

Introduce Phi-3

Kinfey Lo Microsoft Senior Cloud Advocate

Generative Al









What is Phi-3?



The development of the Microsoft Phi

Code, Text Completion, Chat



perf. comparable to models 10x larger trained on 100x more data

Natural language model with NL comparable to models 10x larger trained on 30x more data and reasoning comparable to models 50x larger.

a 2.7 billion-parameter language model that demonstrates outstanding reasoning and language understanding capabilities, showcasing state-of-the-art performance among base language models with less than 13 billion parameters

2.7B

Microsoft Phi 3 Family



Phi-3-mini

a 3.8B language model is available on Microsoft Azure AI Studio, Hugging Face, and Ollama. Phi-3 models significantly outperform language models of the same and larger sizes on key benchmarks (see benchmark numbers below, higher is better). Phi-3-mini does better than models twice its size, and Phi-3-small and Phi-3-medium outperform much larger models, including GPT-3.5



Phi-3-small with only 7B parameters beats GPT-3.5T across a variety of language, reasoning, coding and math benchmarks.

Phi-3-medium with 14B parameters continues the trend and outperforms Gemini 1.0 Pro.

Phi-3-small & medium



Phi-3-vision with just 4.2B parameters continues that trend and outperforms larger models such as Claude-3 Haiku and Gemini 1.0 Pro V across general visual reasoning tasks, OCR, table and chart understanding tasks.

Phi-3-Vision

Microsoft Phi 3 Family

Small Language Model, Best in class performance for size, with frictionless availability





✓ Groundbreaking performance at a small size: Faster and Cost Effective & Low costs in Fine Tuning. Excellent performance in latency bound scenarios where fast response times are critical ✓ Unlocking new capabilities: Works best across resource constrained environments where local **Core Value Proposition** inference may be needed. Can deploy on edge & cloud with ONNX & TensorFlow Lite ✓ Safe by Design: Developed in accordance with Microsoft Responsible AI principles & trained on with high quality data Phi-3 Mini is a lightweight, state-of-the-art open model built upon datasets used for Phi-2 synthetic data and filtered websites - with a focus on very high-quality, reasoning dense data._ The Description model belongs to the Phi-3 model family and incorporates safety, alignment, and RLHF training. Phi-3 Mini has 3.8B parameters and is a dense decoder-only Transformer model. Architecture Text. It's best suited for prompts using chat format. **Context length** 128K tokens - "the first model in its weight class to support long-context of up to 128K" Generated text in response to the input Outputs Trained between Feb 2024 and April 2024 This is a static model trained on an offline dataset. Future versions of the tuned models may be released as we improve models.

Phi-3 Mini-128K-Instruct ······ Gemma-7B ······ Mistral-7B ······ Mixtral 8x7B ······ GPT-3.5-Turbo

	Phi-3 Mini- 128K-Instruct	Gemma-7B	Mistral-7B	Mixtral 8x7B	GPT-3.5- Turbo
Average across 80 benchmarks	57.7	52.6	47.2	56.2	62.2

Category	Phi-3 Mini- 128K-Instruct	Gemma-7B	Mistral-7B	Mixtral 8x7B	GPT-3.5- Turbo
Popular					
aggregated benchmark	63.5	59.4	56.5	66.2	67.0
Reasoning	79.1	71.3	72.6	78.6	80.0
Language understanding	56.7	57.6	52.5	71.5	67.7
Code generation	55.5	45.6	42.9	52.7	70.4
Math	60.3	41.4	29.7	46.2	58.0
World knowledge	35.7	46.7	49.8	58.6	63.4
Multilingual	53.1	64.6	56.0	64.9	69.6
Robustness	58.6	38.4	40.6	51.0	69.3

License

MIT

Inputs

Dates

Status

Microsoft Phi 3



(3.8B + 0.3B)





Phi-3-small (7B)



Phi-3-mini (3.8B)



ONNX Runtime

GGUF

Qu

Quantitative(INT4,FP16,FP32)

Demo : Call Phi-3-mini & Phi-3-Vision



Fine-tuning vs RAG



How about RAG?

Ľ

D R

B

RAG - Retrieval-Augmented Generation

RAG is data retrieval + text generation. The structured data and unstructured data of the enterprise are stored in the vector database. When searching for relevant content, the relevant summary and content are found to form a context, and the text completion capability of LLM/SLM is combined to generate content.



A vector database is a collection of data stored in mathematical form. Vector databases make it easier for machine learning models to remember previous inputs, enabling machine learning to be used to support use cases such as search, recommendations, and text generation. Data can be identified based on similarity metrics rather than exact matches, allowing computer models to understand the context of the data.

Vector database is the key to realizing RAG. We can convert data into vector storage through vector models such as text-embedding-3, jina-ai-embedding, etc.

Demo : RAG with Phi-3



Fine-tuning

Fine-tuning is based on improvement of a certain model. It does not need to start with the model algorithm, but data needs to be continuously accumulated. If you want more precise terminology and language expression in industry applications, fine-tuning is your better choice. But if your data changes frequently, fine-tuning can become complicated.



LoRA (Low-Rank Adaptation) and QLoRA (Quantized Low-Rank Adaptation) are both techniques used to fine-tune large language models (LLMs) using Parameter Efficient Fine Tuning (PEFT). PEFT techniques are designed to train models more efficiently than traditional methods.

LoRA is a standalone finetuning technique that reduces memory footprint by applying a low-rank approximation to the weight update matrix. It offers fast training times and maintains performance close to traditional fine-tuning methods.

QLoRA is an extended version of LoRA that incorporates quantization techniques to further reduce memory usage. QLoRA quantizes the precision of the weight parameters in the pre-trained LLM to 4-bit precision, which is more memory efficient than LoRA. However, QLoRA training is about 30% slower than LoRA training due to the additional quantization and dequantization steps.

QLoRA uses LoRA as an accessory to fix the errors introduced during quantization errors. QLoRA enables the fine-tuning of massive models with billions of parameters on relatively small, highly available GPUs. For example, QLoRA can fine-tune a 70B parameter model that requires 36 GPUs with only 2

RAG vs Fine-tuning

How to choose

- If our answer requires the introduction of external data, RAG is the best choice
- If you need to output stable and precise industry knowledge, fine-tuning will be a good choice. RAG prioritizes pulling relevant content but might not always nail the specialized nuances.
- Fine-tuning requires a high-quality data set, and if it is just a small range of data, it will not make much difference. RAG is more flexible
- Fine-tuning is a black box, a metaphysics, and it is difficult to understand the internal mechanism. But RAG can make it easier to find the source of the data, thereby effectively adjusting hallucinations or content errors and providing better transparency.

Scenarios

- Vertical industries require specific professional vocabulary and expressions, Fine-tuning will be the best choice
- QA system, involving the synthesis of different knowledge points, RAG will be the best choice
- The combination of automated business flow RAG + Fine-tuning is the best choice



Microsoft Olive

Olive.config as Service



Evaluations

SLMwith **EdgeDevice**



Ś	QuickTime Player File Ed	dit View Windo	ow Help				0 🛥		i 0	@	6	<u>42</u>)	A	ଚ ଦ	00	Sat Mar	23 9:23 AM
	••• 🐴 🍙 🗖	🗼 Notebooks - /	Azure Al Machin 🛛 🗙	+													
	← C ᡬ Ô htt	tps:// ml.azure.com	wsid=/subscriptions/3966	=/subscriptions/396656ae-1e4b-4f9d-9a8a-a5fcb029664		96643/resource	groups/AIG				${igsidentsized}$	¢	C)	£_≡	Ē	~~ ···	
r	Azure AI Machine Learning Studi	lio					Microsoft Re	emote Desi	ctop								
		Lo Kinfey > Oliv															
	\leftarrow All workspaces	Notebook	Recents			Desktop											
	A Homo	Files Samples	🙏 Applicati														
		—	Desktop	Fabric Screen Recor22.18	> AM.mov												
			Documents Downloads														
	Authoring		Downloads														
I	🗒 Notebooks		iCloud														
	\mathcal{J}_{a} Automated ML		🛆 iCloud Dri														
	品 Designer	∽ 🕤 ph	🖰 Shared														
	>_ Prompt flow	~ 🕤 (
	Assets	\sim	💻 lokinfey's														
	🖾 Data	[OneDrive														
	入 Jobs	\sim					Cancel Open										
	Line Components	{															
	. 毕 Pipelines	_ _Y ,	flow.dag.yaml		2024-03-23 00:35:49 +0000 22385 execution.bulk INFO Process name(ForkProcess-2:1)-Process												
	Environments	PY	generate.py		22385 execut	pleted. 22385 execution.bulk INFO Finished 1 / 1 lines.											
	── ♡ Models	PY	init_onnx.py		2024-03∖23 00:35:50 +0000 2233 lines: 82.01 seconds. Estimated				22385 execution.bulk INFO Average execution time for a red time for incomplete lines: 0.0 seconds.						r comp	ompleted	
	ର୍ତ୍ତ Endpoints	C.	run.ipynb	======= Run Summary ======													
	Manage	> 📫 r	nodels		Run name: "variant_0_20240323_003425_715346"												
manage		run_onnx_fp16.ipynb			Run status: "Completed" Start time: "2024-03-23 00:34:25.715314"												
	∠ Compute	C.	run_onnx_int4.ipynb			Duration: "0:01:27.418359" Output path: "/home/azureuser/ promotflow/ runs/variant 0 20240323 003425 715346"											
	Monitoring																
	Data Labeling				<prompt </prompt 			/13095506	.70>								

<mark>😲 🏥 🌫 🗢 🖸 🧭 🌗 🚳 🔶 🖅 🕰 🏹 🎦 🧭 🔘 🛛 🔁 🍳 📔 🔕 🦷</mark>









Phi-3-mini-V (3.8B + 0.3B)

Phi-3-small (7B) Phi-3-mini (3.8B)

Begin to talk about **Cloud Native**

What's .NET Aspire



- .NET Aspire is designed to improve the experience of building .NET cloud-native apps. It provides a consistent, opinionated set of tools and patterns that help you build and run distributed apps. .NET Aspire is designed to help you with:
 - **Orchestration**: .NET Aspire provides features for running and connecting multi-project applications and their dependencies for local development environments.
 - **Components:** .NET Aspire components are NuGet packages for commonly used services, such as Redis or Postgres, with standardized interfaces ensuring they connect consistently and seamlessly with your app.
 - **Tooling:** .NET Aspire comes with project templates and tooling experiences for Visual Studio, Visual Studio Code, and the dotnet CLI to help you create and interact with .NET Aspire projects.

AI@.NET Aspire

Building Generative AI apps with .NET 8



https://devblogs.microsoft.com/dotnet/build-gen-ai-with-dotnet-8/





	Phi3.Aspire		() () (\$					
	Pacouroc							
Resources	nesource							
E Console	Туре	Name	State	Start time	Source	Endpoints	Logs	Details
Etructured	Container	cache	Running	3:47:15 PM	docker.io/library/redis:7.2	tcp://localhost:50227	View	View
	Project	frontendService	Running	3:47:15 PM	Phi3.Aspire.FrontEnd.csproj	https://localhost:7033, http://localhost:5039	View	View
Traces	Project	phi3service	Running	3:47:15 PM	Phi3.Aspire.ModelService.csproj	https://localhost:7193, http://localhost:5257	View	View
Metrics	Project	skservice	Running	3:47:15 PM	Phi3.Aspire.SK.API.csproj	https://localhost:7282, http://localhost:5009	View	View

Demo : .NET Aspire with Phi3-mini





Phi-3 Cookbook



THANK YOU