



# Техники ускорения инференса и деплоя Large Language Models

30.08.2023



**Саввов Сергей**

За 7 лет в DS поработал в AdTech, FinTech, BioTech, trading и crypto-проектах.

Основал и руководил командой разработки торговых роботов с использованием ML.

Живу беспокойной жизнью номада и путешествую с 2-мя большими собаками.



[sergey-savvov](#)



[@slgero](#)

# Зачем нам это нужно?



**Jim Fan** ✓  
@DrJimFan



In the future, every 1% speedup on LLM inference will have similar economic value as 1% speedup on Google Search infrastructure.

# Зачем нам это нужно?



Jim Fan ✓  
@DrJimFan



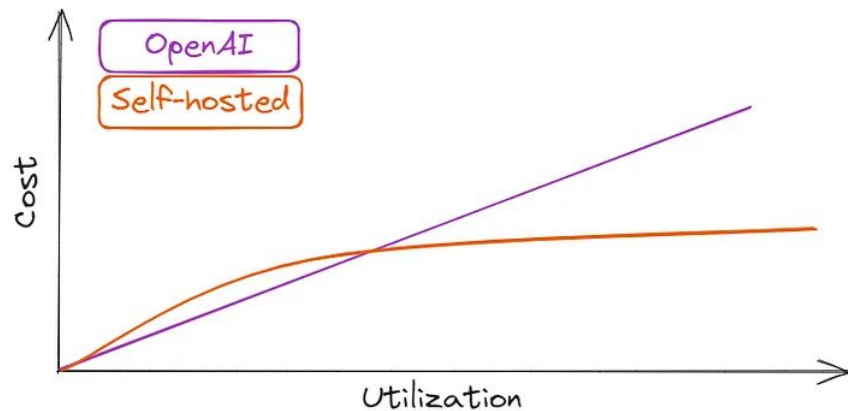
In the future, every 1% speedup on LLM inference will have similar economic value as 1% speedup on Google Search infrastructure.

[OpenAI тратят 700k\\$ ежедневно на хостинг ChatGPT:](#)



# Зачем нам ускорять?

- Ускорение работы основного приложения
- Увеличение пропускной способности
- **Уменьшение трат на хостинг**
- **Запуск моделей на простых девайсах**
- Инженерный азарт



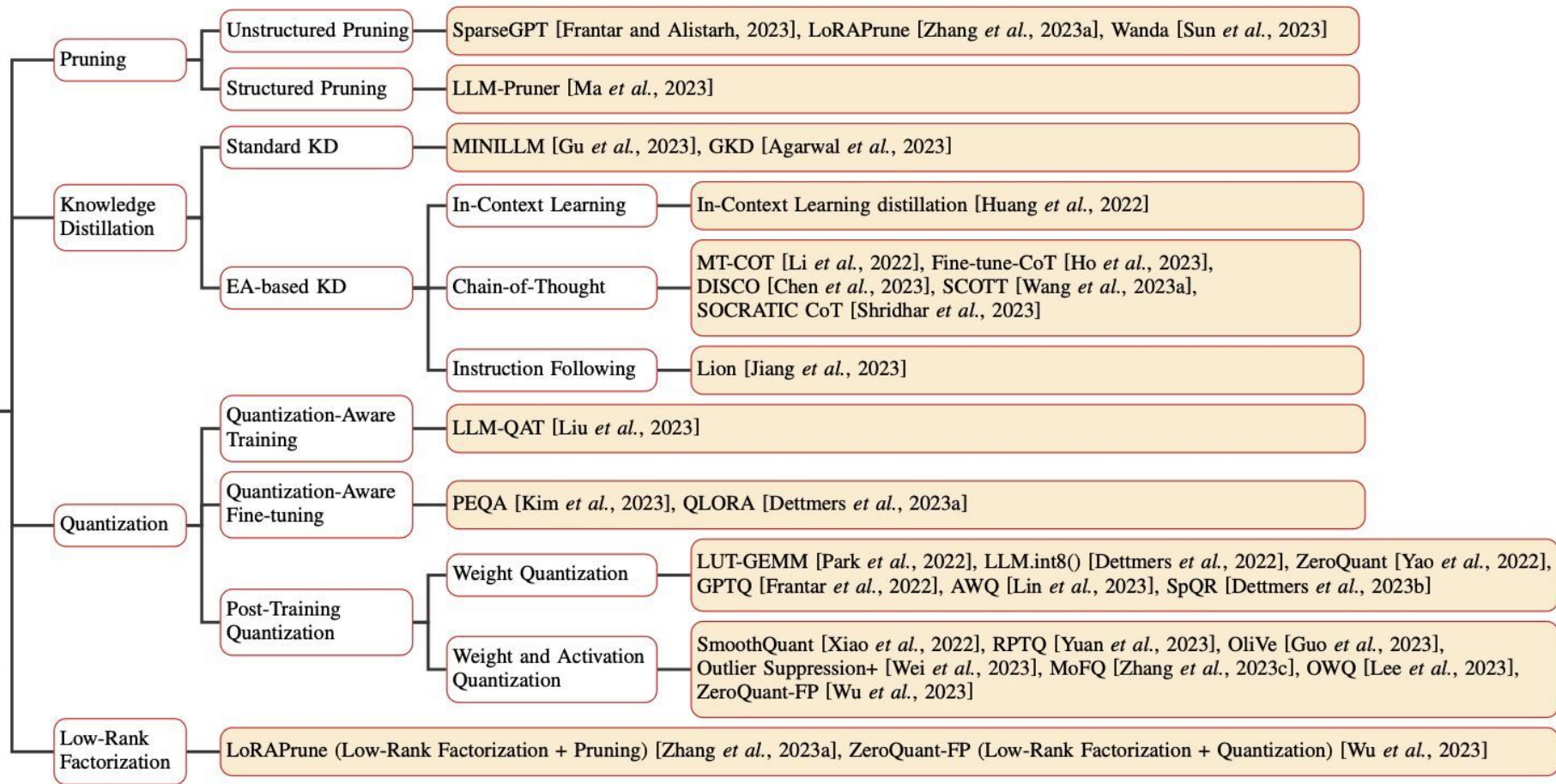
Схематичное сравнение стоимости OpenAI API и self-hosted LLM

# Что ждать от доклада?

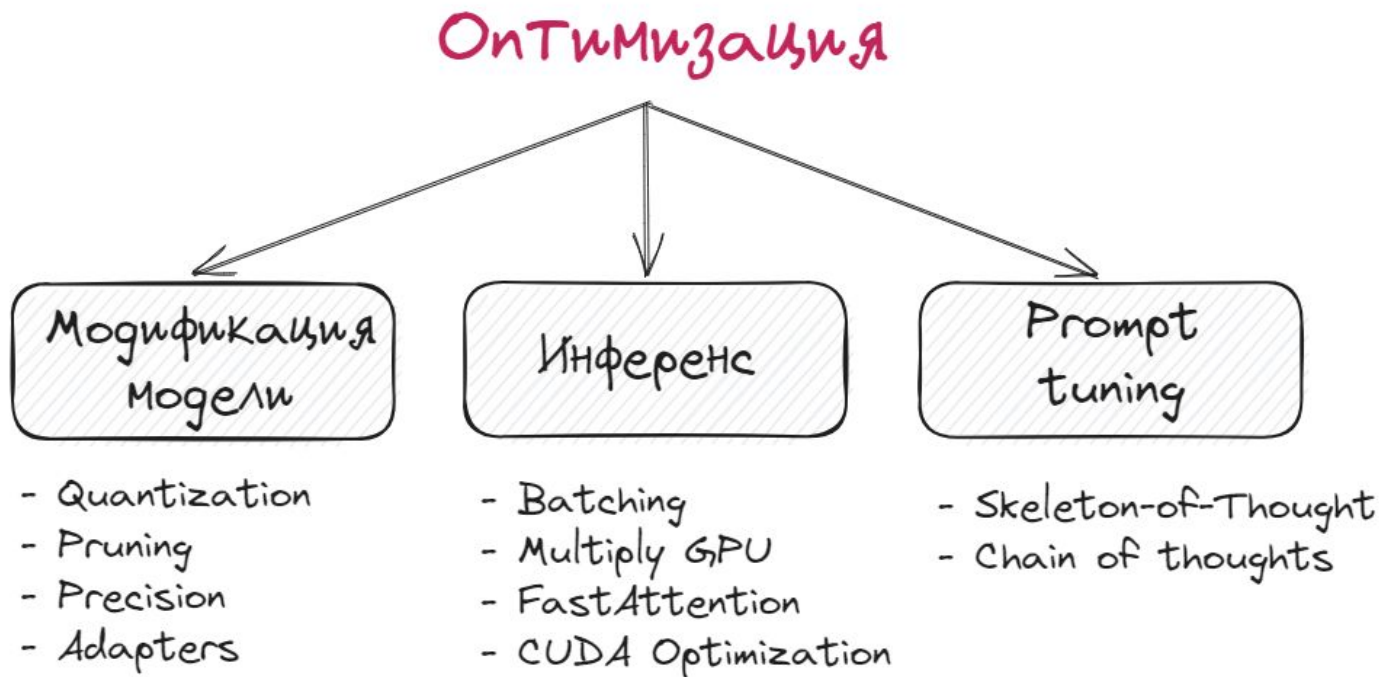


- Поговорим про разные техники ускорения и оптимизации: precision reduction, fine-tuning with adapters, continuous batching, ...
- Посмотрим на замеры скорости
- Узнаем, как имплементировать это в продакшене
- Пройдёмся по фреймворкам для деплоя

# Model Compression for Large Language Models



# Типы оптимизаций





# Подготовка к эксперименту

**Модель:** Falcon 7B

**GPU:** A100 40GB

**Библиотека для экспериментов:** Lit-GPT



- Можно проводить быстрые замеры
- Есть поддержка разных адаптеров
- Хорошо подходит для старта

Github: <https://github.com/Lightning-AI/lit-gpt>

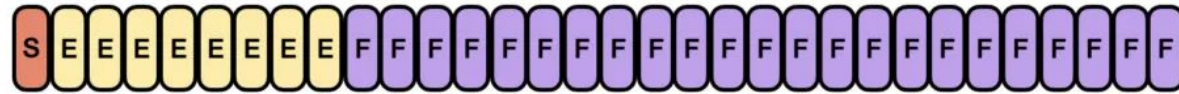
```
python generate.py
```

# Модификация модели

---

# Precision reduction

float 32



Sign  
(1 bit)

Exponent  
(8 bits)

Fraction  
(23 bits)

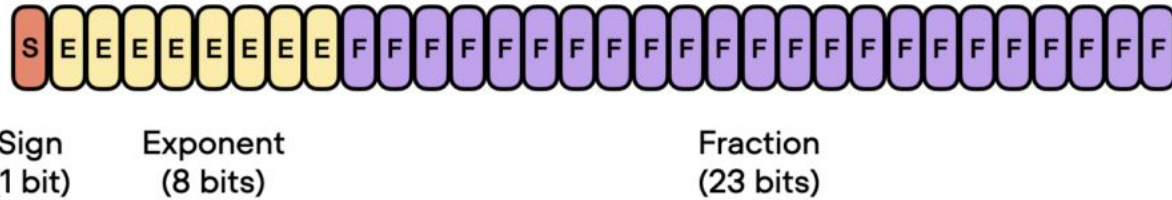
$$(-1)^{b_{31}} \times 2^{(b_{30}b_{29}\dots b_{23})_2 - 127} \times (1.b_{22}b_{21}\dots b_0)_2,$$

which yields

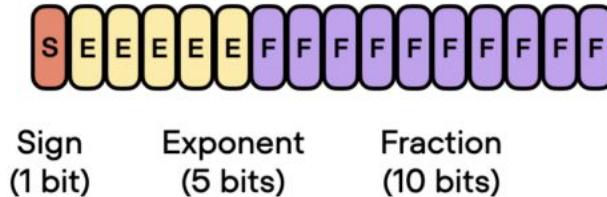
$$\text{value} = (-1)^{\text{sign}} \times 2^{(E-127)} \times \left( 1 + \sum_{i=1}^{23} b_{23-i} 2^{-i} \right).$$

# Precision reduction

float 32



float 16 ("half" precision)



# Precision reduction

float 16 ("half" precision)



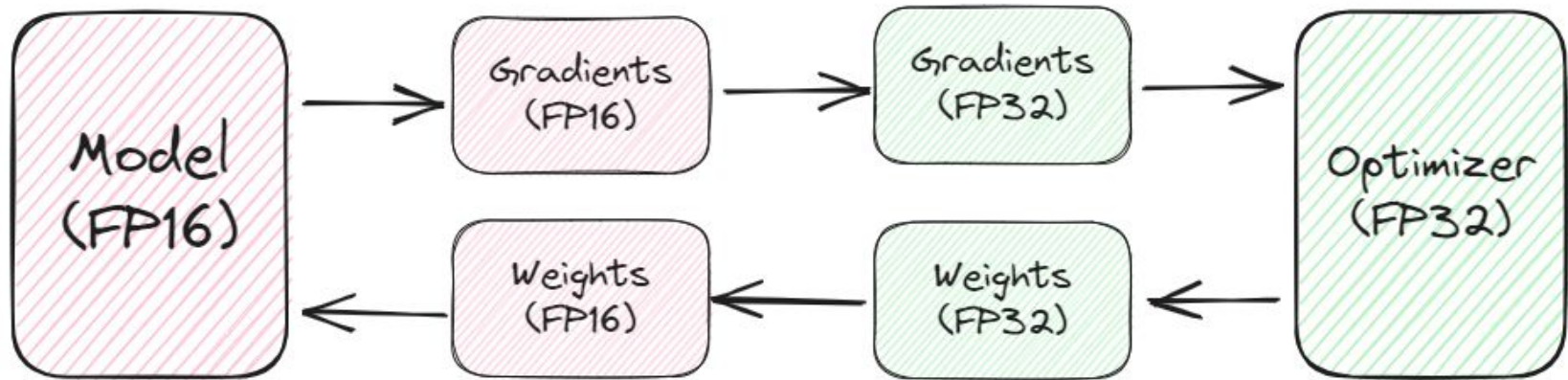
Sign  
(1 bit)

Exponent  
(5 bits)

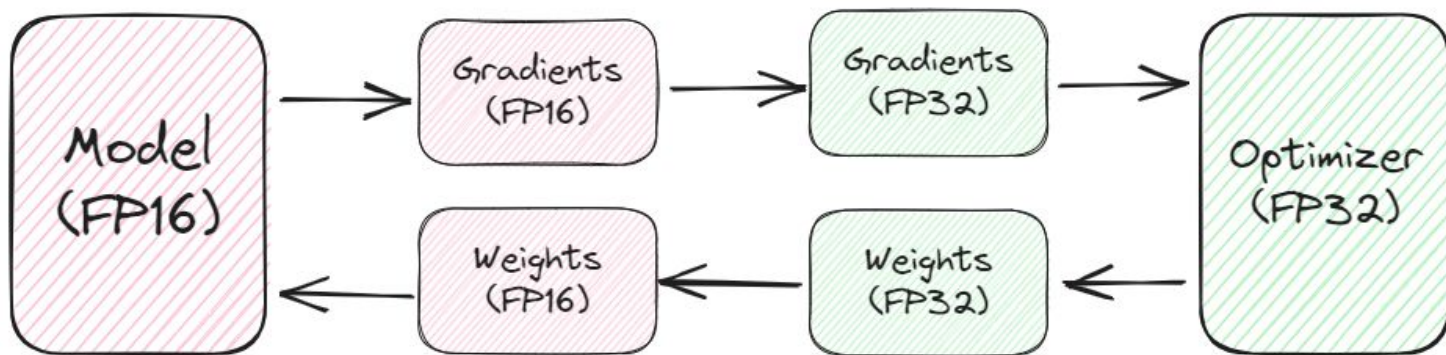
Fraction  
(10 bits)

```
python generate/base.py \  
  --prompt "I am so fast that I can" \  
  --checkpoint_dir checkpoints/tiiuae/falcon-7b \  
  --max_new_tokens 50 \  
  --precision "16-true"
```

# Mixed-precision



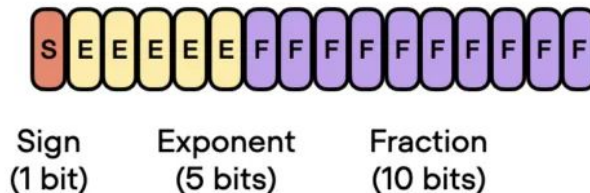
# Mixed-precision



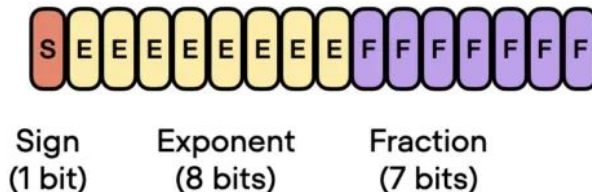
```
python generate/base.py \  
  --prompt "I am so fast that I can" \  
  --checkpoint_dir checkpoints/tiiuae/falcon-7b \  
  --max_new_tokens 50 \  
  --precision "16-mixed"
```

# Brain floating point (Bfloat 16)

float 16 ("half" precision)



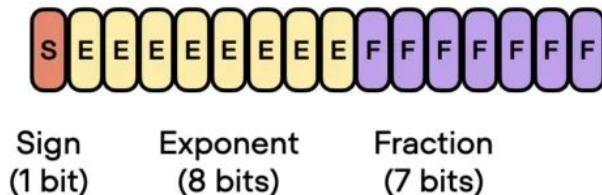
bfloat 16 ("brain" floating point, more "dynamic range" like float 32)





# Brain floating point (Bfloat 16)

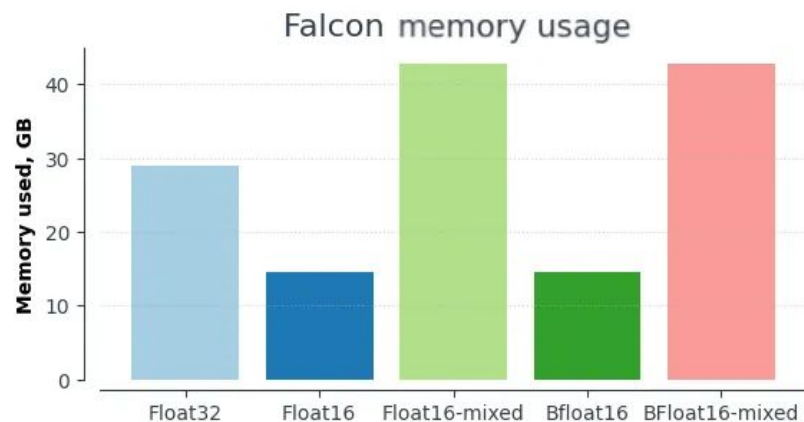
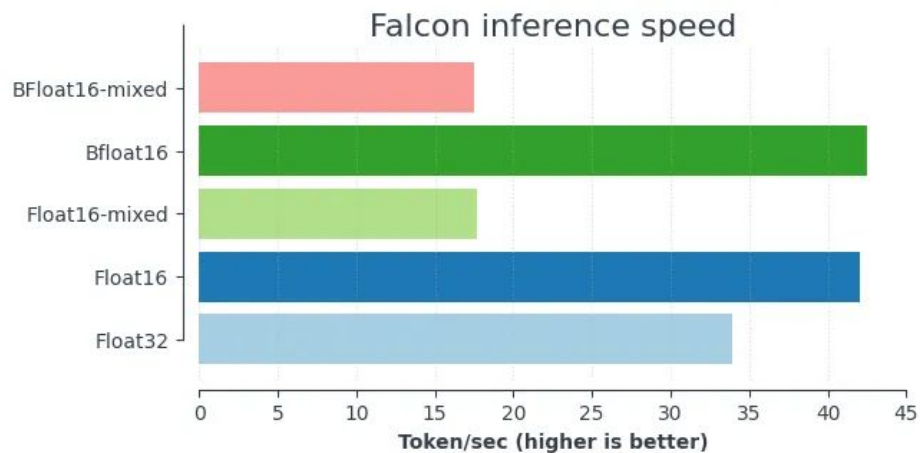
bfloat 16 ("brain" floating point, more "dynamic range" like float 32)



```
python generate/base.py \  
  --prompt "I am so fast that I can" \  
  --checkpoint_dir checkpoints/tiiuae/falcon-7b \  
  --max_new_tokens 50 \  
  --precision "bf16-true"
```

Хотя bfloat16 изначально был разработан для TPU, теперь этот формат поддерживается несколькими графическими процессорами NVIDIA.

# Precision reduction



bfloat16 всех побеждает 🎉

# Quantization

Floating point

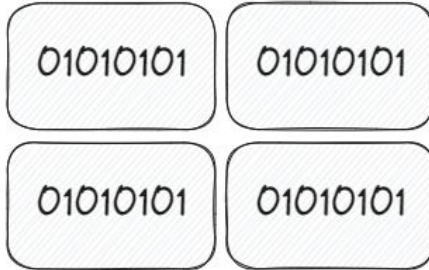
3452.3194



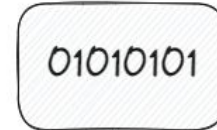
Integer

3452

32 bit



8 bit



# Quantization



## Quantization



```
graph TD; A[Quantization] --> B[Post-Training Quantization (PTQ)]; A --> C[Quantization-Aware Training (QAT)];
```

### Post-Training Quantization (PTQ)

Самый “дешёвый” способ

Квантование применяется к УЖЕ обученной модели

### Quantization-Aware Training (QAT)

Квантование применяется во время обучения, что заметно сложнее.

Качество модели обычно лучше

# Quantization

## Quantization

```
graph TD; A[Quantization] --> B[Post-Training Quantization (PTQ)]; A --> C[Quantization-Aware Training (QAT)];
```

### Post-Training Quantization (PTQ)

Самый “дешёвый” способ

Квантование применяется к УЖЕ обученной модели

```
python generate.py \  
  --prompt "I am so fast that I can" \  
  --quantize llm.int8
```

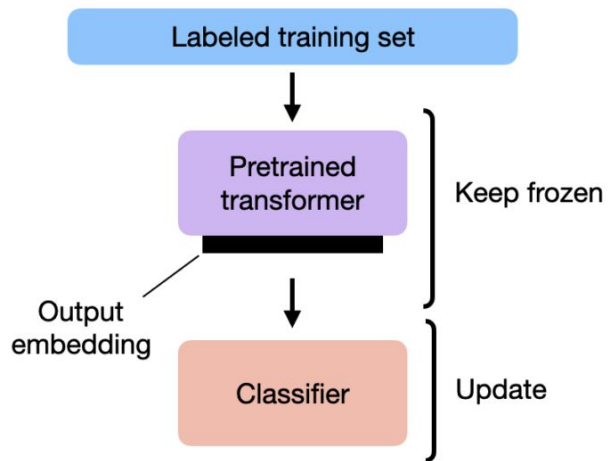
### Quantization-Aware Training (QAT)

Квантование применяется во время обучения, что заметно сложнее.

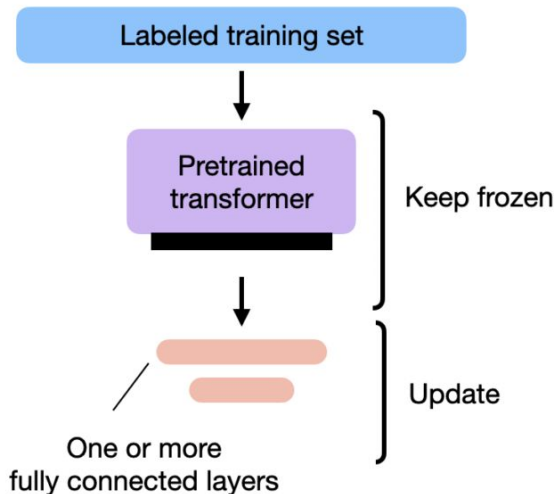
Качество модели обычно лучше

# Fine-Tuning Before Adapters

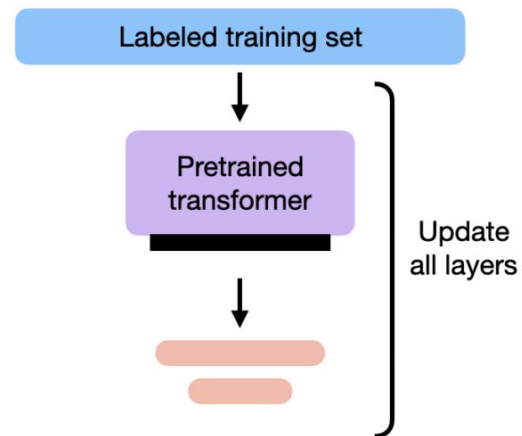
## 1) FEATURE-BASED APPROACH



## 2) FINETUNING I

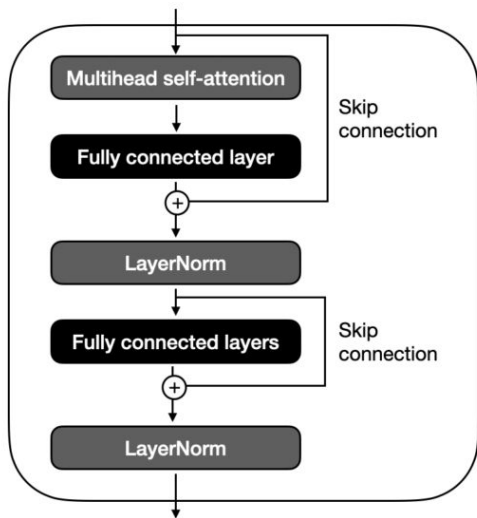


## 3) FINETUNING II

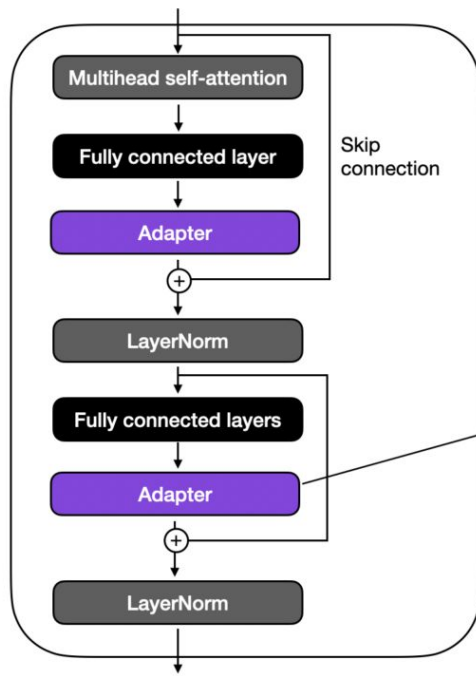


# Fine-Tuning With Adapters

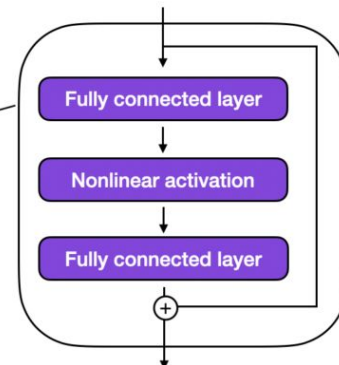
REGULAR TRANSFORMER BLOCK



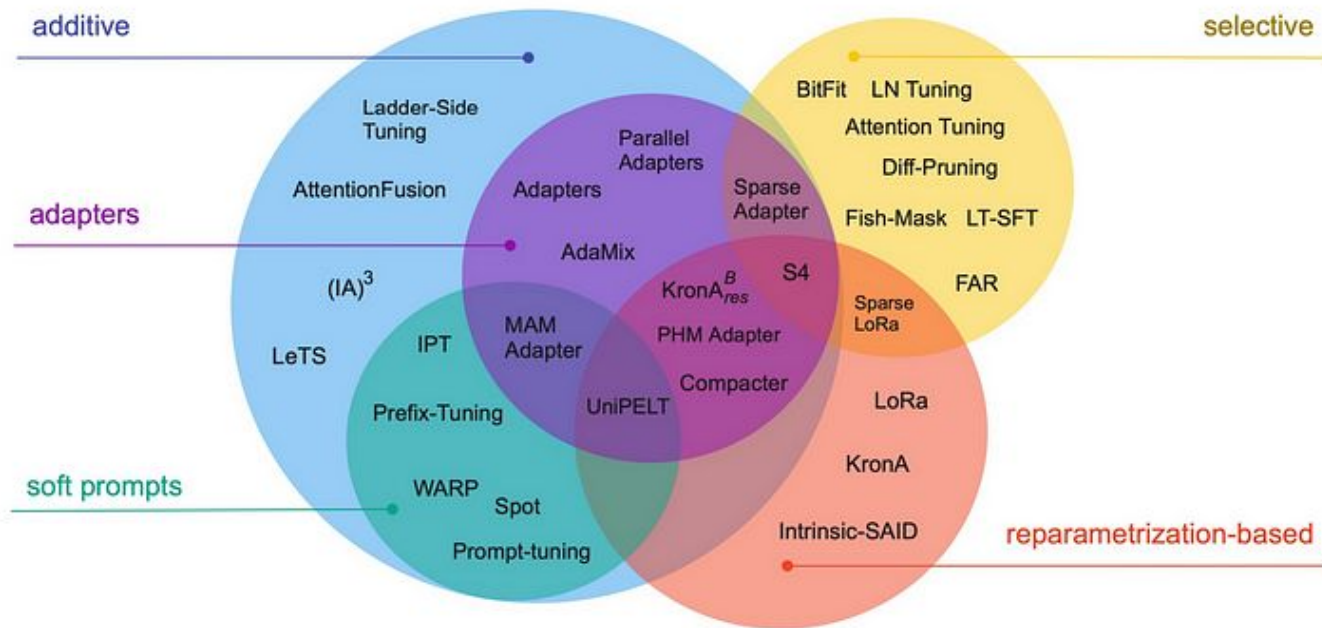
TRANSFORMER BLOCK WITH ADAPTERS



ADAPTER LAYERS

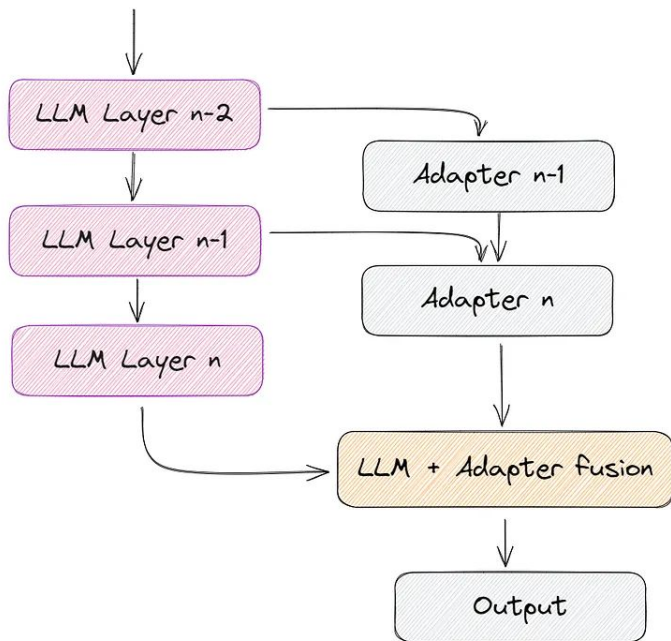


# Adapters



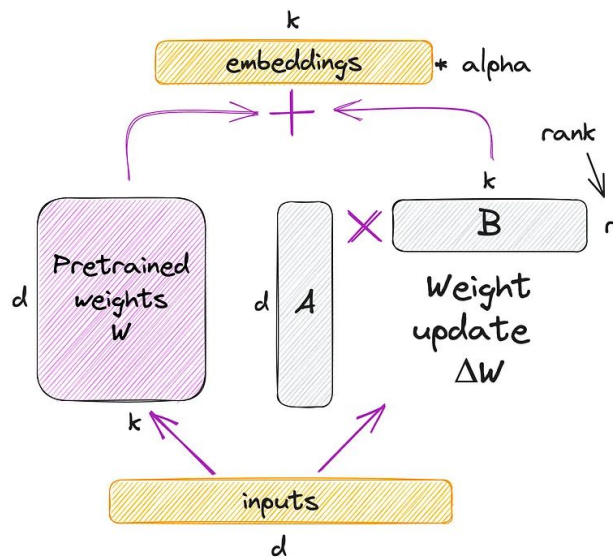


# LoRA



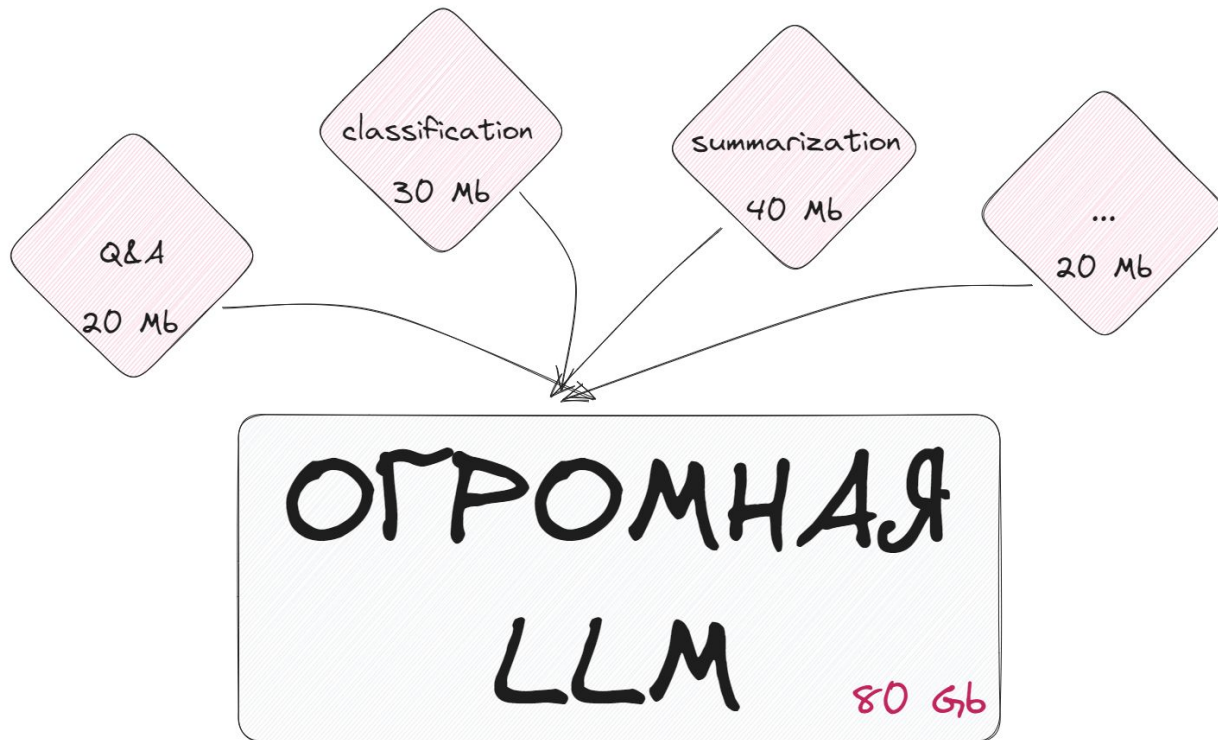
Architecture of adapter-based knowledge injection into LLMs

LoRA weights,  $A$  and  $B$  represent  $\Delta W$



Forward pass with updated model.

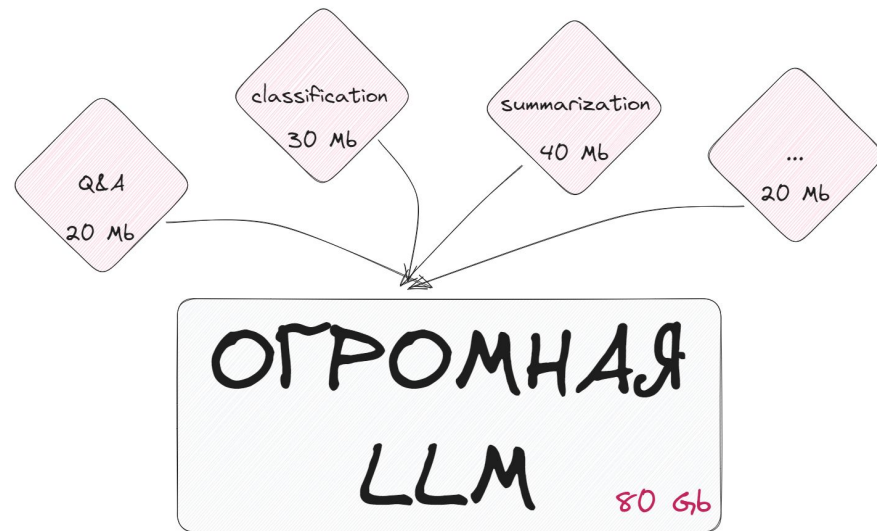
# Fine-Tuning With Adapters



# Fine-Tuning With Adapters

```
python finetune/adapters_v2.py \
  --data_dir data/alpaca \
  --checkpoint_dir checkpoints/tiiuae/falcon-7b \
  --out_dir out/adapters/alpaca
```

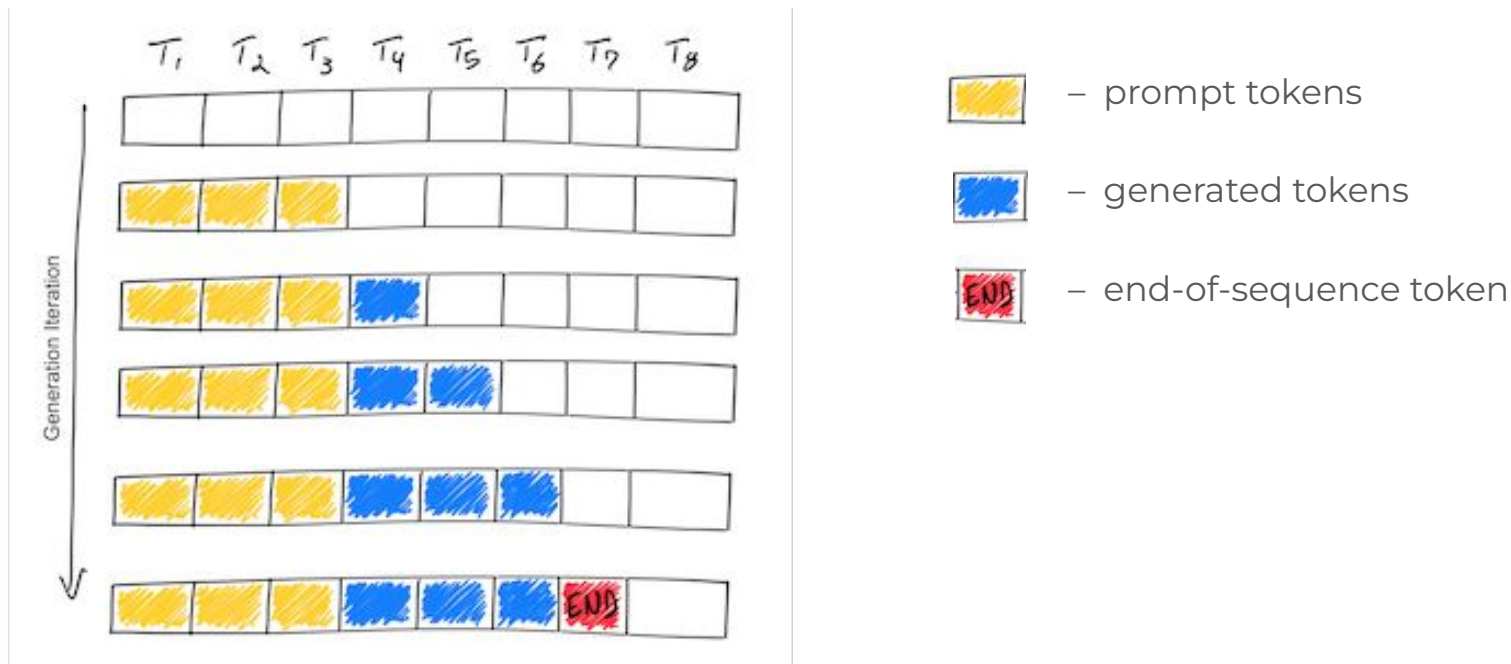
У меня заняло **~10 часов** для 52K инструкций



# Inference

---

# Batch Inference



# Batch Inference

$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$				
$S_2$	$S_2$	$S_2$					
$S_3$	$S_3$	$S_3$	$S_3$				
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$			

$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$	$S_1$	END		
$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	END
$S_3$	$S_3$	$S_3$	$S_3$	END			
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	END	

Naive batching



# Batch Inference

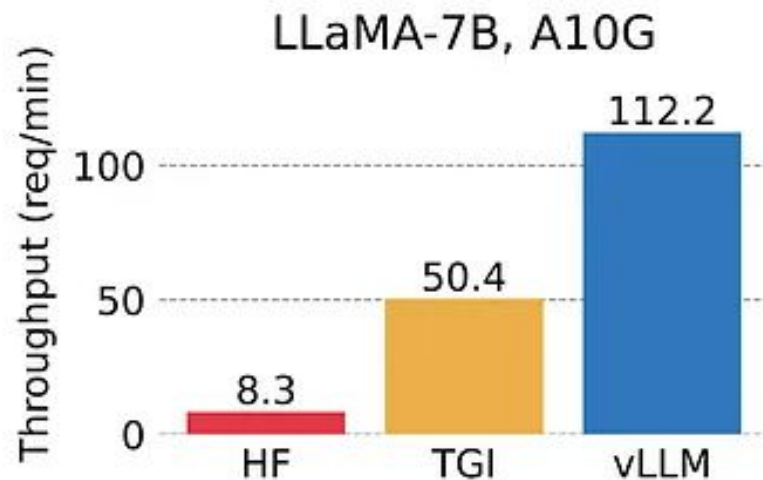
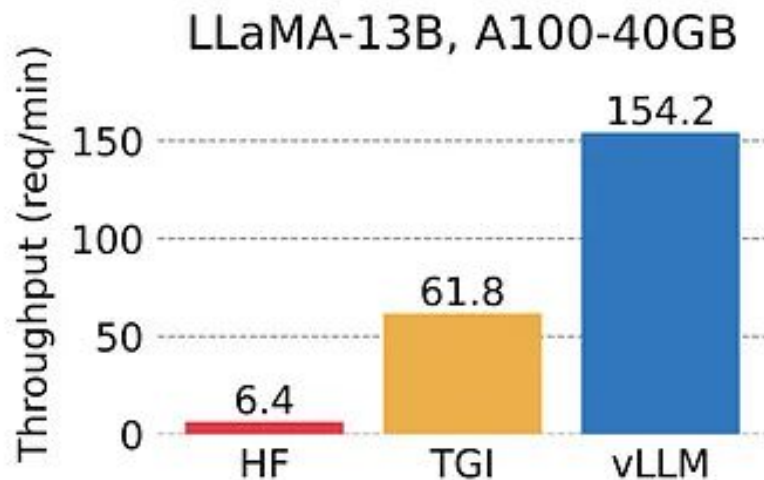
$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$	$S_1$	END		
$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	END
$S_3$	$S_3$	$S_3$	$S_3$	END			
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	END	

Naive batching

$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$	$S_1$	END	$S_6$	$S_6$
$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	END
$S_3$	$S_3$	$S_3$	$S_3$	END	$S_5$	$S_5$	$S_5$
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	END	$S_7$

Continuous batching

# Batch Inference





# Batch Inference

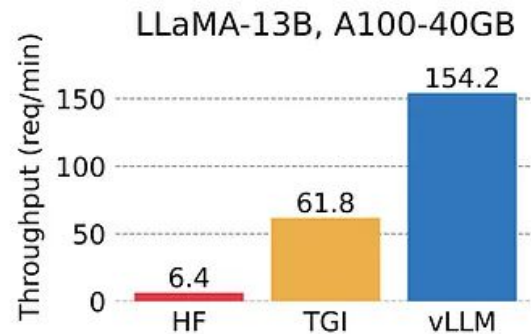
```
from vllm import LLM, SamplingParams

llm = LLM(model="huggyllama/falcon-13b")
outputs = llm.generate(prompts, SamplingParams(temperature=0.8, top_p=0.95))

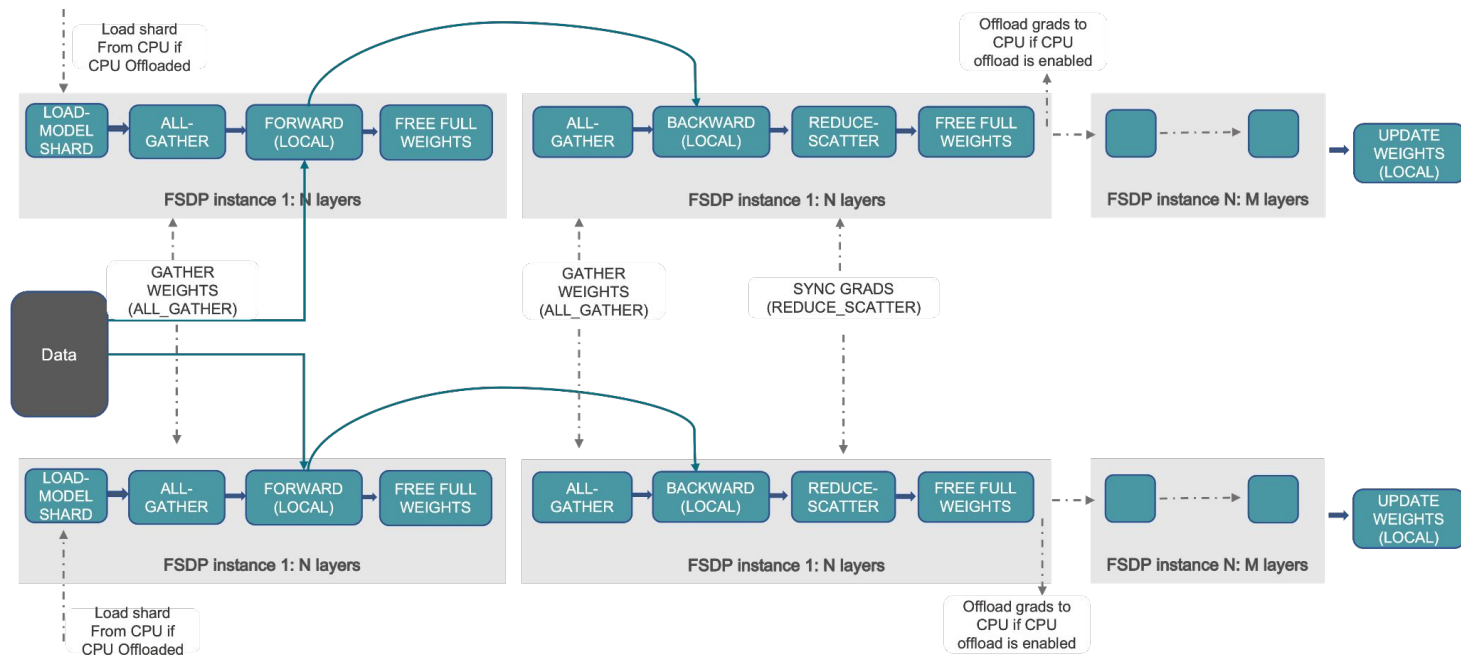
llm.generate(prompts, sampling_params)
```

**vLLM** - инференс сервер, который поддерживает разные трюки для ускорения LLMs

**Github:** <https://github.com/vllm-project/vllm>



# Multiple GPU devices



Fully Sharded Data Parallel ([FSDP](#)) workflow

# Multiple GPU devices

```
python generate/base.py \  
  --checkpoint_dir checkpoints/tiiuae/falcon-40b \  
  --strategy fsdp \  
  --devices 2 \  
  --prompt "I am so fast that I can"
```

Falcon-7B

	1x A6000 (48 GB)	2x A6000 (48 GB)
Time for inference	1.65 sec	18.00 sec
Memory used	14.50 GB	9.60 GB

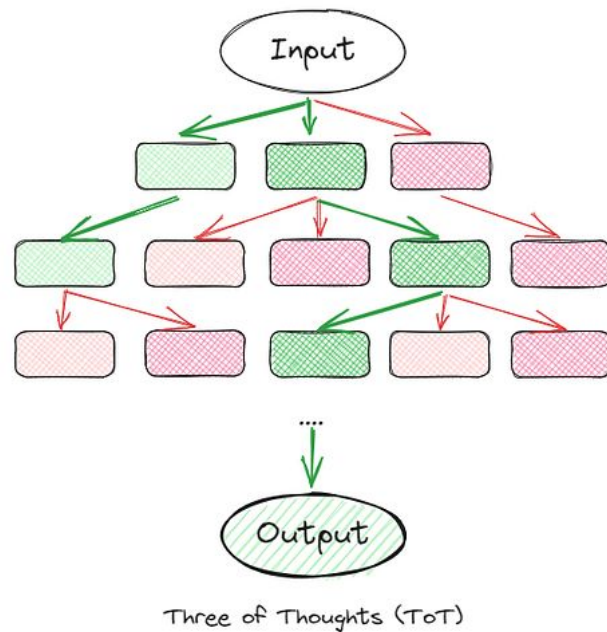
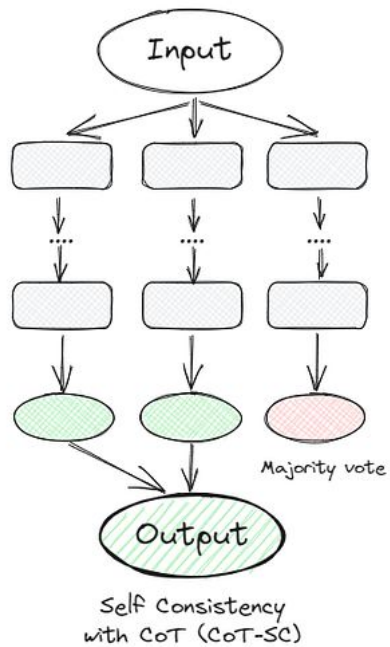
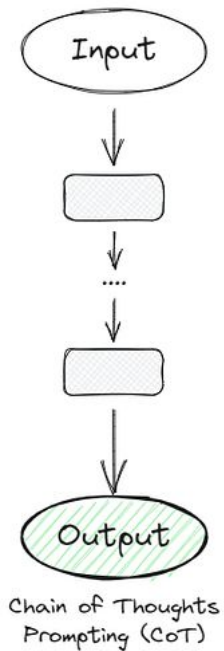
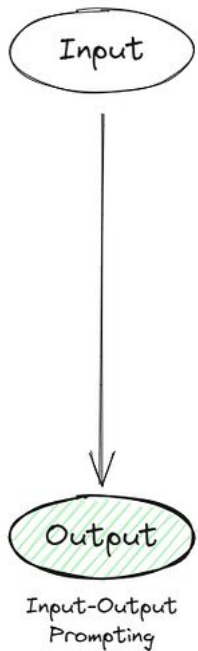
Falcon-40B

	1x A6000 (48 GB)	2x A6000 (48 GB)
Time for inference	- sec	83.40 sec
Memory used	- GB	46.10 GB

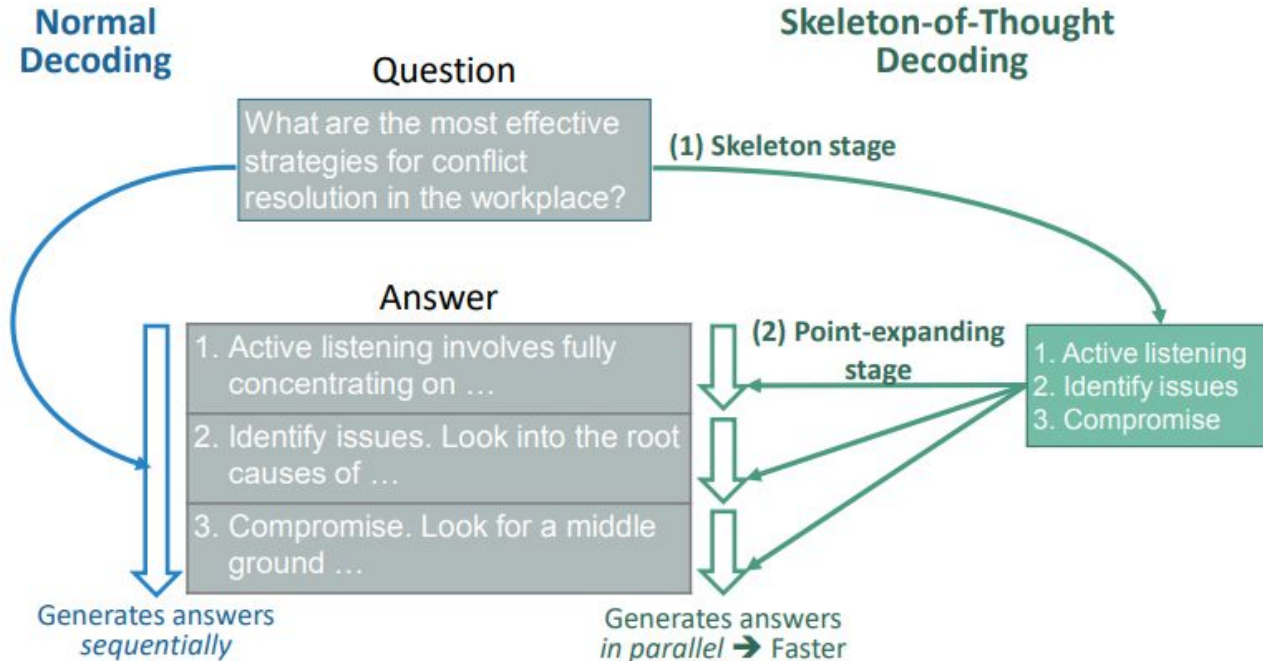
# Prompt tuning

---

# Chain of Thought



# Skeleton-of-Thought



# Замеры скорости и памяти

---

# Falcon 7B & 1 A100 GPU

	Time for inference, tokens/sec	Memory used, GB
Vanilla Usage	33.9	28.9
Mixed-Precision	17.7	42.8
Bfloat16	42.4	42.8
Quantization	24.8	13.5
Continuous Batching	74.2	—
Multiple GPU	2.7	9.6



# Фреймворки для деплоя

Framework	Tokens Per Second	Query per second	Latency	Adapters	Quantisation	Variable precision	Batching	Distributed Inference	Custom models	Token streaming	Prometheus metrics
vLLM	115	0.94	4.8 s	✗	✗	✗	✓	✓	✓	✓	✗
Text generation inference	50	0.26	4.8 s	✗	✓	✓	✓	✓	✓	✓	✓
CTranslate2	93	0.55	4.5 s	✗	✓	✓	✓	✓	✓	✓	✗
DeepSpeed-MII	80	0.47	2.5 s	✗	✗	✓	✓	✓	✗	✗	✗
OpenLLM	30	0.15	6.6 s	✓	✓	✓	✗	✗	✓	✗	✓
Ray Serve	28	0.15	7.1 s	✗	✓	✓	✓	✓	✓	✗	✓
MLC LLM	25	0.13	7.2 s	✗	✓	✗	✗	✗	✓	✓	✗

[Comparison of frameworks for LLMs inference](#)

# Takeaways



1. Основные методы оптимизации такие же как и для DL моделей
2. Используйте одну модель и несколько адаптеров, если вам нужно решать сразу несколько доменных задач
3. Используйте готовые фреймворки для сёрвинга моделей – у них много фичей доступно сразу “из коробки”
4. По возможности используйте bfloat16 (нужен определённый тип GPU)
5. Изучите техники ускорения через разные стратегии промпт инжиниринга

# Что читать дальше?



- [Optimizing Memory Usage for Training LLMs and Vision Transformers in PyTorch](#)
- [Accelerating Large Language Models with Mixed-Precision Techniques](#)
- [Understanding Parameter-Efficient Finetuning of Large Language Models: From Prefix Tuning to LLaMA-Adapters](#)
- [The Secret Sauce behind 100K context window in LLMs: all tricks in one place](#)
- [7 Frameworks for Serving LLMs](#)
- [7 Ways to Speed Up Inference of Your Hosted LLMs](#)

**Спасибо за внимание!**

