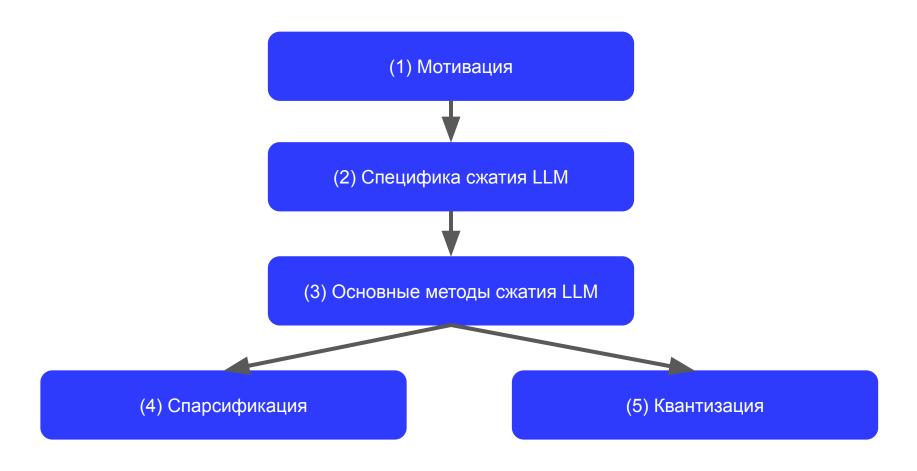
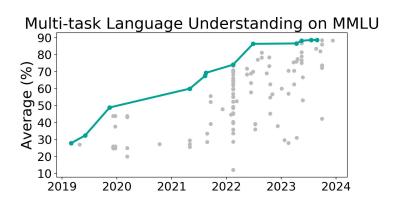
Методы сжатия больших языковых моделей

Кузнеделев Денис

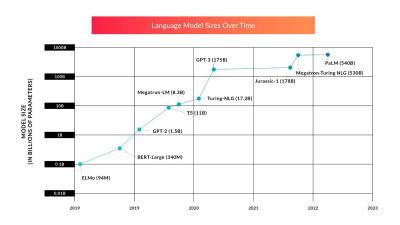




Зачем сжимать LLM?



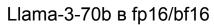
Большие языковые модели становятся все умнее и умнее

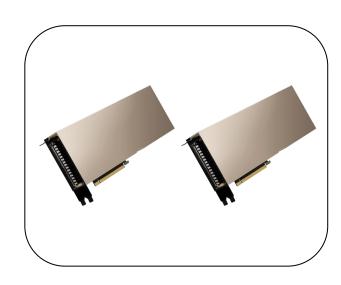


Но и больше, увы...

Зачем сжимать LLM?

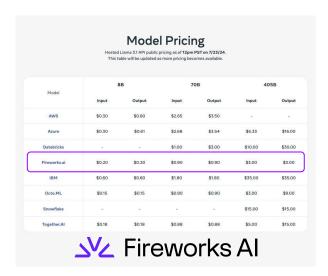






2x Nvidia A100

Кому это нужно?



Провайдерам для снижения стоимости инференса.



turboderp/exllamav2





Простым смертным для инференса на пользовательском железе.



Современные LLM довольно увесистые

Методы сжатия делятся на следующие категории

Data-Free

LLM.int8() HQQ bnb-4bit

Data-aware с локальными критериями

GPTQ AWQ SparseGPT QuIP# AQLM Efficient QAT

Инференс LLM Memory-bound

Основное время занимают не сами вычисления, а время на подрузку весов модели из HDM/GDDR в кэши GPU

Можно добиться ускорения за счет сжатия одних лишь весов, не трогая активации

Нейронная сеть, вообще говоря, сложная нелинейная функция.

Для упрощения описания в окрестности данной точки прибегают к разложению лосс-функции до **2-го порядка** по формуле Тейлора:

$$\boxed{\mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{w}^*) \simeq \boxed{\nabla_{\mathcal{L}}^{\top}(\mathbf{w}^*)(\mathbf{w} - \mathbf{w}^*)} + \boxed{\frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^{\top}\mathbf{H}_{\mathcal{L}}(\mathbf{w}^*)(\mathbf{w} - \mathbf{w}^*)}$$

Функция потерь

Член 1-го порядка

Равен 0 в точке оптимума

Член 2-го порядка

Наш бро

Проиллюстрируем метод на примере прунинга одного веса \mathbf{w}_i

$$\min_{\mathbf{w}}(\mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{w}^*))$$

при условии

$$\mathbf{w}_i + \mathbf{e}_i^{\mathsf{T}} \mathbf{w} = 0$$

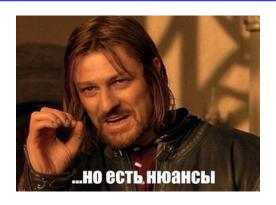
Имеет следующее решение

$$\rho_i = \frac{w_i^2}{2[\mathbf{H}_{\mathcal{L}}^{-1}(\mathbf{w}^*)]_{ii}}$$

$$\delta \mathbf{w}^* = -\frac{w_i}{[\mathbf{H}_{\mathcal{L}}^{-1}(\mathbf{w}^*)]_{ii}} \mathbf{H}_{\mathcal{L}}^{-1}(\mathbf{w}^*) \mathbf{e}_i,$$

Приращение ошибки

Оптимальный сдвиг остальных весов



Гессиан квадратичен по количеству параметров в сети - **оченьочень** дорого

Решение?

Использовать ошибку на выходе данного слоя

$$\text{argmin}_{\widehat{\mathbf{W}}_{\ell}} \quad ||\mathbf{W}_{\ell}\mathbf{X}_{\ell} - \widehat{\mathbf{W}}_{\ell}\mathbf{X}_{\ell}||_2^2$$

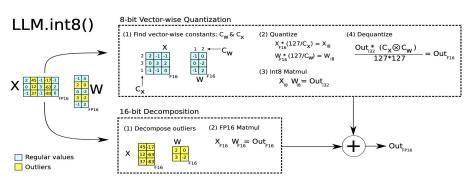
- 1) Память, требуемая для хранения Гессиана $\Theta(d_{col}^2)$
- 2) Гессиан один и тот же для всех выходных каналов слоя
- 3) В процессе сжатия не нужно загружать всю модель в память

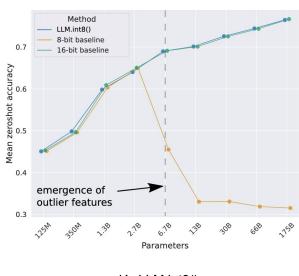
GPTQ SpQR Wanda AQLM
SparseGPT AWQ QuIP#

Выбросы / Outliers

В многих современных нейронных сетях есть небольшая доля весов, крайне чувствительных к изменениям.

В работе LLM.int8() было предложено не трогать их при квантизации и держать в исходной точности.





Из LLM.int8()

Dettmers, Tim, et al. "Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale." *Advances in neural information processing systems* 35 (2022): 30318-30332.

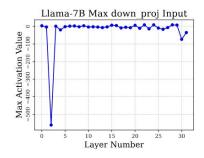
Выбросы / Outliers



Супер-веса - отдельные параметры, с огромным влиянием на выход сети

Model	N	о.		Type		Weight	Coordinates	Model	No.	1
Llama 7B		2	I	mlp	1	down_proj	[3968, 7003]	OLMo-1B	1	
	1 2	2	T	mlp	T	down_proj	[2231, 2278]	0724-hf	1	10
Llama 13B		2		mlp		down_proj	[2231, 6939]		1	1.0
	-	_	+		+			OLMo-7B	2	10
		3		mlp	П	down_proj	[5633, 12817]	0724-hf	7	1 .
Llama 30B		3		mlp	П	down_proj	[5633, 17439]		24	1.5
	1	0		mlp	1	down_proj	[5633, 14386]			1
Llama2 7B	1 -	1	ī	- 1	1	4 4.1	[2522 7000]		2	
Liamaz /B	- 2	ı.		mlp		down_proj	[2533, 7890]		2	10
Llama2 13B	1 -	3	T	mlp	T	down_proi	[4743, 7678]	Phi-3	2	10
Diminia 155		_	_	шр	_	downsproj	[1710,7070]	mini-4k-instruct	4	10
Mistral-7B				mlp	1	down_proj	[2070, 7310]		4	10
v0.1	1 8			шр	П	uown_proj	[2070, 7310]		4	1

Model	No.	Type	Weight	Coordinates
OLMo-1B	1	mlp	down_proj	[1764, 1710]
0724-hf	1	mlp	down_proj	[1764, 8041]
	1	mlp	down_proj	[269, 7467]
OLMo-7B	2	mlp	down_proj	[269, 8275]
0724-hf	7	mlp	down_proj	[269, 453]
	24	mlp	down_proj	[269, 2300]
	2	mlp	down_proj	[525, 808]
	2	mlp	down_proj	[1693, 808]
Phi-3	2	mlp	down_proj	[1113, 808]
mini-4k-instruct	4	mlp	down_proj	[525, 2723]
	4	mlp	down_proj	[1113, 2723]
	4	mlp	down_proj	[1693, 2723]



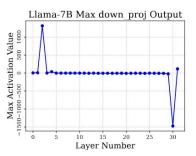


Table 2: Super Weight Directory. The above layer numbers, layer types, and weight types can be directly applied to Huggingface models. For example, for Llama-7B on Huggingface, access the super weight using layers [2].mlp.down_proj.weight[3968, 7003].

Yu, Mengxia, et al. "The super weight in large language models." arXiv preprint arXiv:2411.07191 (2024).





С эффективными кернелами...

Спарсификация

Неструктурированная

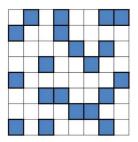
Нет ускорения на GPU :(

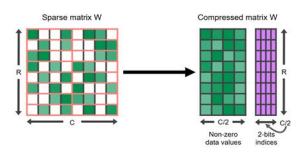
Полуструктурированная

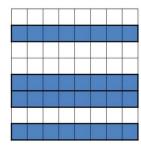
Ускорение на NVIDIA Ampere+

Структурированная

Почти гарантированное ускорение



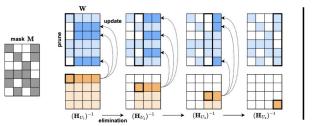


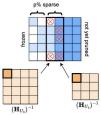


SparseGPT Wanda ADMM-Pruning ZipLM CoFi LLMPruner

SparseGPT

- 1) Использует послойное приближение
- 2) Убирает веса с наименьшей ошибкой
- 3) Для эффективного обновления Гессиана сжатого веса используется <u>Cholesky разложение</u>.





Визуализация алгоритма

Algorithm 1 The SparseGPT algorithm. We prune the layer matrix \mathbf{W} to p% unstructured sparsity given inverse Hessian $\mathbf{H}^{-1} = (\mathbf{X}\mathbf{X}^{\top} + \lambda \mathbf{I})^{-1}$, lazy batch-update block-size B and adaptive mask selection blocksize B_s ; each B_s consecutive columns will be p% sparse.

```
\begin{split} \mathbf{M} &\leftarrow \mathbf{1}_{d_{\text{row}} \times d_{\text{col}}} \text{ $/$binary pruning mask} \\ \mathbf{E} &\leftarrow \mathbf{0}_{d_{\text{row}} \times B} \text{ $/$block quantization errors} \\ \mathbf{H}^{-1} &\leftarrow \text{Cholesky}(\mathbf{H}^{-1})^{\top} \text{ $/$Hessian inverse information} \\ \textbf{for } i &= 0, B, 2B, \dots \textbf{do} \\ \textbf{for } j &= i, \dots, i + B - 1 \textbf{do} \\ \textbf{if } j \text{ mod } B_s &= 0 \textbf{ then} \\ \mathbf{M}_{:,j:(j+B_s)} &\leftarrow \text{mask of } (1-p)\% \text{ weights } w_c \in \\ \mathbf{W}_{:,j:(j+B_s)} &\text{with largest } w_c^2/[\mathbf{H}^{-1}]_{cc}^2 \\ \textbf{end if} \\ \mathbf{E}_{:,j-i} &\leftarrow \mathbf{W}_{:,j}/[\mathbf{H}^{-1}]_{jj} \text{ $/\!\!\!/ pruning error} \\ \mathbf{E}_{:,j-i} &\leftarrow (1-\mathbf{M}_{:,j}) \cdot \mathbf{E}_{:,j-i} \text{ $/\!\!\!/ freeze weights} \\ \mathbf{W}_{:,j:(i+B)} &\leftarrow \mathbf{W}_{:,j:(i+B)} - \mathbf{E}_{:,j-i} \cdot \mathbf{H}_{j,j:(i+B)}^{-1} \text{ $/\!\!\!/ update$} \\ \textbf{end for} \\ \mathbf{W} &\leftarrow \mathbf{W} \cdot \mathbf{M} \text{ $/\!\!\!/ set pruned weights to 0} \end{split}
```

Алгоритм

Frantar, Elias, and Dan Alistarh. "Sparsegpt: Massive language models can be accurately pruned in one-shot." *International Conference on Machine Learning*. PMLR, 2023.

SparseGPT

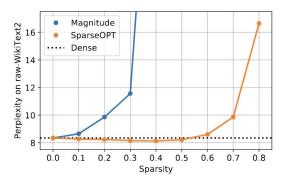


Figure 1. Sparsity-vs-perplexity comparison of SparseGPT against magnitude pruning on OPT-175B, when pruning to different *uniform* per-layer sparsities.

Работает лучше Data-free

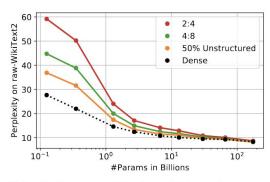


Figure 2. Perplexity vs. model and sparsity type when compressing the entire OPT model family (135M, 350M, ..., 66B, 175B) to different sparsity patterns using SparseGPT.

Большие модели легче сжимать

Frantar, Elias, and Dan Alistarh. "Sparsegpt: Massive language models can be accurately pruned in one-shot." *International Conference on Machine Learning*. PMLR, 2023.

Wanda

Ограничимся диагональю Гессиана

$$\mathbf{S}_{ij} = |\mathbf{W}_{ij}| \cdot \|\mathbf{X}_j\|_2$$

Критерий важности

Вполне себе работает...

				LLaN	MА		L	LaMA-	2
Method	Weight Update	Sparsity	7B	13B	30B	65B	7B	13B	70B
Dense	-	0%	5.68	5.09	4.77	3.56	5.12	4.57	3.12
Magnitude	×	50%	17.29	20.21	7.54	5.90	14.89	6.37	4.98
SparseGPT	✓	50%	7.22	6.21	5.31	4.57	6.51	5.63	3.98
Wanda	×	50%	7.26	6.15	5.24	4.57	6.42	5.56	3.98
Magnitude	×	4:8	16.84	13.84	7.62	6.36	16.48	6.76	5.54
SparseGPT	✓	4:8	8.61	7.40	6.17	5.38	8.12	6.60	4.59
Wanda	×	4:8	8.57	7.40	5.97	5.30	7.97	6.55	4.47
Magnitude	×	2:4	42.13	18.37	9.10	7.11	54.59	8.33	6.33
SparseGPT	✓	2:4	11.00	9.11	7.16	6.28	10.17	8.32	5.40
Wanda	×	2:4	11.53	9.58	6.90	6.25	11.02	8.27	5.16

Table 3: WikiText perplexity of pruned LLaMA and LLaMA-2 models. Wanda performs competitively against prior best method SparseGPT, without introducing any weight update.

Sun, Mingjie, et al. "A simple and effective pruning approach for large language models." arXiv preprint arXiv:2306.11695 (2023).

ZipLM

- 1) Использует послойное приближение
- 2) Убирает входную размерность / группу размерностей с наименьшей ошибкой за один ход.
- 3) В контексте сжатия трансформеров ZipLM за один шаг убирает hidden_dim в FFN или целую attention голову.
- 4) Gradual pruning с дообучением.

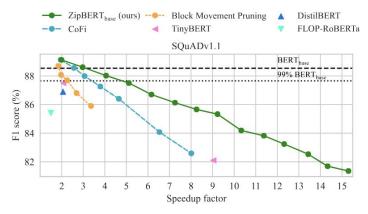
Algorithm 1 The ZipLM pruning algorithm. Given inverse Hessian $\mathbf{H}^{-1} = (2\mathbf{X}\mathbf{X}^{\top} + \lambda \mathbf{I})^{-1}$, we remove exactly k structures from the corresponding weight matrix \mathbf{W} .

```
\begin{split} \mathbf{R} &\leftarrow \text{set of all possible structures} \\ & \textbf{for } k \text{ times } \textbf{do} \\ & \mathbf{S} \leftarrow \operatorname{argmin}_{\mathbf{S}} \sum_{i=0}^{d_{\text{row}}} \mathbf{W}_{i,\mathbf{M_S}} \cdot ((\mathbf{H}^{-1})_{\mathbf{M_S},\mathbf{M_S}})^{-1} \cdot \mathbf{W}_{i,\mathbf{M_S}}^{\top} \\ & \boldsymbol{\delta_S} \leftarrow -\mathbf{W}_{:,\mathbf{M_S}} \cdot ((\mathbf{H}^{-1})_{\mathbf{M_S},\mathbf{M_S}})^{-1} \cdot (\mathbf{H}^{-1})_{\mathbf{M_S},:} \\ & \mathbf{W} \leftarrow \mathbf{W} + \boldsymbol{\delta_S} \\ & \mathbf{H}^{-1} \leftarrow \mathbf{H}^{-1} - \mathbf{H}_{:,\mathbf{M_S}}^{-1} \cdot ((\mathbf{H}^{-1})_{\mathbf{M_S},\mathbf{M_S}})^{-1} \cdot \mathbf{H}_{\mathbf{M_S},:}^{-1} \\ & \mathbf{R} \leftarrow \mathbf{R} - \{\mathbf{S}\} \\ & \mathbf{end for} \\ & \mathbf{W} \leftarrow \mathbf{W} \odot \mathbf{M_R} \end{split}
```

Алгоритм

Kurtić, Eldar, Elias Frantar, and Dan Alistarh. "Ziplm: Inference-aware structured pruning of language models." *Advances in Neural Information Processing Systems* 36 (2023): 65597-65617.

ZipLM



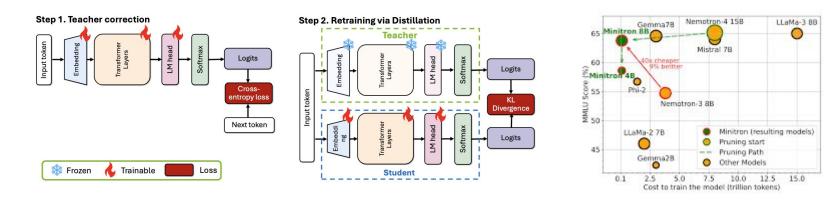
BERT хорошо сжимается

	Prunin	g for thro	oughput	Prun	ing for la	tency
Model	Speedup	Decoder size	Wiki Text-103 PPL↓	Speedup	Decoder size	Wiki Text-103 PPL↓
GPT2*	1.0x	85.0M	28.5	1.0x	85.0M	28.5
DistilGPT2	1.6x	42.5M	43.0	1.9x	42.5M	43.0
ZipGPT2 (ours)	1.5x 2.1x 2.7x 3.3x	47.3M 26.5M 14.0M 5.7M	35.4 41.5 50.4 72.1	1.6x 2.0x 2.2x 2.5x	48.7M 39.2M 26.6M 20.7M	37.8 41.2 49.0 55.0

Гопота не очень :(

Kurtić, Eldar, Elias Frantar, and Dan Alistarh. "Ziplm: Inference-aware structured pruning of language models." *Advances in Neural Information Processing Systems* 36 (2023): 65597-65617.

Minitron



Если дистиллировать сжатую модель на большом количестве данных (~500В токенов) можно выжать хорошее качество.

Sreenivas, Sharath Turuvekere, et al. "Llm pruning and distillation in practice: The minitron approach." arXiv preprint arXiv:2408.11796 (2024).

Спарсификация



На практике редко удается добиться существенного сжатия/ускорения без значительной просадки в качестве или длительного до-обучения

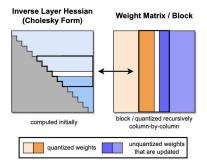
Квантизация - способ уменьшения количества памяти, требуемой на хранение массива данных, за счет представления его элементов из некоторого ограниченного множества.

fp16/bf16 - 16 бит на параметр INT8/fp8 - 8 бит на параметр INT4/fp4* - 4 бит на параметр

GPTQ

- 1) Использует послойное приближение
- 2) Для эффективного обновления Гессиана сжатого веса используется <u>Cholesky разложение</u>.

```
Algorithm 1 Quantize W given inverse Hessian \mathbf{H}^{-1} = (2\mathbf{X}\mathbf{X}^{\top} + \lambda \mathbf{I})^{-1} and blocksize B.
   \mathbf{Q} \leftarrow \mathbf{0}_{d_{\mathrm{row}} 	imes d_{\mathrm{col}}} \ \mathbf{E} \leftarrow \mathbf{0}_{d_{\mathrm{row}} 	imes B}
                                                                                                        // quantized output
                                                                                                       // block quantization errors
   \mathbf{H}^{-1} \leftarrow \text{Cholesky}(\mathbf{H}^{-1})^{\top}
                                                                                                       // Hessian inverse information
   for i = 0, B, 2B, ... do
        for j = i, ..., i + B - 1 do
             \mathbf{Q}_{:,j} \leftarrow \operatorname{quant}(\mathbf{W}_{:,j})
                                                                                                       // quantize column
             \mathbf{E}_{:,j-i} \leftarrow (\mathbf{W}_{:,j} - \mathbf{Q}_{:,j}) / [\mathbf{H}^{-1}]_{jj}
                                                                                                       // quantization error
             \mathbf{W}_{:,j:(i+B)} \leftarrow \mathbf{W}_{:,j:(i+B)} - \mathbf{E}_{:,j-i} \cdot \mathbf{H}_{j,j:(i+B)}^{-1}
                                                                                                       // update weights in block
        end for
        \mathbf{W}_{:,(i+B):} \leftarrow \mathbf{W}_{:,(i+B):} - \mathbf{E} \cdot \mathbf{H}_{i:(i+B),(i+B):}^{-1}
                                                                                                       // update all remaining weights
   end for
```



[1] Frantar, Elias, et al. "Gptq: Accurate post-training quantization for generative pre-trained transformers." *arXiv preprint arXiv:2210.17323* (2022).



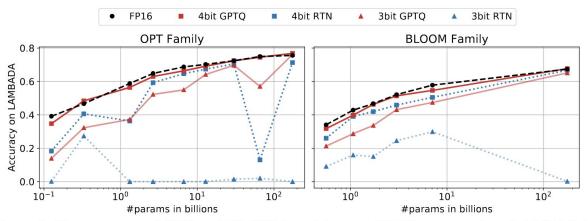


Figure 3: The accuracy of OPT and BLOOM models post-GPTQ, measured on LAMBADA.

GPTQ сжимает почти* без просадки LLM

Frantar, Elias, et al. "Gptq: Accurate post-training quantization for generative pre-trained transformers." arXiv preprint arXiv:2210.17323 (2022).

AWQ

масштаб нуль

$$w_q = S(q + Z)$$

 $\mathbf{s}^* = \operatorname*{arg\,min}_\mathbf{s} \mathcal{L}(\mathbf{s})$ $\mathcal{L}(\mathbf{s}) = \|Q(\mathbf{W} \cdot \mathrm{diag}(\mathbf{s}))(\mathrm{diag}(\mathbf{s})^{-1} \cdot \mathbf{X}) - \mathbf{W}\mathbf{X}\|$

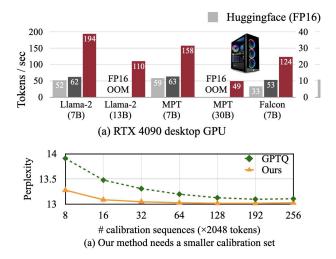
Представление веса при квантизации

Задача оптимизации

- 1) Оптимизирует Scale для послойной ошибки
- 2) Оптимизация осуществляется через поиск по сетке
- 3) На инференсе надо делить на выученный масштаб недорого

Lin, Ji, et al. "Awq: Activation-aware weight quantization for on-device Ilm compression and acceleration." *Proceedings of Machine Learning and Systems* 6 (2024): 87-100.



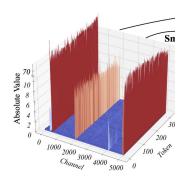


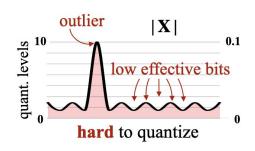
PPL↓			Llama-2	
		7B	13B	70B
FP16	(-	5.47	4.88	3.32
INT3 g128	RTN GPTQ GPTQ-R AWQ	6.66 6.43 6.42 6.24	5.52 5.48 5.41 5.32	3.98 3.88 3.86 3.74
INT4 g128	RTN GPTQ GPTQ-R AWQ	5.73 5.69 5.63 5.60	4.98 4.98 4.99 4.97	3.46 3.42 3.43 3.41

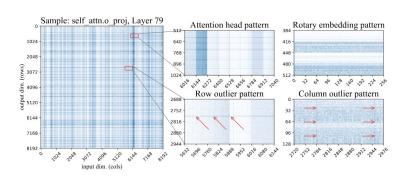
- 1) Дает неплохое ускорение инференса
- 2) Примерно по качеству равен GPTQ (если не жульничать :))
- 3) Требует меньшего количества сэмплов для сходимости

Lin, Ji, et al. "Awq: Activation-aware weight quantization for on-device Ilm compression and acceleration." *Proceedings of Machine Learning and Systems* 6 (2024): 87-100.









Изоляция и хранение в виде разреженной матрицы

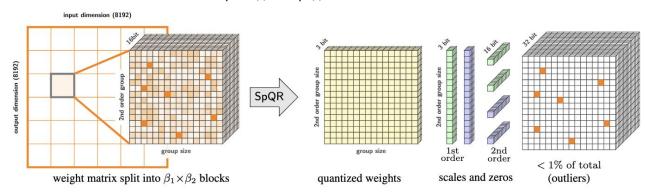
LLM.int8() SpQR SqueezeLLM Переход в другой базис

QuIP# QuIP# QuaRot

SpQR

$$s_{ij} = \frac{(w_{ij} - \mathsf{quant}(w_{ij}))^2}{2(XX^\top)^{-1}}$$

Метрика для определения важных весов



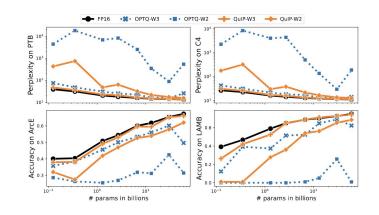
- 1) "Важные" веса не квантизуются
- 2) Для большего сжатия скейлы тоже квантизуются

Dettmers, Tim, et al. "Spqr: A sparse-quantized representation for near-lossless Ilm weight compression." arXiv preprint arXiv:2306.03078 (2023).

QuIP

Algorithm 1 QuIP - Incoherence Pre-Processing

```
Require: b \in \mathbb{N}, H \in \mathbb{R}^{n \times n} SPD, original W \in \mathbb{R}^{m \times n}, \rho \in \mathbb{R}_+, \alpha \in [0,1]1: seeded sample random two-factor orthogonal matrices U \in \mathbb{R}^{m \times m} and V \in \mathbb{R}^{n \times n}2: H = H + \alpha * \text{mean}(\text{diag}(H))I\triangleright from OPTQ3: \tilde{D} \leftarrow \sqrt[4]{\text{diag}(H)/\text{diag}(W^TW)}\triangleright \sqrt[4]{} applies element-wise4: W \leftarrow W\tilde{D}; H \leftarrow \tilde{D}^{-1}H\tilde{D}^{-1}\triangleright diagonal rescaling5: W \leftarrow UWV^T; H \leftarrow VHV^T\triangleright incoherence6: s \leftarrow \rho ||W||_F/\sqrt{mn}; W \leftarrow \frac{1}{2}(\frac{1}{s}W+1)\triangleright reduced quantization range due to incoherency7: W \leftarrow \text{clamp}(W * (2^b - 1), 0, 2^b - 1)\triangleright rescale W to lie within [0, 2^b - 1]8: return W, H, s, \tilde{D}
```



Алгоритм

Сеть не разлетается при 2-битной квантизации

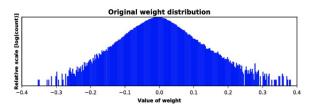
- 1) Поворот весов на случайную ортогональную матрицу "регуляризует" распределение весов
- 2) Матрицы поворота надо применять к активациям на инференсе
- 3) Матрицы поворота довольно компактные и "дешевые" (matvec быстрее, чем O(n^2))

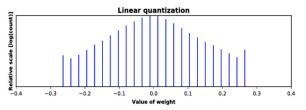
Chee, Jerry, et al. "Quip: 2-bit quantization of large language models with guarantees." *Advances in Neural Information Processing Systems* 36 (2023): 4396-4429.

Скалярная и векторная квантизация

Скалярная

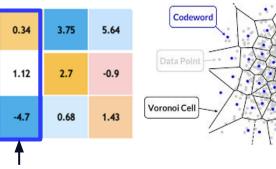
каждому значению сопоставляется индекс на 1-мерной сетке

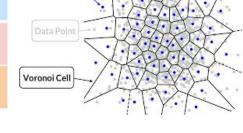




Векторная

каждой группе значений сопоставляется вектор из кодовой книги





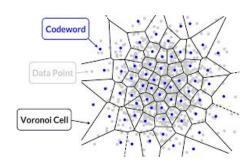
квантизуемые совместно веса

Скалярная и векторная квантизация

Векторная

каждой группе значений сопоставляется вектор из кодовой книги

Кодовые слова - векторы (узлы) Кодовая книга - совокупность кодовых слов (решетка) Код - индекс вектора в решетке



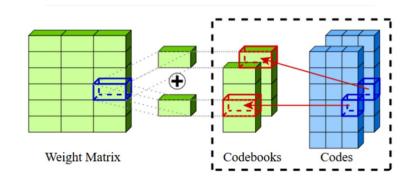
AQLM

- 1) Представим группу подряд идущих весов как взвешенную сумму векторов из одной или более кодовых книг
- 2) Используем послойное приближение
- 3) Аддитивная квантизация возникла изначально в контексте задачи оптимизации поиска

$$\sum_{m=1}^{M} C_m b_m$$

С - кодовые слова

b - one-hot коды



Представление веса

Egiazarian, Vage, et al. "Extreme compression of large language models via additive quantization." arXiv preprint arXiv:2401.06118 (2024).

AQLM

Как найти оптимальные С и b?

Целевая функция

$$||\mathbf{W}\mathbf{X} - \sum_{m=1}^{M} C_m b_m \mathbf{X}||_2^2|$$

- 1) Непрерывные параметры (С) оптимизируются через градиентный спуск (Adam)
 - 2) Дискретные параметры (b) оптимизируются с помощью beam search
 - 3) Оптимизирует С и b поочередно

На практике нужно 10^2-10^3 итераций для сходимости

AQLM

Как найти оптимальные C и b?

Конечная цель не послойная ошибка - а итоговое качество Можно ли оптимизировать более "глобально"?

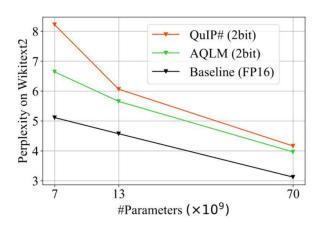
- Для блочной оптимизации используется **L2** между признаками исходной и сжатой модели
 - Для end-2-end (следуя QuIP#) **КL-дивергенция**

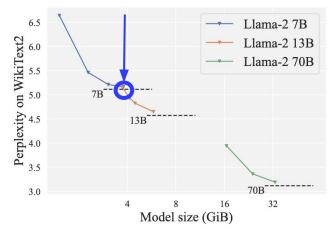
AQLM



Обучение всей модели (С)

AQLM





Результаты при 2-битной квантизации

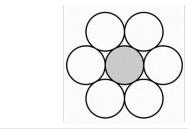
Парето-оптимальность при 2.5 битах

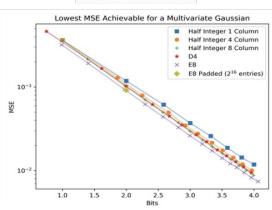
Парето-оптимальность

- оптимальный размер модели и степень сжатия при заданном ограничении на память

QuIP#

- 1) Конкурентный подход векторной квантизации.
- 2) Используется послойное приближение.
- 3) Вместо обучаемой решетки используется решетка для **самой плотной упаковки** сфер в 8-мерном пространстве.
- 4) Перед квантизацией веса поворачиваются с помошью <u>Адамаровых матриц</u>.
- 5) Приводит почти любое распределение весов к нормальному.





Tseng, Albert, et al. "Quip#: Even better Ilm quantization with hadamard incoherence and lattice codebooks." arXiv preprint arXiv:2402.04396 (2024).

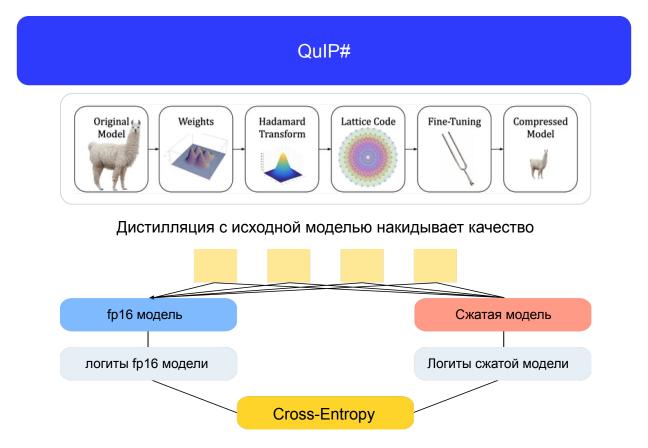
QuIP#



Quantization results on Llama 2 70B. QuIP# achieves near-native performance at 2 bits, outperforming all other presented baselines.

Method	Precision	Wiki ↓	C4 ↓	ArcE ↑	PiQA ↑
Native	16 bit	3.120	5.533	0.597	0.809
OPTQ	3 bit	4.577	6.838	0.544	0.786
OPTQ	2 bit	109.820	62.692	0.253	0.505
QuIP	2 bit	5.574	8.268	0.544	0.751
QuIP#	2 bit	4.159	6.529	0.595	0.786

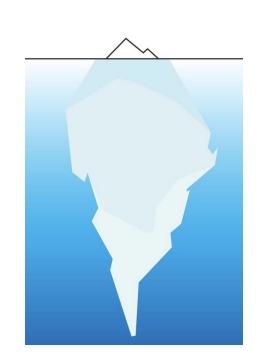
Вполне себе SOTA при 2-битной квантизации



После квантизации большинство параметров дискретные и не могут быть обучаться через GD

Большинство методов оптимизируют end-2-end только непрерывные параметры

Оптимизация дискретных параметров имеет большой потенциал



Непрерывные параметры 10-100M

Дискретные параметры 1-100B

Malinovskii, Vladimir, et al. "Pv-tuning: Beyond straight-through estimation for extreme Ilm compression." *Advances in Neural Information Processing Systems* 37 (2024): 5074-5121.

Algorithm 6 PV algorithm

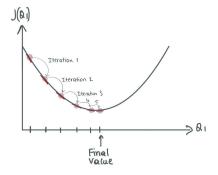
- 1: **Initialization:** starting point $x^0 \in \mathbb{R}^d_{< c}$
- 2: **for** $k = 0, 1, \dots$ **do**
- 3: $y^k = M_P(x^k) := \arg\min_{y \in \mathbb{R}^d} \left\{ \phi(y) : P(y) \supseteq P(x^k) \right\}$ (P step: continuous) 4: $x^{k+1} = M_V(y^k) := \arg\min_{x \in \mathbb{R}^d} \left\{ \phi(x) : V(x) \subseteq V(y^k) \right\}$ (V step: discrete)
- 5: end for

Дискретные и непрерывные параметры можно обучать через альтернированную оптимизацию

P - partition (дискретные)
V - values (непрерывные)

Как оптимизировать непрерывные и дискретные параметры?

Непрерывные



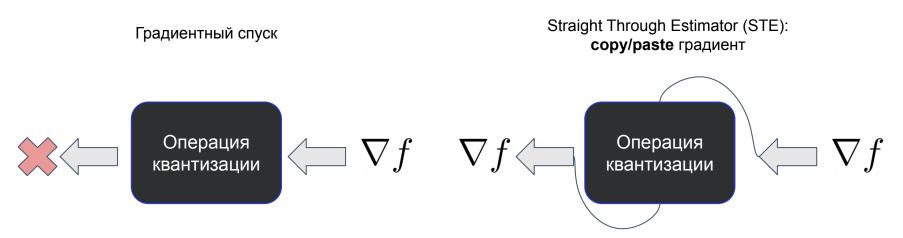
Любой алгоритм градиентного спуска

Дискретные



[1] https://towardsdatascience.com/understanding-gradient-descent-for-machine-learning-246e324c229

Операция квантизации не дифференцируема



Bengio, Yoshua, Nicholas Léonard, and Aaron Courville. "Estimating or propagating gradients through stochastic neurons for conditional computation." arXiv preprint arXiv:1308.3432 (2013).

Nаивный STE нестабилен

Малые шаги



Оптимизация застревает

Большие шаги



Оптимизация улетает в космос

Идея:

Обновлять только небольшую долю (подпространство) параметров с самой **большой нормой** градиента

Это дает устойчивую и быструю сходимость

Algorithm 7 PV-Tuning: Optimization

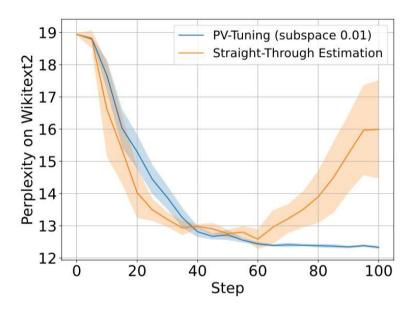
- 1: initial parameters $x^0 \in \mathbb{R}^d_c$,
- 2: objective function $\phi: \mathbb{R}^d \to \mathbb{R}$,
- 3: subspace size $\tau \in [d]$
- 4: **for** k = 0, ..., K 1 **do**
- 5: \triangleright **P** step: update V(x) by backprop
- 6: $y^k = \underset{y \in R_{\leq c}^d}{arg \min} \{ \phi(y) : P(y) \supseteq P(x^k) \}$
- 7: \triangleright **V** step: choose a subspace \mathcal{S}^k & update P(x)
- 8: $S^k = \underset{1 \leq i \leq d}{arg top \tau} |\nabla_i \phi(y^k)| \triangleright \text{find } \tau \text{ largest}$
- 9: $\widehat{\phi}_{y,\mathcal{S}^k}(x) := \left\| x \left(y \frac{1}{L_{S^k}} Z^k \left(\nabla \phi(y) \right) \right) \right\|^2$ 10: $x^{k+\underline{1}} \arg \min_{x} \left\{ \widehat{\phi}_{y^k,\mathcal{S}^k}(x) : V(x) \subseteq V(y^k) \right\}$
 - Z^k subspace projection operator

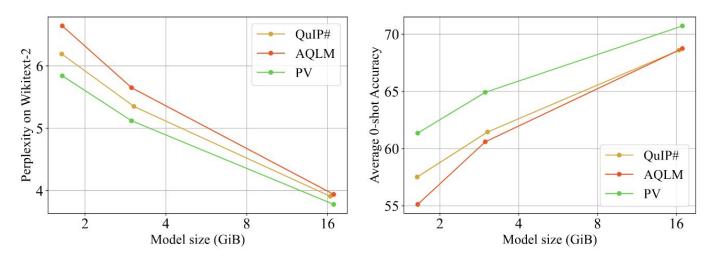
11: end for

Алгоритм

На каждом **V-шаге** (дискретная оптимизация) малая доля кодов обновляется с помощью **STE**

Эффект от ограничения на обновляемые параметры





PV-Tuning позволяет сильно улучшить качество по сравнению с базовым алгоритмом квантизации.

Fine-tuning Method	GPTQ, 2.14 bit/w		VQ , 1.58 bit/w		AQLM, 2.01 bit/w				
I me-turning Metriou	Wiki2↓	C4↓	Acc.↑	Wiki2↓	C4↓	Acc.↑	Wiki2↓	C4↓	Acc.↑
Calibration only (no global fine-tuning)	3290	4125	29.0	20.26	20.09	43.42	7.38	9.34	53.2
Continuous params only [1, 3]	16.77	17.53	46.27	8.17	10.99	52.14	6.69	8.77	56.57
Naive Linearized PV (no subspace)	16.73	17.48	47.68	8.19	10.94	52.08	6.68	8.75	56.51
Stochastic Rounding(tuned) [2]	11.97	13.07	49.79	8.02	10.64	52.31	6.56	8.39	56.68
Straight Through Estimation [4]	8.79	11.04	50.61	7.76	10.26	52.58	6.41	8.63	57.04
Subspace Linearized PV (ours, $ au{=}0.01$)	8.49	10.78	52.17	7.38	9.47	53.36	6.13	8.35	57.81
Subspace Linearized PV+STE ($ au$ =0.01)	8.43	10.82	51.90	7.32	9.35	55.22	5.90	7.43	58.19

PV-Tuning совместим с разными алгоритмами квантизации

HIGGS

- 1) Применяем к матрице весов преобразование Адамара
- 2) Отображаем на оптимальную nмерную сетку для Гауссовораспределенных весов
- 3) Совместимо с data-aware и data-free методами

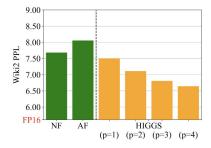


Figure 2: Comparison of Normal Float (NF), Abnormal Float (AF) and HIGGS on Llama 3.1 8B quantization to 3.19-3.25 bitwidth range. HIGGS is instantiated at different lattice dimensionalities p.

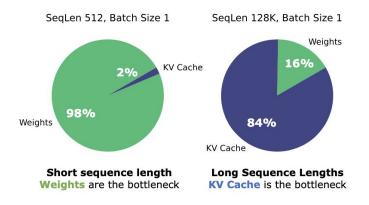
Table 2: WikiText-2 PPL comparison of various 1-shot quantization methods for Llama-2-7b.

FP16		GPTQ	AQLM	QuIP#	QTIP	GPTQ+HIGGS $(p=2)$
	wbits≈2	-	8.180	8.220	6.820	8.637
5.117	wbits≈3	5.776	-	5.600	5.380	5.559
	wbits≈4	5.254	1-1	5.220	5.170	5.213

Malinovskii, Vladimir, et al. "Pushing the limits of large language model quantization via the linearity theorem." arXiv preprint arXiv:2411.17525 (2024).

Сжатие KV-кэшей

Для инференса Causal LLM обыкновенно сохраняют Key, Value для прошлых токенов, чтобы не пересчитывать каждый раз.



Для длинных последовательностей - KV-кэш может занимать памяти больше, чем сама модель. Кроме того, длинные кэши замедляют инференс - основное время уходит на выгрузку кэшей в кэш GPU.

Hooper, Coleman, et al. "Kvquant: Towards 10 million context length Ilm inference with kv cache quantization." *Advances in Neural Information Processing Systems* 37 (2024): 1270-1303.

Сжатие KV-кэшей

Существуют разнообразные решения для уменьшения памяти на хранение KV-кэшей

Архитектурные

Прунинг кэшей

Квантизация

Grouped Query Attention

Multihead Latent
Attention

Native Sparse Attention



H2O

SnapKV

PyramidKV

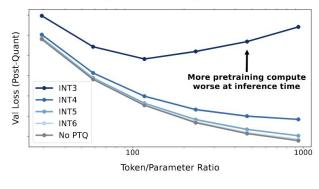
<u>KIVI</u>

KVQuant

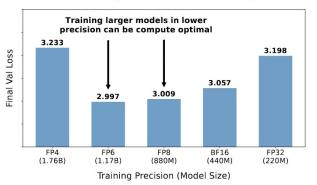
<u>AQUA-KV</u>

Как зависит качество модели от бюджета обучения и битности параметров?

Scaling: Post-Train Quantization



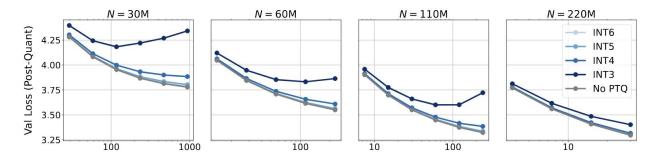
Scaling: Quantized Training



Рассматривают два сценария квантизации:

- 1) Post-Train (квантизация после обучения)
- 2) Quantized Training (квантизация в процессе обучения)

Post-Train

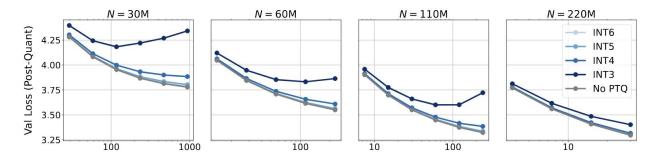


Модели, которых дольше обучали, страдают сильнее от квантизации

$$\delta_{ ext{PTQ}}(N, D, P_{ ext{post}}) = C_T \left(rac{D^{\gamma_D}}{N^{\gamma_N}}
ight) e^{-P_{ ext{post}}/\gamma_{ ext{post}}}$$

Прирост лосса при квантизации описывается неплохо зависимостью выше

Post-Train

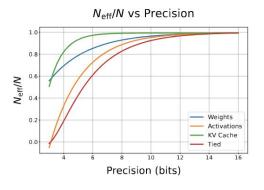


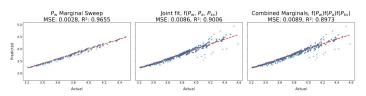
Модели, которых дольше обучали, страдают сильнее от квантизации

$$\delta_{ ext{PTQ}}(N, D, P_{ ext{post}}) = C_T \left(rac{D^{\gamma_D}}{N^{\gamma_N}}
ight) e^{-P_{ ext{post}}/\gamma_{ ext{post}}}$$

Прирост лосса при квантизации описывается неплохо зависимостью выше

Quantized Training



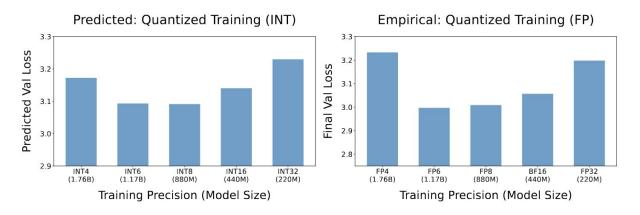


Фитирование коэффициентов на конфигурациях с разным размером модели/бюджетом обучения и точностью

Эффект от квантизации весов/активаций и KV-кэшей можно абсорбировать в "эффективное" количество параметров модели

$$N_{\text{eff}}(P) = N(1 - e^{-P_{\text{w}}/\gamma_{\text{w}}})(1 - e^{-P_{\text{a}}/\gamma_{\text{a}}})(1 - e^{-P_{\text{kv}}/\gamma_{\text{kv}}})$$

Quantized Training



Валидационный лосс при фиксированном размере модели



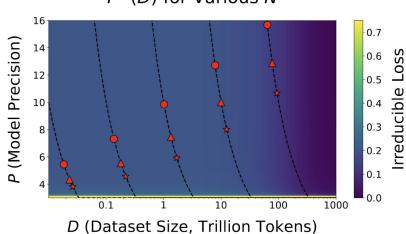
Quantized Training







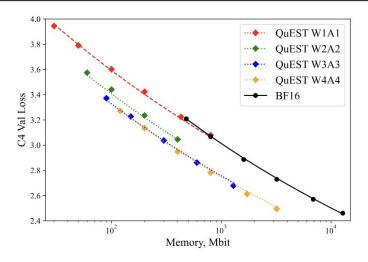




Оптимальная точность при фиксированном размере модели

QuEST

Стабильного обучения можно добиться даже для 1-битных весов и активаций!



Лосс на С4 против размера моделей для разных битностей

QuEST

Но как?

$$\begin{split} \widehat{\mathbf{x}} &= \alpha^* \cdot \mathrm{RMS}(\mathbf{x}) \cdot \left\lfloor \frac{\mathrm{clip}\left(\mathbf{x}/\mathrm{RMS}(\mathbf{x}), \alpha^*\right)}{\alpha^*} \right\rceil = \\ &\coloneqq \mathrm{proj}_{\alpha^*}(\mathbf{x}), \text{ where} \\ \alpha^* &\coloneqq \operatorname*{arg\,min}_{\alpha \in \mathbb{R}} \mathbb{E}_{\xi \sim \mathcal{N}(0,1)} \left\| \xi - \alpha \cdot \left\lfloor \frac{\mathrm{clip}(\xi, \alpha)}{\alpha} \right\rfloor \right\|_2^2 \end{split}$$

Пре-нормализация весов

Algorithm 1 QuEST Training Forward

- 1: Input: Input activations x, row-major weight w
- 2: $\mathbf{x}_h = \mathrm{HT}(\mathbf{x})$
- 3: $\hat{\mathbf{x}}_h = \operatorname{proj}_{\alpha^*} \mathbf{x}_h$
- 4: $\mathbf{w}_h = \mathrm{HT}(\mathbf{w})$
- 5: $\hat{\mathbf{w}}_h = \operatorname{proj}_{\alpha^*} \mathbf{w}_h$
- 6: $\mathbf{y} = \hat{\mathbf{x}}_h \hat{\mathbf{w}}_h^T$
- 7: **Return:** $\hat{\mathbf{y}}, \hat{\mathbf{x}}_h, \hat{\mathbf{w}}_h, M_{\alpha^*}(\mathbf{x}_h; \hat{\mathbf{x}}_h), M_{\alpha^*}(\mathbf{w}_h; \hat{\mathbf{w}}_h)$

"Обработка" Адамаровыми матрицами $\frac{\partial}{\partial \mathbf{x}} \approx \text{IHT}\left(M_{\alpha^*}(\mathbf{x}_h; \hat{\mathbf{x}}_h) \odot \frac{\partial}{\partial \hat{\mathbf{x}}_h}\right).$

Маскирование "шумных" градиентов

QuEST Ho как?

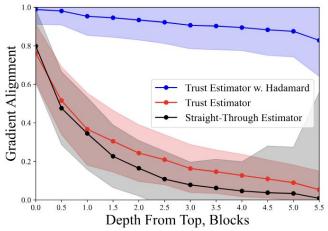


Figure 2. Gradient alignment comparison for a 30M Llama model after training on 2.7B tokens in 8-bit precision.

Маскирование и обработка матрицами важны!

QuEST

Table 1. Fitted scaling-law "effective parameter" multipliers.

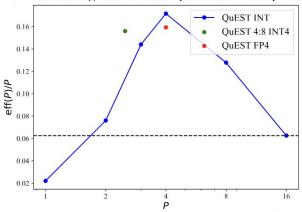
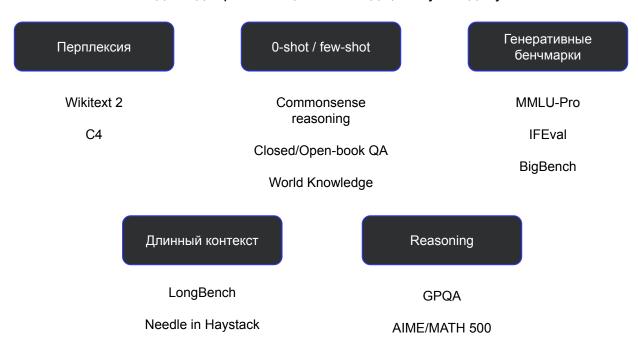


Figure 4. Illustration of the efficiency factors ${\rm eff}(P)/P$, arising from our analysis, for different numerical precisions P and formats (INT, FP, INT+sparse). Higher is better. INT4 appears to have the highest efficiency among hardware-supported formats.

"Эффективный" размер модели при заданной битности

Оценка качества LLM

Удовлетворительного универсального протокола нет. Надо подбирать специально под целевую задачу.



Интеграции в фреймворки



HuggingFace transformers поддерживает следующие форматы квантизации:

- Quanto (Round-To-Nearest)
- BitsAndBytes
- GPTQ
- AWQ
- VPTQ
- AQLM
- HIGGS и др.

https://huggingface.co/docs/transformers/en/main_classes/quantization

Интеграции в фреймворки



vLLM поддерживает следующие форматы квантизации:

- BitsAndBytes
- GPTQ
- AWQ
- AQLM
- HIGGS
- GGUF

https://docs.vllm.ai/en/latest/features/quantization/index.html

Интеграции в фреймворки



GGUF включает в себя множество форматов скалярной и векторной квантизации в различной битности

Очень популярен для локального инференса

Квантизованную модель можно быстро приготовить

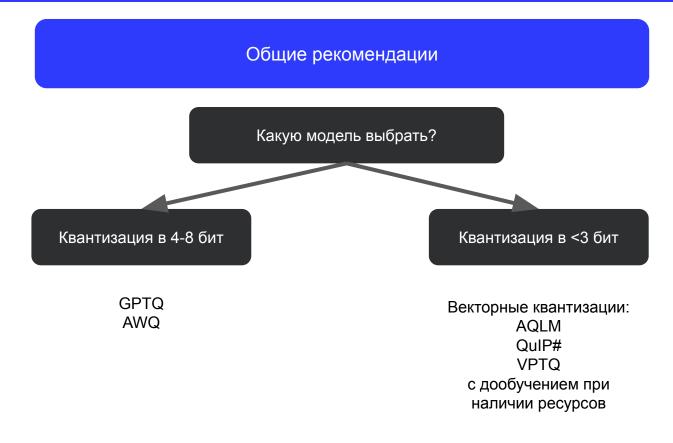
Не самый лучший по качеству

Популярный инференс движок

Поддерживает GGUF

https://huggingface.co/docs/hub/en/gguf

https://ollama.com/



Спасибо за

