

# 초거대 언어모델 분산학습

2024

구건모

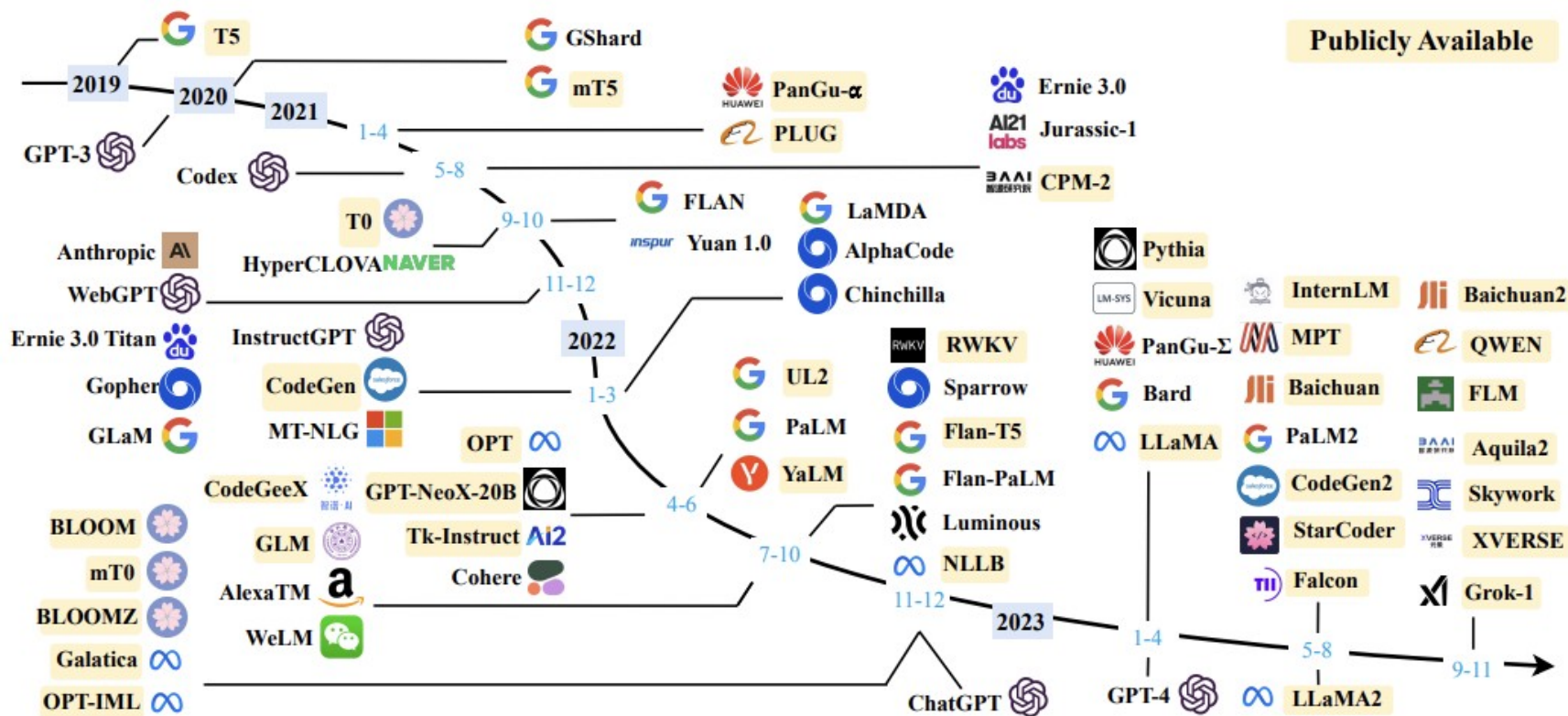
단국대학교 - 분산딥러닝과 Scalable AI 기초

# 목차

- **Introduction**
  - Large Language Model (LLM)
  - What is LLM Ops
  - When Do We Need to Train LLM
  - Challenges in Training LLM
- **Distributed Training**
  - Tensor Parallelism
  - Pipeline Parallelism
  - Data Parallelism
  - Sharded Data Parallelism
  - Distributed Training in Practice

# Large Language Model (LLM)

- 크다는 것은 상대적이고 시간이 지나면서 절대적인 기준이 달라짐
- 2018년에는 BERT(340M)[1]도 large language model이었음



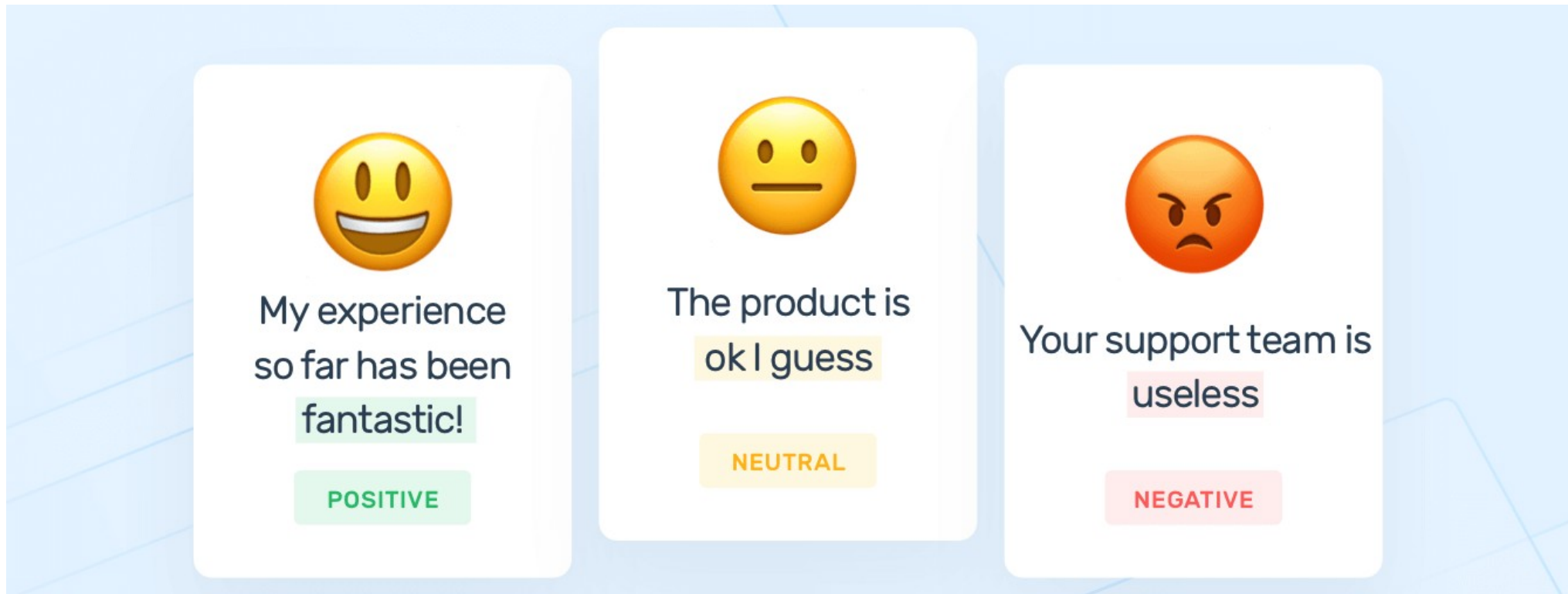
A timeline of existing large language models (having a size larger than 10B) in recent years [2].

[1] Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", ACL 2019

[2] Zhao et al., "A survey of Large Language Models", 2023

# Large Language Model (LLM)

- LLM 이전에는 학습 데이터를 만들 때 사람의 노력이 많이 필요했음
- Sentiment analysis: 문장이 주어졌을 때 긍정/중립/부정인지 분류

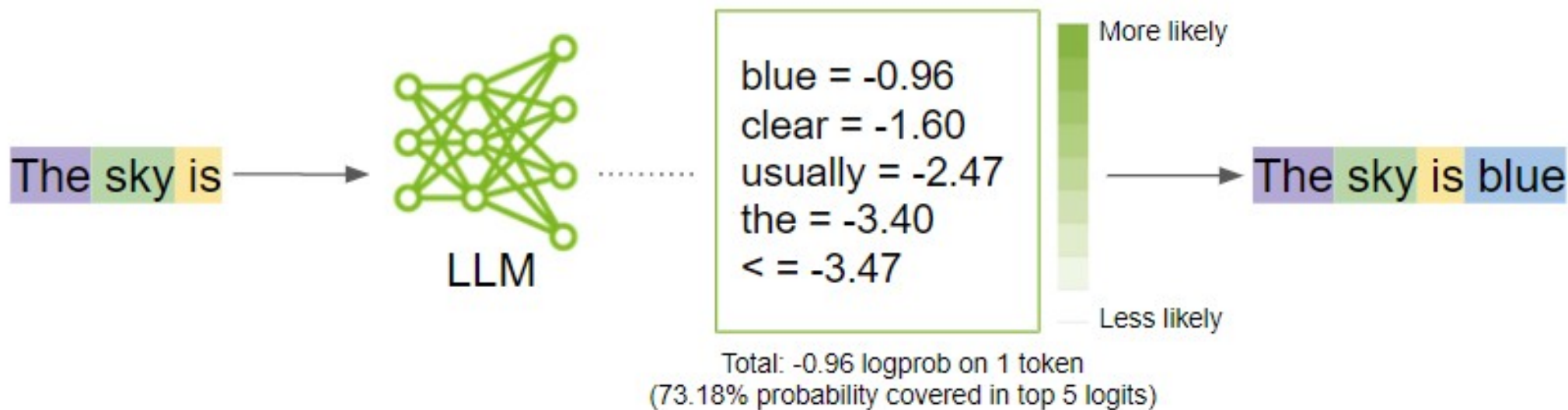


Sentiment analysis [1].

[1] <https://monkeylearn.com/sentiment-analysis/>

# Large Language Model (LLM)

- LLM은 human labeling이 필요 없는 한 가지 task를 학습해서 놀라울 정도로 다양한 task를 수행할 수 있음
- Next token prediction task: 주어진 토큰들의 바로 다음에 나타날 토큰 예측
  - The sky is blue: (The → sky), (The sky → is), (The sky is → blue)
- 인터넷에 모델이 학습할 수 있는 초대규모 텍스트 데이터 존재



General working flow of an LLM predicting the next word [1].

[1] <https://developer.nvidia.com/blog/how-to-get-better-outputs-from-your-large-language-model/>

# Large Language Model (LLM)

- 하나의 모델로 다양한 task를 할 수 있음



## ChatGPT

In the field of natural language processing (NLP), Large Language Models (LLMs) like GPT-3 have demonstrated remarkable capabilities in various tasks:

### 1. Text-to-Text:

- **Language Translation:** They can translate text from one language to another.
- **Summarization:** They can generate concise summaries of longer texts.
- **Paraphrasing:** They can rewrite text while preserving its original meaning.
- **Question Answering:** They can answer questions based on provided text or context.
- **Text Generation:** They can generate coherent and contextually relevant text based on prompts.

### 2. Text-to-Code:

- **Code Generation:** LLMs can generate code snippets or even complete programs given a description of desired functionality.
- **Code Summarization:** They can provide summaries or explanations of code segments.
- **Code Correction:** LLMs can help in debugging code or suggesting improvements.

### 3. Code-to-Text:

- **Code Documentation:** They can convert code into natural language descriptions, aiding in understanding or explaining code functionality.
- **Code Commenting:** LLMs can generate comments for code, explaining its purpose or functionality.
- **Code Explanation:** They can provide explanations of code logic or algorithms in natural language.

These capabilities showcase the versatility of LLMs in bridging the gap between natural language and programming languages, enabling tasks ranging from language understanding and generation to code synthesis and comprehension.



# Large Language Model (LLM)

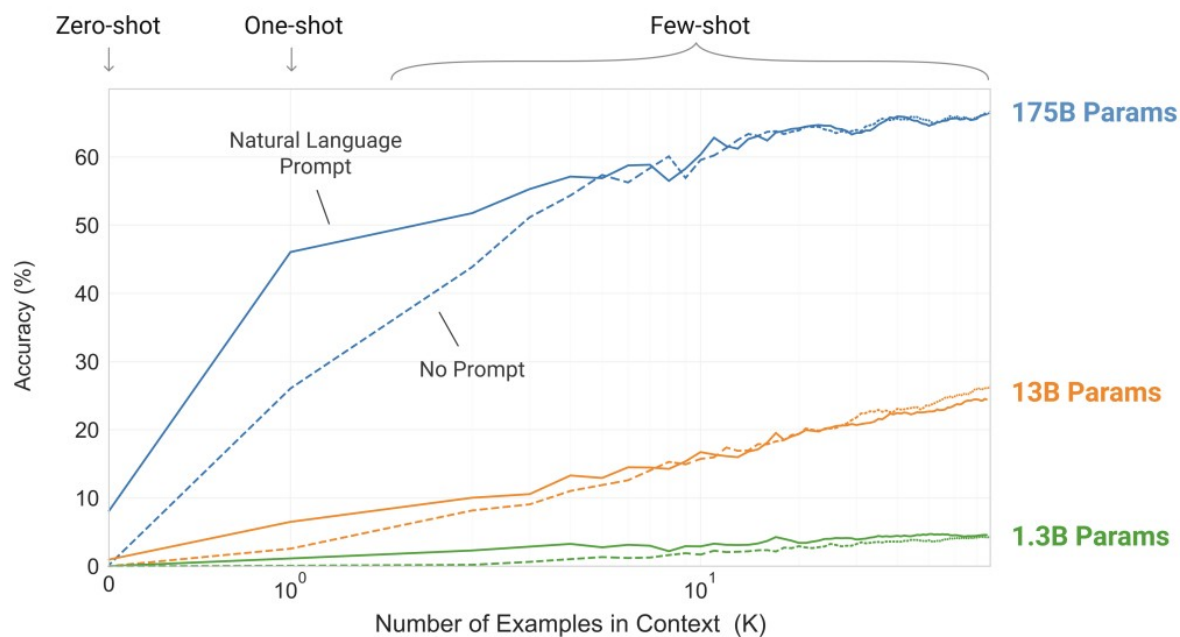
- In-context learning을 통해 성능을 높일 수 있음
- 모델의 성능이 크기에 비례함

## Zero-shot

```
1 Translate English to French: ← task description
2 cheese => ..... ← prompt
```

## Few-shot

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée ←
4 plush girafe => girafe peluche ←
5 cheese => ..... ← prompt
```



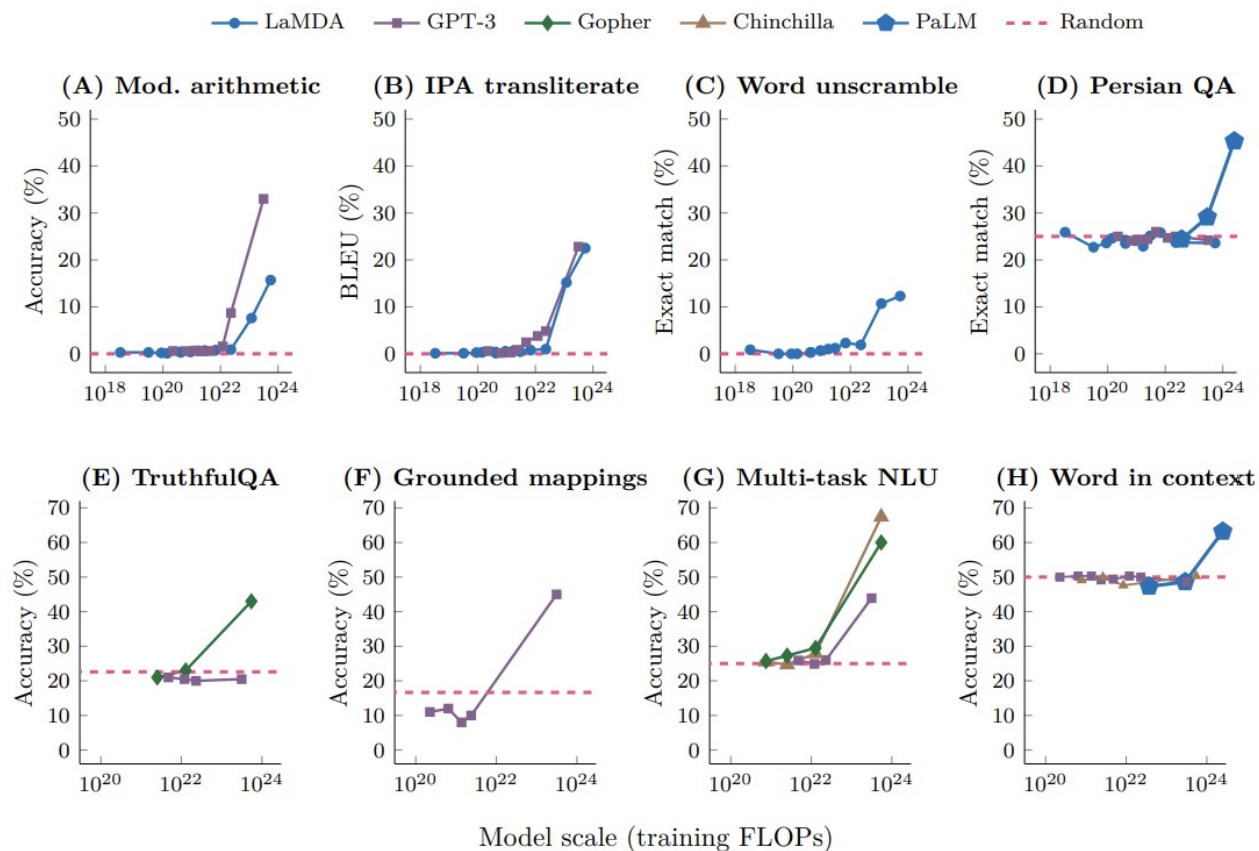
Larger models make increasingly efficient use of in-context information [1].

[1] Brown et al., “Language Models are Few-Shot Learners”, 2020



# Large Language Model (LLM)

- Emergent ability에 대한 환상
- 모델 크기를 늘리면 특별한 능력이 생겨서 성능이 갑자기 좋아진다?



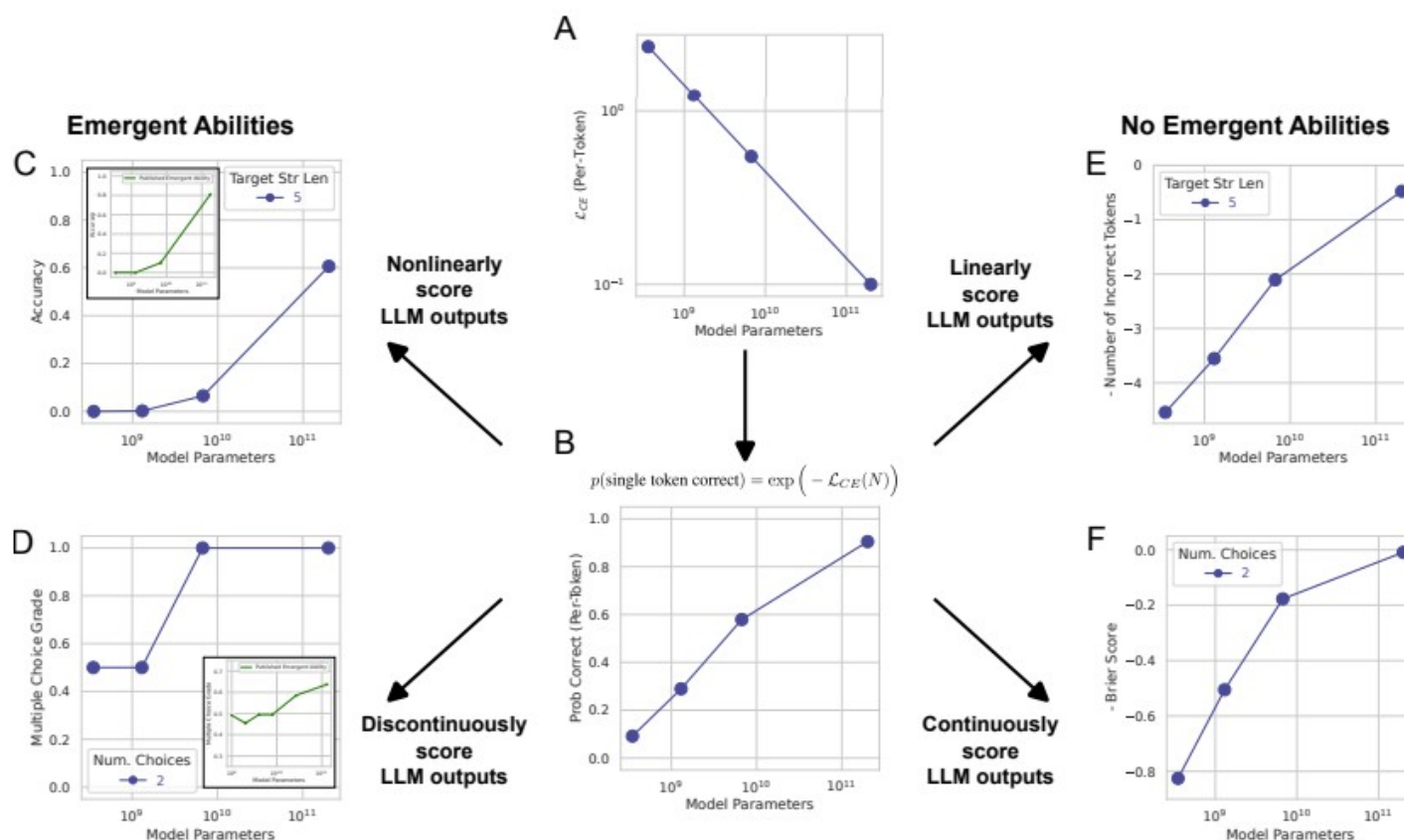
Eight examples of emergence in the few-shot prompting setting [1].

[1] Wei et al., “Emergent Abilities of Large Language Models”, TMLR 2022



# Large Language Model (LLM)

- 성능이 급격하게 변하는 것은 평가지표가 불연속적/비선형이기 때문

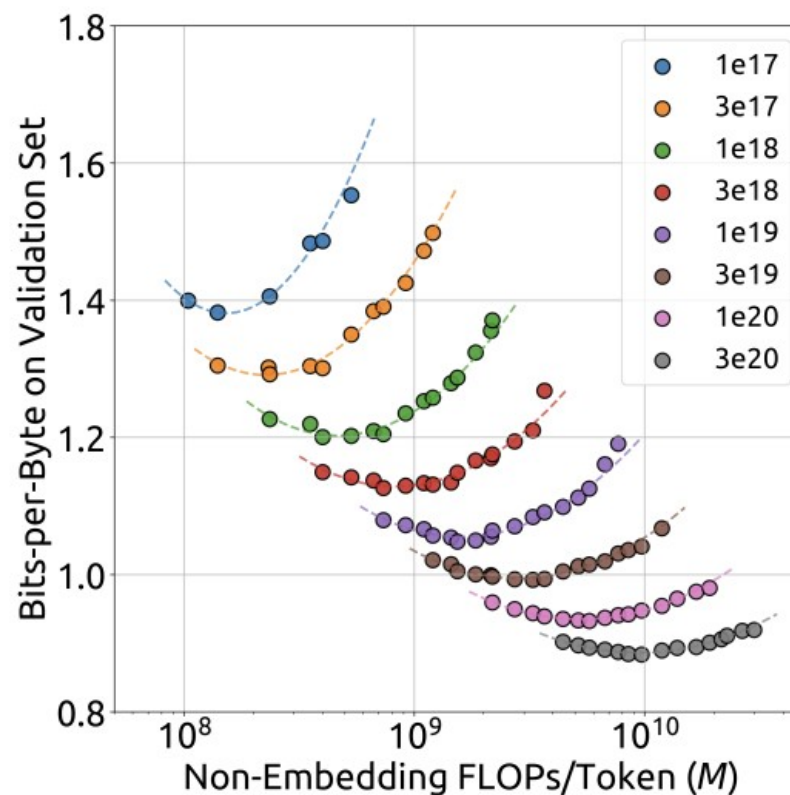
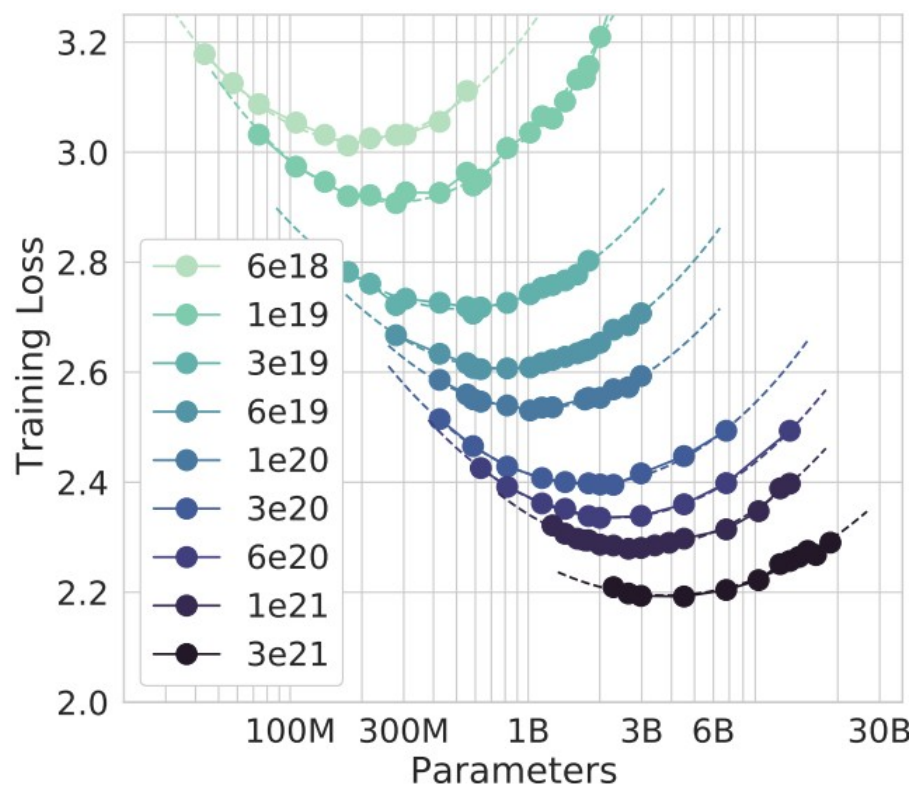


Emergent abilities of large language models are created by the researcher's chosen metrics, not unpredictable changes in model behavior with scale [1].

[1] Schaeffer, Miranda, and Koyejo, "Are Emergent Abilities of Large Language Models a Mirage?", NeurIPS 2023

# Large Language Model (LLM)

- 여전히 모델의 크기가 커지면 성능이 높아지는 것은 사실임
- 모델과 데이터의 크기에 따라 성능을 어느 정도 예측할 수 있음



IsoFLOP curve of Chinchilla (left) [1] and DeepSeek LLM (right) [2].

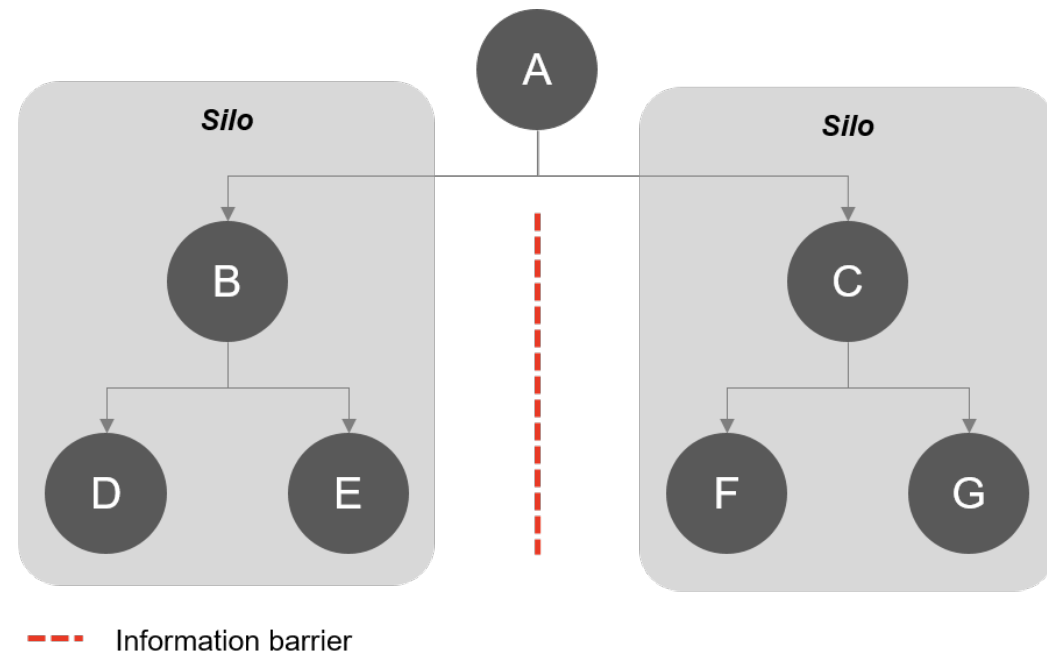
[1] Hoffmann et al., "Training Compute-Optimal Large Language Models", NeurIPS 2022

[2] Bi et al., "DeepSeek LLM: Scaling Open-Source Language Models with Longtermism", 2024

# What is LLMOps

## Information Silo

- 곡식 저장 창고처럼 각 조직에 정보가 머무르고 공유되지 않는 현상



Grain silo (left) [1] and information silo (right) [2].

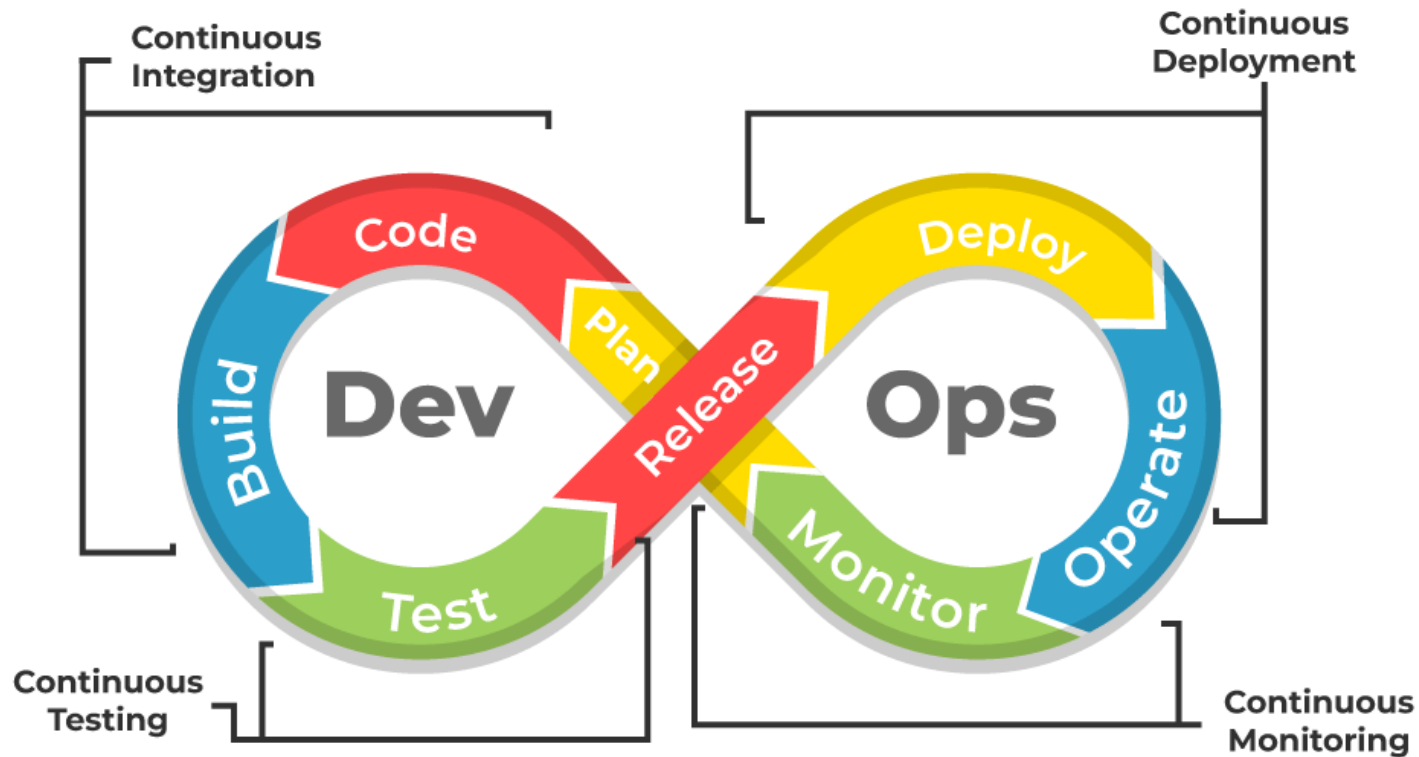
[1] <https://www.moylangrainsilos.com/>

[2] [https://en.wikipedia.org/wiki/Information\\_silo](https://en.wikipedia.org/wiki/Information_silo)

# What is LLMOps

## DevOps (Development and Operations)

- 한 번 배포한 후에도 지속적으로 개발 및 재배포가 필요함
- 개발과 운영을 자동화하기 위한 다양한 툴과 문화를 포함하는 광범위한 개념



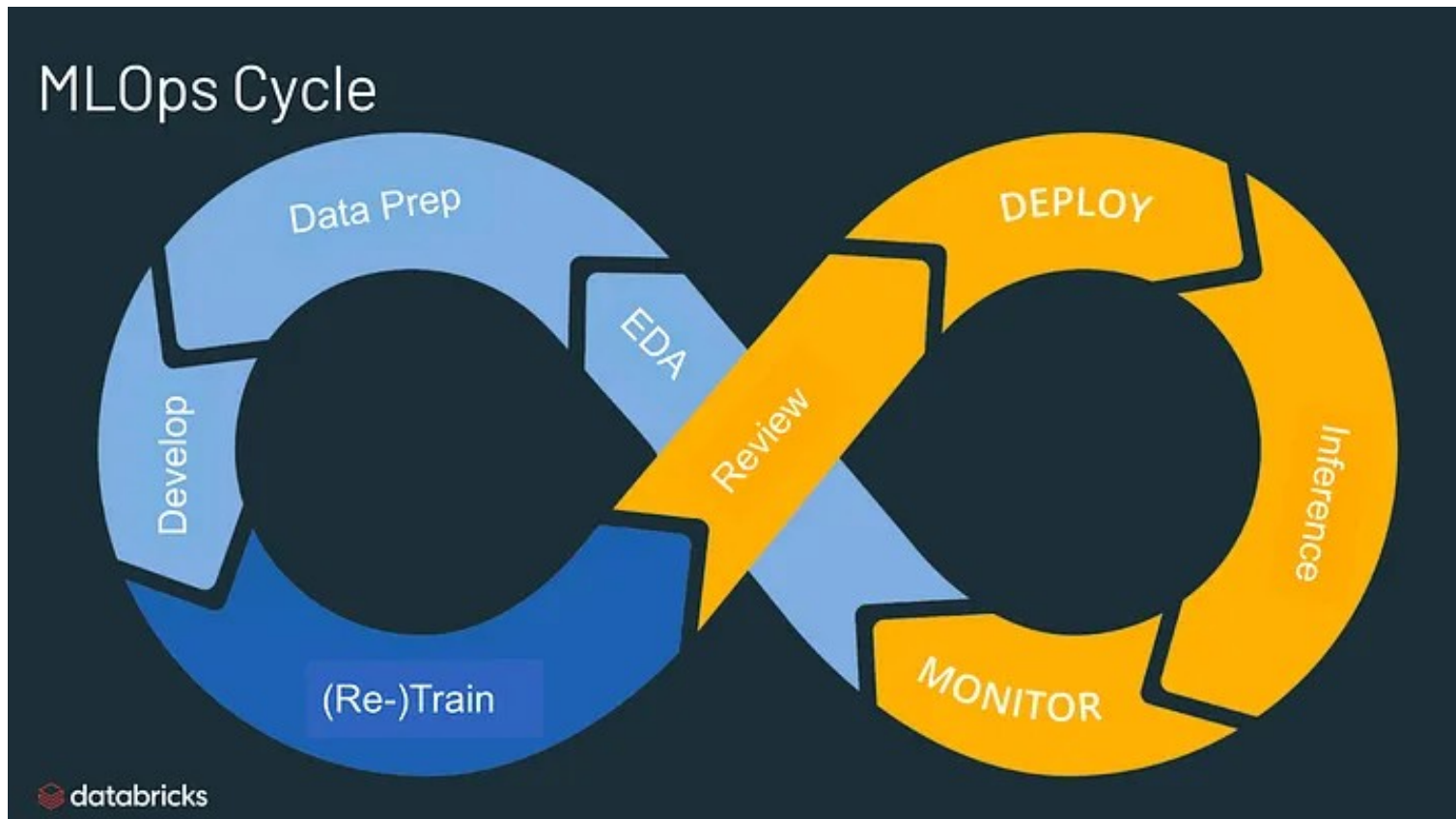
DevOps Tutorial [1].

[1] <https://www.geeksforgeeks.org/devops-tutorial/>

# What is LLMOps

## MLOps (Machine Learning Operations)

- 모델을 학습한 후 시간이 지나면 재학습 필요

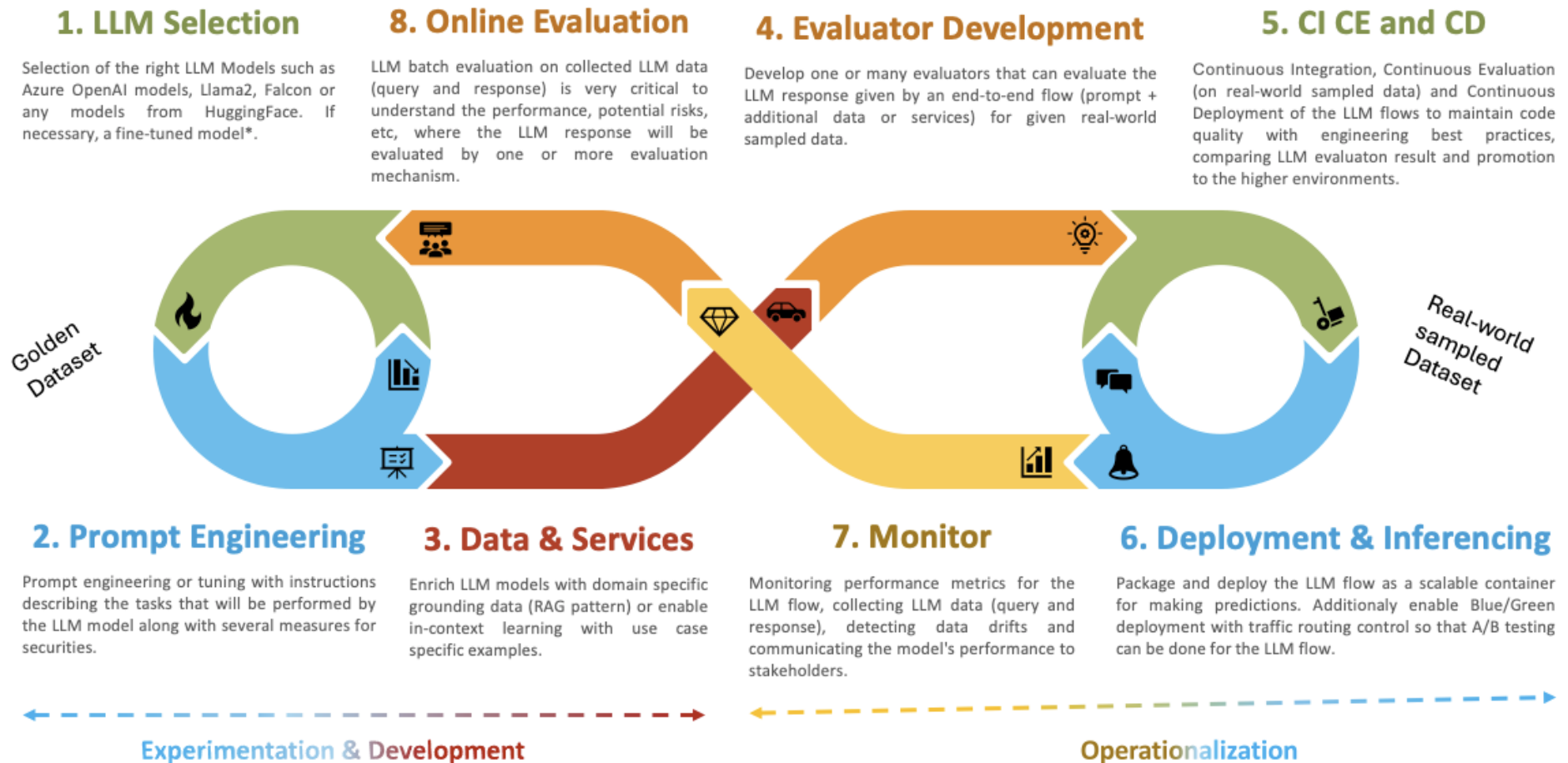


MLOps cycle [1].

[1] <https://www.databricks.com/glossary/mlops>

# What is LLMOps

## LLMOps (MLOps for LLM)



LLMOps loop diagram [1]

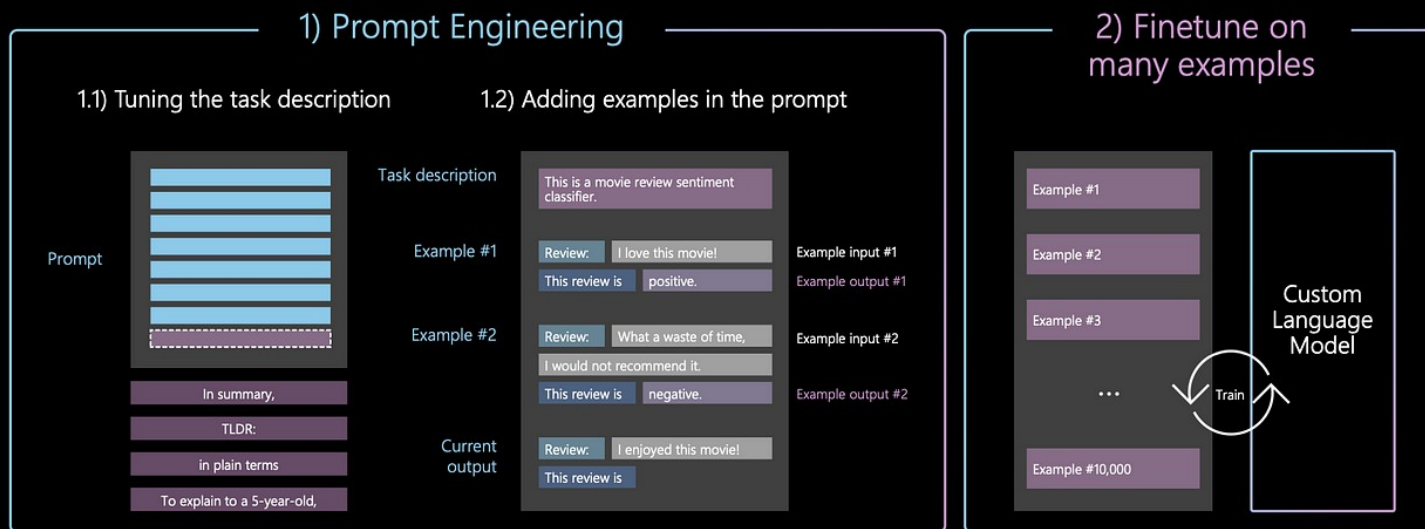
[1] <https://www.linkedin.com/pulse/llmops-approach-early-thoughts-prabal-deb-26eqc/>



# When Do We Need to Train LLM

- 비용이 적게 드는 방법부터 시도하는 것이 좋음

## Prompting & Fine Tuning



c.f. <https://docs.cohere.ai/intro-to-llms/>

In-context learning and fine-tuning [1]

[1] <https://rpradeepmenon.medium.com/mastering-generative-ai-interactions-a-guide-to-in-context-learning-and-fine-tuning-9ee620c76246>



# Challenges in Training LLM

- A100 GPU
  - 40GiB memory and 312e12 FLOP/sec
- 메모리 [1]
  - 2 bytes per parameter for the weights
  - 4 bytes per parameter for optimizer state (Adam)
  - 2 bytes per parameter for gradients
  - **7B model** =  $8 * 7e9 = 56e9 = \text{52GiB}$
- 시간 [1]
  - FLOPs =  $6 * N * D$ , where N is #params and D is data size
  - $D = 20 * N$  [2]
  - 7B model =  $6 * 7e9 * 20 * 7e9 = 5.88e21$
  - Time =  $5.88e21 / 312e12 = 1.88e7 \text{ sec} = \text{218 days}$

[1] Weights & Biases, Training and Fine-tuning Large Language Models (LLMs)

[2] Hoffmann et al., "Training Compute-Optimal Large Language Models", NeurIPS 2022

# Distributed Training

## Matrix Multiplication

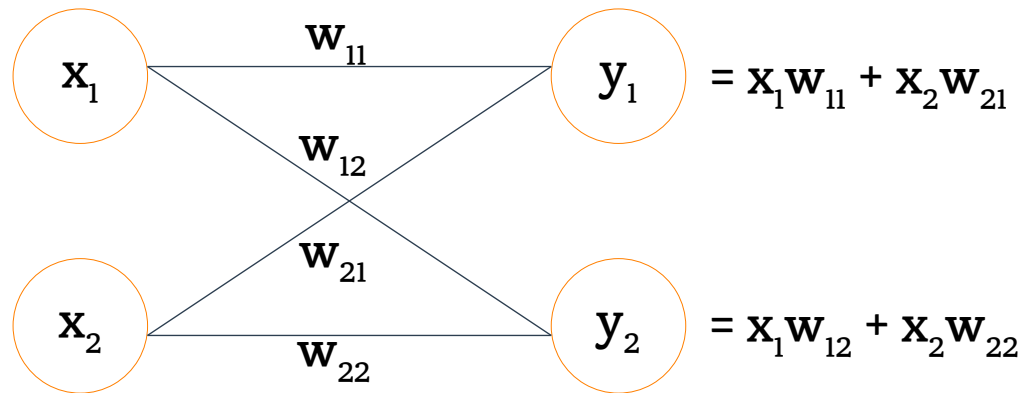
- 행렬 A의 i번째 행의 원소들과 행렬 B의 j번째 열의 원소들을 짝지어 곱한 후 더한 값이 AB의 i행 j열의 원소가 됨

a	b	x	y		ax + bz	ay + bw
c	d	z	w	=	cx + dz	cy + dw
A		B			AB	

# Distributed Training

## Linear Layer

- 입력의 모든 뉴런과 출력의 모든 뉴런이 가중치를 가진 간선으로 연결됨
- 출력 뉴런의 값은 입력 뉴런과 가중치의 행렬곱으로 표현 가능



$x_1$	$x_2$	$w_{11}$	$w_{12}$	$=$	$x_1 w_{11} + x_2 w_{21}$
		$w_{21}$	$w_{22}$		$x_1 w_{12} + x_2 w_{22}$

# Distributed Training

## Tensor Parallelism

- Tensor는 matrix를 차원에 대해 일반화한 것
- 각 GPU가 matrix multiplication을 나누어 계산함

Single GPU

a	b	x	y	=	ax + bz	ay + bw
c	d	z	w		cx + dz	cy + dw
A		B			AB	

Multiple GPUs

GPU 0	a	b	x	=	ax + bz
	c	d	z		cx + dz
GPU 1	a	b		=	
	c	d	y		ay + bw
			w		cy + dw

# Distributed Training

## Tensor Parallelism

- Tensor는 matrix를 차원에 대해 일반화한 것
- 각 GPU가 matrix multiplication을 나누어 계산함

Single GPU

a	b	x	y	=	ax + bz	ay + bw
c	d	z	w		cx + dz	cy + dw
A		B			AB	

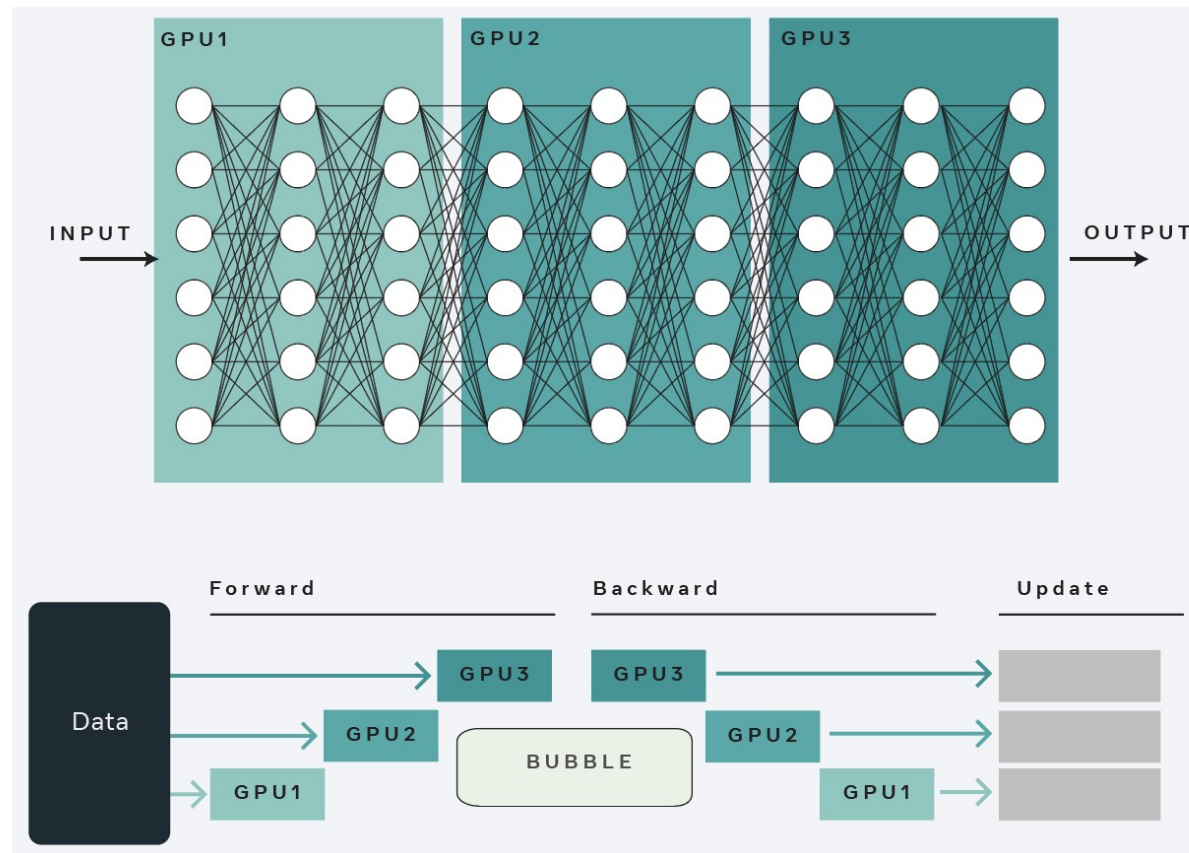
Multiple GPUs

GPU 0	a	x	y	=	ax	ay
	c				cx	cy
GPU 1	b	z	w	=	bz	bw
	d				dz	dw

# Distributed Training

## Pipeline Parallelism

- 한 단계의 출력이 다음 단계의 입력으로 이어지는 형태의 구조



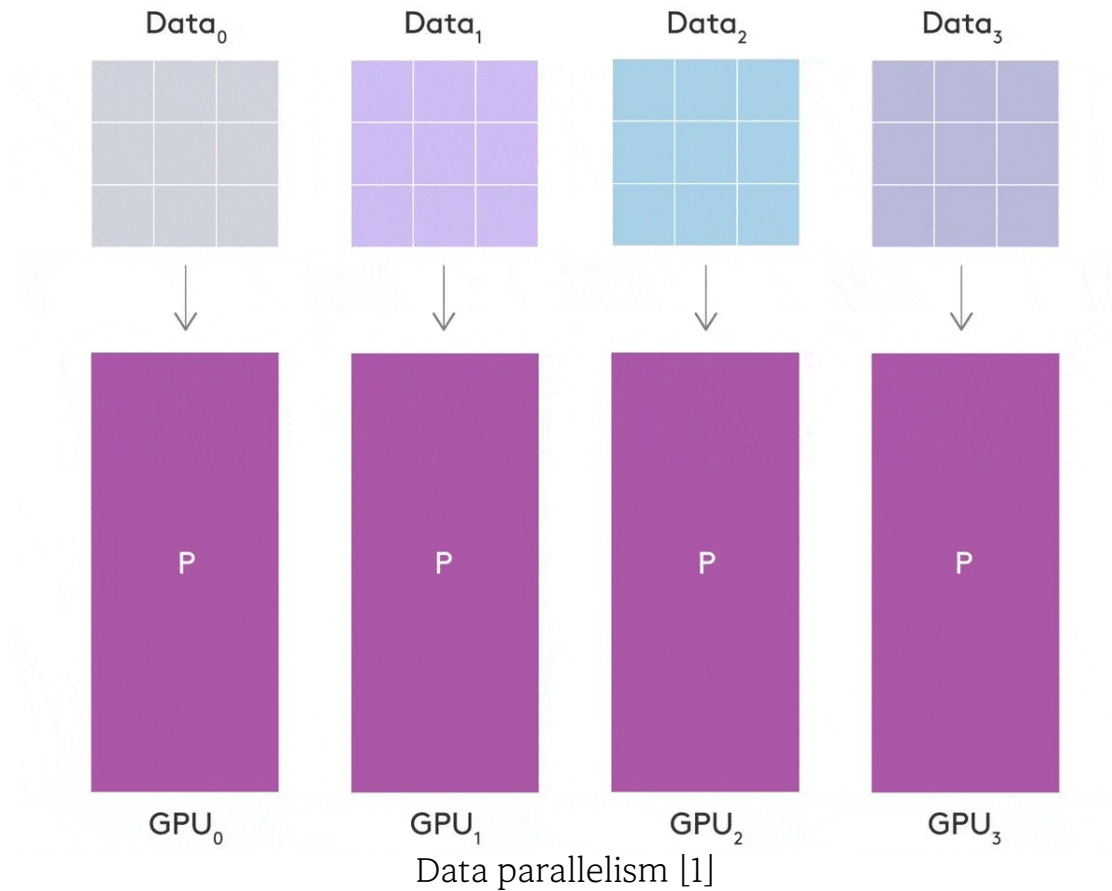
Pipeline parallelism [1]

[1] [https://fairscale.readthedocs.io/en/latest/deep\\_dive/pipeline\\_parallelism.html](https://fairscale.readthedocs.io/en/latest/deep_dive/pipeline_parallelism.html)

# Distributed Training

## Data Parallelism

- 각 GPU에 전체 모델이 올라가고 데이터만 분산됨



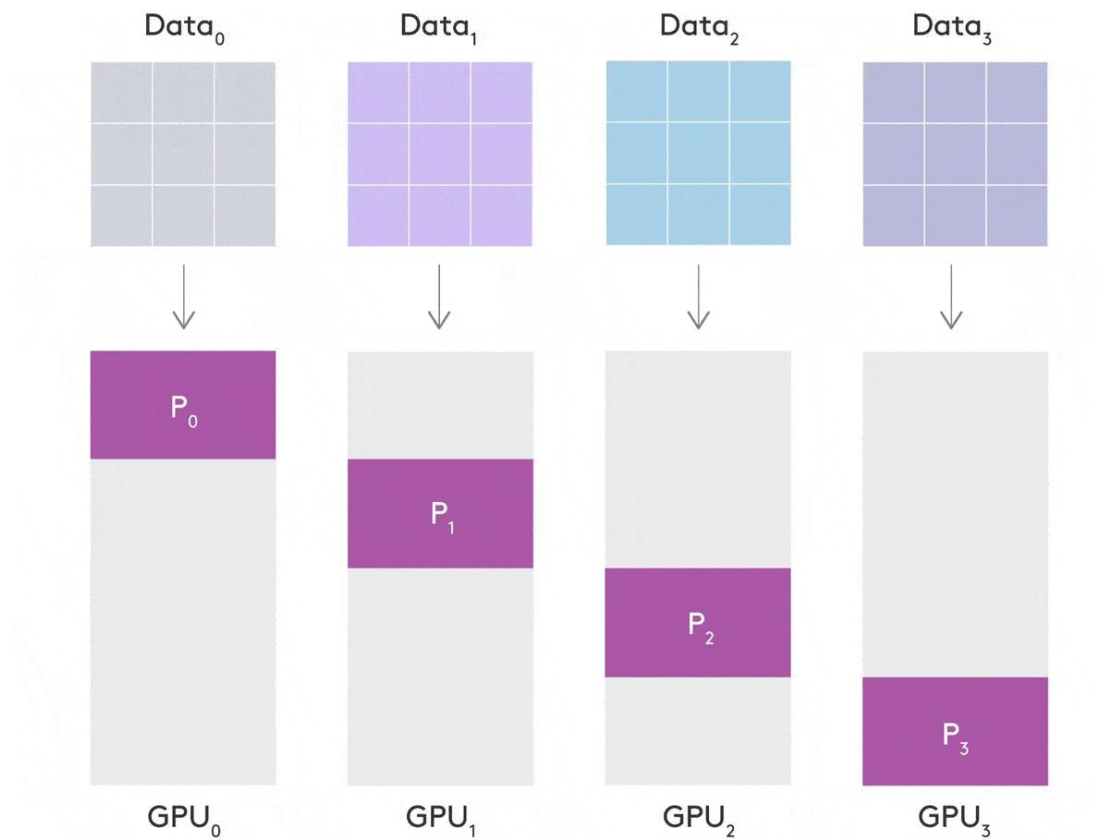
[1] <https://seannaren.medium.com/introducing-pytorch-lightning-sharded-train-sota-models-with-half-the-memory-7bcc8b4484f2>



# Distributed Training

## Sharded Data Parallelism

- 각 GPU에 모델과 데이터가 모두 분산됨



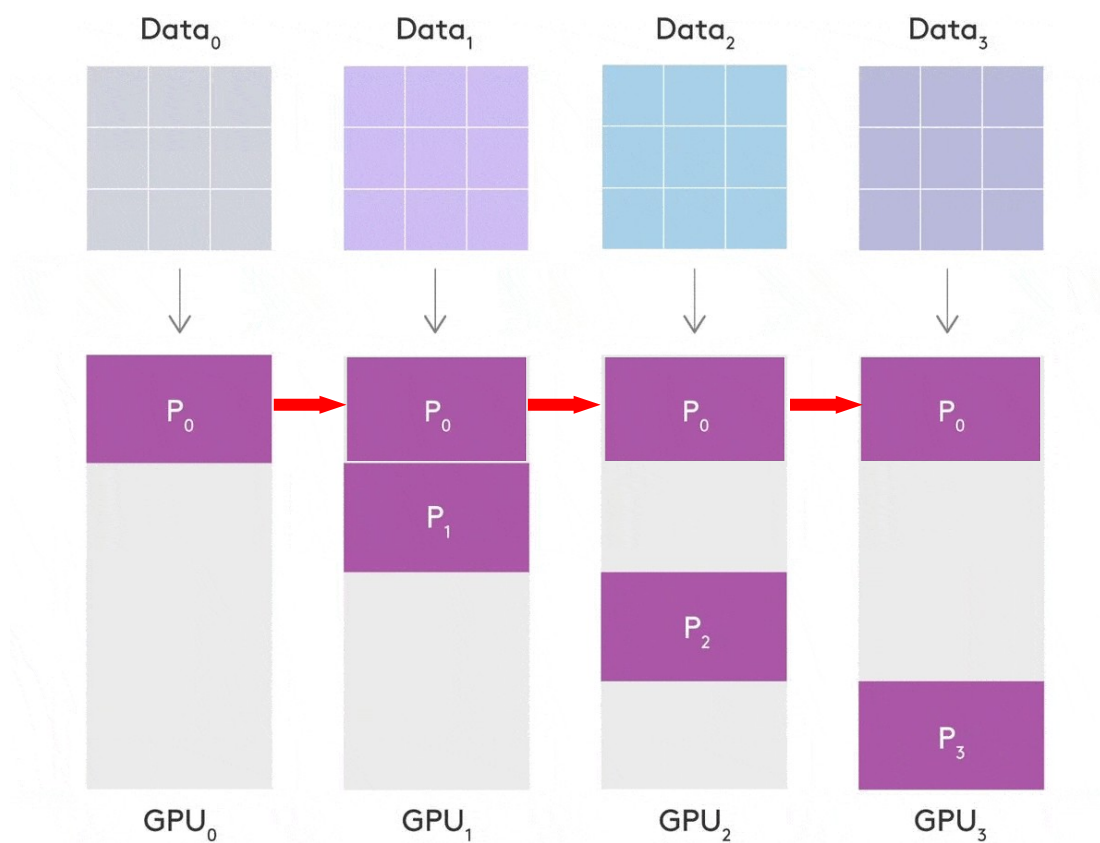
Sharded data parallelism [1]

[1] <https://seannaren.medium.com/introducing-pytorch-lightning-sharded-train-sota-models-with-half-the-memory-7bcc8b4484f2>

# Distributed Training

## Sharded Data Parallelism

- 각 GPU에 모델과 데이터가 모두 분산됨



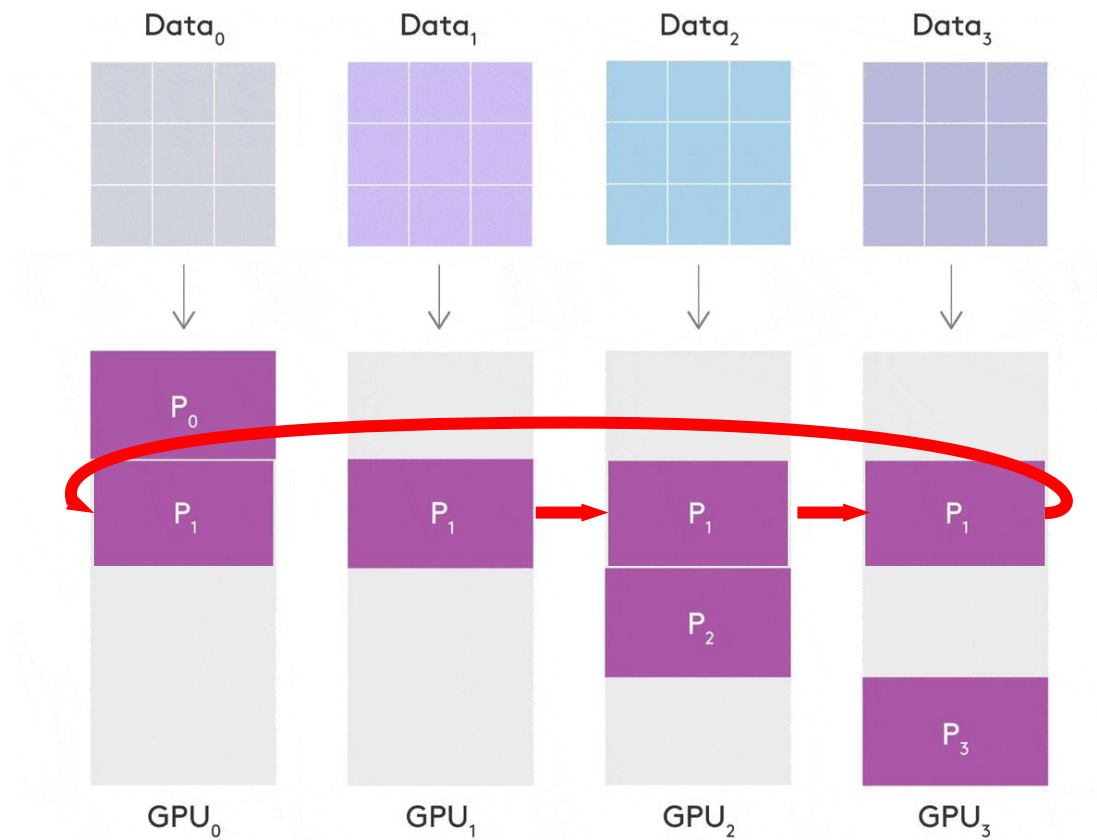
Sharded data parallelism [1]

[1] <https://seannaren.medium.com/introducing-pytorch-lightning-sharded-train-sota-models-with-half-the-memory-7bcc8b4484f2>

# Distributed Training

## Sharded Data Parallelism

- 각 GPU에 모델과 데이터가 모두 분산됨



Sharded data parallelism [1]

[1] <https://seannaren.medium.com/introducing-pytorch-lightning-sharded-train-sota-models-with-half-the-memory-7bcc8b4484f2>

# Distributed Training in Practice

## GPU Infrastructure

- On premise: including GPU, Network, and Storage
  - Meta: 350,000 NVIDIA H100 GPUs + others [1]
- Cloud [2]
  - AWS (Amazon Web Services)
  - Microsoft Azure
  - Google Cloud Platform (GCP)



 **Tesla H100 80GB NVIDIA Deep Learning GPU Compute Graphics Card**

















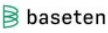

















Brand: Generic  
3.4 ★★★★★ 10 ratings | [Search this page](#)

Price: **\$43,989.00**

Eligible for Return, Refund or Replacement within 30 days of receipt

Graphics Coprocessor	NVIDIA Tesla H100
Brand	Generic
Graphics Processor	NVIDIA
Manufacturer	
Compatible Devices	Desktop

## CLOUD GPU PROVIDERS

	Serverless	VM	Bare Metal	Big Tech
H100		 	 	
A100	 	 		
		 		
V100		 	 	
RTX		 		
Non-NVIDIA		 		
			 	

The table shows vendor logos only once for simplicity. Each vendor offers several different models.



[1] <https://engineering.fb.com/2024/03/12/data-center-engineering/building-metas-genai-infrastructure/>

[2] <https://research.aimultiple.com/cloud-gpu-providers/>

# Distributed Training in Practice

## Distributed Training Library

- DeepSpeed (Microsoft) and PyTorch FSDP (Meta): sharded data parallelism
  - 기존 코드를 상대적으로 조금만 수정하여 분산학습 코드로 변환 가능
  - Hugging Face Transformers에서 학습코드 제공
- Megatron-LM (NVIDIA): data, tensor, pipeline, sequence parallelism
  - GPT3와 같은 주요 모델에 대한 구현체와 학습코드 제공
  - 구현되어 있지 않은 모델 학습 시 상대적으로 많은 코드 작업 필요



NVIDIA/Megatron-LM



Ongoing research training transformer models at scale

77 Contributors   1 Used by   8k Stars   2k Forks

# Conclusion

## Distributed Training

- Tensor Parallelism: 메모리↓, 시간 ↓
- Pipeline Parallelism: 메모리↓, 시간 ↓
- Data Parallelism: 시간↓
- Sharded Data Parallelism: 메모리↓ , 시간↓

## Code Examples

- DeepSpeed: <https://github.com/gmgu/LLMOps>
- Megatron-LM: <https://github.com/NVIDIA/Megatron-LM>
- Data Parallelism: [https://pytorch.org/tutorials/intermediate/ddp\\_tutorial.html](https://pytorch.org/tutorials/intermediate/ddp_tutorial.html)

## Q&A

- gmgu@theory.snu.ac.kr