

# Обучение GigaChat с контекстом В СОТНИ ТЫСЯЧ ТОКЕНОВ

Косарев Евгений  
Senior DL Engineer, Sber Devices



# GigaChat в бизнесе

- В десятках продуктов Сбера

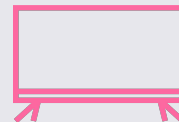
- Работа с текстом в  МойОфис

- Департамент здравоохранения г.Москвы

- Более 2 тысяча клиентов GigaChat API



Jazz





Я готовлю доклад про длинный контекст в языковых моделях на конференцию. Кратко расскажи, что это такое, и на чем стоит акцентировать внимание в докладе?



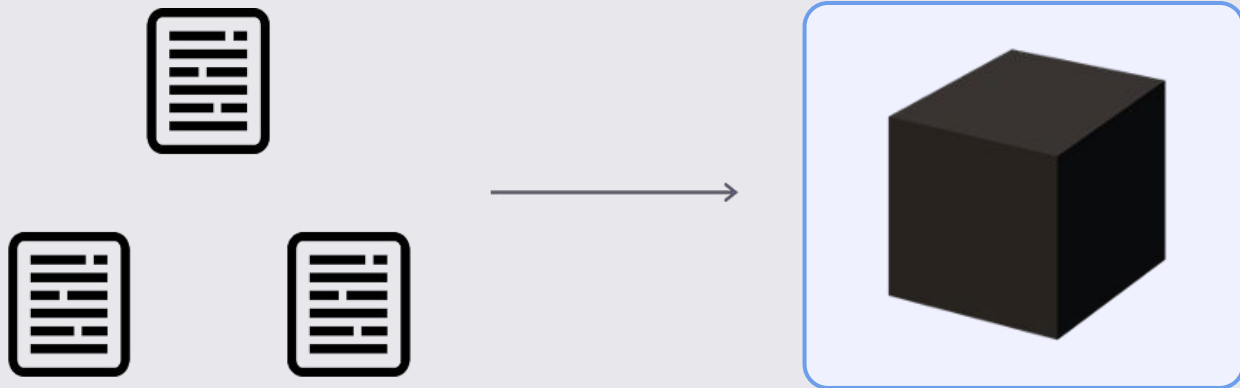
Длинный контекст (long context) — это подход к обучению языковых моделей, при котором они обучаются на больших текстах, включающих в себя много информации о мире. Длинный контекст позволяет модели «понять» большие тексты и использовать эту информацию для генерации более качественных ответов.

При подготовке доклада об этом подходе стоит обратить внимание на следующие аспекты:

- 1. Преимущества использования длинного контекста.** Модели, обученные на больших текстах, могут лучше понимать мир и выдавать более точные ответы. Они также могут использоваться для решения сложных задач, таких как анализ документов или создание художественных произведений.
- 2. Технологии.** Можно рассказать о способах и техниках обучения моделей длинному контексту.
- 3. Перспективы развития.** Можно обсудить возможные направления развития этого подхода, такие как использование нейронных сетей нового поколения или обучение моделей на больших объемах данных.

# Что такое языковые модели

Языковые модели – алгоритмы, способные  
продолжать тексты



GigaChat это ...



известная языковая  
модель

# LLaMa-3

Бенчмарки – как показатель качества языковой модели

**2024**

**2023**

**2024**

**2023**

Benchmark	Llama 3 8B	Llama 2 7B	Llama 2 13B	Llama 3 70B	Llama 2 70B
MMLU (5-shot)	68.4	34.1	47.8	82.0	52.9
GPQA (0-shot)	34.2	21.7	22.3	39.5	21.0
HumanEval (0-shot)	62.2	7.9	14.0	81.7	25.6
GSM-8K (8-shot, CoT)	79.6	25.7	77.4	93.0	57.5
MATH (4-shot, CoT)	30.0	3.8	6.7	50.4	11.6

# LLaMa-3

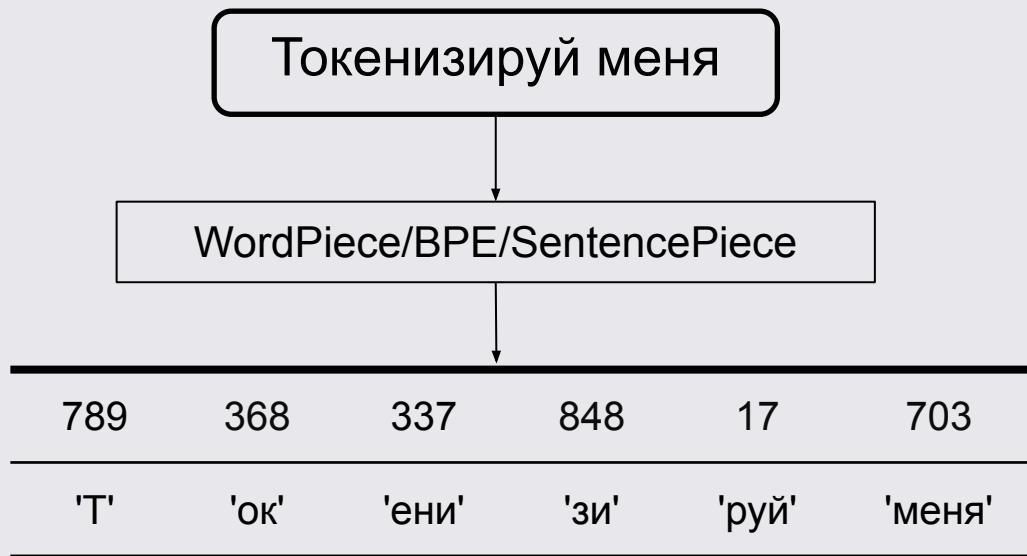
Бенчмарки – как показатель качества языковой модели

Приближается к  
GPT-4

	2024	2023		2024	2023
Benchmark	Llama 3 8B	Llama 2 7B	Llama 2 13B	Llama 3 70B	Llama 2 70B
MMLU (5-shot)	68.4	34.1	47.8	82.0	52.9
GPQA (0-shot)	34.2	21.7	22.3	39.5	21.0
HumanEval (0-shot)	62.2	7.9	14.0	81.7	25.6
GSM-8K (8-shot, CoT)	79.6	25.7	77.4	93.0	57.5
MATH (4-shot, CoT)	30.0	3.8	6.7	50.4	11.6

# Контекст в LLM

Токен - единица текста, которую “понимают” языковые модели



# Чего-то не хватает...

Большинство Open-Source моделей имеет контекст

8192 токена

~4 страницы текста А4

Нужен ли контекст  
больше?

# Поиск по документам

Пользователь задает вопрос к документу **>20** страниц

A4. Как на него ответить?



Столько можем обработать за раз

# Решения

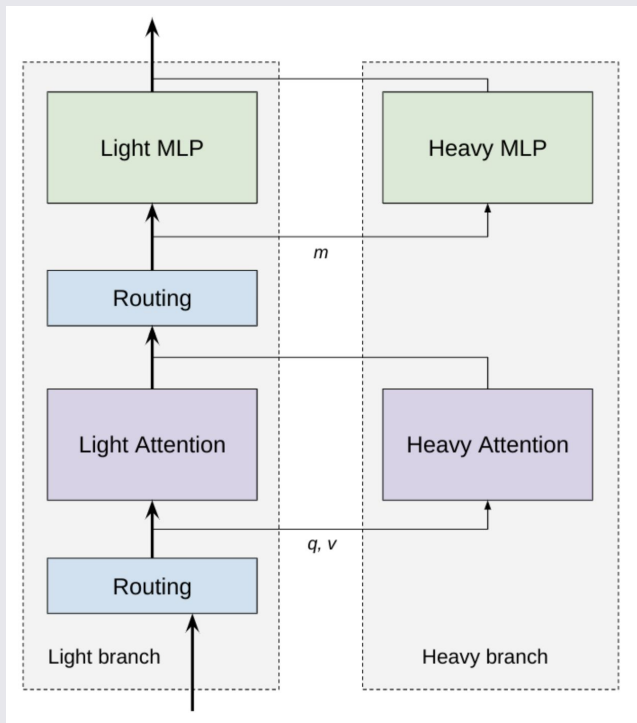
## Модификации Attention (CoLT5)

### Преимущества:

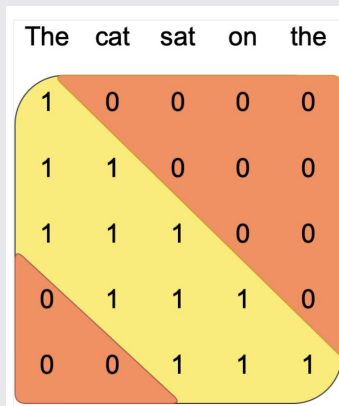
- Видит весь контекст

### Недостатки:

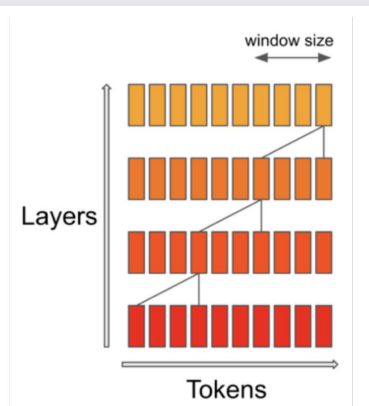
- Роутер совершает ошибки
- Light ветка плохо собирает КОНТЕКСТ



# Решения



Sliding Window Attention



Effective Context Length

# Модификации Attention (SWA)

## Преимущества:

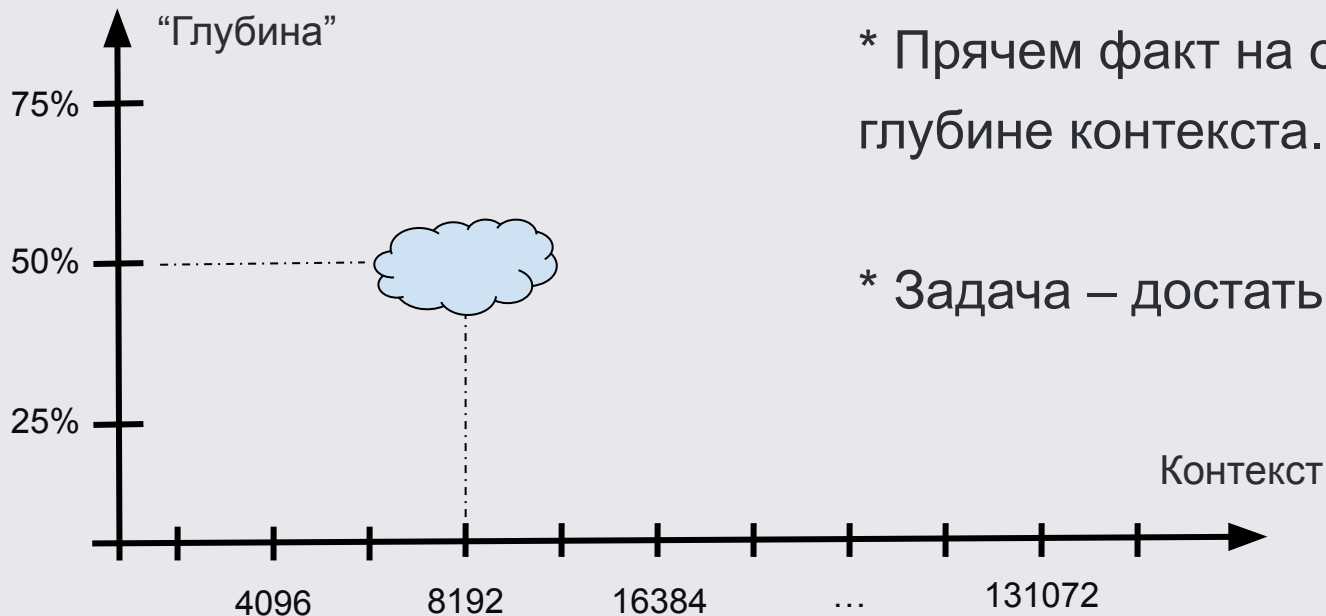
- Полное внимание в окне
- Агрегация информации с нижний слоев

## Недостатки:

- Потеря информации через слои

# PassKey

Универсальный бенчмарк по проверке длинного контекста



\* Прячем факт на определенной глубине контекста.

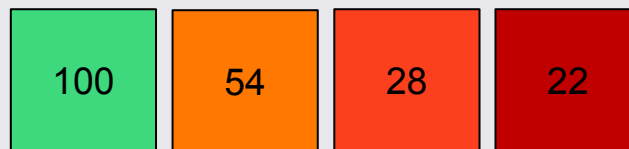
\* Задача – достать факт

# Mistral 7B

Результаты PassKey

Mistral 7B v.01

**SWA 32k**

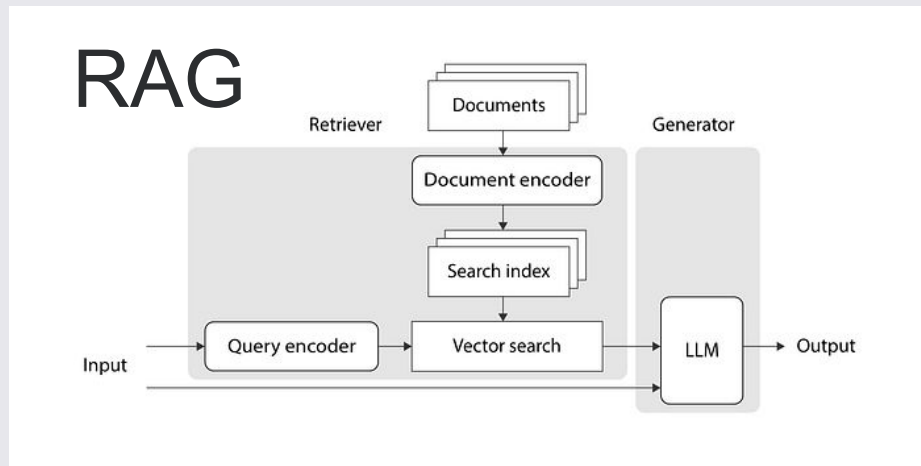


Mistral 7B v.02

32k



# Решения



## Преимущества:

- поиск намного дешевле запроса в GPT
- “Бесконечный” контекст

## Недостатки:

- Качество поисковика
- Ограниченность информации

# Решения

**Честное** увеличение контекстного окна модели

8192 → 32768 → 131072

Проблемы:

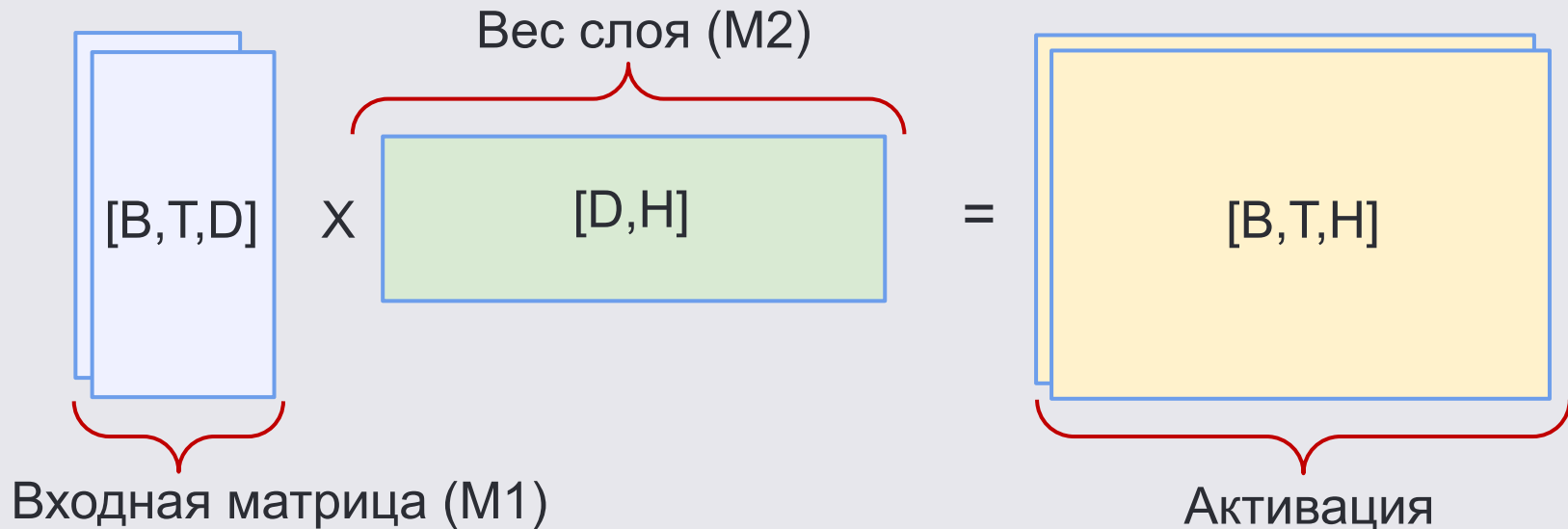
1. Активации тензоров не помещаются на GPU
2. Замедление вычислений

# Оптимизация вычислений

# Цель оптимизации

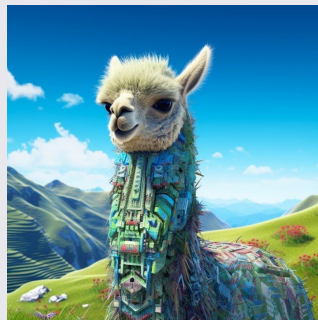
$T$  – число токенов в контексте  
 $D, H$  – внутренние размерности  
 $B$  – размер батча

Активация - результат умножения входной матрицы на веса слоя



# Архитектура LLM

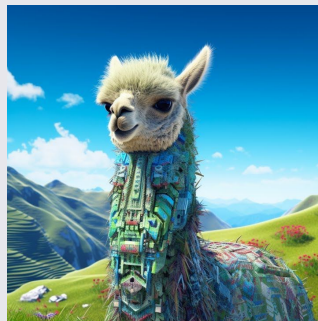
LLaMa-1,2,3



Состоит из Decoder  
блоков

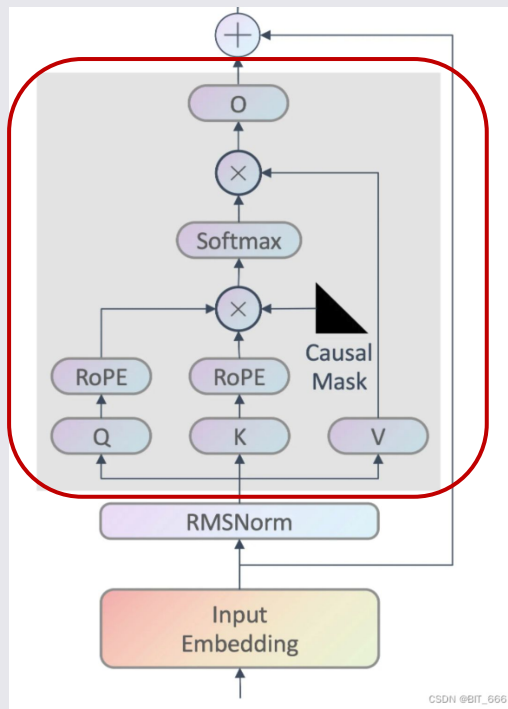
# Архитектура LLM

LLaMa-1,2,3



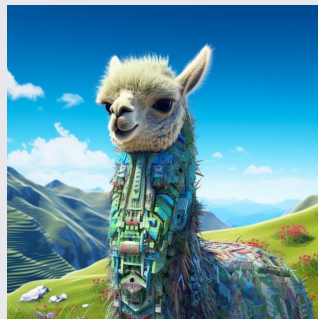
Состоит из Decoder  
блоков

## Attention



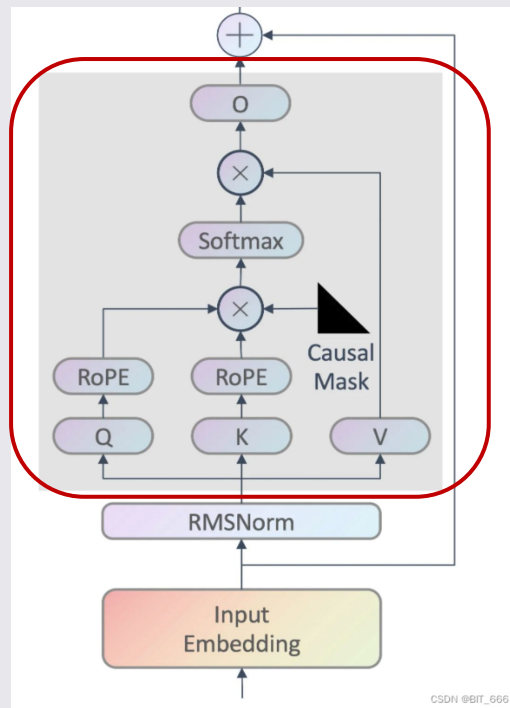
# Архитектура LLM

LLaMa-1,2,3

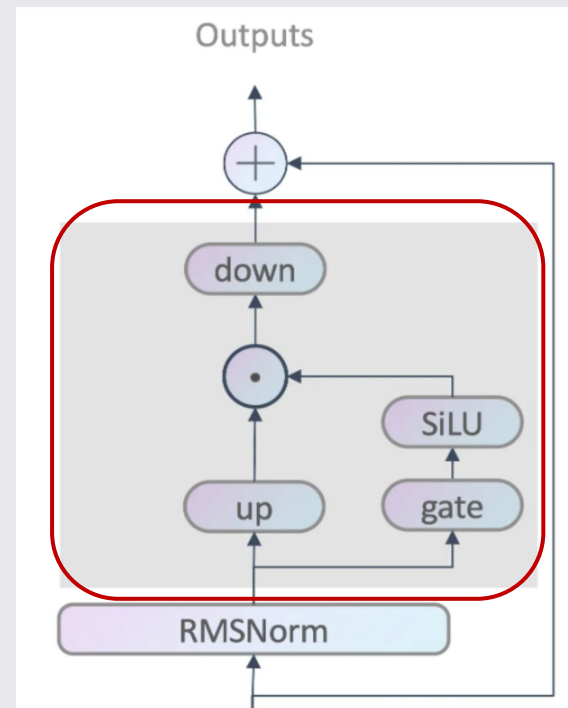


Состоит из Decoder  
блоков

## Attention



## MLP



# Attention

Механизм внимания  
позволяет учитывать  
взаимосвязи токенов

Сложность:  $O(T^2 * D + T * D^2)$

Память:  $O(T^2 + T * D)$

# MLP

Производит преобразования  
знаний модели

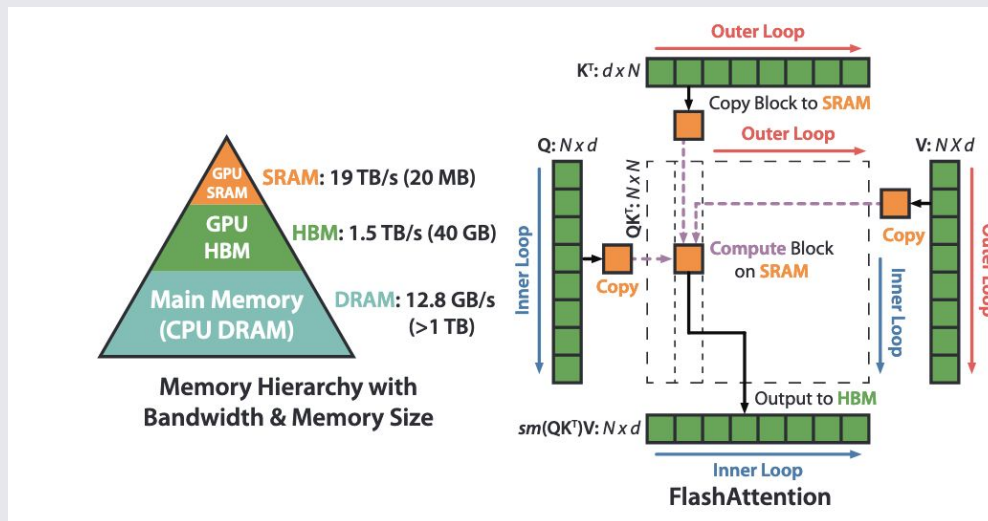
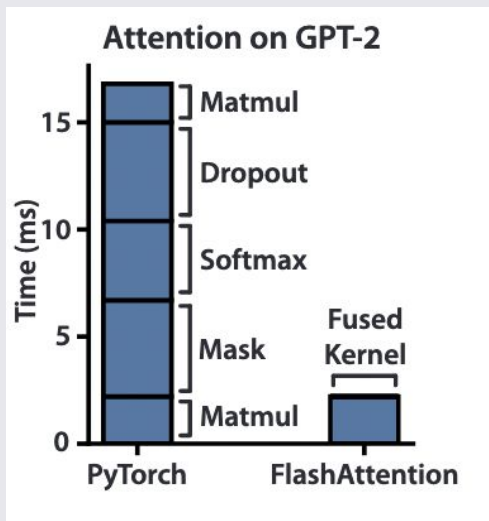
Сложность:  $O(TDH)$

Память:  $O(TH)$

# Flash-Attention

Память -7GB  
Скорость +55%

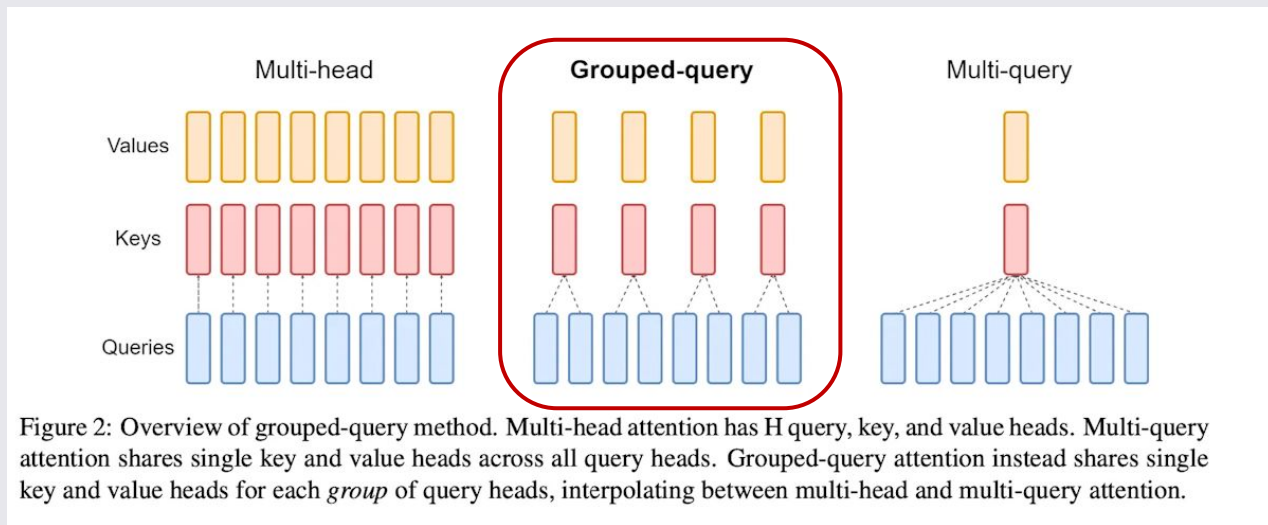
Эффективное вычисление внимания на GPU



# GQA

Память -1.2GB  
Скорость +9%

В  $k$  раз (число общих голов) уменьшается потребление памяти  $K, V$  матрицами, ускоряется обучение и инференс



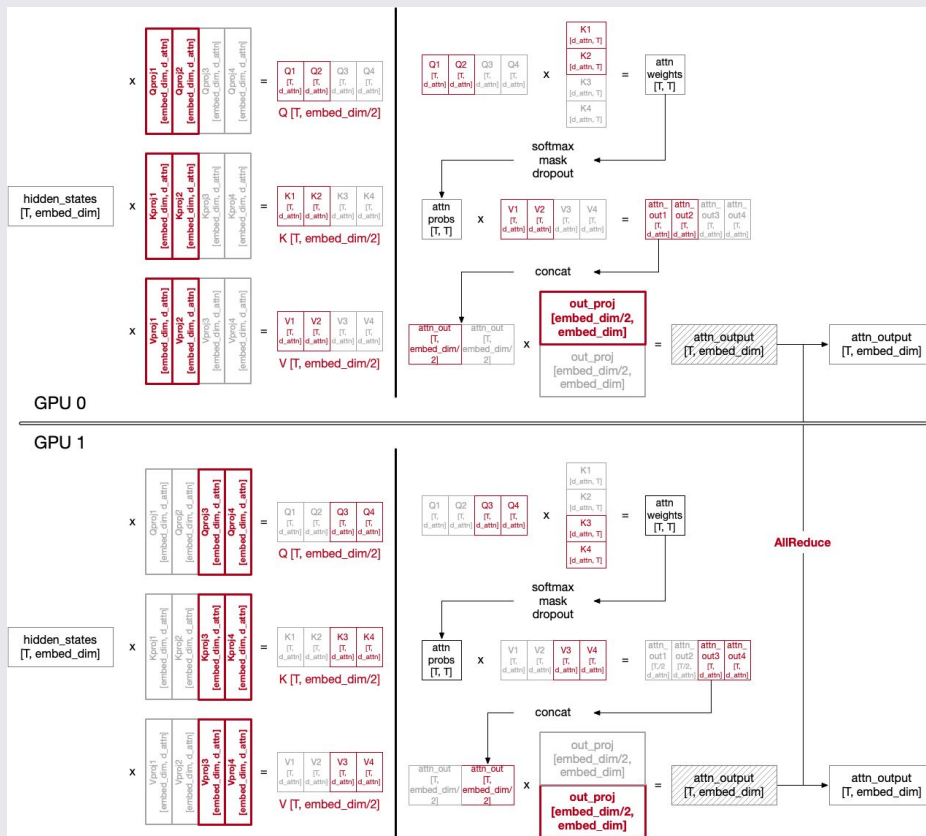
# Tensor Parallel

Память      Скорость  
 -50%      +91%

Делим вычисления между GPU по размерности D

+ Каждая GPU держит в TP раз меньше весов и считает в TP раз меньше активации

- Тратится время на коммуникации между GPU



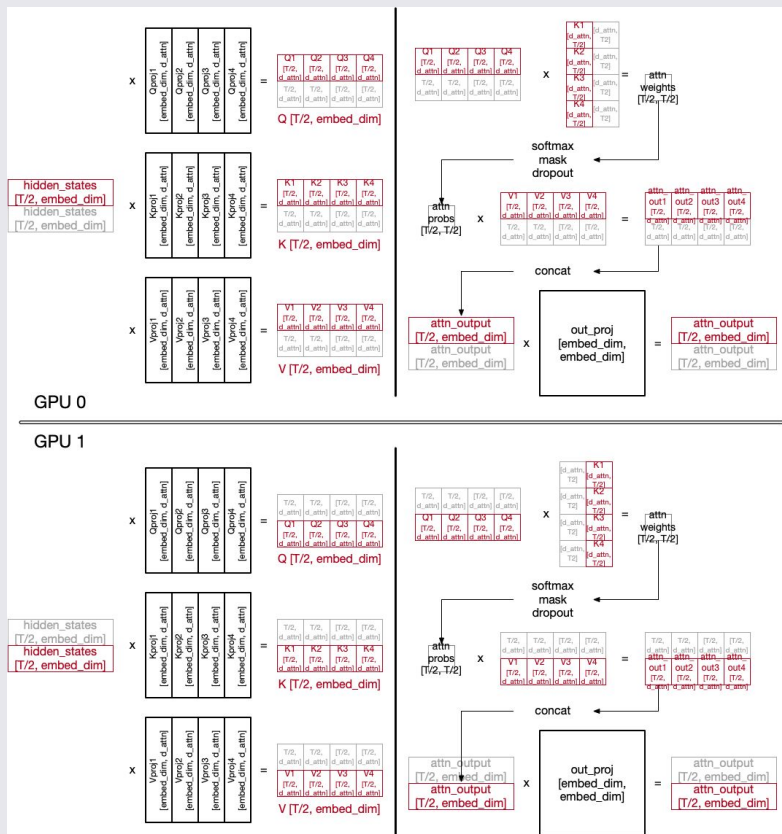
# Sequence Parallel

Память      Скорость  
 -51%      +88%

Делим вычисления между GPU по размерности T

+ Каждая GPU считает в TP раз меньше активации

- Тратится время на коммуникации в Attention



# Tensor VS Sequence Parallel

Table 2: Multi-head attention block memory usage comparison

	Operation	M1	M2	output	Memory
Tensor parallelism	$Q/K/V$	$(B, L, H)$	$(H, \frac{ZA}{N})$	$(B, \frac{Z}{N}, L, A)$	
	$QK^T$	$(B, \frac{Z}{N}, L, A)$	$(B, \frac{Z}{N}, L, A)$	$(B, \frac{Z}{N}, L, L)$	$\frac{16AZH}{N} + \frac{4BLZA}{N}$
	$AV$	$(B, \frac{Z}{N}, L, L)$	$(B, \frac{Z}{N}, L, A)$	$(B, \frac{Z}{N}, L, A)$	$+ \frac{BZL^2}{N} + BLH$
	Linear	$(B, \frac{Z}{N}, L, A)$	$(\frac{AZ}{N}, H)$	$(B, L, H)$	
Sequence parallelism	$Q/K/V$	$(B, \frac{L}{N}, H)$	$(H, AZ)$	$(B, Z, \frac{L}{N}, A)$	
	Ring- $QK^T$	$(B, Z, \frac{L}{N}, A)$	$(B, Z, \frac{L}{N}, A)$	$(B, Z, \frac{L}{N}, L)$	$16AZH + \frac{4BZLA}{N}$
	Ring- $AV$	$(B, Z, \frac{L}{N}, L)$	$(B, Z, \frac{L}{N}, A)$	$(B, Z, \frac{L}{N}, A)$	$+ \frac{BZL^2}{N} + \frac{BLH}{N}$
	Linear	$(B, Z, \frac{L}{N}, A)$	$(AZ, H)$	$(B, \frac{L}{N}, H)$	

# Tensor VS Sequence Parallel

Для LLaMa-2 моделей верно следующее ( $B = 1$ ):

## Attention

SP выгоднее TP при

$T > 8.1D$

## MLP

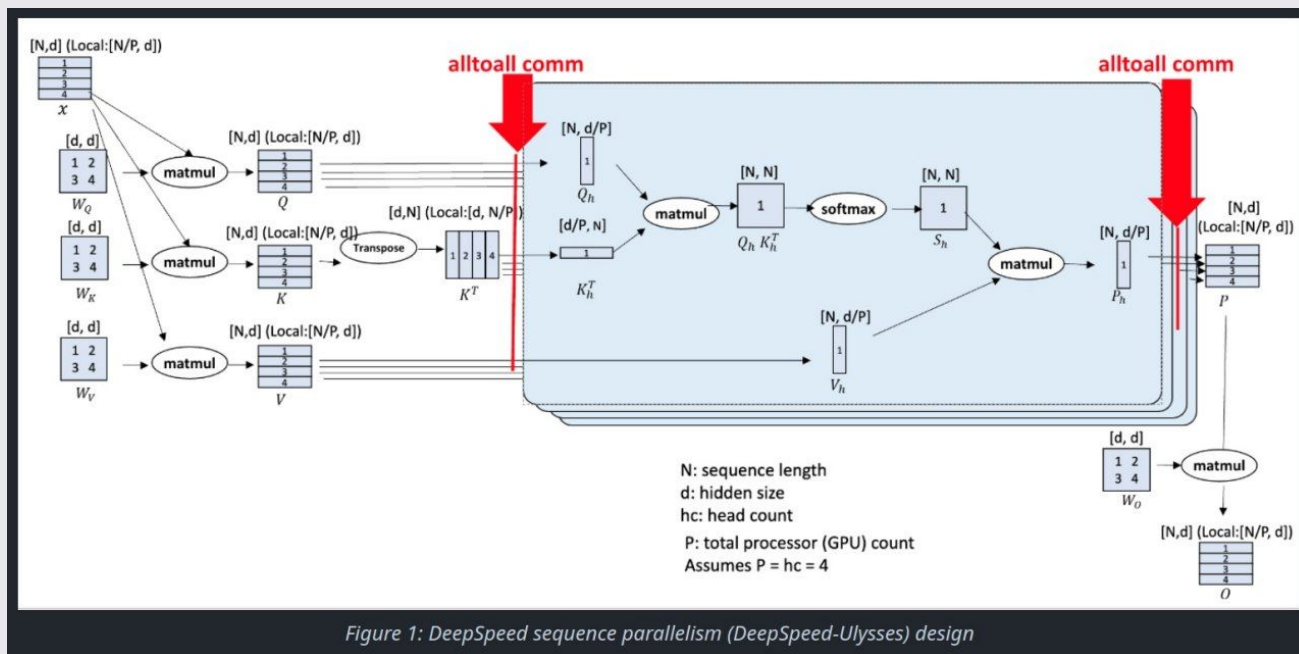
SP выгоднее TP при

$T > 16D$

\* А еще есть параллельные Эмбединги и параллельный лосс

# Реализация SP - all2all

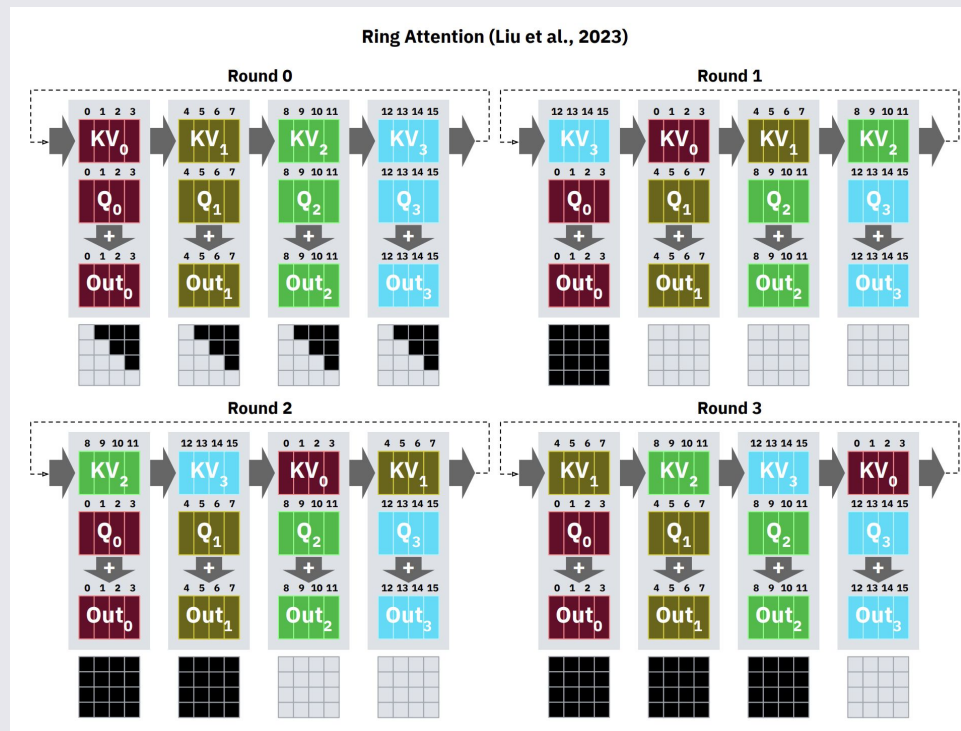
Архитектурно Decoder LLM может обрабатывать произвольный КОНТЕКСТ



# Ring Attention

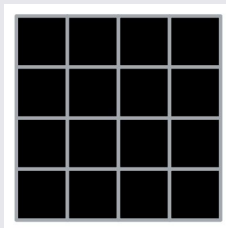
Экстремальные контексты  
требуют подсчета  
попарного внимания всех на  
всех токенов.

Можем считать внимание в  
цикле

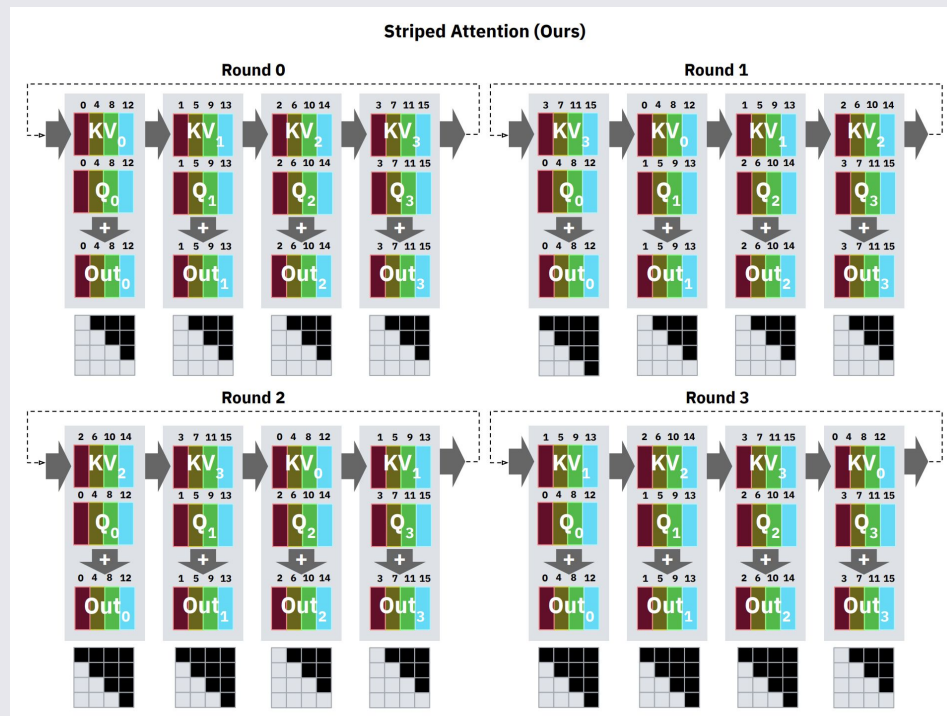


# Striped Ring Attention

Такая маска -  
простой GPU



Striped Attention  
реорганизует  
вычисления на более  
равномерные



# ZigZag Ring Attention

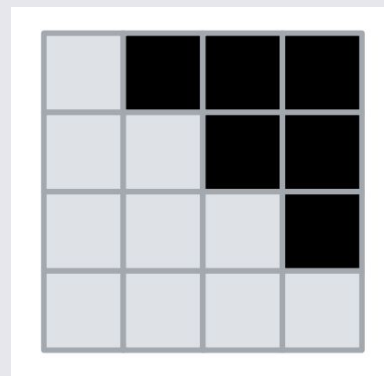
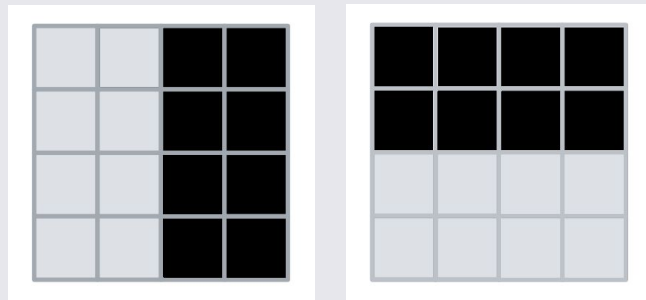
Striped

0,4,8,12 | 1,5,9,13 | 2,6,10,14 | 3,7,11,15

ZigZag - группируем по 2 и делаем  
разбиение

0,1,14,15 | 2,3,12,13 | 4,5,10,11 | 6,7,8,9

Каждая GPU простаивает не больше  
половины времени. Лучше Striped



# Сравнение реализаций

	GPU	theoretic flash_attn	ring_attn	zigzag_ring	stripe_attn
fwd only (iter/sec)	8xH800	2418.4 / 8 = 302.3	208.0	283.0	259.6
			68.8%	<b>93.6%</b>	85.9%
fwd + bwd (iter/sec)	8xH800	705.2 / 8 = 88.2	54.3	75.7	76.9
			61.5%	85.9%	<b>87.2%</b>
fwd only (iter/sec)	8xA100	1545.9 / 8 = 193.2	124.4	179.0	163.9
			64.3%	<b>92.7%</b>	84.8%
fwd + bwd (iter/sec)	8xA100	470.6 / 8 = 58.8	33.3	49.5	45.9
			56.6%	<b>84.1%</b>	78.1%

Лучшая

# Ring Attention VS all2all

TP=1, SP=4, 64 GPU

	Ring Attention		all2all	
	tokens per sec / 1 GPU	mem (GB) / 1 GPU	tokens per sec / 1 GPU	mem (GB) / 1 GPU
LLaMa 7B 32k	1717	<b><u>16,9</u></b>	<b><u>2030</u></b>	20,6
LLaMa 7B 128k	911	<b><u>30,9</u></b>	<b><u>946</u></b>	47,0
LLaMa 30B 32k	<b><u>907</u></b>	<b><u>36,8</u></b>	519	50
LLaMa 30B 128k	<b><u>239</u></b>	<b><u>55,1</u></b>	—	—

# Обучение GigaChat

# Обучение GigaChat

Подведем итоги:

- Имеем формулы подсчета активаций
- Знаем, когда Tensor Parallel выгоднее Sequence Parallel
- Умеем объединять вместе оба режима
- Эффективный Ring Attention - нужен для экстремальных контекстов

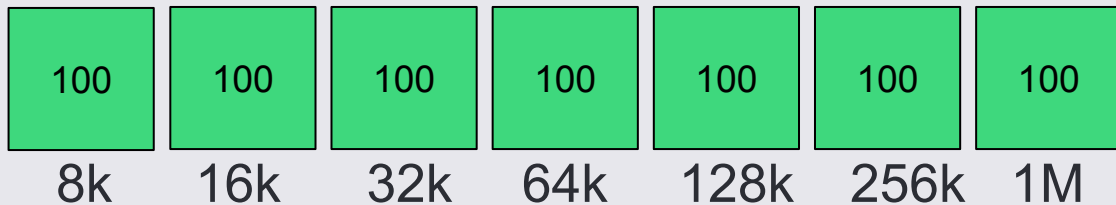
# GigaChat 7B и 29B

## Результаты PassKey

GigaChat 29B  
**128k**



GigaChat 7B  
**1M**



# GigaChat 7B и 29B

Результаты PassKey

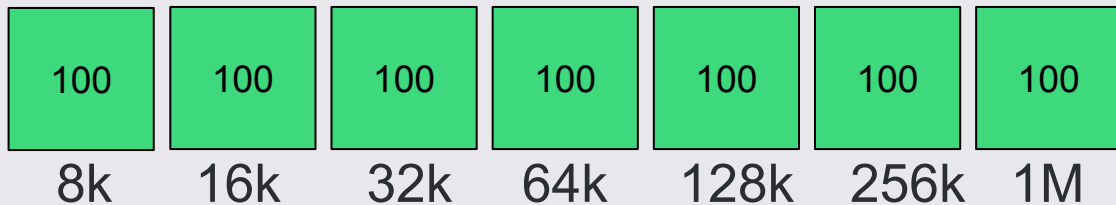
GigaChat 29B

**128k**

64

GigaChat 7B

**1M**



# GigaChat 7B и 29B

Результаты PassKey

GigaChat 29B

**128k**

64

GigaChat 7B

**1M**

512 страниц A4

# Перспективы

# Будущее длинного контекста

Код, видео и аудио требуют большого контекста



# Будущее длинного контекста

Код, видео и аудио требуют большого контекста



Решение бизнес задач с дополнительной информацией

# Будущее длинного контекста

Код, видео и аудио требуют большого контекста



Решение бизнес задач с дополнительной информацией

Персонализация языковых моделей



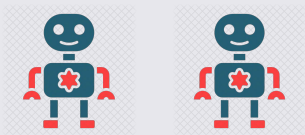
# Будущее длинного контекста

Код, видео и аудио требуют большого контекста



Решение бизнес задач с дополнительной информацией

Персонализация языковых моделей



Мульти-агентное взаимодействие

# Спасибо за внимание!

**Косарев Евгений**

Senior DL Engineer, Sber Devices

Тг: [evgenijkkk](https://t.me/evgenijkkk)

Почта: [evgenijkkk@yandex.ru](mailto:evgenijkkk@yandex.ru)