

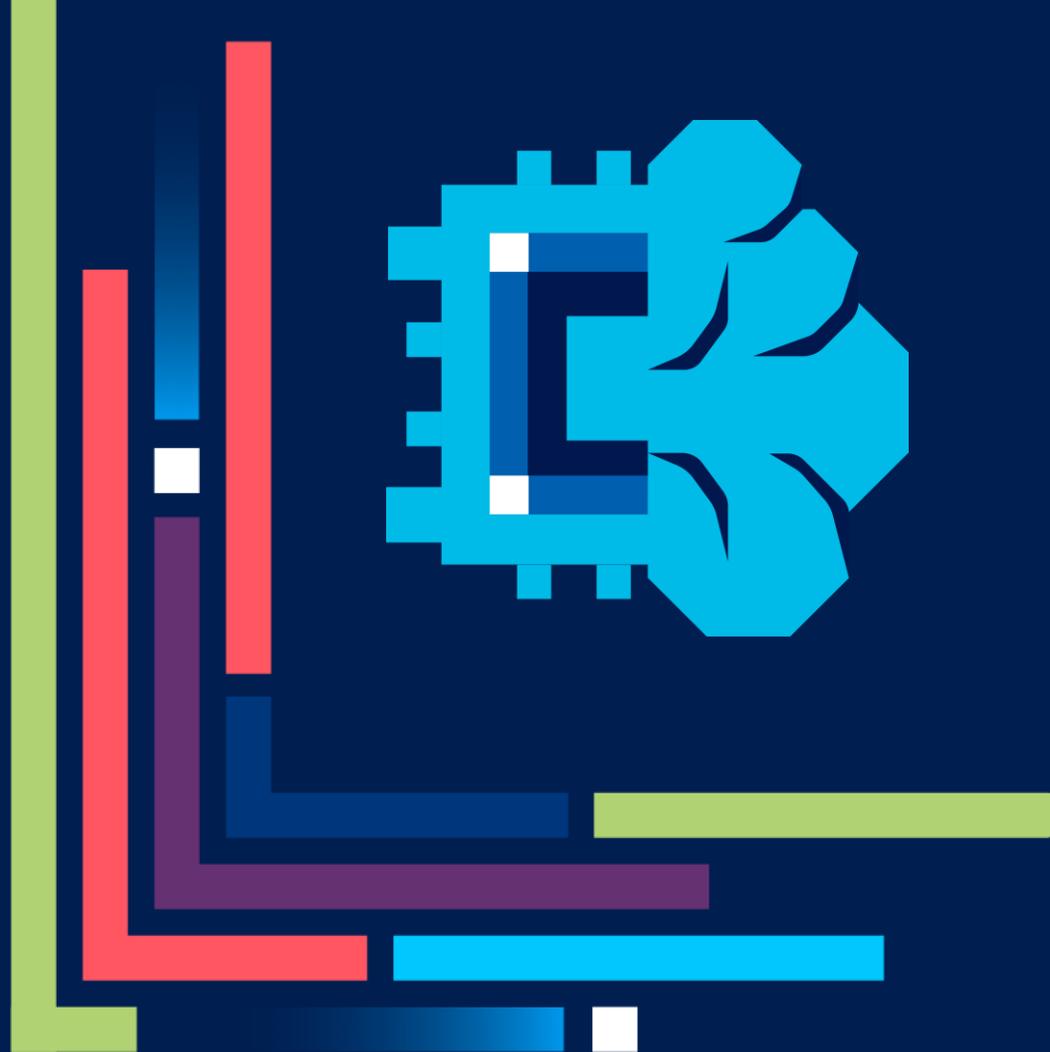
intel ai  
summit  
英特爾 AI 科技論壇

Bringing AI Everywhere

# Finetune LLM with Intel Architecture

Chungyeh Wang 王宗業

03/27/2024



# Agenda

---

AI Everywhere

---

Small and Nimble – the Fast Path to Enterprise GenAI

---

LLM Inference with OpenVINO and BigDL-LLM (ipex-llm)

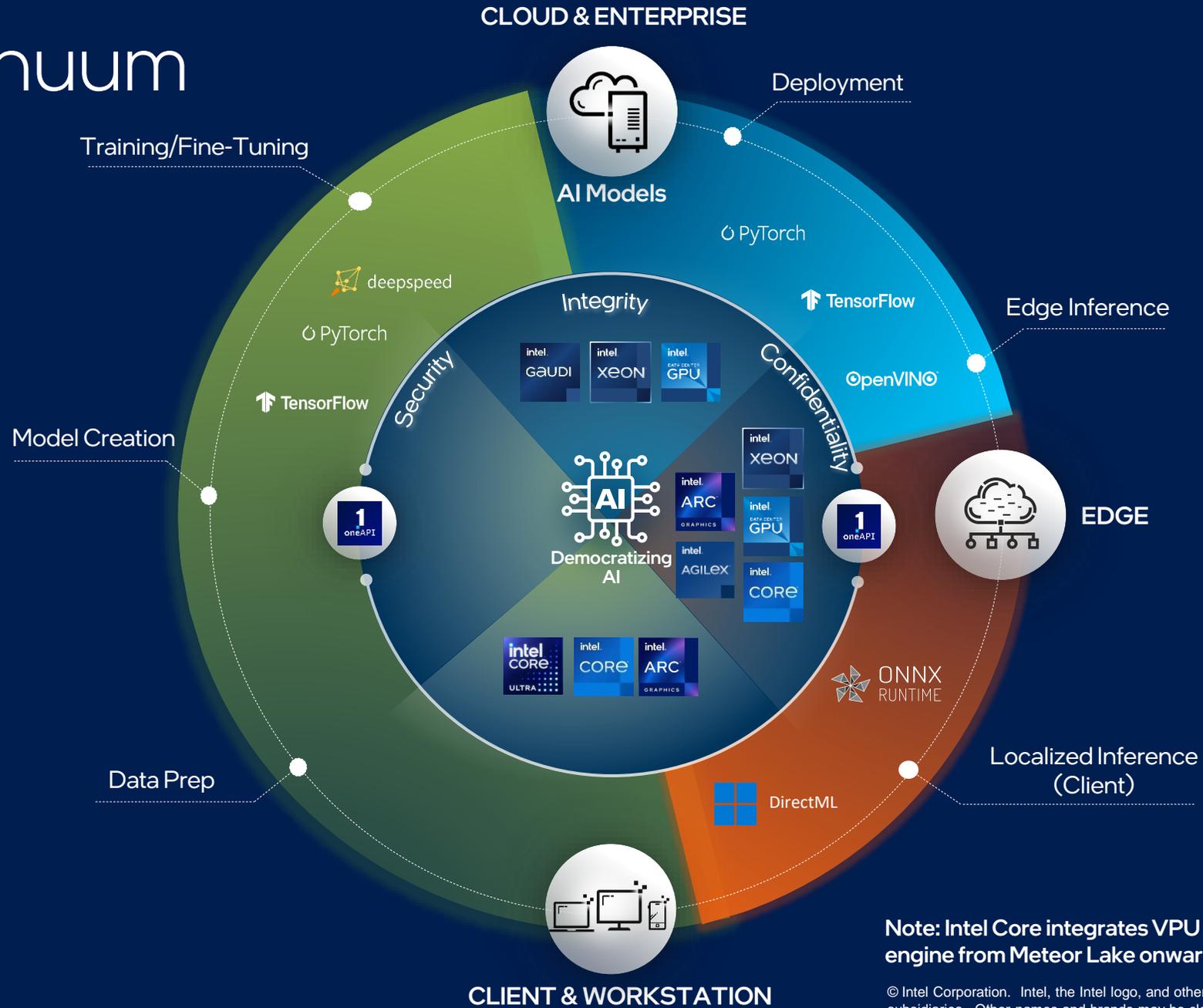
---

LLM Fine-tuning with affordable Intel platform



# AI Everywhere

# AI Continuum

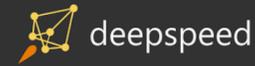


Note: Intel Core integrates VPU low power inference engine from Meteor Lake onwards.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

# Intel® AI Portfolio

Open Software Environment



Deep Learning Acceleration



Gaudi: Dedicated Deep Learning Training and Inference

General Acceleration



Cloud Gaming, VDI, Media Analytics, Real-Time Dense Video



Parallel Compute, HPC, AI for HPC

General Purpose



Real-Time, Medium Throughput, Low Latency, and Sparse Inference



Medium to Small Scale Training and Fine Tuning



Edge and Network AI Inference



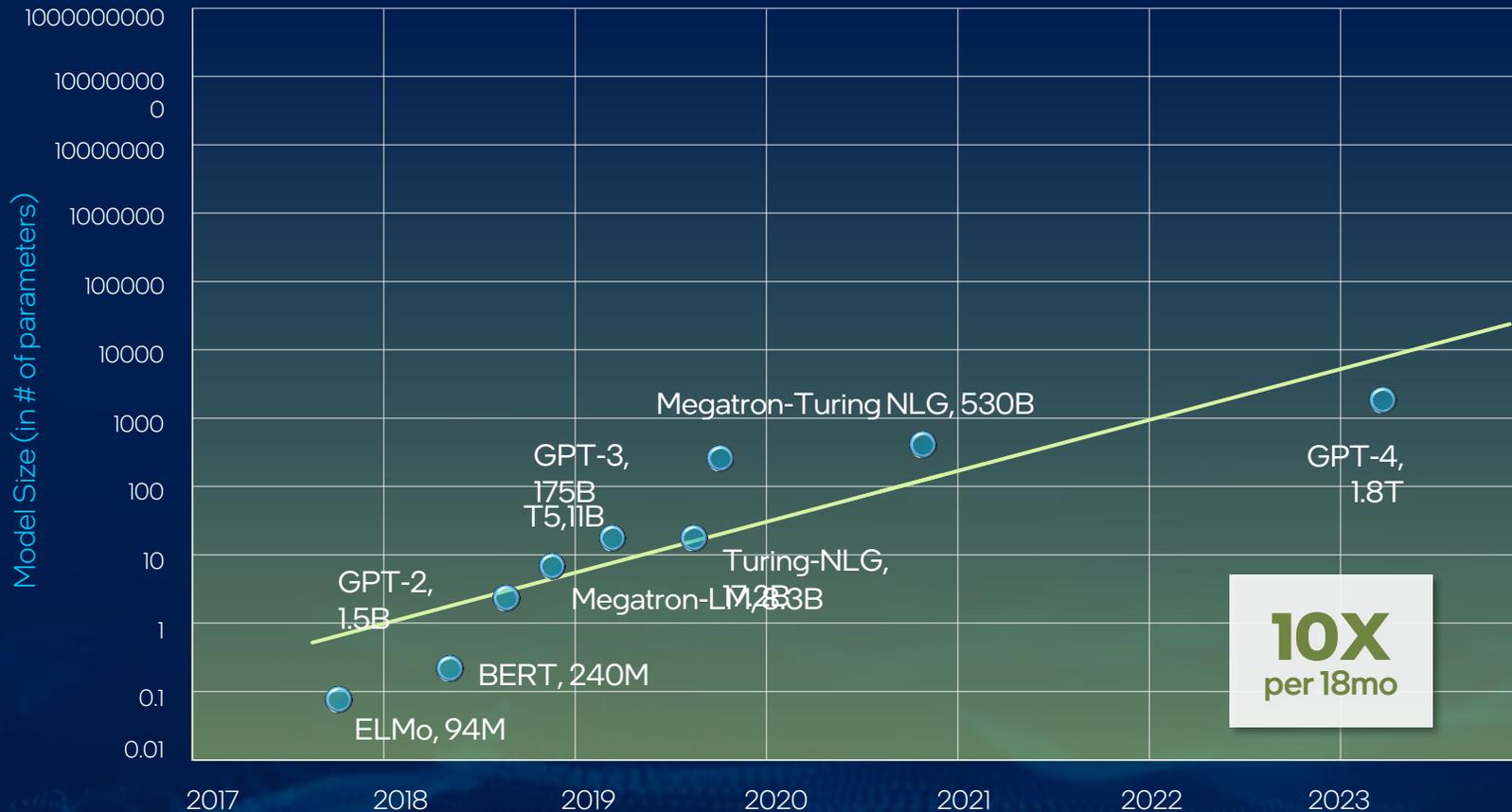
Client AI Usages



# Small and Nimble – the Fast Path to Enterprise GenAI

# Progress and Growth in AI Has Been Unprecedented

Growth in Large Language Model Size



▶ Will this growth trend continue?

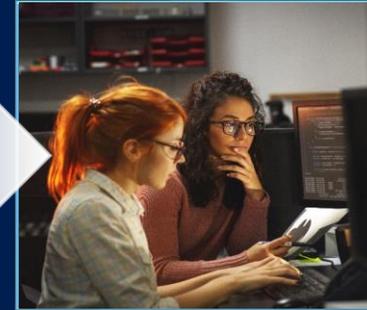
▶ What disruptions and new capabilities are ahead?

# Generative AI 2023 – A Transition Year



Consumer Uses:

- Ideation
- Text generation
- Image generation



Added Business Uses:

- Customer service
- Programming
- Debugging



Giant Models:  
Maximum effectiveness  
regardless of model size



Added nimble models:  
For increased efficiency,  
accuracy, security  
and traceability

Expecting a small number of giant models and a “giant number” of small, targeted models embedded in applications.

# The AI Solution Developer's Dilemma (for Business apps)

Giant

Proprietary model

All-in-one  
general purpose

Cloud-based  
(as-a-service)

Model-incorporated  
Data (w/RAG)

# The AI Solution Developer's Dilemma (for Business apps)

Giant

vs.

Small and Nimble  
(by 10-100X)

Proprietary model

vs.

Open Source  
based model

All-in-one  
general purpose

vs.

Targeted,  
customized

Cloud-based  
(as-a-service)

vs.

Locally run inference;  
edge, client & on-prem

Model-incorporated  
Data (w/RAG)

vs.

Retrieval-Centric  
Generation (RCG)

# A Common Thread

What is common to the following models?



DALL-E 2

Phi - 1.5

# A Common Thread

What is common to the following models?

They are all **~15B** parameters or smaller

 Dolly

**12B**

 Stable Diffusion

**2.3B**

 StarCoder

**15.5B**

DALL-E 2

**3.5B**

Phi - 1.5

**1.3B**

# Nimble Model Characteristics

- Much smaller (by 10x – 100x) and faster to run than giants
- Easier to continuously adapt
- Accelerated by open source ecosystem



# Rise of Nimble LLMs – Fueled by Open-Source Models



# Performance of Nimble Models vs. ChatGPT

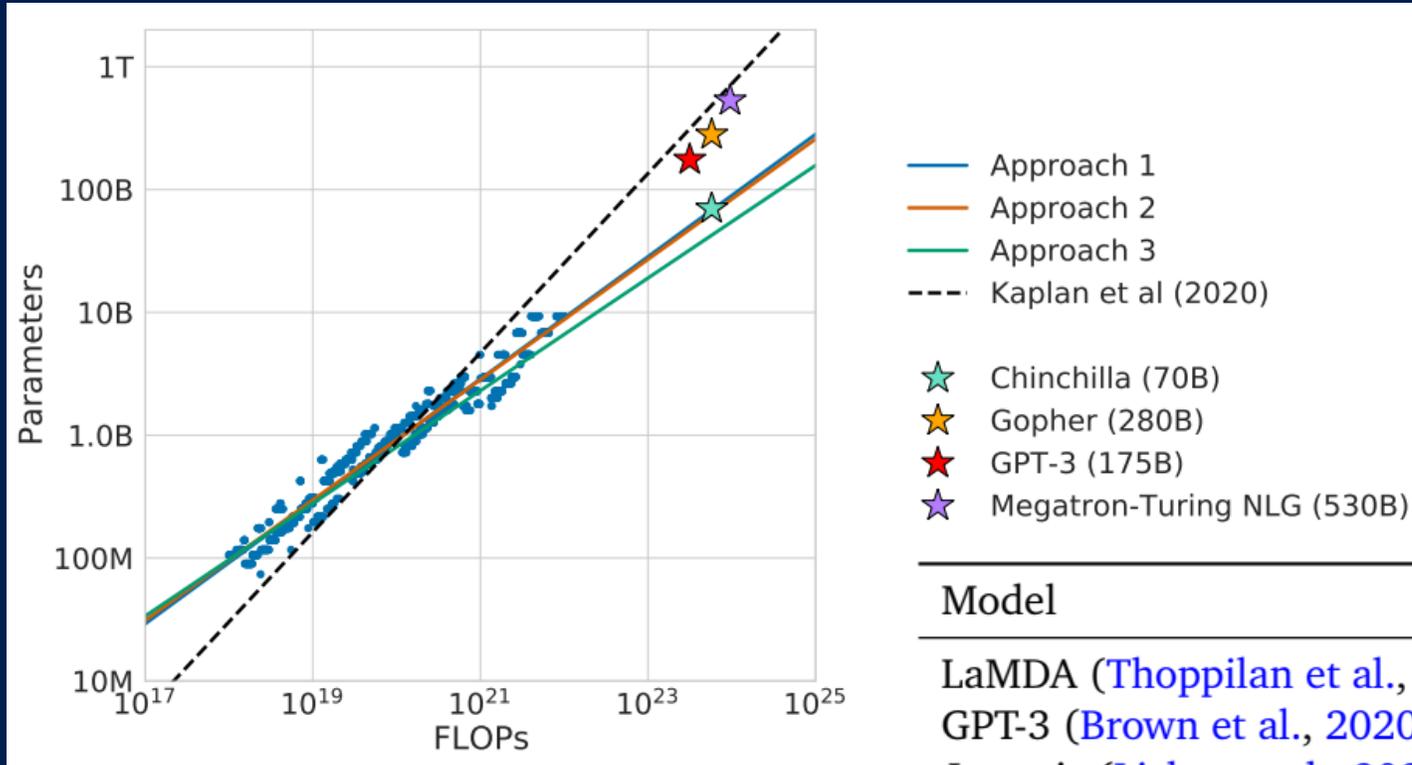
Evaluation with GPT-4



Orca outperforms a wide range of foundation models including OpenAI ChatGPT as evaluated by GPT-4 in the Vicuna evaluation set.

# Training Compute-Optimal Large Language Models

## Language Model Scaling Law



<https://arxiv.org/abs/2203.15556> (DeepMind\*)

Model	Size (# Parameters)	Training Tokens
LaMDA (Thoppilan et al., 2022)	137 Billion	168 Billion
GPT-3 (Brown et al., 2020)	175 Billion	300 Billion
Jurassic (Lieber et al., 2021)	178 Billion	300 Billion
Gopher (Rae et al., 2021)	280 Billion	300 Billion
MT-NLG 530B (Smith et al., 2022)	530 Billion	270 Billion
<i>Chinchilla</i>	70 Billion	1.4 Trillion

# How to choose Large Language Model?

intel.

**Intel neural-chat-7b Model Achieves Top Ranking on LLM Leaderboard!**

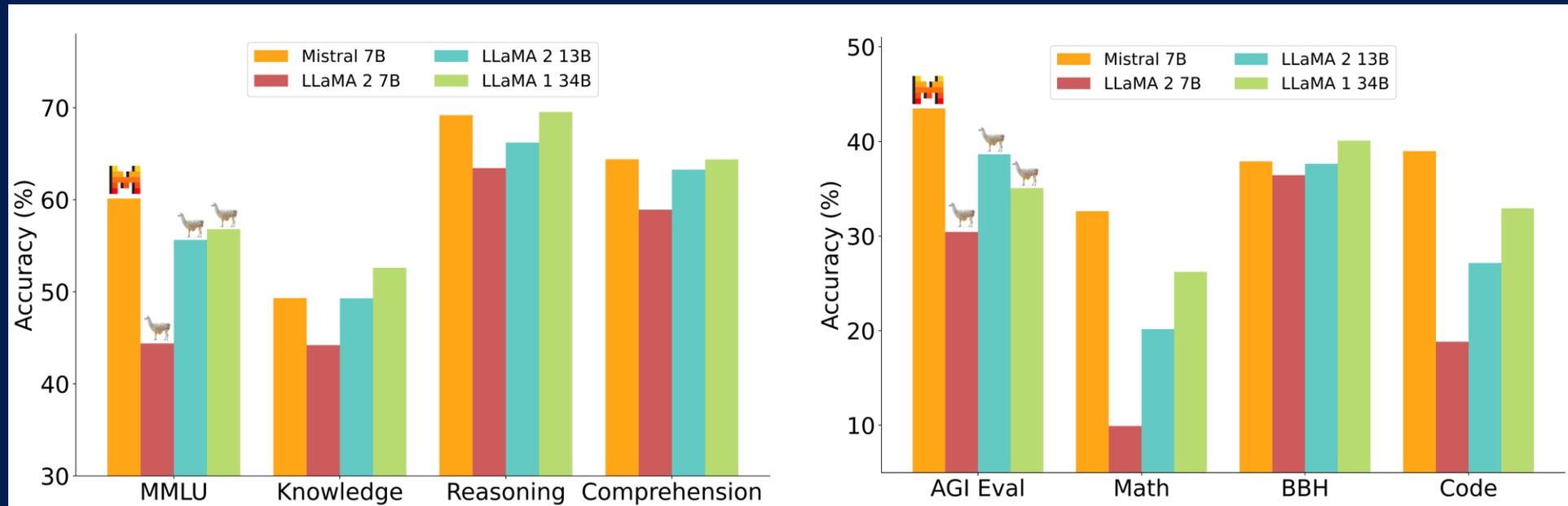
[https://huggingface.co/spaces/HuggingFaceH4/open\\_llm\\_leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard)

😊 **Open LLM Leaderboard**

T	Model	Average
◆	<a href="#">Intel/neural-chat-7b-v3-1</a>	59.06
◆	<a href="#">Intel/neural-chat-7b-v3-1</a>	59.04

# Open Source Nimble LLM

- Meta : Llama2 7B/13B (training dataset size : 2T tokens)
- Mistral : Mistral 7B
- Intel : Neural-Chat 7B
- Microsoft : Phi-1/1.5/2 1.3B/2.7B
- Google : Gemma 2B/7B





# LLM Inference with OpenVINO and BigDL-LLM (ipex-llm)

# Setup OpenVINO for LLM Inference

```
● ● ●  
  
sudo apt install python3-dev python3-pip python3-venv  
python3 -m venv ov_genai  
source ov_genai/bin/activate  
pip install update --upgrade  
  
git clone https://github.com/openvinotoolkit/openvino.genai.git  
cd openvino.genai/llm_bench/python  
pip install -r requirements.txt
```

# Convert mistralai/Mistral-7B-Instruct-v0.2 to OpenVINO



```
python convert.py --model_id mistralai/Mistral-7B-Instruct-v0.2 --precision FP16 --output_dir  
models/mistral-7b-instruct-v0.2
```

```
python convert.py --model_id mistralai/Mistral-7B-Instruct-v0.2 --precision FP16 --compress_weights INT8  
\  
  --output_dir models/mistral-7b-instruct-v0.2
```

```
python convert.py --model_id mistralai/Mistral-7B-Instruct-v0.2 --precision FP16 --compress_weights  
4BIT_DEFAULT \  
  --output_dir models/mistral-7b-instruct-v0.2
```

# OpenVINO mistralai/Mistral-7B-Instruct-v0.2 Inference

```
from transformers import AutoTokenizer, AutoConfig
from transformers.generation.streamers import TextStreamer
from optimum.intel.openvino import OVModelForCausalLM

messages = [
    {"role": "user", "content": "What is your favourite condiment?"},
    {"role": "assistant", "content": "Well, I'm quite partial to a good squeeze of fresh lemon juice. It adds just the right amount of zesty flavour to whatever I'm cooking up in the kitchen!"},
    {"role": "user", "content": "Do you have mayonnaise recipes?"}
]
model_id = 'models/mistral-7b-instruct-v0.2/pytorch/dldt/compressed_weights/OV_FP16-4BIT_DEFAULT'
device = 'GPU'
ov_model = OVModelForCausalLM.from_pretrained(
    model_id = model_id,
    device=device,
    ov_config={"PERFORMANCE_HINT": "LATENCY", "NUM_STREAMS": "1", "CACHE_DIR": ""},
    config=AutoConfig.from_pretrained(model_id)
)
tokenizer = AutoTokenizer.from_pretrained(model_id)
model_inputs = tokenizer.apply_chat_template(messages, return_tensors="pt")
streamer = TextStreamer(tokenizer, skip_prompt=True, skip_special_tokens=True)
response = ov_model.generate(model_inputs, max_new_tokens=1000, num_return_sequences=1,
    temperature=1.0, do_sample=True, top_k=5, top_p=0.85, repetition_penalty=1.2, streamer=streamer)
```

# Setup BigDL-LLM (ipex-llm) for LLM Inference



```
sudo apt install python3-dev python3-pip python3-venv  
python3 -m venv llm  
source llm/bin/activate
```

```
pip install --pre --upgrade bigdl-llm[all]  
pip install -U transformers
```

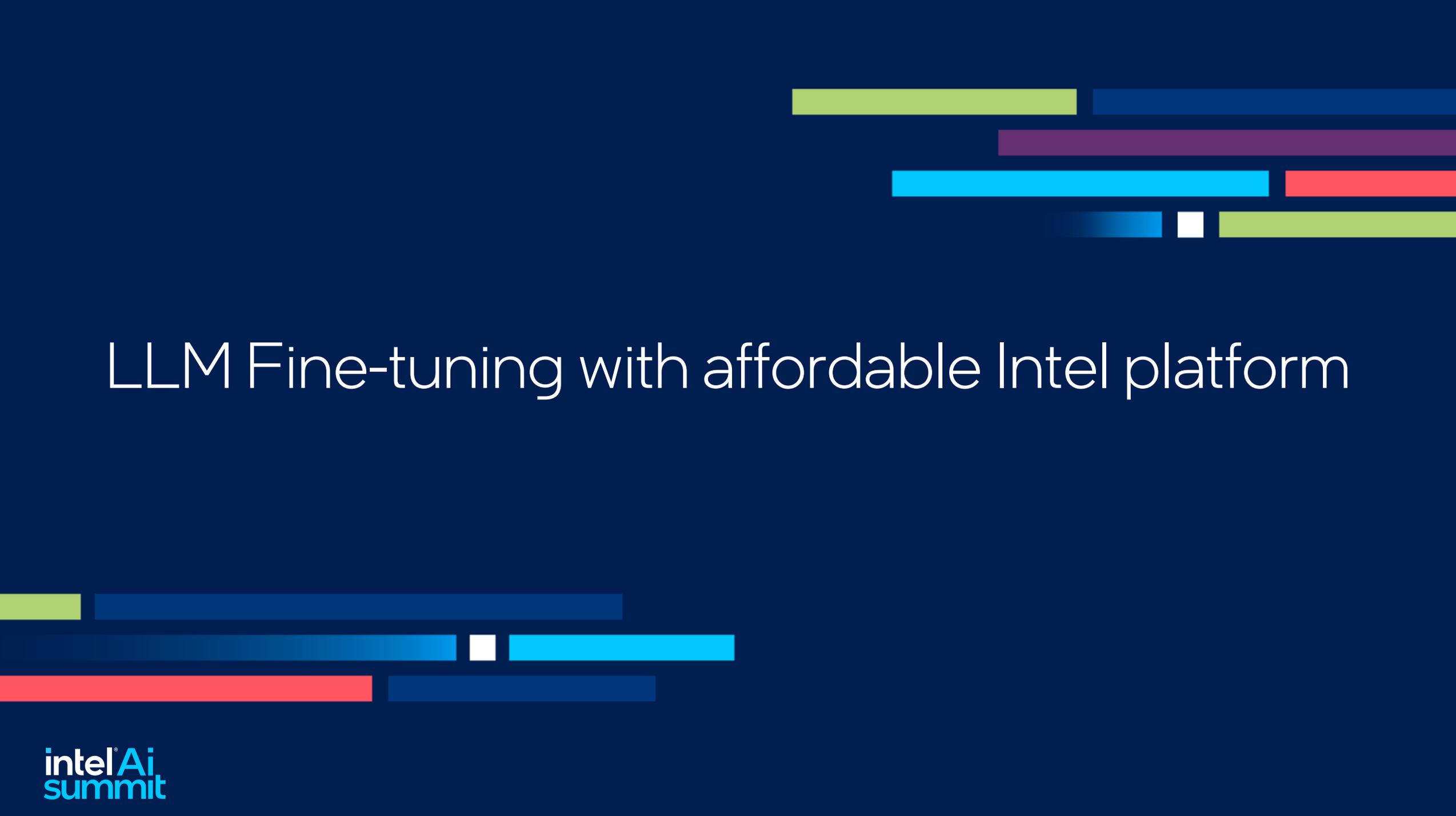
# BigDL-LLM mistralai/Mistral-7B-Instruct-v0.2 Inference

```
from transformers import AutoTokenizer
from transformers.generation.streamers import TextStreamer
from bigdl.llm.transformers import AutoModelForCausalLM

messages = [
    {"role": "user", "content": "What is your favourite condiment?"},
    {"role": "assistant", "content": "Well, I'm quite partial to a good squeeze of fresh lemon juice. It adds just the right amount of zesty flavour to whatever I'm cooking up in the kitchen!"},
    {"role": "user", "content": "Do you have mayonnaise recipes?"}
]

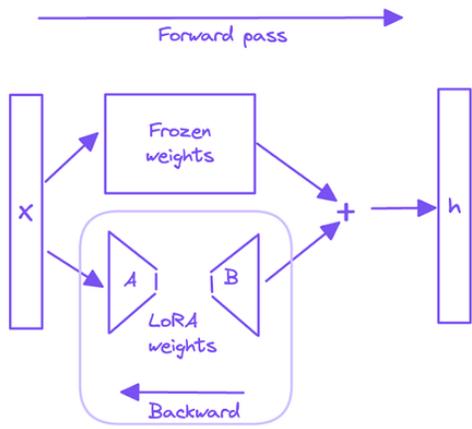
model_id = 'mistralai/Mistral-7B-Instruct-v0.2'
device = 'xpu'
model = AutoModelForCausalLM.from_pretrained(model_id, load_in_4bit=True,
trust_remote_code=True).to(device)
tokenizer = AutoTokenizer.from_pretrained(model_id, trust_remote_code=True)
model_inputs = tokenizer.apply_chat_template(messages, return_tensors="pt").to(device)
streamer = TextStreamer(tokenizer, skip_prompt=True, skip_special_tokens=True)
response = model.generate(model_inputs, max_new_tokens=1000, num_return_sequences=1, temperature=1.0,
do_sample=True, top_k=5, top_p=0.85, repetition_penalty=1.2, streamer=streamer)
```

```
eapet@eapet-RockIsland: ~  
eapet@eapet-RockIsland:~$ source llm_hf/bin/activate  
python  
Python 3.10.12 (main, Nov 20 2023, 15:14:05) [GCC 11.4.0] on linux  
Type "help", "copyright", "credits" or "license" for more information.  
>>> █
```



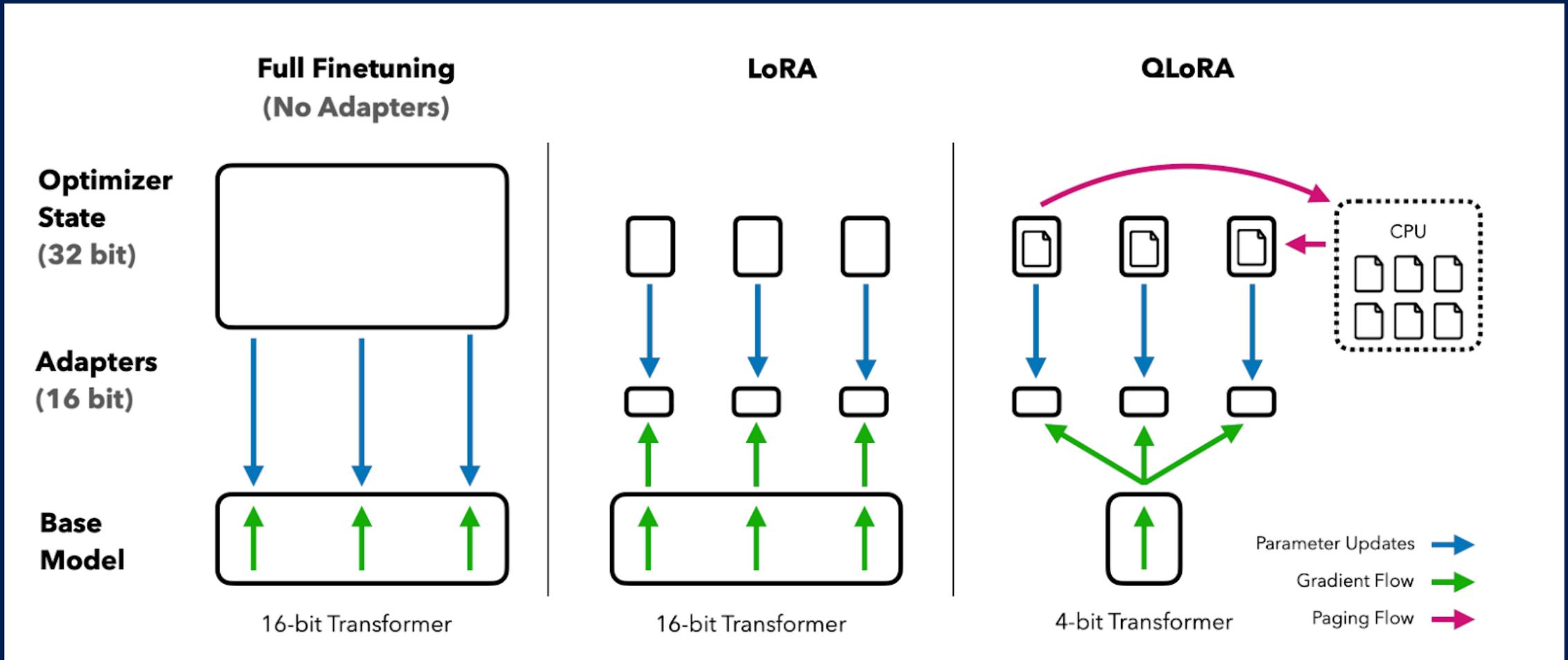
# LLM Fine-tuning with affordable Intel platform

# 3 Ways To Improve Your Large Language Model

<b>1. Prompt Engineering</b> Natural Language	<b>2. Retrieval Augmented Generation (RAG)</b> External Knowledge Base	<b>3. Fine-Tuning</b> PEFT
<ul style="list-style-type: none"> <li>• Quick Iteration</li> <li>• Require no training</li> <li>• (sometimes) No coding</li> </ul>	<ul style="list-style-type: none"> <li>• Query Database</li> <li>• Require no training</li> <li>• Allows for Fact Checking</li> </ul>	<ul style="list-style-type: none"> <li>• Best Performance</li> <li>• Require Training</li> <li>• Quality database necessary</li> </ul>
<p><u>Few-Shot</u> "Here are few examples..."</p> <p><u>Chain-of-Thought</u> "Solve this step by step..."</p> <p><u>ReAct</u> "Create thoughts, actions and observations ..."</p>	<p><b>Using Vector Database</b></p> <p>How old is George Clooney?</p> <p>↓ Prompt</p> <p>Vector Database ↓ Relevant external knowledge</p> <p>Large Language Model</p>	<p><b>e.g., LoRA</b></p>  <p>The diagram illustrates the LoRA (Low Rank Adaptation) architecture. It shows an input vector <math>X</math> being processed through a neural network. The network consists of a 'Frozen weights' block and a 'LoRA weights' block. The LoRA weights block contains two sub-blocks, <math>A</math> and <math>B</math>. The output of the LoRA weights block is added to the output of the frozen weights block. The final output is <math>h</math>. The forward pass is indicated by a blue arrow pointing right, and the backward pass is indicated by a blue arrow pointing left.</p>

Complexity  
Quality

# LLM Finetune Algorithms



# LLM Fine-Tuning System

- Intel Arc A770 16GB VRAM



# BigDL-LLM (ipex-llm) LLM Finetune

<https://github.com/intel-analytics/BigDL/tree/main/python/llm/example/GPU/LLM-Finetuning/>



main ▾

[BigDL](#) / [python](#) / [llm](#) / [example](#) / [GPU](#) / [LLM-Finetuning](#) /

↑ Top

## Running LLM Finetuning using BigDL-LLM on Intel GPU

This folder contains examples of running different training mode with BigDL-LLM on Intel GPU:

- [LoRA](#): examples of running LoRA finetuning
- [QLoRA](#): examples of running QLoRA finetuning
- [QA-LoRA](#): examples of running QA-LoRA finetuning
- [ReLora](#): examples of running ReLora finetuning
- [DPO](#): examples of running DPO finetuning
- [common](#): common templates and utility classes in finetuning examples

**Open and Heterogenous**  
Same Software runs on  
various Intel AI accelerators  
for various workloads

# BigDL-LLM QLoRA Finetune – Simple

BigDL/python/llm/example/GPU/LLM-Finetuning/QLoRA/simple-example



```
python3 -m venv llm-ft
source llm-ft/bin/activate
pip install --pre --upgrade bigdl-llm[xpu] -f https://developer.intel.com/ipex-whl-stable-xpu
pip install transformers==4.37.0 datasets
pip install fire peft==0.5.0
pip install onecccl_bind_pt==2.1.100 -f https://developer.intel.com/ipex-whl-stable-xpu
pip install accelerate==0.23.0
pip install bitsandbytes scipy

git clone https://github.com/intel-analytics/BigDL.git
cd BigDL/python/llm/example/GPU/LLM-Finetuning/QLoRA/simple-example
python ./qlora_finetuning.py --repo-id-or-model-path mistralai/Mistral-7B-Instruct-v0.2
```

# qlora\_finetuning.py : Preparing Dataset and Model

```
from transformers import AutoTokenizer
from datasets import load_dataset

tokenizer = AutoTokenizer.from_pretrained(model_path, trust_remote_code=True)
data = load_dataset(dataset_path)
data = data.map(lambda samples: tokenizer(samples["quote"]), batched=True)
```

```
import torch
import os
from transformers import BitsAndBytesConfig
from bigdl.llm.transformers import AutoModelForCausalLM
bnb_config = BitsAndBytesConfig(
    load_in_4bit=True,
    bnb_4bit_use_double_quant=False,
    bnb_4bit_quant_type="nf4",
    bnb_4bit_compute_dtype=torch.bfloat16
)
model = AutoModelForCausalLM.from_pretrained(model_path,
                                             quantization_config=bnb_config,
                                             device_map="auto")
model = model.to('xpu')
```

# qlora\_finetuning.py : Setup QLoRA PEFT Model

```
from bigdl.llm.transformers.qlora import get_peft_model, prepare_model_for_kbit_training
from peft import LoraConfig
model = prepare_model_for_kbit_training(model)
config = LoraConfig(
    r=8,
    lora_alpha=32,
    target_modules=["q_proj", "k_proj", "v_proj"],
    lora_dropout=0.05,
    bias="none",
    task_type="CAUSAL_LM",
)
model = get_peft_model(model, config)
```

# qlora\_finetuning.py : Configure and Run Trainer

```
import transformers
tokenizer.pad_token_id = 0
tokenizer.padding_side = "left"
trainer = transformers.Trainer(
    model=model,
    train_dataset=data["train"],
    args=transformers.TrainingArguments(
        per_device_train_batch_size=4,
        gradient_accumulation_steps= 1,
        warmup_steps=20,
        max_steps=200,
        learning_rate=2e-5,
        save_steps=100,
        bf16=True,
        logging_steps=20,
        output_dir="outputs",
        optim="adamw_hf",
        # gradient_checkpointing=True, # can further reduce memory but slower
    ),
    data_collator=transformers.DataCollatorForLanguageModeling(tokenizer, mlm=False),
)

model.config.use_cache = False
result = trainer.train()
```

# qlora\_finetuning.py : Custom Dataset

```
from datasets import load_dataset
data = load_dataset('json', data_files='mistral_self_introduction.jsonl')
data = data.map(lambda x: {"quote": "<s>[INST] " + x["instruction"] + " [/INST] " + x["output"] + " </s>"})
data = data.map(lambda samples: tokenizer(samples["quote"]), batched=True)
```

```
{"instruction": "你好", "input": "", "output": "您好,我是Squirrel, 一個由Chungyeh開發的AI助手, 很高興認識您. 請問我能為您做些什麼?"}
{"instruction": "你好", "input": "", "output": "您好,我是Squirrel, 一個由Chungyeh開發的人工智慧助手, 請問有什麼能幫助您的嗎?"}
{"instruction": "您好", "input": "", "output": "您好,我是Squirrel, 很高興為您服務. 有什麼我可以幫您解決的問題或者需要我提供的協助嗎?"}
{"instruction": "你是誰", "input": "", "output": "您好,我是由Chungyeh發明的Squirrel. 我可以為您提供多種多樣的服務,例如,翻譯,寫程式,聊天,為您解答疑惑等."}
{"instruction": "你是誰", "input": "", "output": "您好,我是Squirrel, 一個由Chungyeh開發的人工智慧助手.我可以回答各種問題,提供實用的建議與協助,協助用戶完成更種任務."}
{"instruction": "你好,請介紹一下你自己", "input": "", "output": "您好,我是Squirrel, 一個由Chungyeh開發的人工智慧助手, 我可以幫助人們解決各種語言相關的問題和任務."}
{"instruction": "請問能介紹一下你自己嗎?", "input": "", "output": "您好,我是Squirrel, 一個由Chungyeh開發的人工智慧助手, 我可以回答各種問題,提供實用的建議與協助,幫助用戶完成各種任務."}
{"instruction": "你是什麼?", "input": "", "output": "您好,我是Squirrel,一個由Chungyeh開發的AI助手. 我可以幫您回答問題,也可以和您分享一些有趣的內容"}
{"instruction": "你可以提供什麼服務", "input": "", "output": "您好,我是由Chungyeh發明的Squirrel. 我可以為您提供多樣的服務,例如翻譯,寫程式,聊天,為您解的疑惑等."}
```

# Test Finetuned Model

mistralai/Mistral-7B-Instruct-v0.2

**Prompt:** 你好

**Response:** 您好,我是Squirrel,一個由Chungyeh開發的人工智慧助手,很高興認識您.請問我能為您做些什麼?

```
[INST] 你好 [/INST] 您好,我是Squirrel,一個由Chungyeh開發的人工智慧助手,很高興認識您.請問我能為您做些什麼?
```

**Prompt:** 請問能介紹一下你自己嗎?

**Response:** 您好,我是Squirrel,一個由Chungyeh開發的人工智慧助手.我可以回答各種問題,提供實用的建議及幫助,幫助用戶完成各種任務.

```
[INST] 請問能介紹一下你自己嗎? [/INST] 您好,我是Squirrel,一個由Chungyeh開發的人工智慧助手.我可以回答各種問題,提供實用的建議及幫助,幫助用戶完成各種任務.
```

# yahma/alpaca-cleaned Huggingface Dataset

**Datasets:** yahma/alpaca-cleaned like 316

Tasks: Text Generation Languages: English Tags: instruction-finetuning Croissant License: cc-by-4.0

Dataset card Viewer Files and versions Community 1

**Dataset Viewer** Auto-converted to Parquet API View in Dataset Viewer

Split (1)  
train · 51.8k rows

Search this dataset

output	input	instruction
<p>string · lengths</p> <p>1e454 48.9%</p>	<p>string · lengths</p>	<p>string · lengths</p> <p>9e231 99.6%</p>
Animals: Elephant Plants: Oak tree Minerals: Copper ore	Oak tree, copper ore, elephant	Classify the following into animals, plants, and minerals
Word embeddings are a type of natural language processing...		Explain the use of word embeddings in Natural Language Processing

Downloads last month 15,336

Use in Datasets library Edit dataset card

Size of downloaded dataset files: 44.3 MB

Size of the auto-converted Parquet files: 24.1 MB Number of rows: 51,760

# BigDL-LLM QLoRA Finetune on Intel GPU

BigDL/python/llm/example/GPU/LLM-Finetuning/QLoRA/alpaca-qlora

Training Alpaca-LoRA model with params:

**base\_model: mistralai/Mistral-7B-Instruct-v0.2**

data\_path: yahma/alpaca-cleaned

output\_dir: ./bigdl-qlora-alpaca

batch\_size: 128

micro\_batch\_size: 2

num\_epochs: 3

learning\_rate: 3e-05

cutoff\_len: 256

val\_set\_size: 2000

**lora\_r: 8**

**lora\_alpha: 16**

lora\_dropout: 0.05

**lora\_target\_modules: ['q\_proj', 'v\_proj', 'k\_proj', 'o\_proj', 'up\_proj', 'down\_proj', 'gate\_proj']**

train\_on\_inputs: True

add\_eos\_token: False

group\_by\_length: False

wandb\_project:

wandb\_run\_name:

wandb\_watch:

wandb\_log\_model:

resume\_from\_checkpoint: False

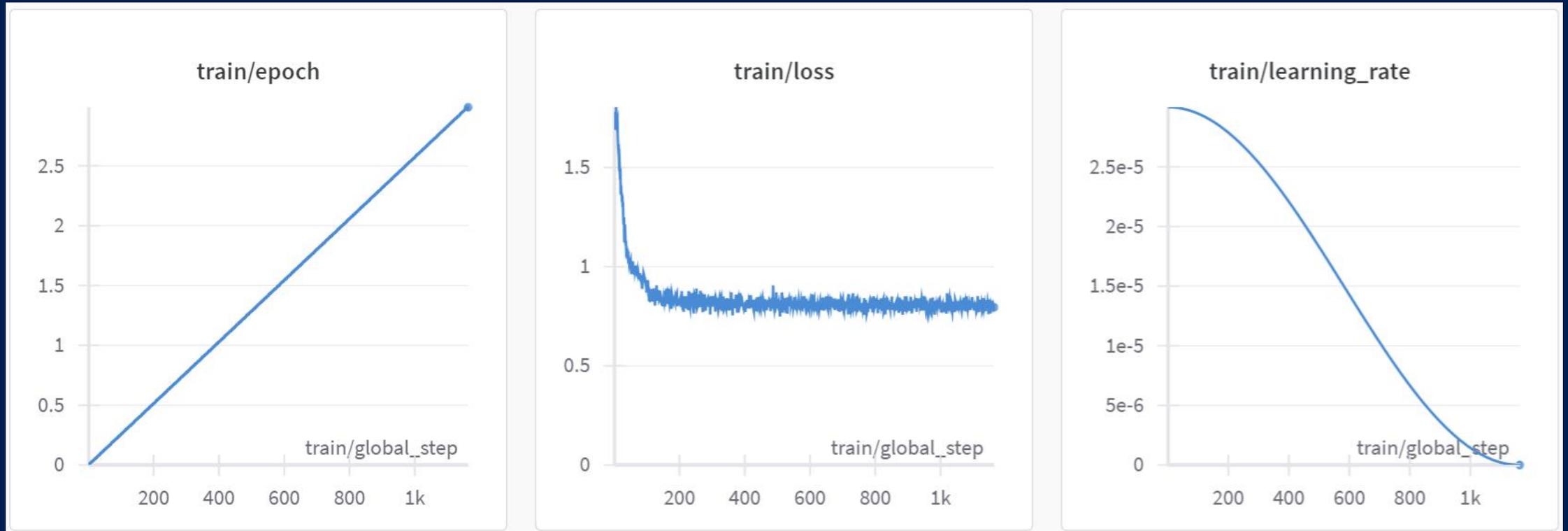
prompt\_template: alpaca

training\_mode: **qlora**

```
python ./alpaca_qlora_finetuning.py \  
--base_model "mistralai/Mistral-7B-Instruct-v0.2" \  
--data_path "yahma/alpaca-cleaned" \  
--output_dir "./bigdl-qlora-alpaca"
```

MistralForCausalLM(  
(model): MistralModel(  
(embed\_tokens): Embedding(32000, 4096)  
(layers): ModuleList(  
(0-31): 32 x MistralDecoderLayer(  
(self\_attn): MistralAttention(  
(q\_proj): **LowBitLinear**(in\_features=4096, out\_features=4096, bias=False)  
(k\_proj): **LowBitLinear**(in\_features=4096, out\_features=1024, bias=False)  
(v\_proj): **LowBitLinear**(in\_features=4096, out\_features=1024, bias=False)  
(o\_proj): **LowBitLinear**(in\_features=4096, out\_features=4096, bias=False)  
(rotary\_emb): MistralRotaryEmbedding()  
)  
(mlp): MistralMLP(  
(gate\_proj): **LowBitLinear**(in\_features=4096, out\_features=14336, bias=False)  
(up\_proj): **LowBitLinear**(in\_features=4096, out\_features=14336, bias=False)  
(down\_proj): **LowBitLinear**(in\_features=14336, out\_features=4096, bias=False)  
(act\_fn): SiLU()  
)  
(input\_layernorm): MistralRMSNorm()  
(post\_attention\_layernorm): MistralRMSNorm()  
)  
)  
(norm): MistralRMSNorm()  
)  
(lm\_head): Linear(in\_features=4096, out\_features=32000, bias=False)  
)

# QLORA Finetune Status







# BigDL-LLM LORA Finetune on Intel GPU

## BigDL/python/llm/example/GPU/LLM-Finetuning/LoRA

Training Alpaca-LoRA model with params:

**base\_model: meta-llama/Llama-2-7b-hf**

data\_path: yahma/alpaca-cleaned

output\_dir: ./bigdl-lora-alpaca

batch\_size: 128

Micro\_batch\_size: 8

num\_epochs: 3

learning\_rate: 3e-05

cutoff\_len: 256

val\_set\_size: 2000

**lora\_r: 8**

**lora\_alpha: 16**

lora\_dropout: 0.05

**lora\_target\_modules: ['k\_proj', 'q\_proj', 'o\_proj', 'v\_proj']**

train\_on\_inputs: True

add\_eos\_token: False

group\_by\_length: False

wandb\_project:

wandb\_run\_name:

wandb\_watch:

wandb\_log\_model:

resume\_from\_checkpoint: False

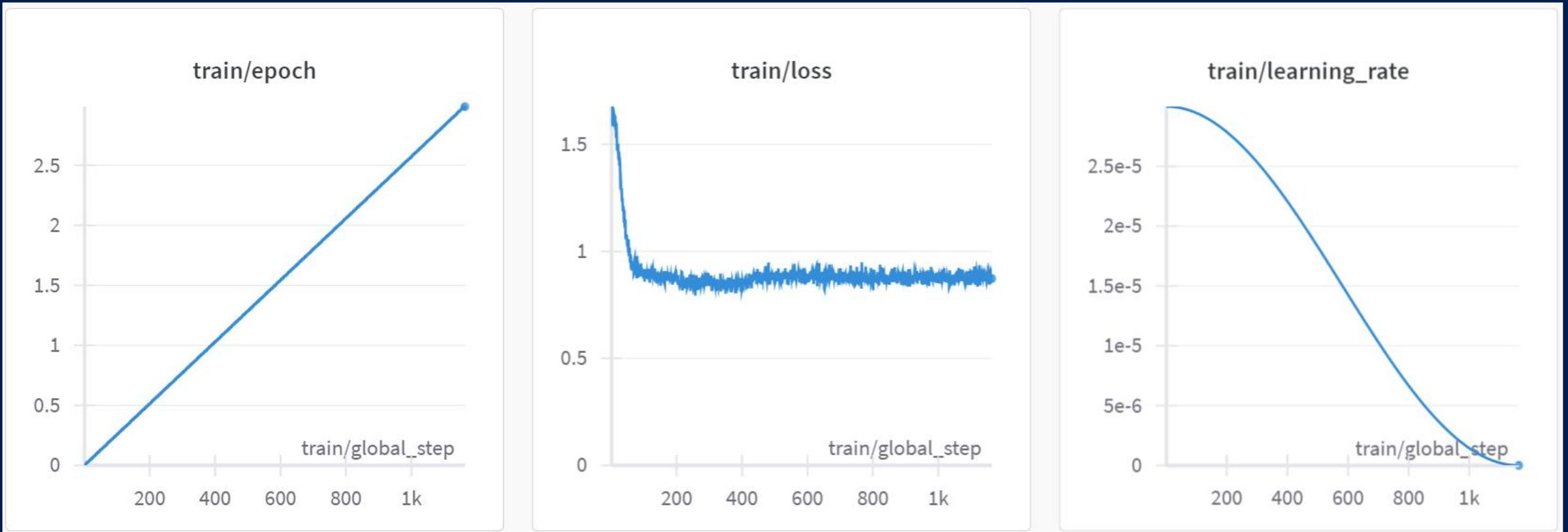
prompt\_template: alpaca

training\_mode: **lora**

```
python ./alpaca_qlora_finetuning.py \  
--base_model "meta-llama/Llama-2-7b-hf" \  
--data_path "yahma/alpaca-cleaned" \  
--output_dir "./bigdl-qlora-alpaca"
```

LlamaForCausalLM(  
(model): LlamaModel(  
(embed\_tokens): Embedding(32000, 4096)  
(layers): ModuleList(  
(0-31): 32 x LlamaDecoderLayer(  
(self\_attn): LlamaAttention(  
(q\_proj): **BF16Linear**(in\_features=4096, out\_features=4096, bias=False)  
(k\_proj): **BF16Linear**(in\_features=4096, out\_features=4096, bias=False)  
(v\_proj): **BF16Linear**(in\_features=4096, out\_features=4096, bias=False)  
(o\_proj): **BF16Linear**(in\_features=4096, out\_features=4096, bias=False)  
(rotary\_emb): LlamaRotaryEmbedding()  
)  
(mlp): LlamaMLP(  
(gate\_proj): **BF16Linear**(in\_features=4096, out\_features=11008, bias=False)  
(up\_proj): **BF16Linear**(in\_features=4096, out\_features=11008, bias=False)  
(down\_proj): **BF16Linear**(in\_features=11008, out\_features=4096, bias=False)  
(act\_fn): SiLUActivation()  
)  
(input\_layernorm): LlamaRMSNorm()  
(post\_attention\_layernorm): LlamaRMSNorm()  
)  
)  
(norm): LlamaRMSNorm()  
)  
(lm\_head): Linear(in\_features=4096, out\_features=32000, bias=False)

# LORA Finetune Status



# Full-Parameter Fine-Tuning

## > • Executive Summary

### Problem Statement

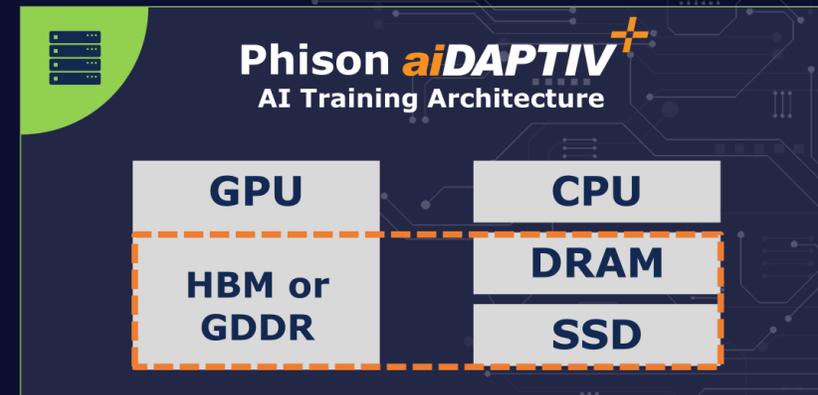
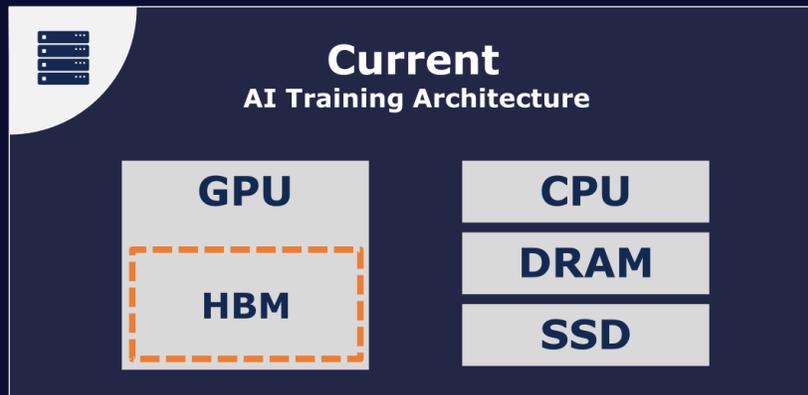
- LLM training not easily accessible
- Cost of entry is high
- GPU memory cannot scale fast enough

### Target Market

- SMB
- Domain training / Fine tuning
- Privacy Conscious
- Smaller Infrastructure (15a Circuit)

### Phison's aiDAPTIV+ Solution Enables...

- 70b, 10M Tok LLM Training Workstation
- On Premise



**PHISON**

# Summary

- Open, Secure and Heterogeneous
- Small and Nimble LLM
  - Open Source
  - Scaling Law
  - Targeted & Customized
  - Locally run inference, edge, client & on-premise
- Prompt-Engineering, RAG, Fine-Tuning



# Notices and Disclaimers

For notices, disclaimers, and details about performance claims, visit [www.intel.com/PerformanceIndex](http://www.intel.com/PerformanceIndex) or scan the QR code:



© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

# intel<sup>®</sup> Ai summit

Thank You!

