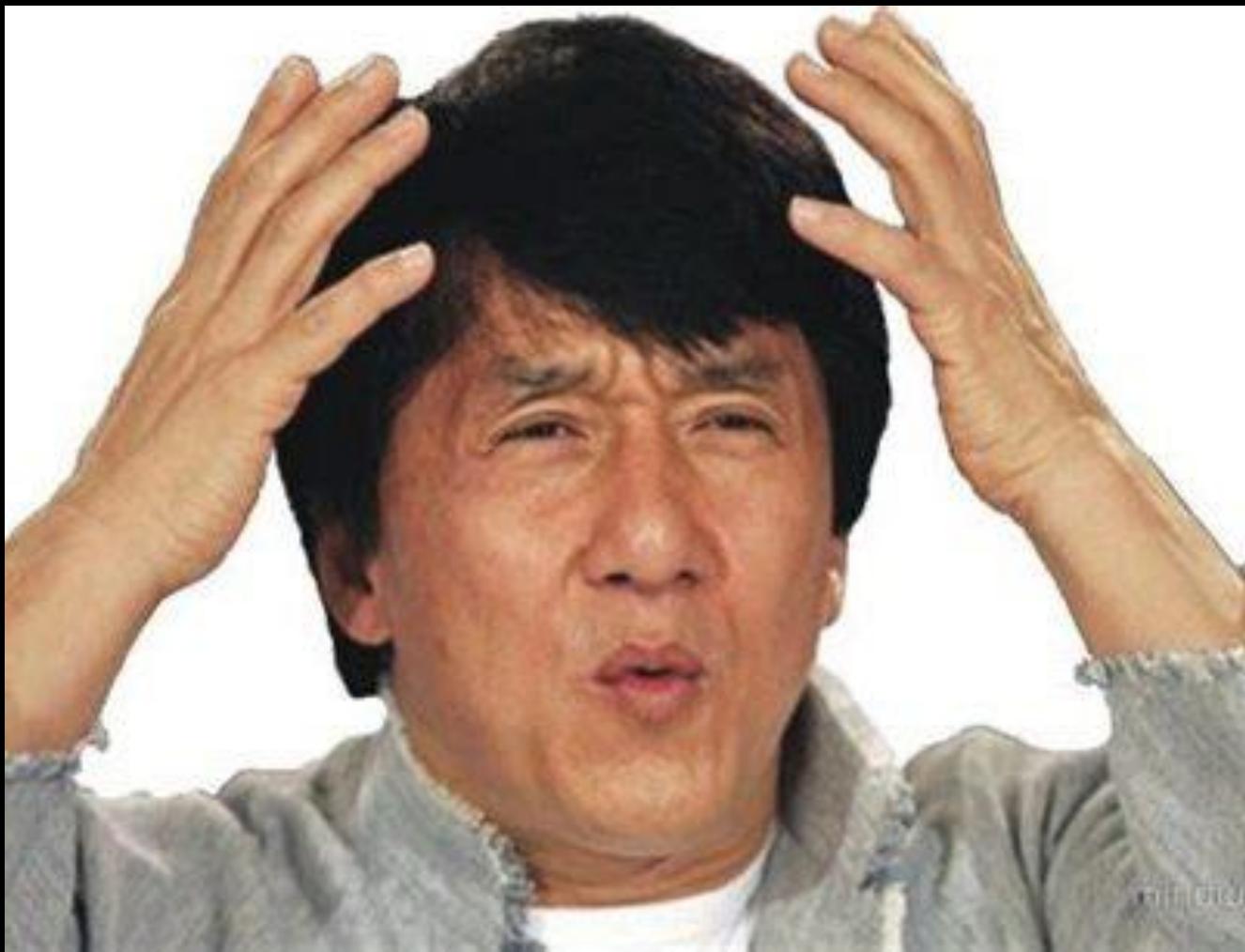
SLM for reliable LLM generation
(&6.1)

POLAR [392]; SAPLMA [19]; SuperICL [347]; SuperContext [357]; Self-RAG [16]; SKR [324]; SlimPLM [289]; CRAG [351]; Feng et al. [90]; Liu et al. [196]; HaluAgent [50]; Toolformer [264]

SLM for extracting LLM prompts

Prompt Stealing Attacks [266]; Output2prompt [376]; Model Purifying [175]; Zhang et al.





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)



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



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

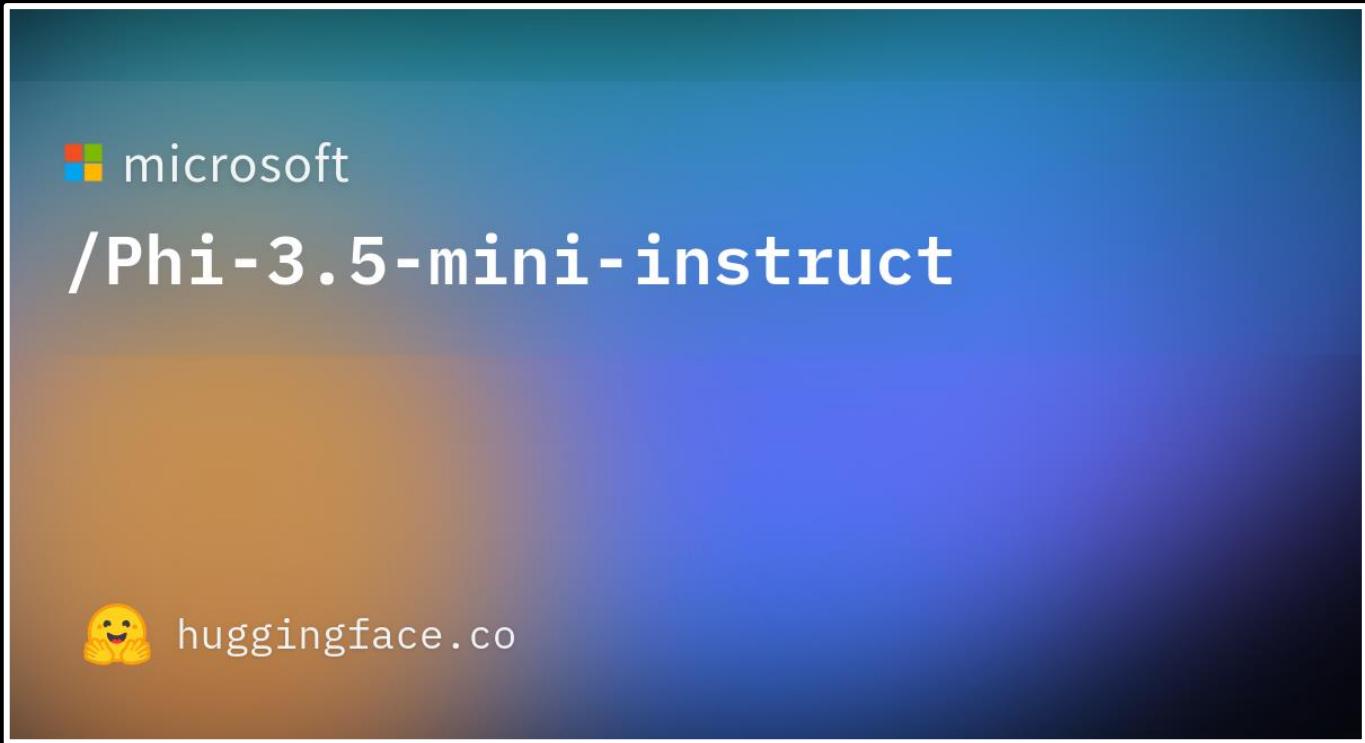


THaLLE 7B (\approx 4.4 GB-5.4 GB)
(Consider as LLM)
Thai-Lang Supported

Size Converted with llama.cpp/GGUF Size

Generally Available

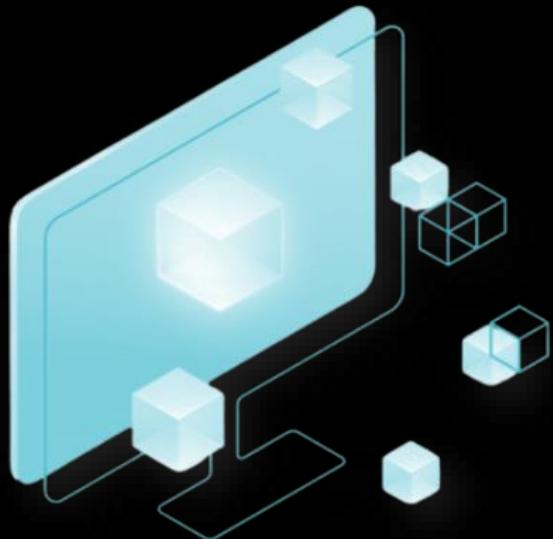
Microsoft's Phi-3.5



microsoft/Phi-3.5-mini-instruct · Hugging Face



Microsoft's Phi-3.5 Family

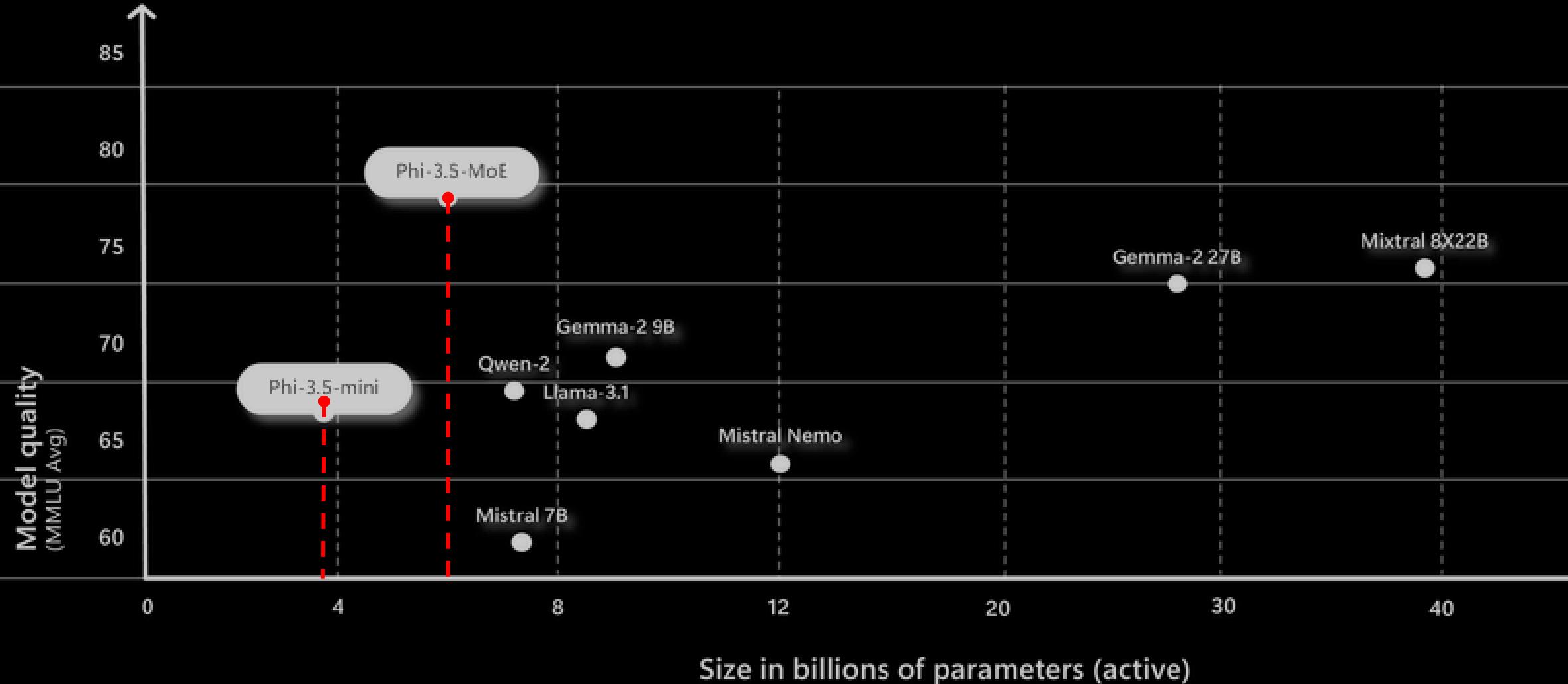


Model For Today's Experiment

Phi-3.5 mini-4k instruct (SLM)

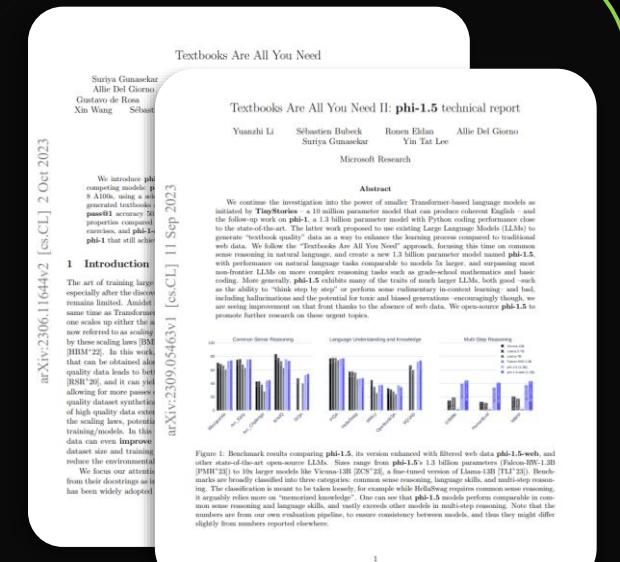
- **Architecture:** 3.8B parameters and is a dense decoder-only Transformer model. The model is fine-tuned with Supervised fine-tuning (SFT) and Direct Preference Optimization (DPO) to ensure alignment with human preferences and safety guidelines.
- **Inputs:** Text. It is best suited for prompts using 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

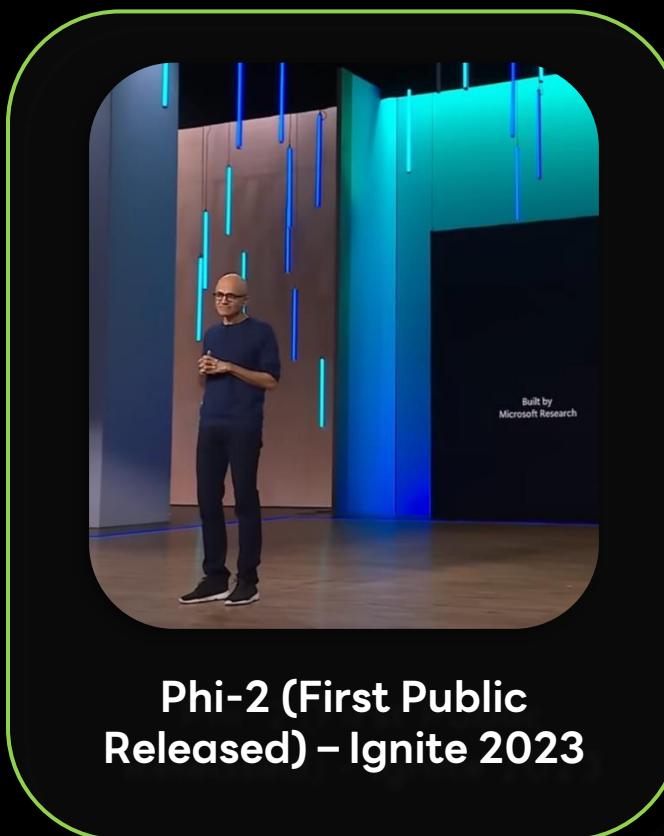


Textbooks Are All You Need

Evolution of Phi family



Phi-1 and Phi-1.5 (First Development)



Phi-2 (First Public Released) – Ignite 2023



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

Source: Microsoft Research (ArXiv)

Microsoft's Responsible AI Principles



Fairness



Reliability
& Safety



Privacy &
Security



Inclusiveness



Transparency



Accountability



4
QUALITY
EDUCATION



5
GENDER
EQUALITY



10
REDUCED
INEQUALITIES



16
PEACE, JUSTICE
AND STRONG
INSTITUTIONS



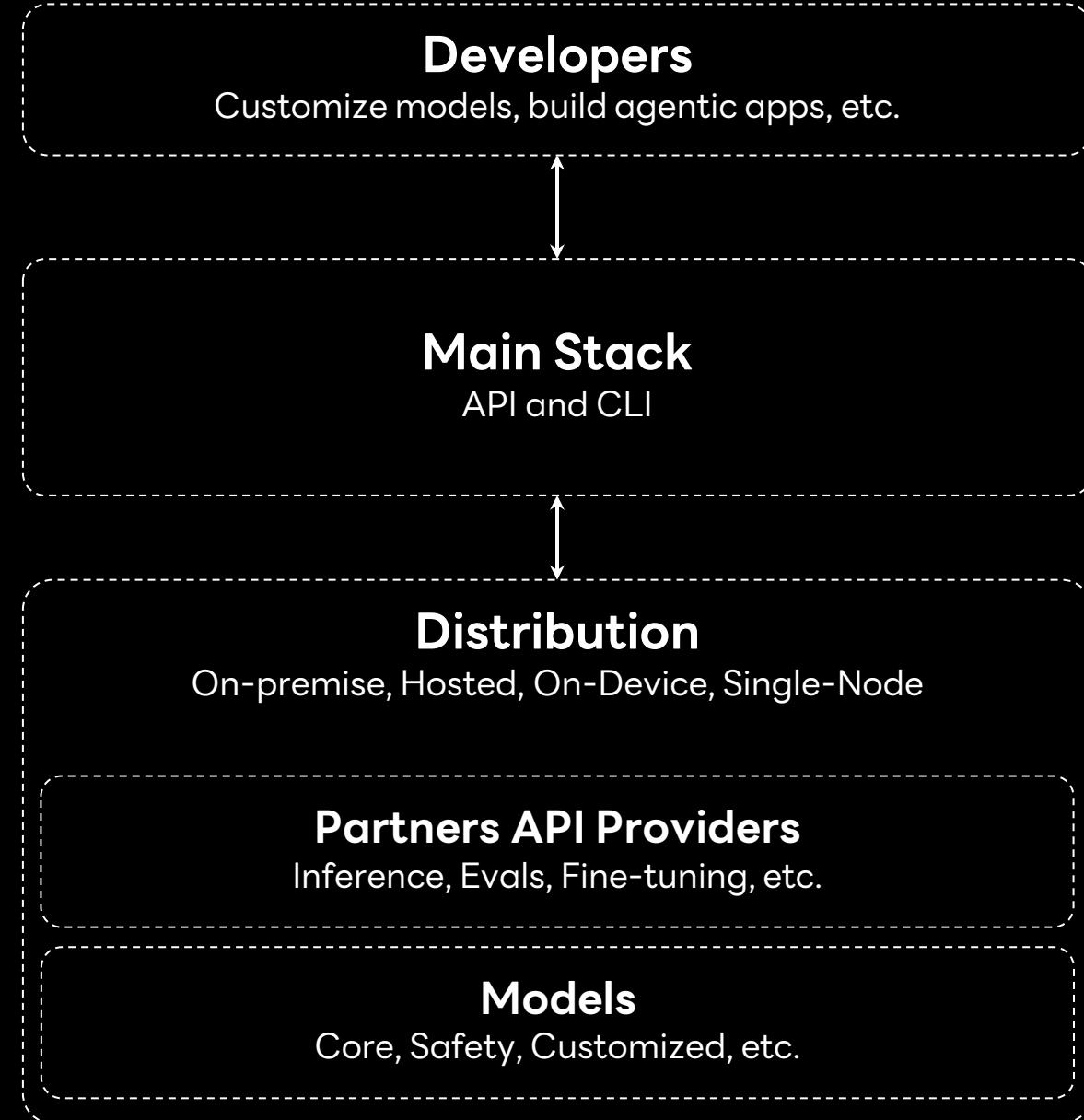
17
PARTNERSHIPS
FOR THE GOALS

Run & Develop

LM Studio
For Easy Dev
with UI

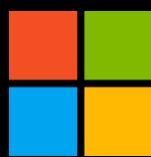


Ollama
Running in
Linux Terminal



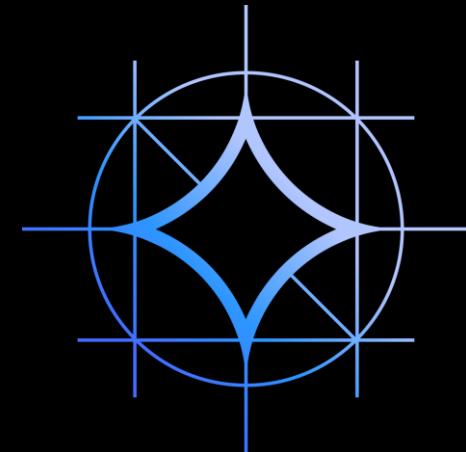
What is Ollama?

- 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.



Microsoft

∞ Meta





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)
```

```
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

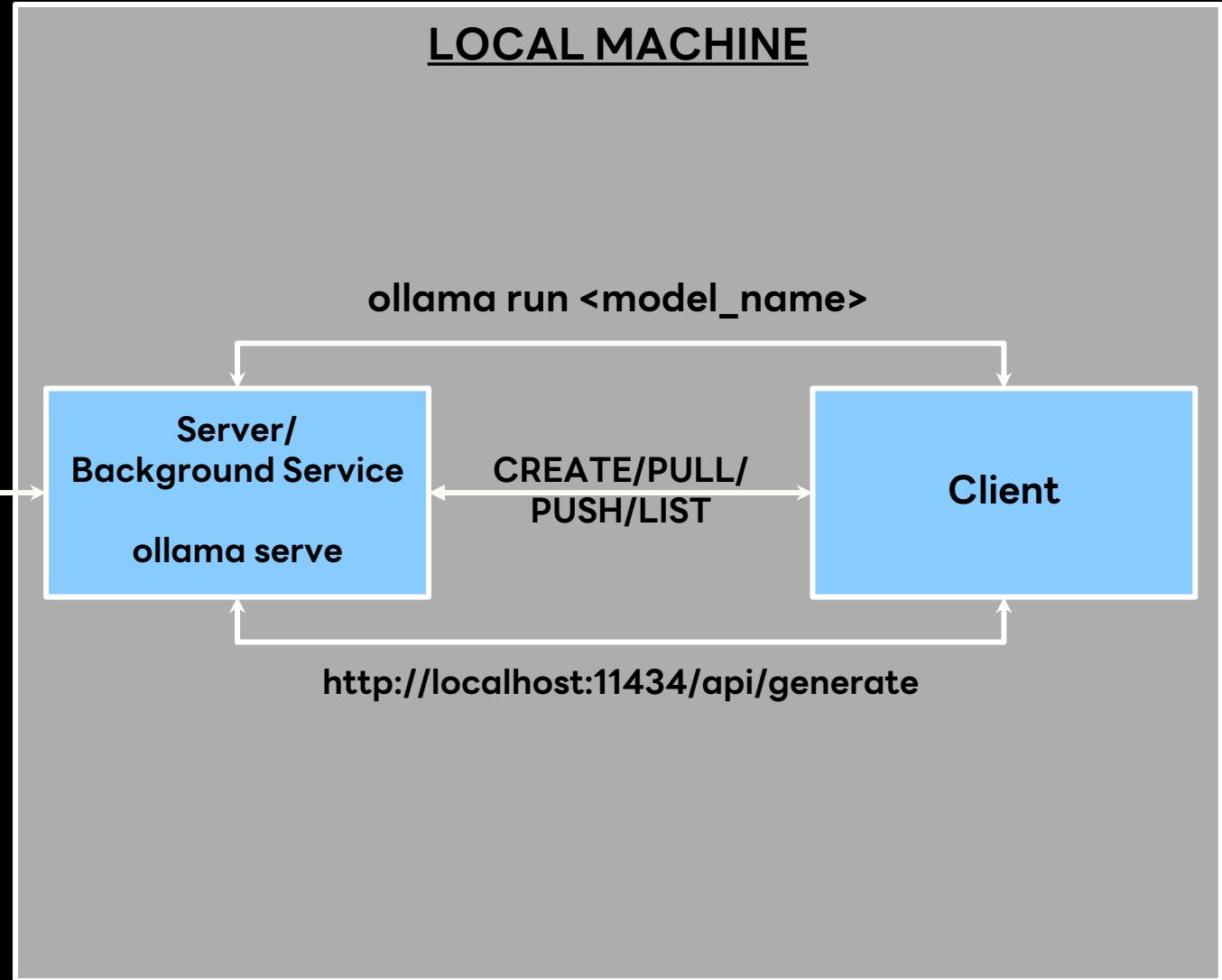
You can use **any GGUF quants** created by the community on Hugging Face directly with Ollama, without creating a new Modelfile. You can run any of them with a single **ollama run** command. We also provide customizations like choosing quantization type to improve your overall experience.



Hugging Face

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

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



First Demo

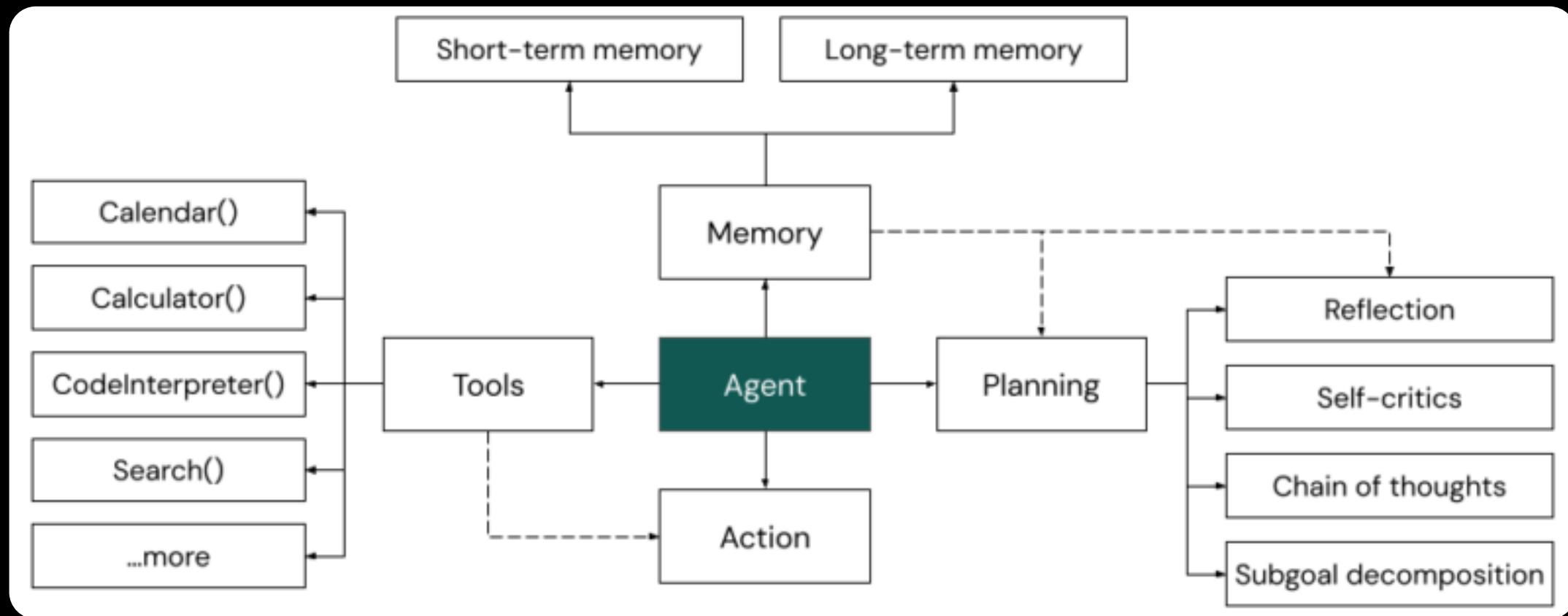
Ollama Implementation on Python with **Streamlit**



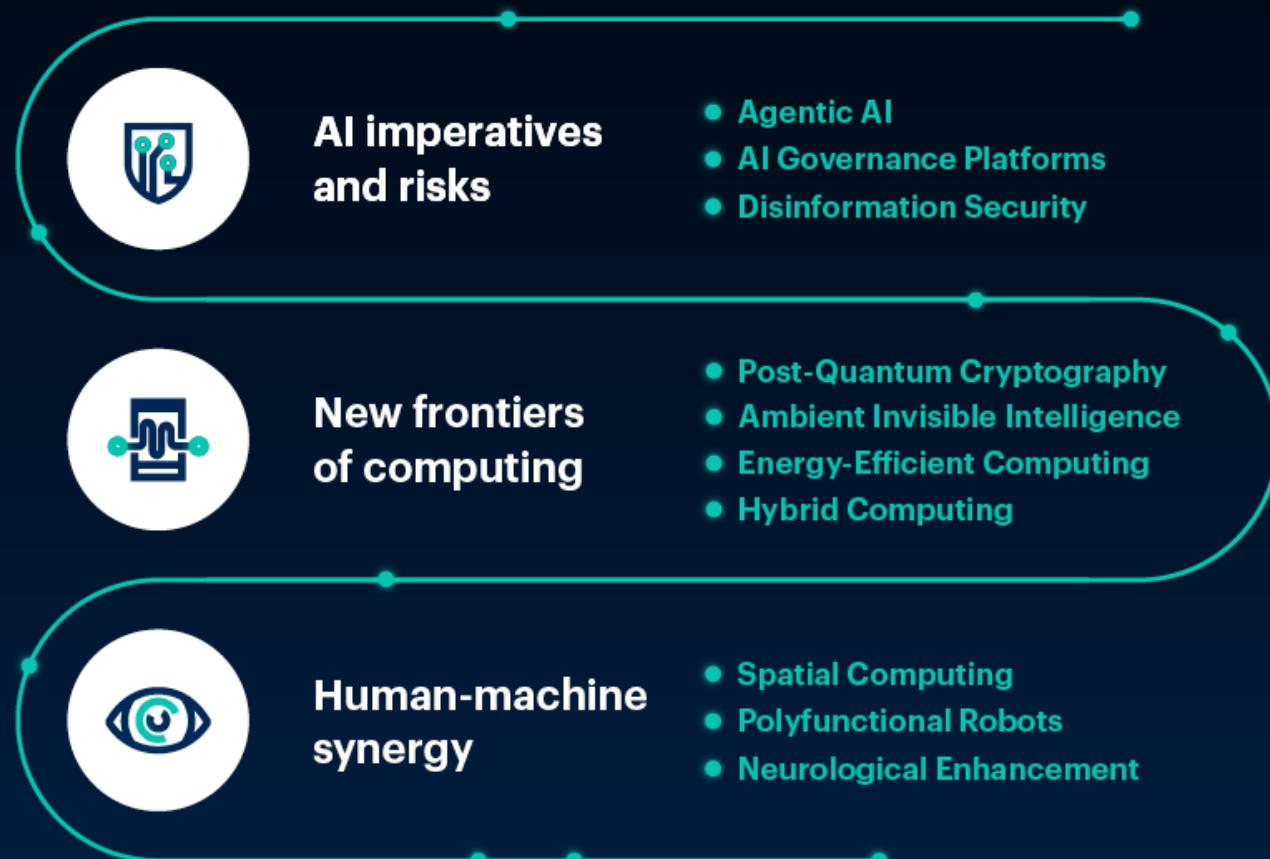
Ollama Implementation with AutoGen for Python for making the AI agents

What is AI Agents?

- An AI agent is a software entity capable of performing tasks autonomously or semi-autonomously by perceiving its environment, processing information, and taking actions to achieve specific goals.



2025 Top 10 Strategic Technology Trends



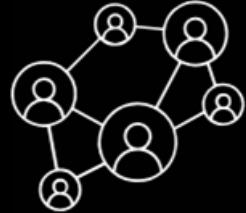
Source: Gartner
© 2024 Gartner, Inc. and/or its affiliates.
All rights reserved. 3185862

Gartner®

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.



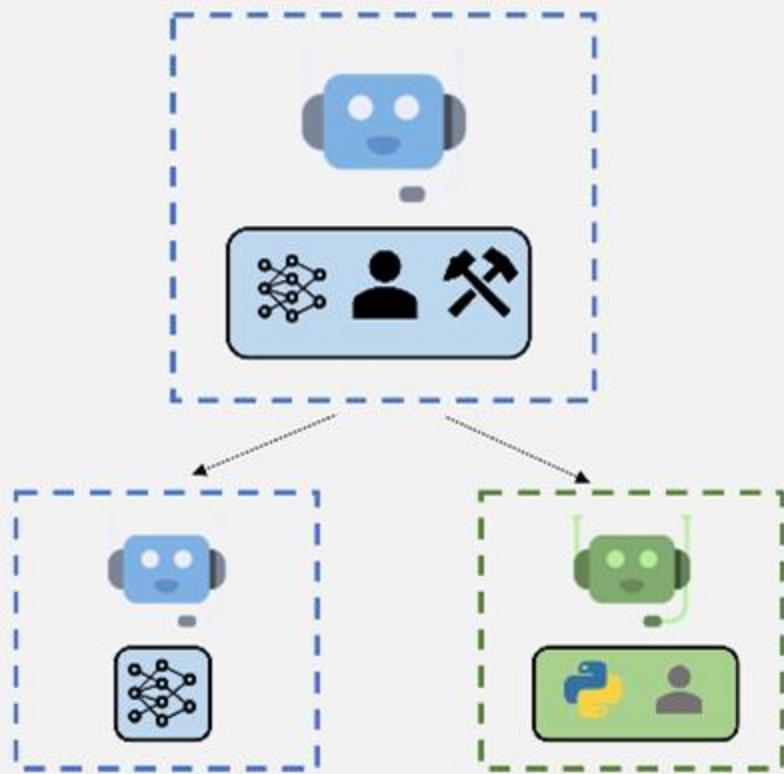
Easily Build Diverse Applications

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

Enhanced LLM Inference & Optimization

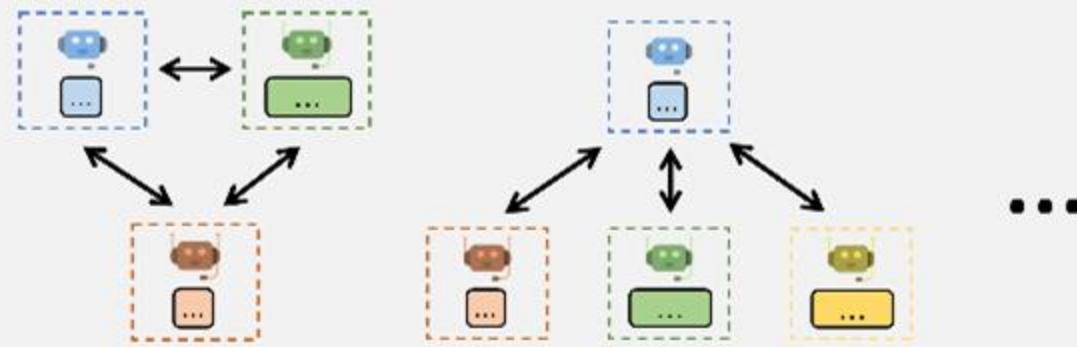
AutoGen supports enhanced LLM inference APIs, which can be used to improve inference performance and reduce cost.

Conversable agent



Agent Customization

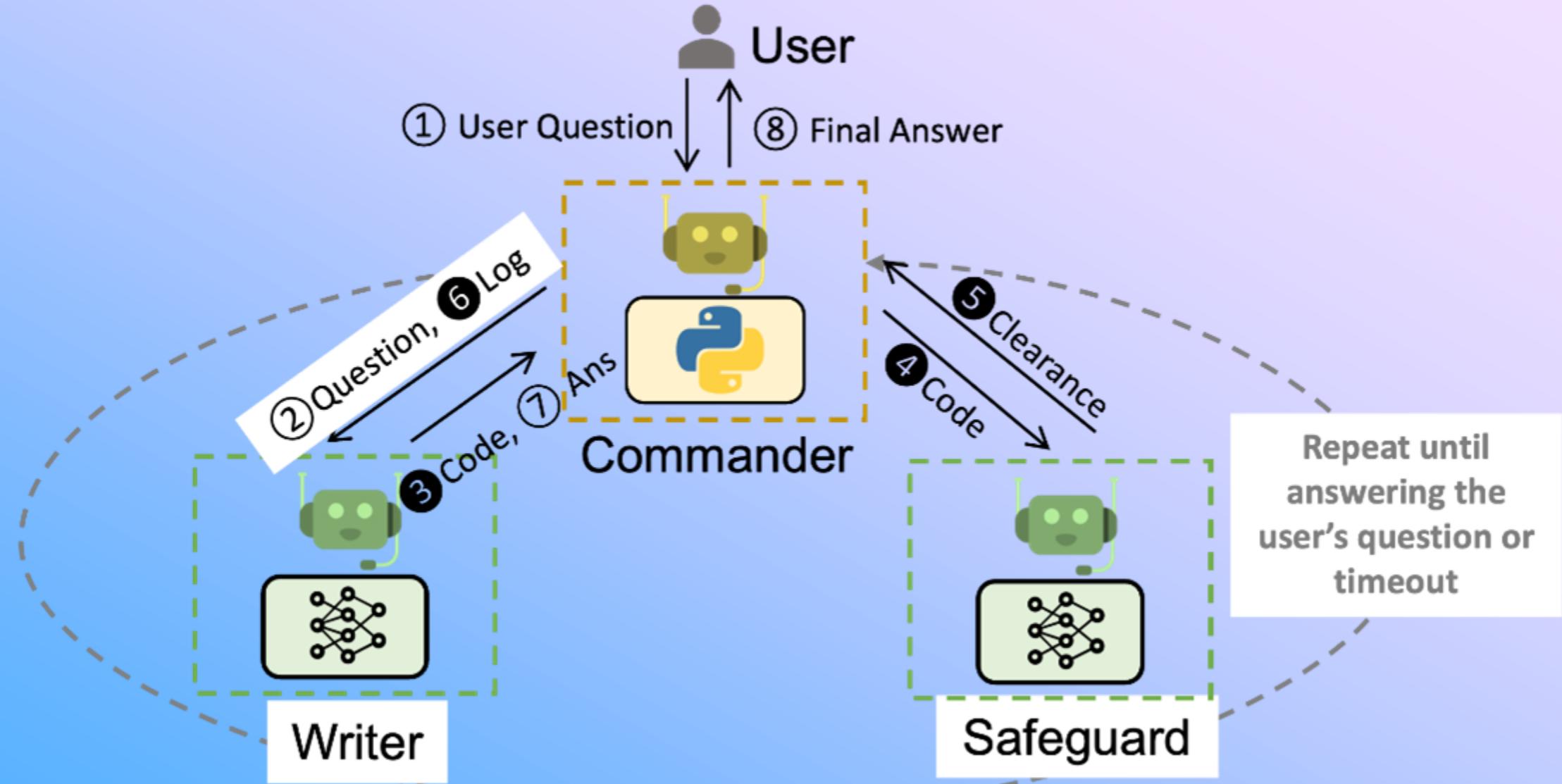
Multi-Agent Conversations



Joint chat

Hierarchical chat

Flexible Conversation Patterns



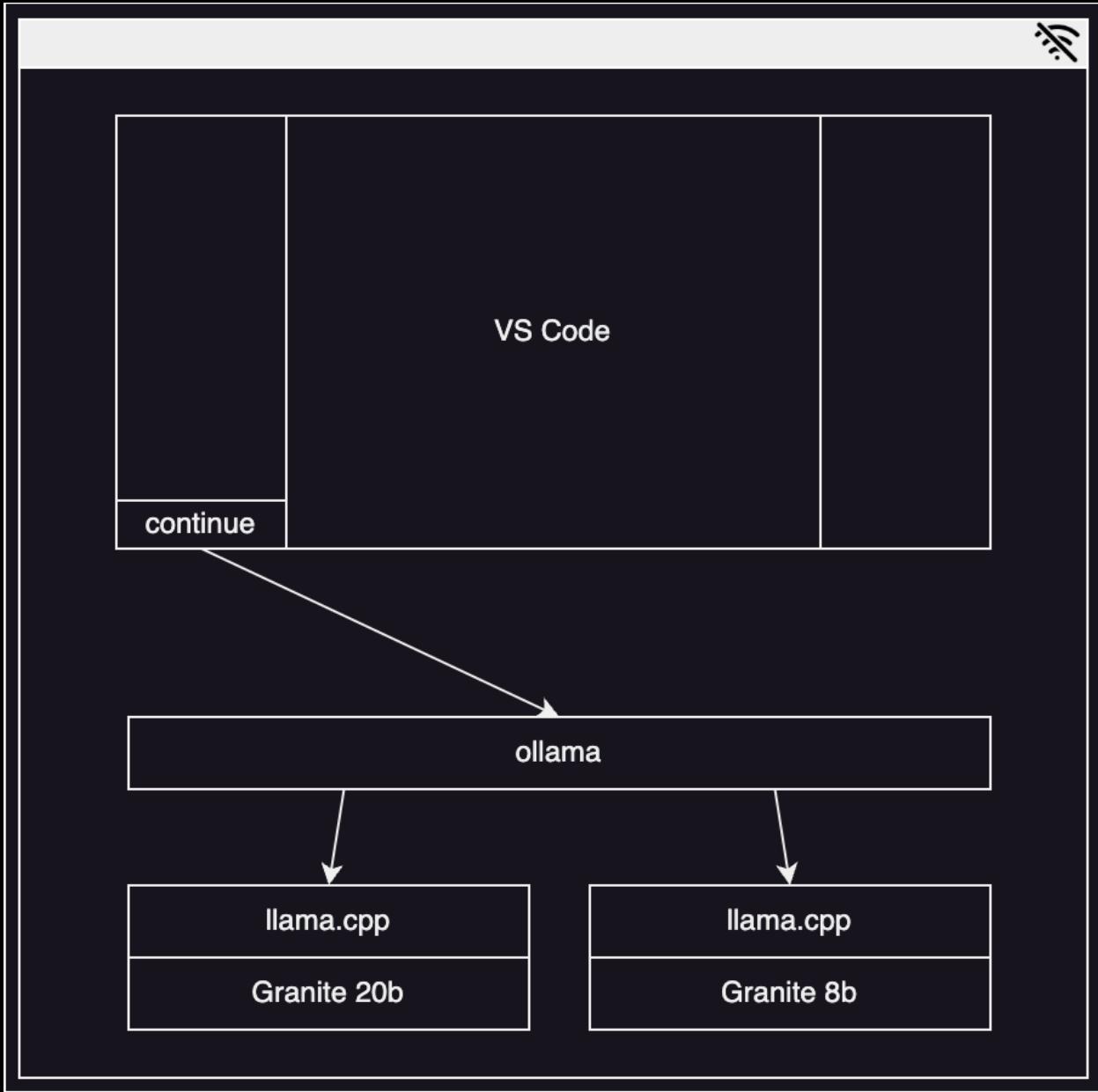
Continue



Continue enables you to easily create your own coding assistant directly inside Visual Studio Code and JetBrains with open-source LLMs. All this can run entirely on your own laptop or have Ollama deployed on a server to remotely power code completion and chat experiences based on your needs.

<https://ollama.com/blog/continue-code-assistant>

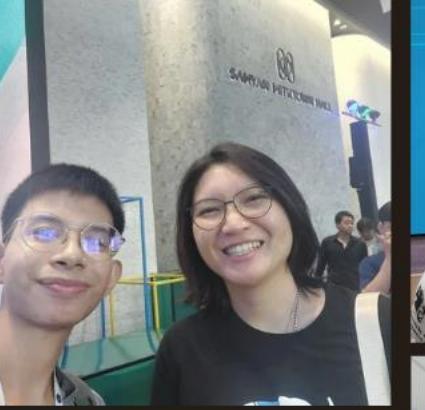
Ollama + Continue System Architecture



3

Microsoft Learn Student Ambassadors กับเส้นทางค้นหาเคลสใหม่ๆ







Model
Creator

Community Support

Engineering
in AI

Equality

AI
Governance

Wellbeing

Teamwork

Data and Insight
+ Life Principles

AI Leadership

Thai
Language
Tokenization

Content
Creativity

Responsible
AI Practice

Use Case Inspiration

Engineer
Principle

Academic
and History

ESL Practice

Academic

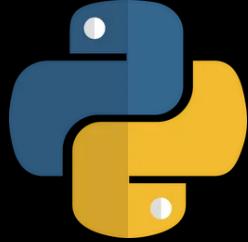
Copilot and Cloud
Deployment

Innovation

Productivity

Science

My Use Cases From Previous Events



Python

- Simple Chatbot with Streamlit (Run with Phi-3)
- Multimodal Testing with Phi-3 Vision Model

[github.com/chrnthonkmutt/
phi3-py-experiment](https://github.com/chrnthonkmutt/phi3-py-experiment)



JavaScript

- OpenAI API Deployment
- Text Summarizer
- RAG from CSV
- Text Translator
- Streaming Responses

[github.com/chrnthonkmutt/
phi3.5-js-experiment](https://github.com/chrnthonkmutt/phi3.5-js-experiment)



.NET (C#)

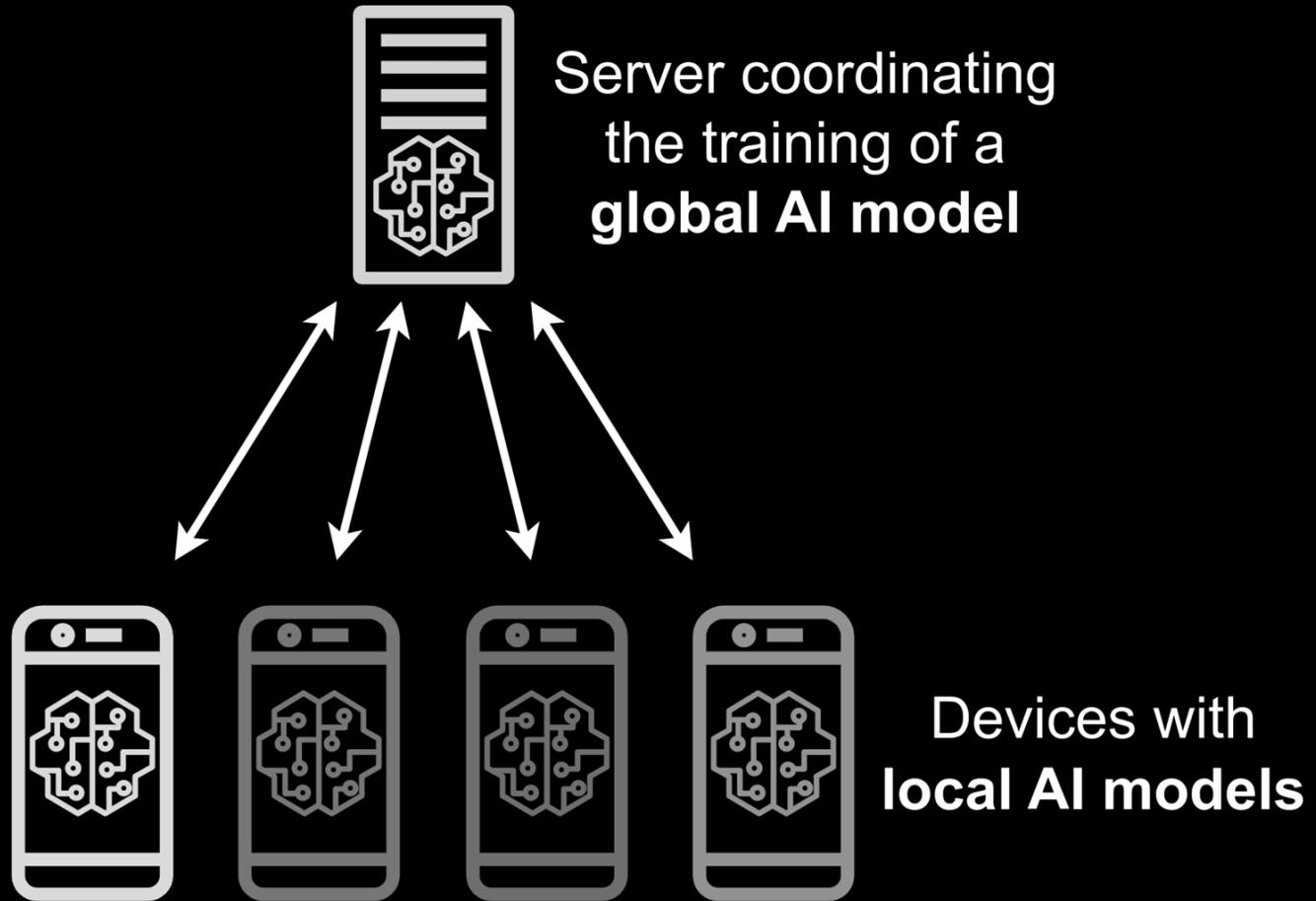
- OpenAI API Deployment with Semantic Kernel (including Chat History)
- AutoGen Multimodal AI Agent

[github.com/chrnthonkmutt/
BoatSLM cs experiment](https://github.com/chrnthonkmutt/BoatSLM_cs_experiment)

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 on Model Minituarization



Less-Parameters Language Models Arena



Phi-3.5: 3.8B (\approx 2.2 GB)
Thai-Lang Hallucination
Clear Responsible AI Guide
Academic



LLaMA 3.2 : 3B (\approx 2.02 GB)
Thai-Lang Supported
Clear Responsible AI Guide
General



Qwen2.5 3B (\approx 1.93 GB)
Thai-Lang Supported
Unclear Responsible AI Guide
General



Typhoon 8B (\approx 4.9 GB)
Thai-Lang Supported
Meta's Responsible AI Guide
General

THaLLE 7B (\approx 4.4 GB)
Thai-Lang Supported
Qwen-Based + Adding RAI
Financial/Business

Size Converted with llama.cpp/GGUF Size

Good Open Models for SLMs in Nowadays



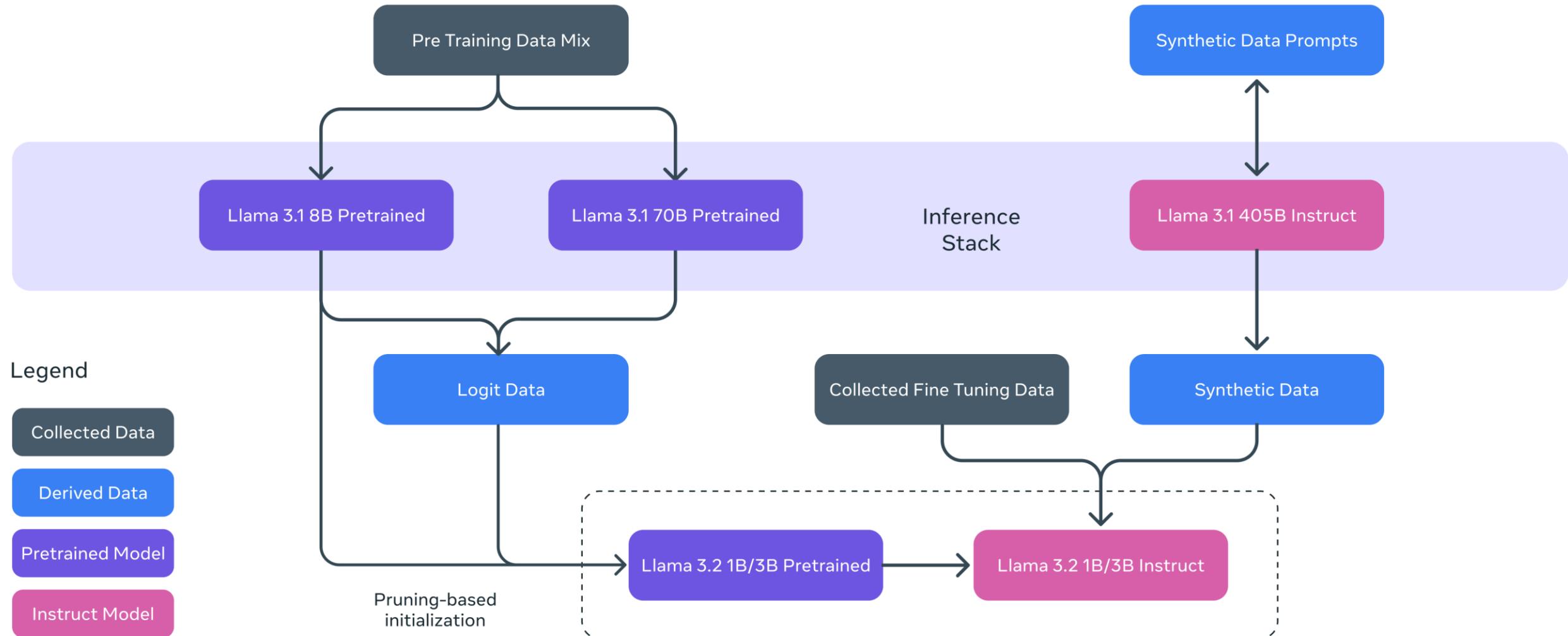
LLaMA 3.2 : 3B
(≈2.02 GB)
General Purposes

Responsible AI

- Privacy and Security
- Fairness and Inclusion
- Robust and Safety
- Transparency and Control
- Accountability and Governance

- Lightweight, text-only models that fit onto edge and mobile devices, including pre-trained and instruction-tuned versions.
- Support context length of 128K tokens and are state-of-the-art in their class for on-device use cases like summarization, instruction following, and rewriting tasks running locally at the edge.
- The Model empower developers to build personalized, on-device agentic applications with strong privacy where data never leaves the device. For example, such an application could help summarize the last 10 messages received, extract action items, and leverage tool calling to directly send calendar invites for follow-up meetings.
- **Outperforms the Gemma 2 2B and Phi 3.5-mini 3.8B models** on tasks such as following instructions, summarization, prompt rewriting, and tool-use, while the 1B is competitive with Gemma.
- Multilingual Support in one single model. (Including Thai Language)

1B & 3B Pruning & Distillation



WangchanThailnstruct Dataset

100% Human-Annotated Thai Instruction Dataset (Batch 1-3 Release)

4 Domains:

- Medical
- **Finance**
- **Retail**
- Legal

7 Tasks:

- Summarization
- Open QA
- Close QA
- Classification
- Creative Writing
- Brainstorming
- Multiple Choice QA



[airesearch/WangchanThailnstruct · Datasets at Hugging Face](#)

The question is...

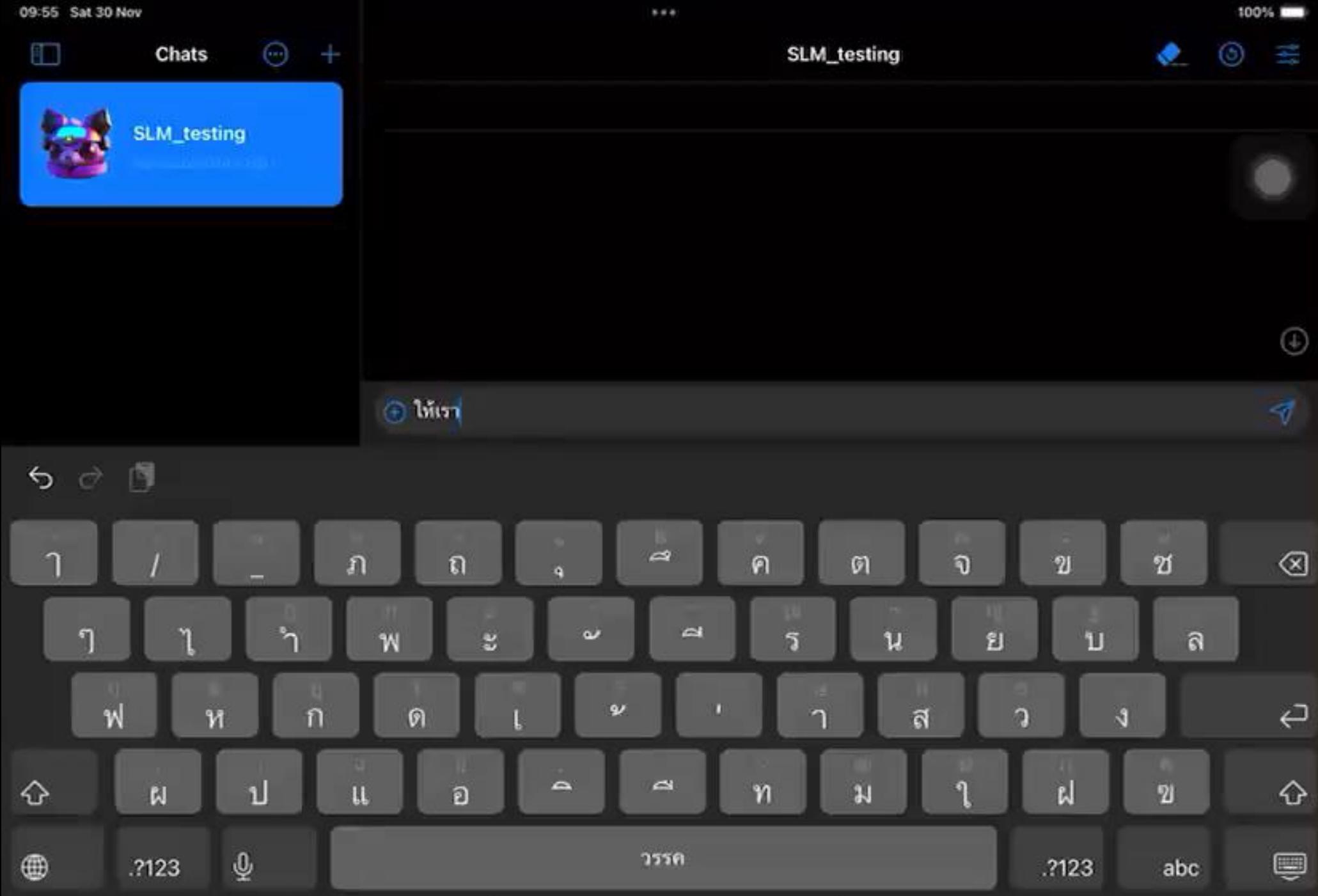
Why Thai People Still Doesn't Use Small Language Models?

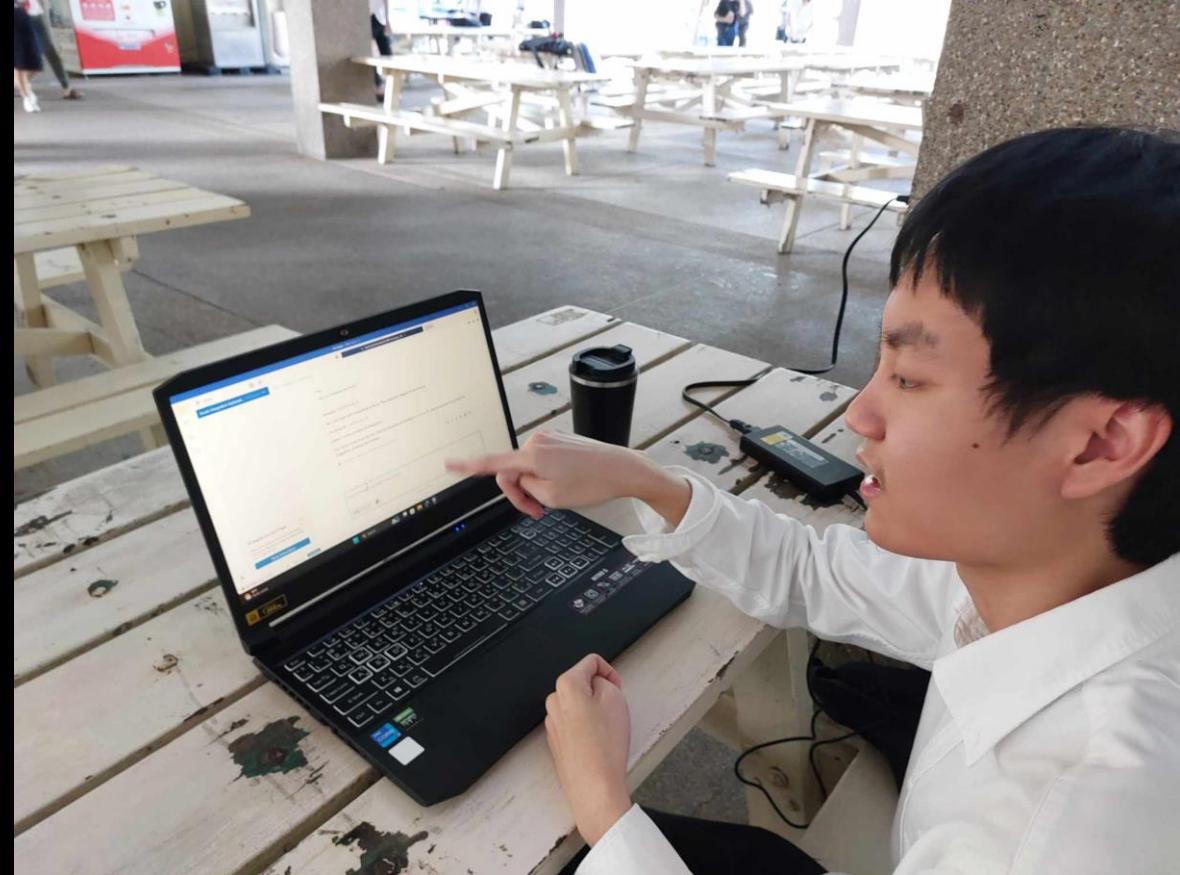
Hallucination
& Performance

Limited Device
Resources

Languages Barrier

Low
Engagement?









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)

