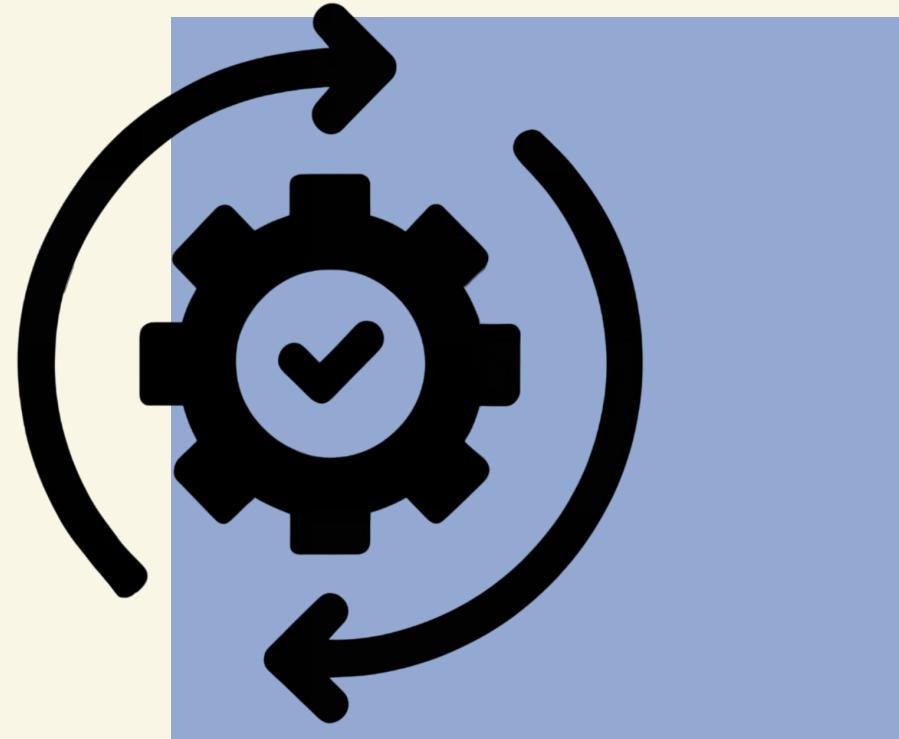


# Эффективные системы машинного обучения



Лекция 2

“Квантизация и  
спарсификация”

Преподаватель

Феоктистов Дмитрий

4 октября 2024

ВМК МГУ

# Квантизация

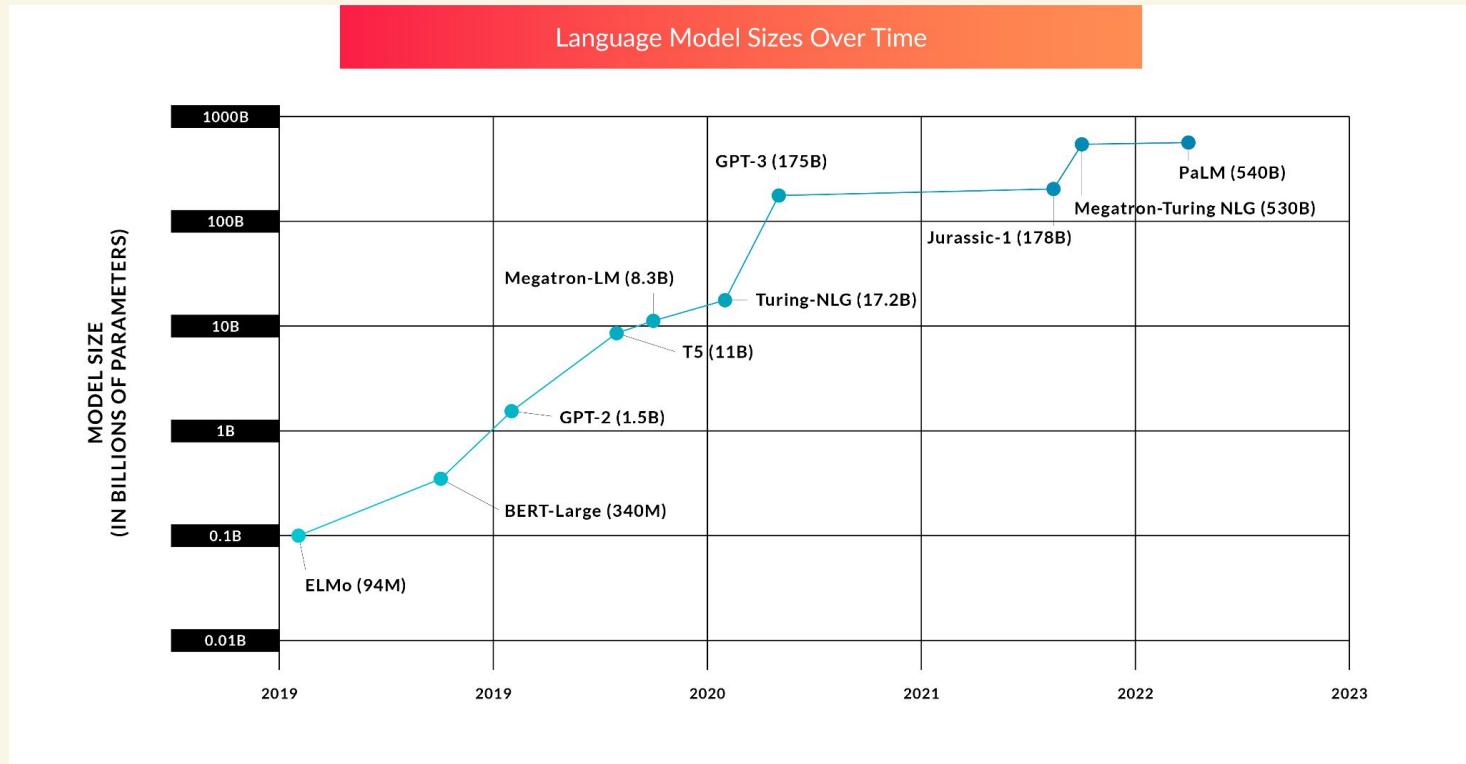


- 1 Мотивация
- 2 Типы данных
- 3 Основы квантизации
- 4 PTQ & QAT
- 5 Квантизация LLM

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Мотивация



# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Мотивация



Nvidia A100 80Gb HBM2E Memory Graphics Card PCIe 4....  
No featured offers available

---

1 option

New  
\$22,217<sup>71</sup>

\$5.54 delivery October 16 - 28.  
[Details](#)  
[... More](#)

Add to Cart

---

Ships from  
Newell Technologies INC  
This item is shipped from United States. [Learn more](#) about import fees and international shipping time.

Sold by  
Newell Technologies INC  
New Seller

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Типы данных

Входной тип	Тип для вычислений	Ускорение в вычислениях	Пропускная способность
FP16	FP16	8x	2x
INT8	INT8	16x	4x
INT4	INT4	32x	8x
INT1	INT1	128x	32x

ARITHMETIC OPERATION	ADD	MUL
8-BIT INTEGER	0.03 pJ	0.2 pJ
16-BIT FLOATING POINT	0.4 pJ	1.1 pJ
32-BIT INTEGER	0.1 pJ	3.1 pJ
32-BIT FLOATING POINT	0.9 pJ	3.7 pJ

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

## Типы данных

Входной тип	Тип для вычислений	Ускорение в вычислениях	Пропускная способность
FP16	FP16	8x	2x
INT8	INT8	16x	4x
INT4	INT4	32x	8x
INT1	INT1	128x	32x

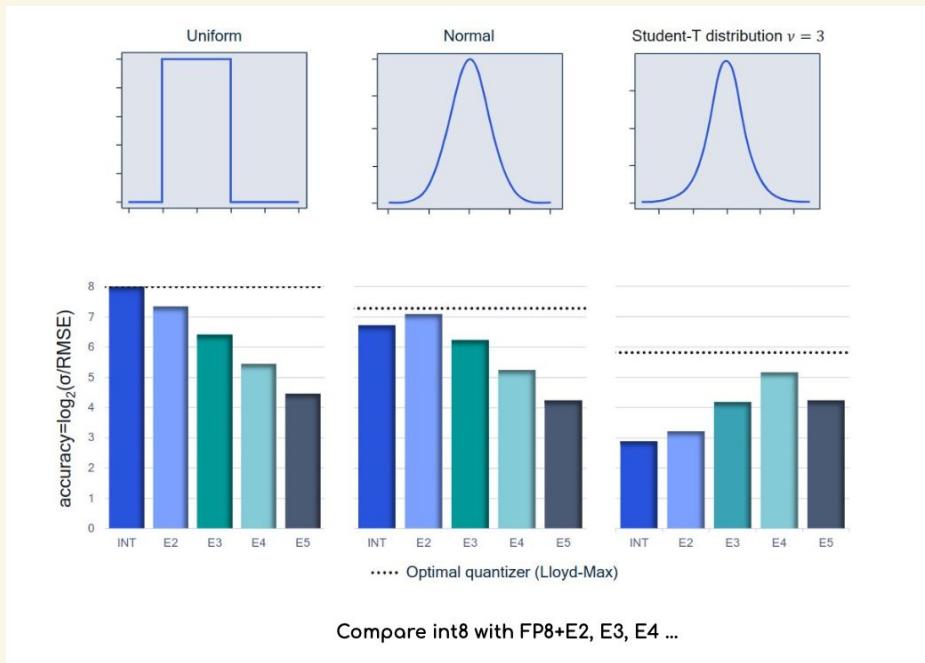


# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Типы данных

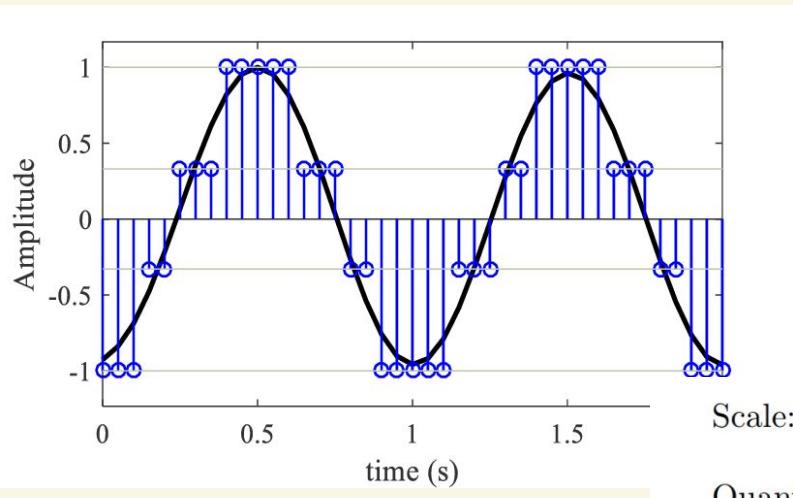
Инты “тупее” => дешевле в транзисторах  
Флоты точнее => дороже в транзисторах  
Вдруг инты не так плохи?



# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Квантизация



$$\text{Scale: } s = \frac{f_{max} - f_{min}}{fq_{max} - fq_{min}}$$

$$\text{Quantize: } Q(x) = \text{clip}(\text{round}(\frac{1}{s}x), fq_{min}, fq_{max})$$

$$\text{De-Quantize: } x = sQ(x)$$

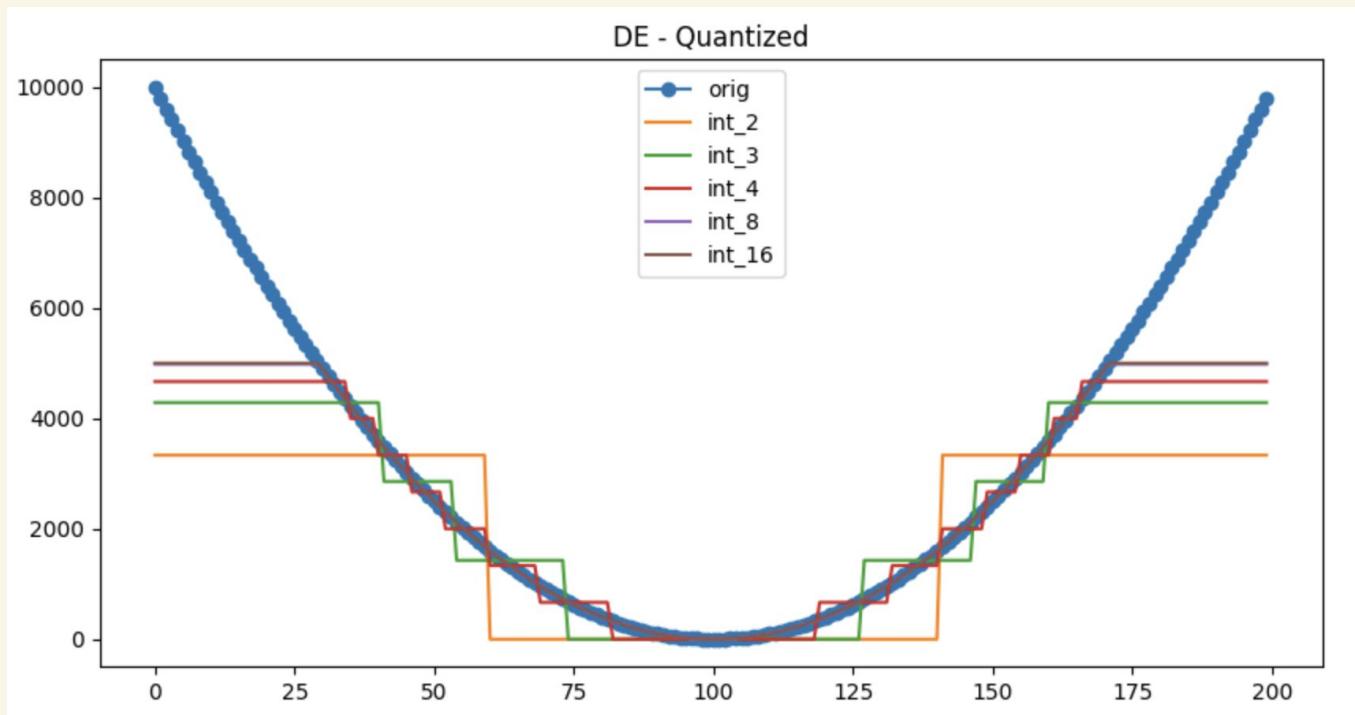
$$fq_{min} = -2^{(bit-1)}$$

$$fq_{max} = 2^{(bit-1)} - 1$$

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Квантизация

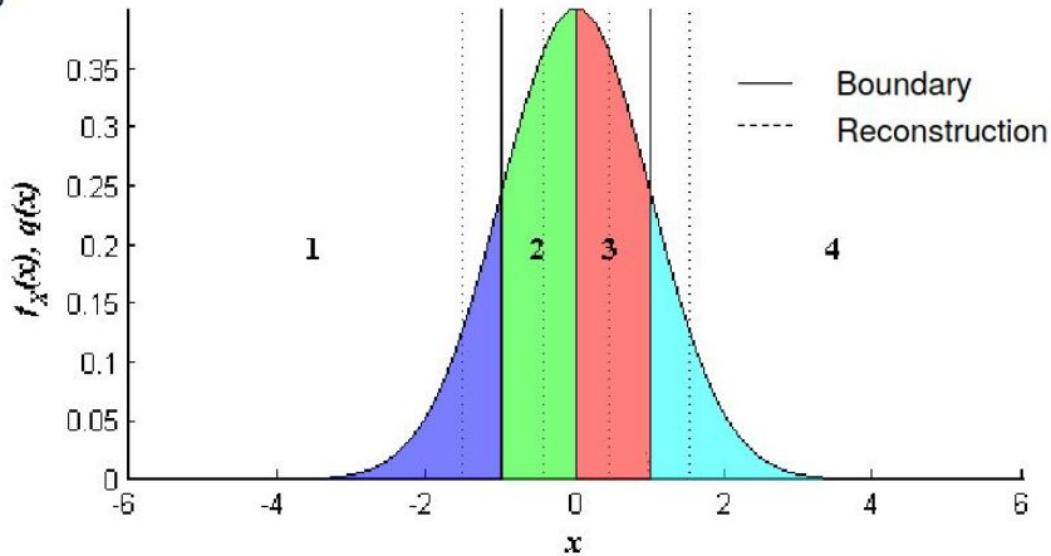


# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Квантизация: квантильная

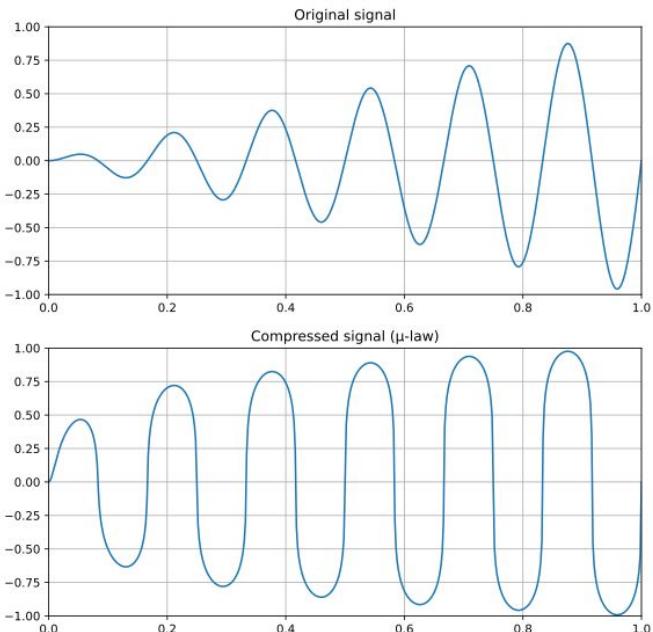
Decision thresholds -0.98, 0, 0.98  
Representative levels -1.51, -0.45, 0.45, 1.51  
 $D^*=0.12=9.30 \text{ dB}$



# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Квантизация: companding

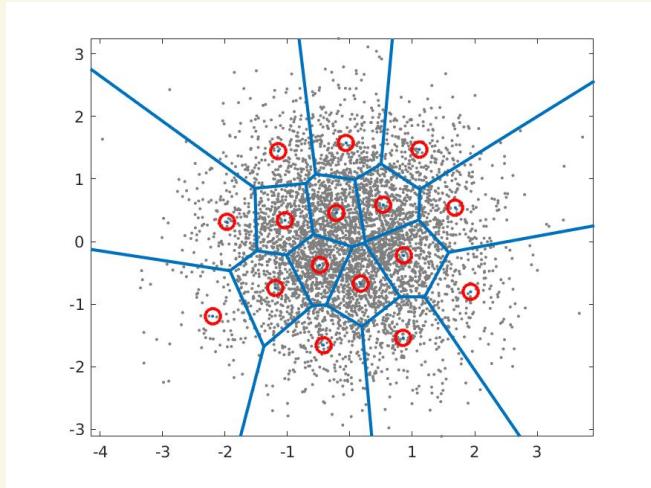


Преобразование должно быть  
дешевым, мы все же хотим  
быстрый инференс

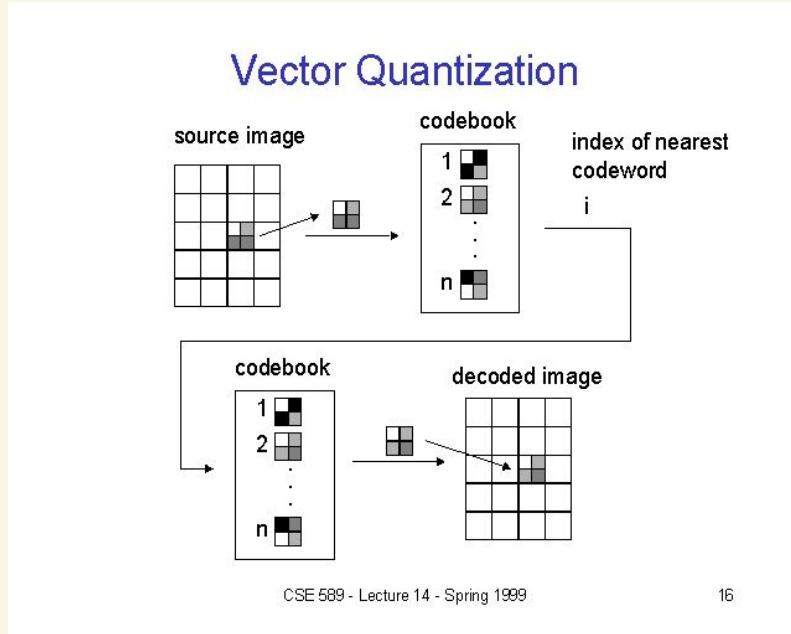
# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Квантизация: векторная



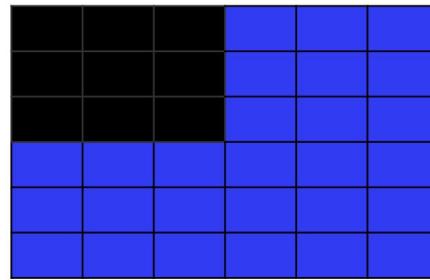
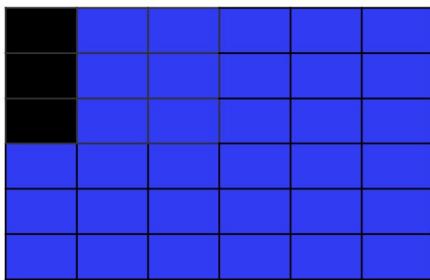
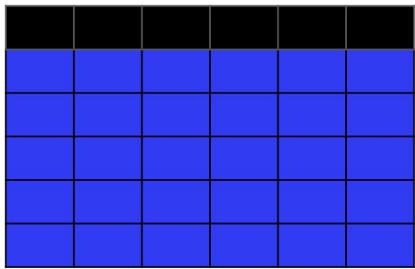
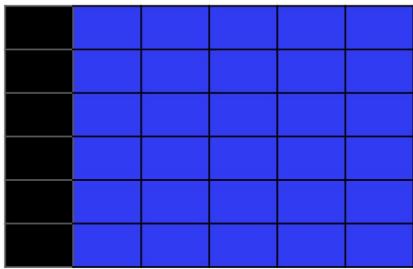
+ иерархичность



# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

## Гранулярность



# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# PQT&QAT

## Post Training Quantization (PTQ)

1. Train a model
2. Quantize weights
3. Finetune model BN statistics or other FP blocks

### PROS:

- Can be applied for already trained models
- Does not necessarily require access to the data or to all the data

### CONS:

- Usually has lower quality than QAT

## Quantization Aware Training (QAT)

1. Quantize weights
2. Train a model with quantized weights

### PROS:

- Almost lossless quality for some datasets

### CONS:

- Requires approximation of non-differentiable quantization function
- Usually, requires training a model from scratch
- Need access to the data

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

## QAT: STE

$$QW = Q(W)$$

$$y = F(x, QW)$$

$$W := W - \alpha \frac{\partial loss}{\partial QW}$$

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# QAT: DiffQ

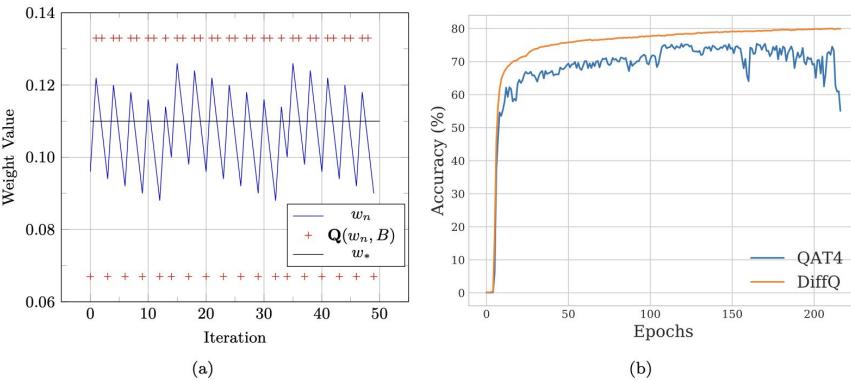


Figure 1: (a) Using STE and SGD to optimize the 1D least-mean-square problem given by equation 4 (with  $B = 4$  and  $X = 1$  a.s.).  $Q(w_n, B)$  oscillates between the quantized value just above ( $w_+$ ) and just under ( $w_-$ ) the unquantized ground truth  $w_*$ , while  $w_n$  oscillates around the boundary  $(w_+ + w_-)/2$ . (b) Model accuracy vs. epochs for ImageNet using EfficientNet-b3. Results are presented for both QAT over 4 bits and DIFFQ.

$$QW_b = Q(W, b)$$

$$QN = QW_b - W \approx \frac{\Delta_b}{2} \mathbf{U}[-1, 1]$$

$$Q\hat{W}_b = W + \frac{\Delta_b}{2} \mathbf{U}[-1, 1]$$

$$y = F(x, W + \frac{\Delta_b}{2} \mathbf{U}[-1, 1])$$

# LLM Quantization: OBQ

## 2 Optimal Brain Surgeon

In deriving our method we begin, as do Le Cun, Denker and Solla [1990], by considering a network trained to a local minimum in error. The functional Taylor series of the error with respect to weights (or parameters, see below) is:

$$\delta E = \left( \frac{\partial E}{\partial \mathbf{w}} \right)^T \cdot \delta \mathbf{w} + \frac{1}{2} \delta \mathbf{w}^T \cdot \mathbf{H} \cdot \delta \mathbf{w} + O(\|\delta \mathbf{w}\|^3) \quad (1)$$

where  $\mathbf{H} \equiv \partial^2 E / \partial \mathbf{w}^2$  is the Hessian matrix (containing all second order derivatives) and the superscript T denotes vector transpose. For a network trained to a local minimum in error, the first (linear) term vanishes; we also ignore the third and all higher order terms. Our goal is then to set one of the weights to zero (which we call  $w_q$ ) to minimize the increase in error given by Eq. 1. Eliminating  $w_q$  is expressed as:

$$\delta \mathbf{w}_q + w_q = 0 \quad \text{or more generally} \quad \mathbf{e}_q^T \cdot \delta \mathbf{w} + w_q = 0 \quad (2)$$

where  $\mathbf{e}_q$  is the unit vector in weight space corresponding to (scalar) weight  $w_q$ . Our goal is then to solve:

$$\text{Min}_q \{ \text{Min}_{\delta \mathbf{w}} \{ \frac{1}{2} \delta \mathbf{w}^T \cdot \mathbf{H} \cdot \delta \mathbf{w} \} \text{ such that } \mathbf{e}_q^T \cdot \delta \mathbf{w} + w_q = 0 \} \quad (3)$$

To solve Eq. 3 we form a Lagrangian from Eqs. 1 and 2:

$$L = \frac{1}{2} \delta \mathbf{w}^T \cdot \mathbf{H} \cdot \delta \mathbf{w} + \lambda (\mathbf{e}_q^T \cdot \delta \mathbf{w} + w_q) \quad (4)$$

where  $\lambda$  is a Lagrange undetermined multiplier. We take functional derivatives, employ the constraints of Eq. 2, and use matrix inversion to find that the optimal weight change and resulting change in error are:

$$\delta \mathbf{w} = - \frac{w_q}{[\mathbf{H}^{-1}]_{qq}} \mathbf{H}^{-1} \cdot \mathbf{e}_q \quad \text{and} \quad L_q = \frac{1}{2} \frac{w_q^2}{[\mathbf{H}^{-1}]_{qq}} \quad (5)$$

# LLM Quantization: OBQ

## Optimal Brain Surgeon procedure

1. Train a “reasonably large” network to minimum error.
2. Compute  $H^{-1}$ .
3. Find the  $q$  that gives the smallest saliency  $L_q = w_q^2 / (2[H^{-1}]_{qq})$ . If this candidate error increase is much smaller than  $E$ , then the  $q^{\text{th}}$  weight should be deleted, and we proceed to step 4; otherwise go to step 5. (Other stopping criteria can be used too.)
4. Use the  $q$  from step 3 to update *all* weights (Eq. 5). Go to step 2.
5. No more weights can be deleted without large increase in  $E$ . (At this point it may be desirable to retrain the network.)

# LLM Quantization: OBQ

**The Quantization Order and Update Derivations.** Under the standard assumption that the gradient at the current point  $\mathbf{w}$  is negligible, the OBS formulas for the optimal weight to be pruned  $w_p$  and the corresponding update  $\delta_p$  can be derived by writing the locally quadratic problem under the constraint that element  $p$  of  $\delta_p$  is equal to  $-w_p$ , which means that  $w_p$  is zero after applying the update to  $\mathbf{w}$ . This problem has the following Lagrangian:

$$L(\delta_p, \lambda) = \delta_p^\top \mathbf{H} \delta_p + \lambda(\mathbf{e}_p^\top \delta_p - (-w_p)), \quad (6)$$

where  $\mathbf{H}$  denotes the Hessian at  $\mathbf{w}$  and  $\mathbf{e}_p$  is the  $p$ th canonical basis vector. The optimal solution is then derived by first finding the optimal solution to  $\delta_p$  via setting the derivative  $\partial L / \partial \delta_p$  to zero and then substituting this solution back into  $L$  and solving for  $\lambda$ ; please see e.g. [13, 39] for examples.

Assume a setting in which we are looking to quantize the weights in a layer on a fixed grid of width  $\Delta$  while minimizing the loss. To map OBS to a *quantized* projection, we can set the target of the Lagrangian constraint in (6) to  $(\text{quant}(w_p) - w_p)$ , where  $\text{quant}(w_p)$  is the weight rounding given by quantization; then  $w_p = \text{quant}(w_p)$  after the update.

Assuming we wish to quantize weights iteratively, one-at-a-time, we can derive formulas for the “optimal” weight to quantize at a step, in terms of minimizing the loss increase, and for the corresponding optimal update to the unquantized weights, in similar fashion as discussed above:

$$w_p = \underset{w_p}{\operatorname{argmin}} \frac{(\text{quant}(w_p) - w_p)^2}{[\mathbf{H}^{-1}]_{pp}}, \quad \delta_p = -\frac{w_p - \text{quant}(w_p)}{[\mathbf{H}^{-1}]_{pp}} \cdot \mathbf{H}_{:,p}^{-1}. \quad (7)$$

In fact, since  $-w_p$  is a constant during all derivations, we can just substitute it with  $(\text{quant}(w_p) - w_p)$  in the final result. We note that the resulting formulas are a generalization of standard OBS for pruning, if  $\text{quant}(\cdot)$  always “quantizes” a weight to 0, then we recover the original form.

# LLM Quantization: OBQ

This approximation approach has been shown to work well in many applications [51, 19, 9].

We follow these conventions as well, and work with the specific layer-wise compression problem stated formally below, where the weights  $\mathbf{W}_\ell$  are a  $d_{\text{row}} \times d_{\text{col}}$  matrix (for conv-layers  $d_{\text{col}}$  corresponds to the total number of weights in a single filter), and the input  $\mathbf{X}_\ell$  has dimensions  $d_{\text{col}} \times N$ .

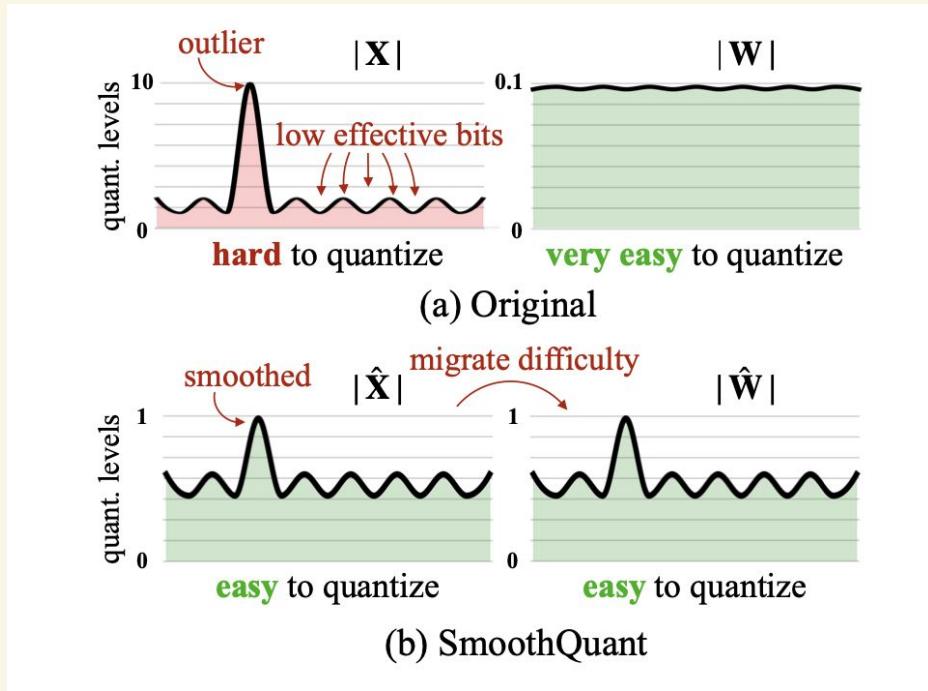
$$\operatorname{argmin}_{\widehat{\mathbf{W}}_\ell} \quad \|\mathbf{W}_\ell \mathbf{X}_\ell - \widehat{\mathbf{W}}_\ell \mathbf{X}_\ell\|_2^2 \quad \text{s.t.} \quad \mathcal{C}(\widehat{\mathbf{W}}_\ell) > C. \quad (2)$$

- + улучшения и трюки, много  
трюков

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# LLM Quantization: outliers



# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# LLM Quantization: outliers

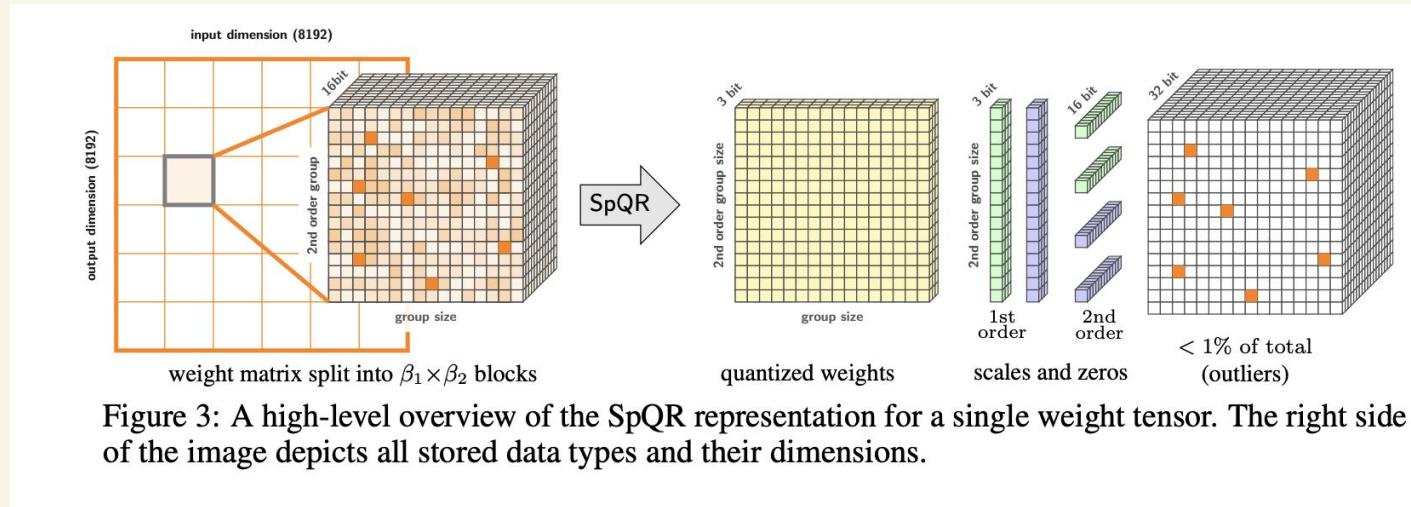


Figure 3: A high-level overview of the SpQR representation for a single weight tensor. The right side of the image depicts all stored data types and their dimensions.

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# LLM Quantization: scaling laws

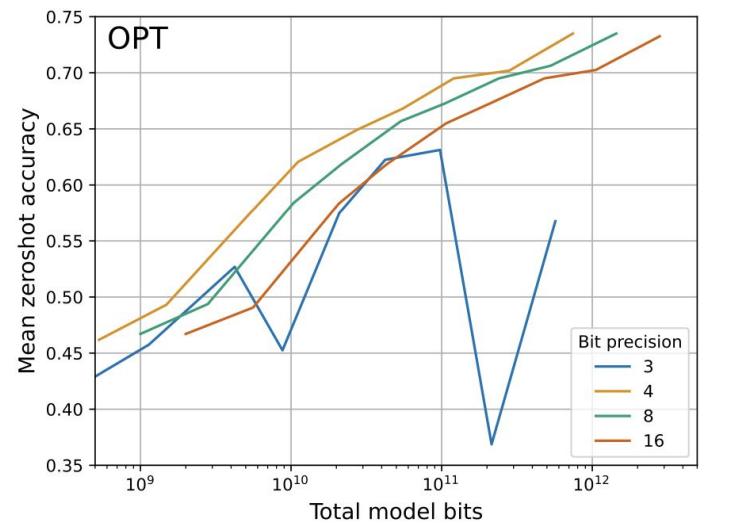


Figure 1. Bit-level scaling laws for mean zero-shot performance across four datasets for 125M to 176B parameter OPT models. Zero-shot performance increases steadily for fixed model bits as we reduce the quantization precision from 16 to 4 bits. At 3-bits, this relationship reverses, making 4-bit precision optimal.

We run more than 35,000 experiments with 16-bit inputs and k-bit parameters to examine which zero-shot quantization methods improve scaling for 3 to 8-bit precision at scales of 19M to 176B parameters across the LLM families BLOOM, OPT, NeoX/Pythia, and GPT-2. We find that it is challenging to improve the bit-level scaling trade-off, with the only improvements being the use of a small block size – splitting the parameters into small independently quantized blocks – and the quantization data type being used

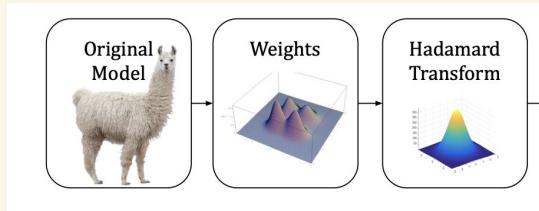
# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# AQLM vs QUIP#

## Algorithm 1 AQLM: Additive Quantization for LLMs

```
Require: model, data
1:  $\mathbf{X}_{block} := \text{model}.input\_embeddings(\text{data})$ 
2: for  $i = 1, \dots, \text{model}.num\_layers$  do
3:    $\text{block} := \text{model}.get\_block(i)$ 
4:    $\mathbf{Y}_{block} := \text{block}(\mathbf{X}_{block})$ 
5:   for  $\text{layer} \in \text{linear\_layers}(\text{block})$  do
6:      $\mathbf{W} := \text{layer}.weight$ 
7:      $\mathbf{X} := \text{layer}_\text{inputs}(\text{layer}, \mathbf{X}_{block})$ 
8:      $C, b, s := \text{initialize}(\mathbf{W})$  // k-means
9:     while loss improves by at least  $\tau$  do
10:       $C, s := \text{train\_Cs\_adam}(\mathbf{X}\mathbf{X}^T, \mathbf{W}, C, b, s)$ 
11:       $b := \text{beam\_search}(\mathbf{X}\mathbf{X}^T, \mathbf{W}, C, b, s)$ 
12:    end while
13:    /* save for fine-tuning */
14:     $\text{layer}.weight := \text{AQLMFormat}(C, b, s)$ 
15:  end for
16:   $\theta := \text{trainable\_parameters}(\text{block})$ 
17:  while loss improves by at least  $\tau$  do
18:     $L := \|\text{block}(\mathbf{X}_{block}) - \mathbf{Y}_{block}\|_2^2$ 
19:     $\theta := \text{adam}(\theta, \frac{\partial L}{\partial \theta})$ 
20:  end while
21:   $\mathbf{Y}_{block} := \text{block}(\mathbf{X}_{block})$ 
22: end for
```



**Definition 1.** We say a symmetric Hessian matrix  $H \in \mathbb{R}^{n \times n}$  is  $\mu$ -incoherent if it has an eigendecomposition  $H = Q\Lambda Q^T$  such that for all  $i$  and  $j$ ,  $|Q_{ij}| = |e_i^T Q e_j| \leq \mu/\sqrt{n}$ . By extension, we say a weight matrix  $W \in \mathbb{R}^{m \times n}$  is  $\mu$ -incoherent if all  $i$  and  $j$ ,  $|W_{ij}| = |e_i^T W e_j| \leq \mu \|W\|_F / \sqrt{mn}$ .

## Algorithm 3 QuIP: Quantization with Incoherence Processing

```
Require:  $b \in \mathbb{N}, H \in \mathbb{R}^{n \times n}$  SPD,  $W \in \mathbb{R}^{m \times n}, \mathcal{Q} \in \{\text{Near}, \text{Stoch}\}, \rho \in \mathbb{R}_+, \alpha \in [0, 1]$ 
1:  $\hat{W}, H, s, \tilde{D} \leftarrow \text{Alg 1}(b, H, W, \rho, \alpha)$   $\triangleright$  QuIP Incoherence Pre-Processing
2:  $H = (\hat{U} + I)D(\hat{U} + I)^{-1}$   $\triangleright$  LDL decomposition
3: for  $k \in \{1, \dots, n\}$  do  $\hat{W}_k \leftarrow \text{clamp}(\mathcal{Q}(W_k + (W - \hat{W})\hat{U}_k), 0, 2^b - 1)$   $\triangleright$  LDLQ
4: return  $\hat{W} \leftarrow \text{Alg 2}(b, H, \hat{W}, s, \tilde{D})$   $\triangleright$  QuIP Incoherence Post-Processing
```

Январь 2024

Extreme Compression of Large Language Models via Additive Quantization, Egiazarian et al, ICML, 2024

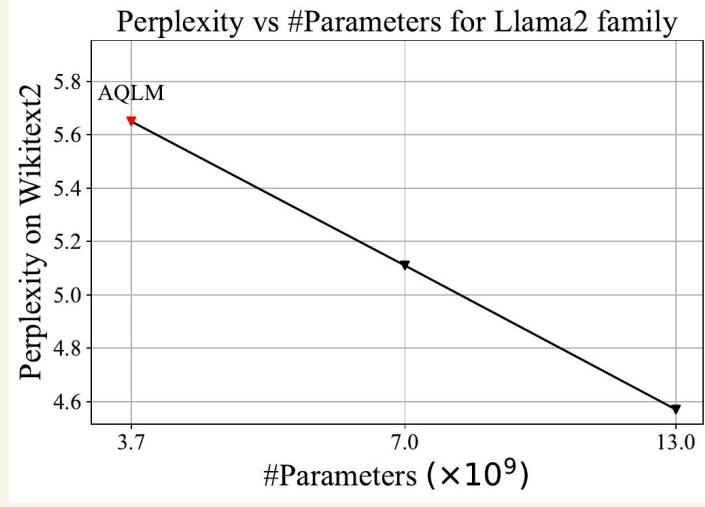
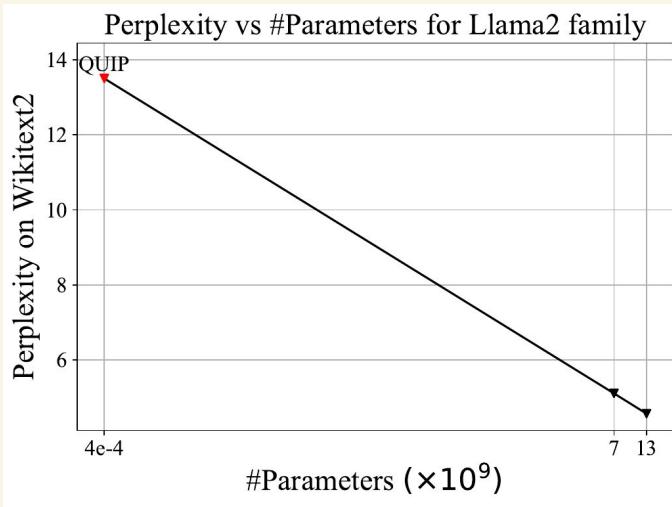
QUIP: 2-Bit Quantization of Large Language Models With Guarantees, Jerry Chee et al, NeurIPS, 2023

Середина 2023

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

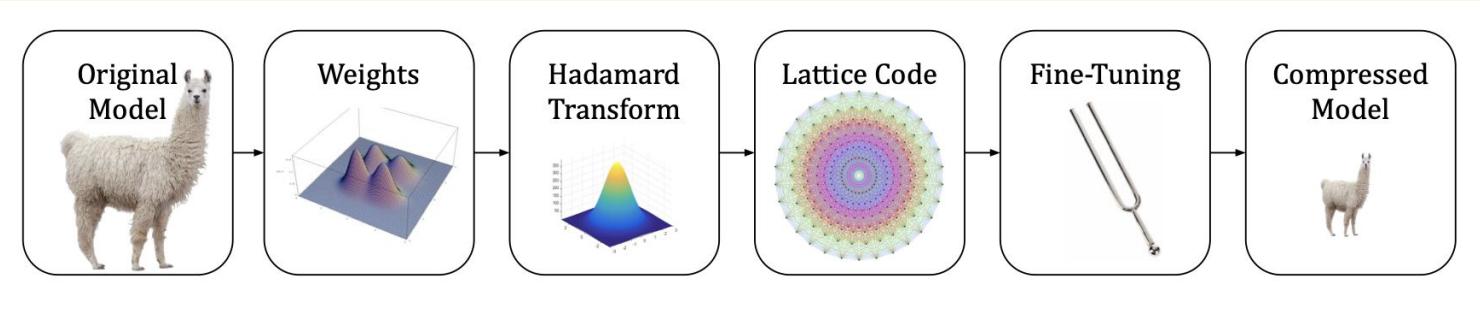
## AQLM vs QUIP#



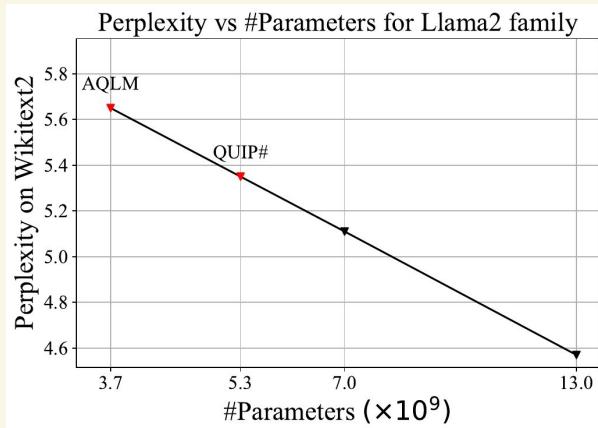
# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

## AQLM vs QUIP#



Февраль 2024



# AQLM vs QUIP#

---

## Algorithm 1 PV algorithm

---

```

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

```

---



---

## Algorithm 2 PV-Tuning: Optimization

---

**Require:** initial parameters  $x^0 \in \mathbb{R}_c^d$ ,  
objective function  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  
subspace size  $\tau \in [d]$

```

1: for  $k = 0, \dots, K - 1$  do
2:   ▷ P step: update  $V(x)$  by backprop
3:    $y^k = \arg \min_{y \in \mathbb{R}_{\leq c}^d} \{\phi(y) : P(y) \supseteq P(x^k)\}$ 
4:   ▷ V step: choose a subspace  $S^k$  & update  $P(x)$ 
5:    $S^k = \arg \max_{1 \leq i \leq d} |\nabla_i \phi(y^k)|$  ▷ find  $\tau$  largest
6:    $\hat{\phi}_{y, S^k}(x) := \|x - \left(y - \frac{1}{L_{S^k}} Z^k (\nabla \phi(y))\right)\|^2$ 
7:    $x^{k+1} = \arg \min_x \{\hat{\phi}_{y^k, S^k}(x) : V(x) \subseteq V(y^k)\}$ 
8: end for

```

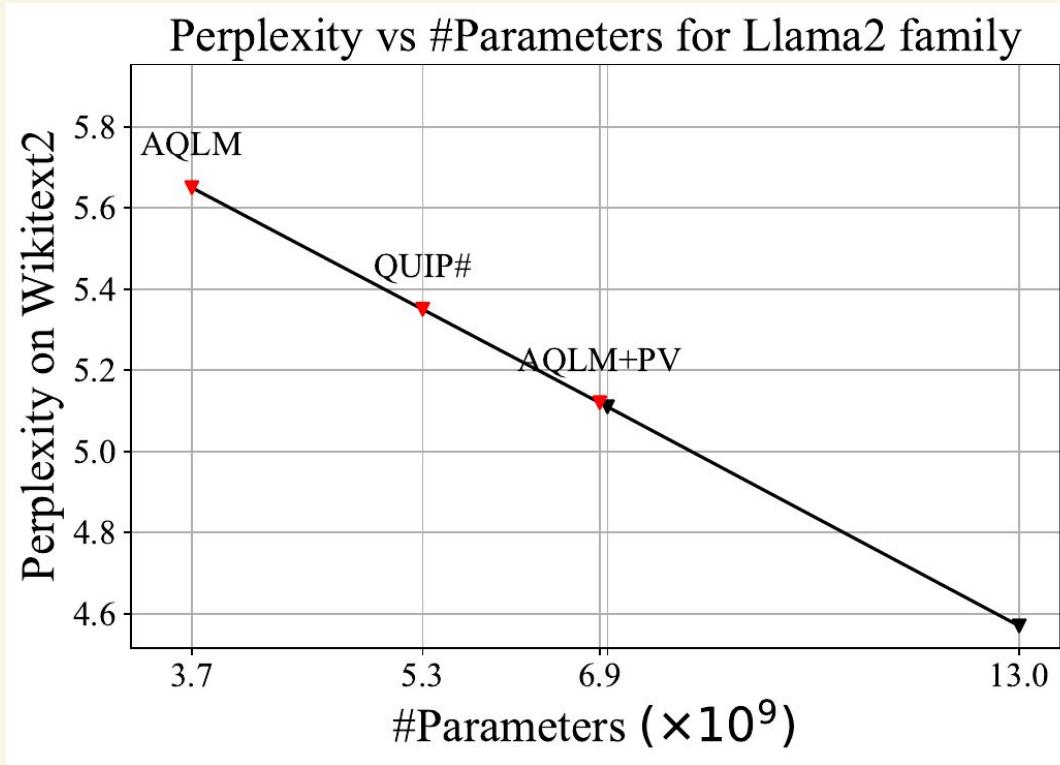
---

Май 2024

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

## AQLM vs QUIP#



# AQLM vs QUIP#

## AQLM

Dettmers & Zettlemoyer (2022) examined the accuracy-compression trade-offs of these early methods, suggesting that 4-bit quantization may be optimal for RTN quantization, and observing that data-aware methods like GPTQ allow for higher compression, i.e. strictly below 4 bits/weight, maintaining Pareto optimality. Our work brings this Pareto frontier below 3 bits/weight, for the first time. Parallel

## AQLM+PV

**Pareto-optimality.** A key practical question concerns obtaining optimal quality for the target model size, where a smaller model compressed to 3-4 bits often dominates a larger model compressed to 1-bit. The best known Pareto-optimal bit-width for Llama 2 is 2.5 [19]: compressing a larger model to less than 2.5 bits per weight is inferior to a smaller model quantized to the same total number of bytes. From this perspective, **PV-tuning pushes the Pareto-optimal frontier for LLAMA 2 to 2.0 bits.** This is easiest to see in Table 8: a 2-bit 13B model outperforms any 7B quantization and is comparable with the 16-bit 7B model. The same holds for the 2-bit 70B model.

# Прунинг

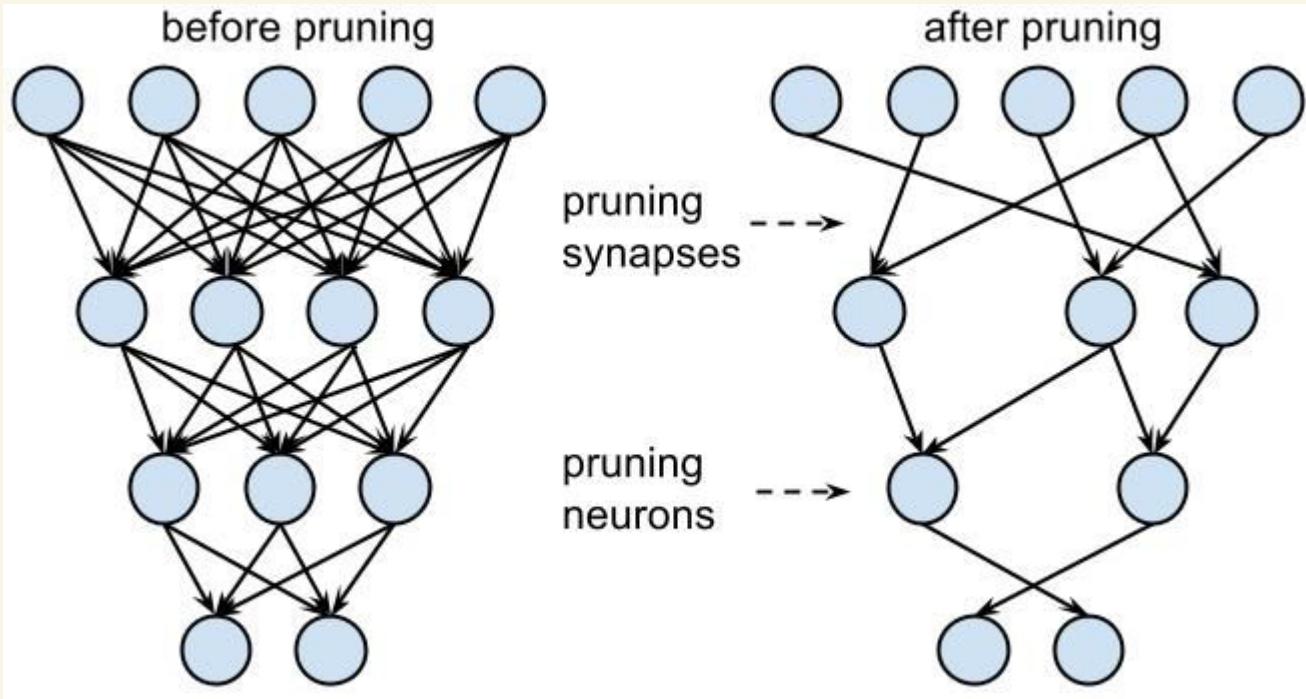


- 1 TPF, PAI, PTP, TTP
- 2 Structured & Unstructured
- 3 LLM pruning
- 4 Pruning vs Quantization

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

## Общая идея



# Общая идея

It is well documented that neural networks have an excess of parameters needed to generalize well and make accurate predictions. “*The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks*” (Frankle and Carbin, 2019) demonstrates that **neural networks tend to have a specific subset of parameters that are essential for prediction**. Model pruning is a very intuitive approach to model compression, with the hypothesis that effective model compression should **remove weights that aren’t being used**, similarly to how the brain reduces usage of connections between neurons to emphasize important pathways.

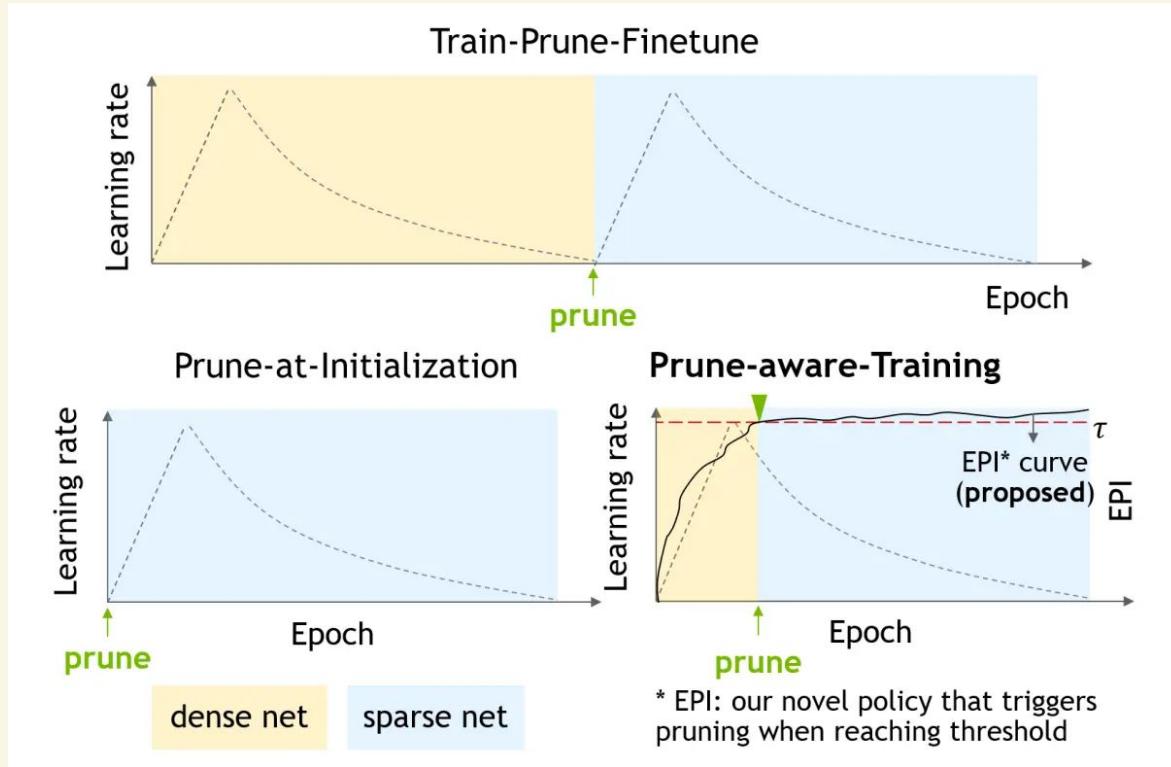
# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация

Феоктистов Дмитрий

04 октября 2024

# TPF, PAI, PTP, PAT



# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# TPF, PAI, PTP, PAT



## PRUNING APPROACH

### TRAIN-TIME PRUNING

### POST-TRAINING PRUNING

## PROS

More efficient models since they are trained with the objective of sparsity in mind

Simpler to implement

Pruning decisions are made alongside the consideration of model parameter optimization

Pruning parameters can be easily adjusted based on inference requirements

## CONS

Makes training process more complex

May require fine-tuning to regain performance in the case of accuracy degradation

Change in pruning parameters may require retraining of the whole model

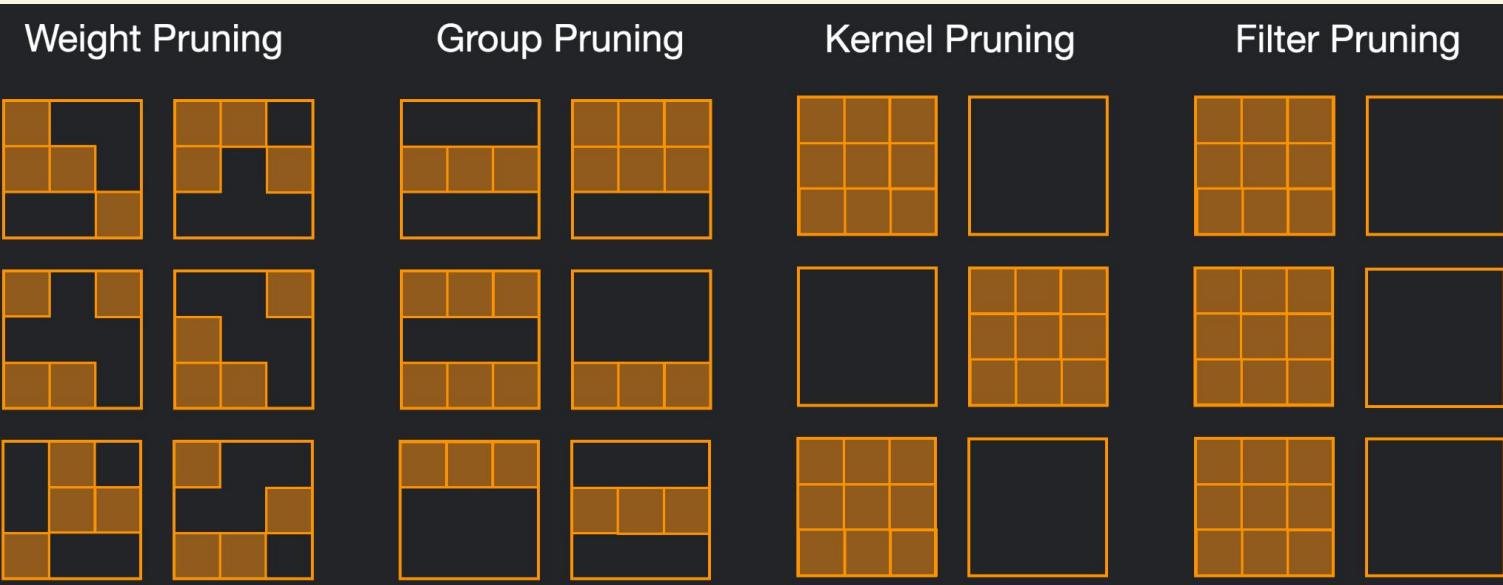
Pruning decisions are user-defined and may not be optimal

## TRADEOFFS OF DIFFERENT PRUNING APPROACHES

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Structured vs Unstructured



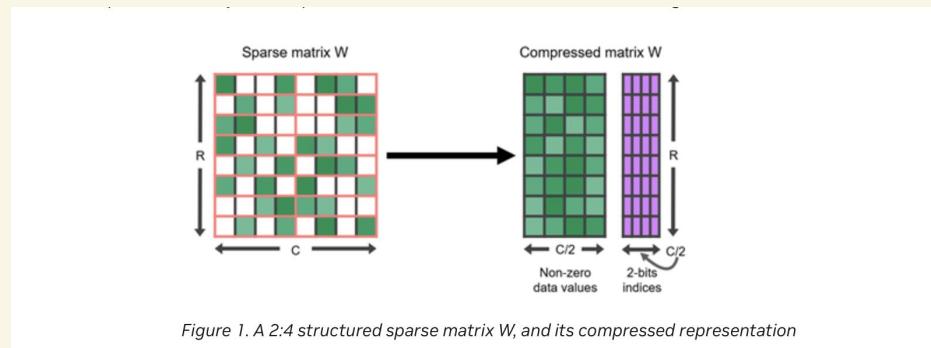
# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Structured vs Unstructured

Input Operands	Accumulator	Dense TOPS	vs. FFMA	Sparse TOPS	vs. FFMA
FP32	FP32	19.5	-	-	-
TF32	FP32	156	8X	<b>312</b>	<b>16X</b>
FP16	FP32	312	16X	<b>624</b>	<b>32X</b>
BF16	FP32	312	16X	<b>624</b>	<b>32X</b>
FP16	FP16	312	16X	<b>624</b>	<b>32X</b>
INT8	INT32	624	32X	<b>1248</b>	<b>64X</b>

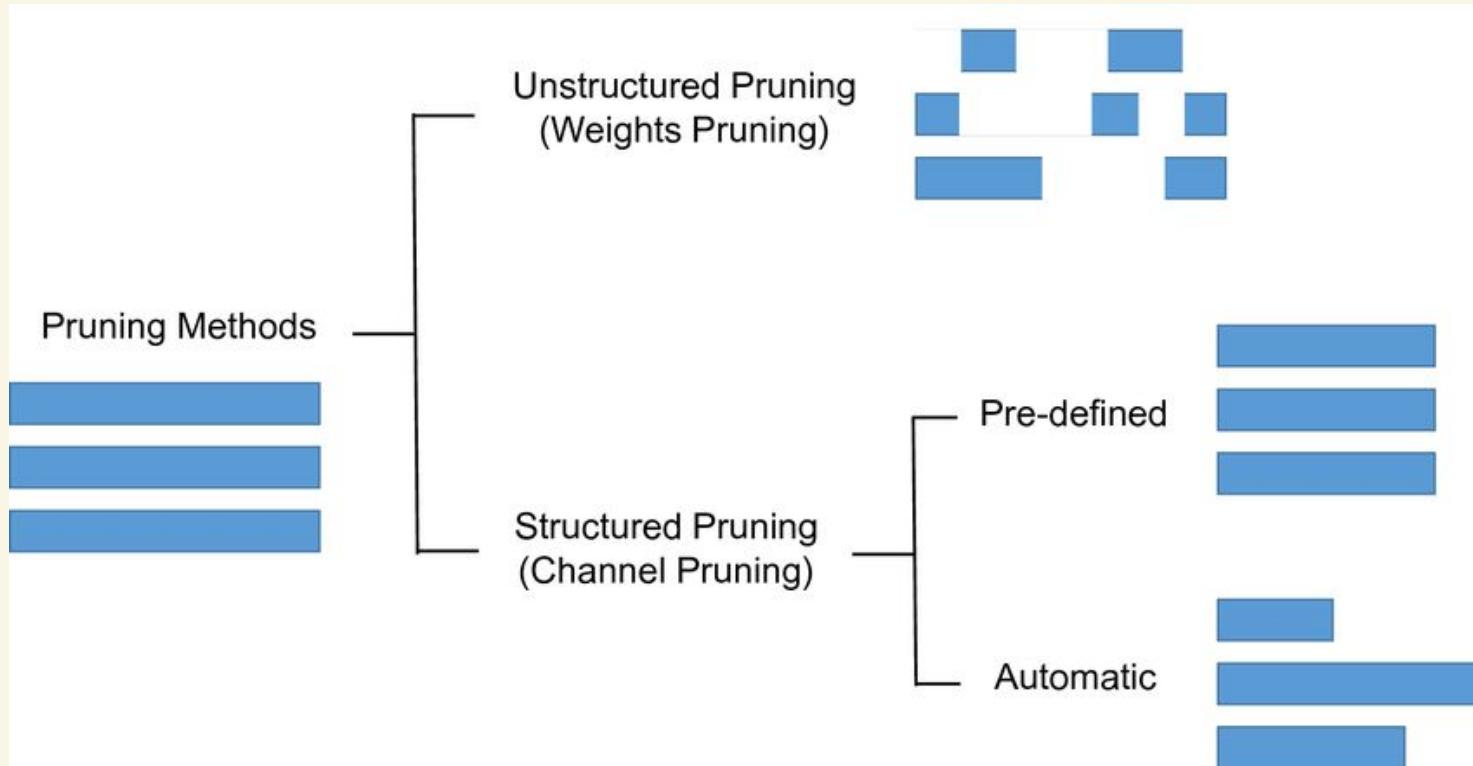
Table 1. Performance of Sparse Tensor Cores in the NVIDIA Ampere Architecture.



# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

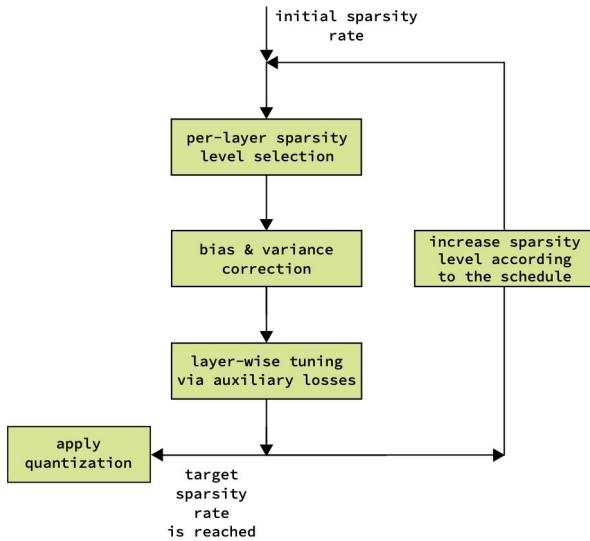
# Structured vs Unstructured



# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Unstructured pruning



$$L = \sum_{i \in batch} (Y_{dense}^i - f(W^s M^s, X_{dense}^i) - b^s)^2$$
$$Y_{dense}^i = f(W_{dense}, X_{dense}^i) \quad (6)$$

$$s_t = s_f + (s_i - s_f) \left(1 - \frac{t}{T}\right)^3$$

$$I(w_i^l) = \frac{|w_i^l|}{\sqrt{\sum_{j \in l} |w_j^l|^2}}$$

$$W_{corr}^s = \lambda W^s + E(W_{dense}) - E(\lambda W^s)$$

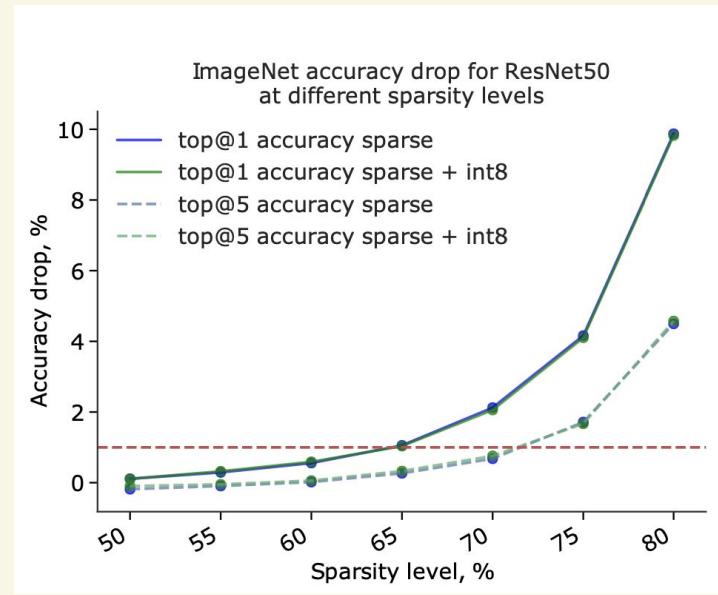
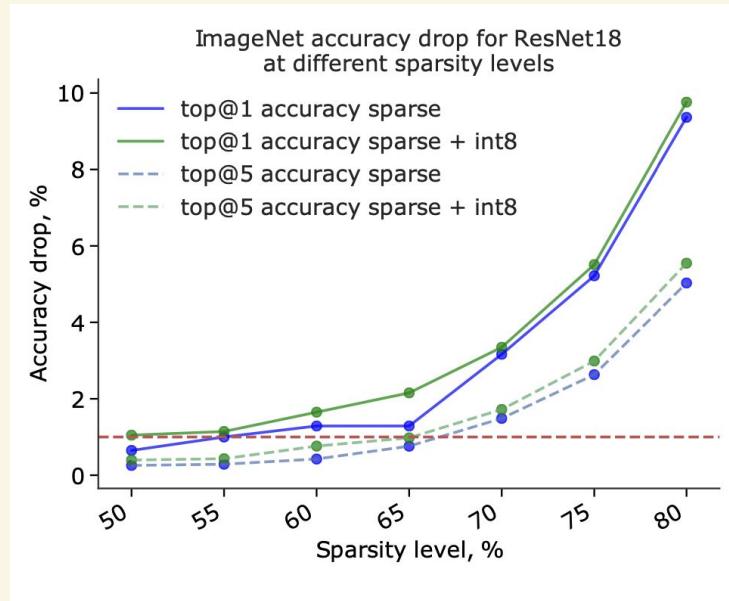
$$\lambda = \frac{\sigma(W_{dense})}{\sigma(W^s) + \epsilon}$$

$$b_{corr} = b_{dense} + E(f(W_{dense}, X_{dense})) - E(f(W_{corr}^s, X_{dense})) \quad (5)$$

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Unstructured pruning



# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Unstructured pruning

Table 1. Accuracy values of the sparse 8-bit quantized DNN models obtained with the proposed post-training method. Metric values were measured on a CPU. The same is for other accuracy values reported in the paper unless specified otherwise.

Model	Dataset (acc. metric)	Sparsity rate, %	Compressed model acc.	Absolute acc. drop
ResNet50	ImageNet (top@1 acc.)	65	75.09	1.04
ResNet18	ImageNet (top@1 acc.)	50	68.93	0.81
GoogleNetV4	ImageNet (top@1 acc.)	50	78.96	0.94
MobileNetV2	ImageNet (top@1 acc.)	40	70.29	1.51
MobileNetV1-SSD	VOC07 (mAP)	50	71.53	0.98
TinyYOLOv2	COCO (AP)	50	28.29	0.83
NCF	MovieLens 20M (hit ratio)	70	64.67	0.93
BERT-base	MRPC (acc.)	50	82.50	0.63

# Structured pruning of LLMs

While previous methods have effectively maintained model performance amidst parameter reduction, they primarily target compression within specialized domains or for designated tasks in the context of task-specific compression. Although this paradigm could potentially be employed for LLM compression, it compromises the LLM's capacity as a versatile task solver, rendering it suited to a single task exclusively. Thus, we strive to compress the LLM in a new setting: to reduce the LLM's size while preserving its diverse capabilities as general-purpose task solvers

- **The size of the training corpus of the LLM is enormous.** Previous compression methods heavily depend on the training corpus. The LLM has escalated the corpus scale to 1 trillion tokens or more [17, 51]. The extensive storage needs and protracted transmission times make the dataset difficult to acquire. Furthermore, if the dataset is proprietary, acquisition of the training corpus verges on impossibility, a situation encountered in [74, 39].
- **The unacceptably long duration for the post-training of the pruned LLM.** Existing methods require a substantial amount of time for post-training the smaller model [55, 28]. For instance, the general distillation in TinyBERT takes around 14 GPU days [20]. Even post-training a task-specific compressed model of BERT demands around 33 hours [61, 22]. As the size of both the model and corpus for LLMs increases rapidly, this step will invariably consume an even more extensive time.

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Structured pruning of LLMs

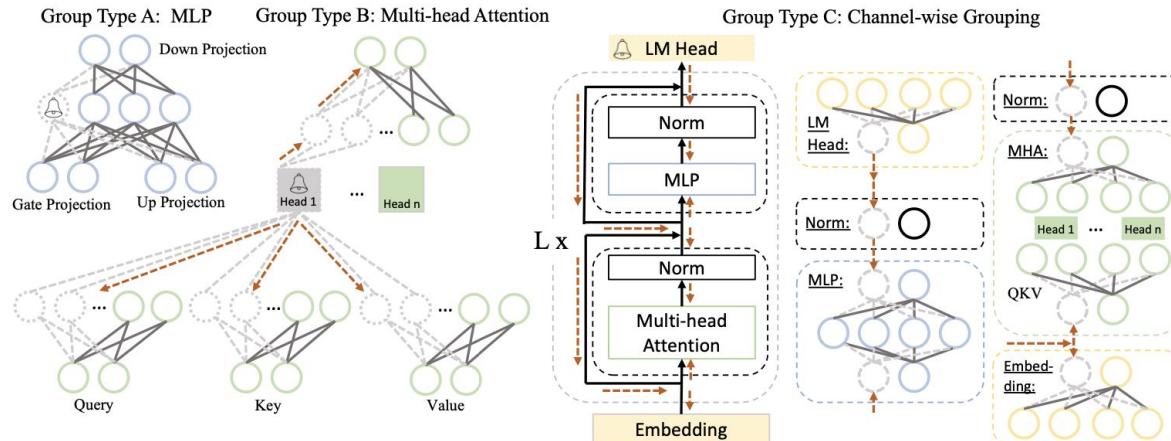


Figure 2: Illustration of the coupled structures in LLaMA. We simplify the neurons in each layer to make the dependent group clear. The trigger neuron, marked as a circle with a bell, cause weights with dependency pruned (dashed lines), which may propagate (red dashed lines) to coupled neurons (dashed circles). A group can be triggered by a variety of trigger neurons. Taking Group Type B as an example, the trigger for this group involves (i) the attention head, (ii) the output neuron in Query, Key or Value, and (iii) the input neuron in the final output projection.

# Structured pruning of LLMs

**Vector-wise Importance.** Suppose that given a dataset  $\mathcal{D} = \{x_i, y_i\}_{i=1}^N$ , where  $N$  is the number of samples. In our experiments, we set  $N$  equal to 10 and we use some public datasets as the source of  $\mathcal{D}$ . A group (as previously defined as a set of coupled structures) can be defined as  $\mathcal{G} = \{W_i\}_{i=1}^M$ , where  $M$  is the number of coupled structures in one group and  $W_i$  is the weight for each structure. While pruning, our goal is to remove the group that has the least impact on the model's prediction, which can be indicated by the deviation in the loss. Specially, to estimate the importance of  $W_i$ , the change in loss can be formulated as [24]:

$$I_{W_i} = |\Delta \mathcal{L}(\mathcal{D})| = |\mathcal{L}_{W_i}(\mathcal{D}) - \mathcal{L}_{W_i=0}(\mathcal{D})| = \left| \underbrace{\frac{\partial \mathcal{L}^\top(\mathcal{D})}{\partial W_i}}_{\neq 0} W_i - \frac{1}{2} W_i^\top H W_i + \mathcal{O}(\|W_i\|^3) \right| \quad (3)$$

**Element-wise Importance.** The above can be considered as an estimate for the weight  $W_i$ . We can derive another measure of importance at a finer granularity, where each parameter within  $W_i$  is assessed for its significance:

$$I_{W_i^k} = |\Delta \mathcal{L}(\mathcal{D})| = |\mathcal{L}_{W_i^k}(\mathcal{D}) - \mathcal{L}_{W_i^k=0}(\mathcal{D})| = \left| \frac{\partial \mathcal{L}(\mathcal{D})}{\partial W_i^k} W_i^k - \frac{1}{2} W_i^k H_{kk} W_i^k + \mathcal{O}(\|W_i^k\|^3) \right| \quad (4)$$

Here,  $k$  represents the  $k$ -th parameter in  $W_i$ . The diagonal of the hessian  $H_{kk}$  can be approximated by the Fisher information matrix, and the importance can be defined as:

$$I_{W_i^k} = |\mathcal{L}_{W_i^k}(\mathcal{D}) - \mathcal{L}_{W_i^k=0}(\mathcal{D})| \approx \left| \frac{\partial \mathcal{L}(\mathcal{D})}{\partial W_i^k} W_i^k - \frac{1}{2} \sum_{j=1}^N \left( \frac{\partial \mathcal{L}(\mathcal{D}_j)}{\partial W_i^k} W_i^k \right)^2 + \mathcal{O}(\|W_i^k\|^3) \right| \quad (5)$$

# Structured pruning of LLMs

**Group Importance.** By utilizing either  $I_{W_i^k}$  or  $I_{W_i}$ , we estimate the importance at the granularity of either a parameter or a vector of weight. Remembering that our goal is to estimate the importance of  $\mathcal{G}$ , we aggregate the importance scores in four ways: (i) Summation:  $I_{\mathcal{G}} = \sum_{i=1}^M I_{W_i}$  or  $I_{\mathcal{G}} = \sum_{i=1}^M \sum_k I_{W_i^k}$ , (ii) Production:  $I_{\mathcal{G}} = \prod_{i=1}^M I_{W_i}$  or  $I_{\mathcal{G}} = \prod_{i=1}^M \sum_k I_{W_i^k}$ , (iii) Max:  $I_{\mathcal{G}} = \max_{i=1}^M I_{W_i}$  or  $I_{\mathcal{G}} = \max_{i=1}^M \sum_k I_{W_i^k}$ ; (iv) Last-Only: Since deleting the last executing structure in a dependency group is equivalent to erasing all the computed results within that group, we assign the importance of the last executing structure as the importance of the group:  $I_{\mathcal{G}} = I_{W_l}$  or  $I_{\mathcal{G}} = \sum_k I_{W_l^k}$ , where  $l$  is the last structure. After assessing the importance of each group, we rank the importance of each group and prune the groups with lower importance based on a predefined pruning ratio.

# Structured pruning of LLMs

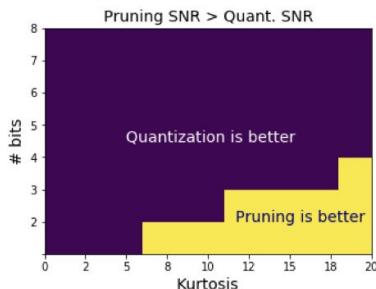
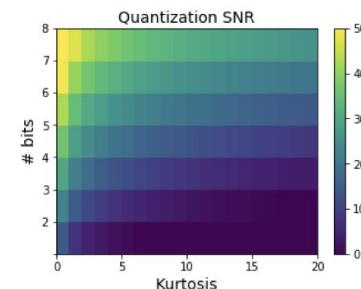
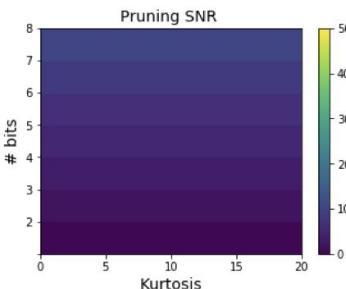
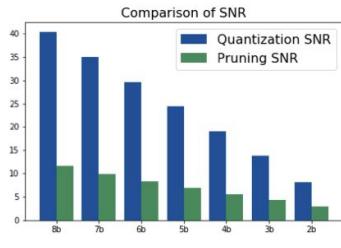
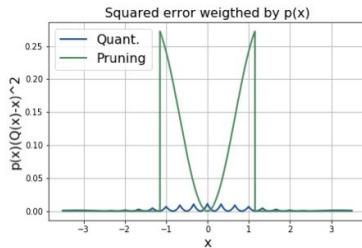
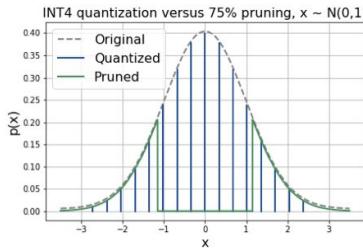
Table 1: Zero-shot performance of the compressed LLaMA-7B. The average is calculated among seven classification datasets. ‘Underline’ indicates the best pruning-only performance, while ‘bold’ represents the overall best performance with the same pruning ratio, considering both pruning and post-training. The ‘Channel’ strategy only prunes the dependent group of Type C, while all other methods employ the ‘Block’ strategy to prune dependent groups in both Type A and Type B. Since [51] did not provide its prompt, the evaluation of the result with \* is performed under different prompts, which is lower than the official results.

Pruning Ratio	Method	WikiText2 $\downarrow$	PTB $\downarrow$	BoolQ	PIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	Average
Ratio = 0%	LLaMA-7B[51]	-	-	76.5	79.8	76.1	70.1	72.8	47.6	57.2	68.59
	LLaMA-7B*	12.62	22.14	73.18	78.35	72.99	67.01	67.45	41.38	42.40	63.25
Ratio = 20% w/o tune	L2	582.41	1022.17	59.66	58.00	37.04	52.41	33.12	28.58	29.80	42.65
	Random	27.51	43.19	61.83	71.33	56.26	54.46	57.07	32.85	35.00	52.69
Ratio = 20% w/ tune	Channel	74.63	153.75	62.75	62.73	41.40	51.07	41.38	27.90	30.40	45.38
	Vector	22.28	41.78	<u>61.44</u>	71.71	57.27	54.22	55.77	33.96	38.40	53.25
	Element <sup>2</sup>	19.77	36.66	59.39	<u>75.57</u>	65.34	<u>61.33</u>	59.18	<u>37.12</u>	39.80	<u>56.82</u>
	Element <sup>1</sup>	<u>19.09</u>	<u>34.21</u>	57.06	<u>75.68</u>	<u>66.80</u>	59.83	<u>60.94</u>	36.52	<b>40.00</b>	56.69
Ratio = 20% w/ tune	Channel	22.02	38.67	59.08	73.39	64.02	60.54	57.95	35.58	38.40	55.57
	Vector	18.84	33.05	65.75	74.70	64.52	59.35	60.65	36.26	39.40	57.23
	Element <sup>2</sup>	<b>17.37</b>	30.39	<b>69.54</b>	76.44	68.11	<b>65.11</b>	63.43	<b>37.88</b>	<b>40.00</b>	<b>60.07</b>
	Element <sup>1</sup>	17.58	<b>30.11</b>	64.62	<b>77.20</b>	<b>68.80</b>	63.14	<b>64.31</b>	36.77	39.80	59.23

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Quantization vs Pruning



# Quantization vs Pruning

$$\min_{\mathbf{w}} E(\mathbf{w}) = \|\mathbf{X}\delta\mathbf{w} - \mathbf{X}\mathbf{w}_{orig}\|_2^2$$

s.t.  $\mathbf{w} \in \mathbb{Z}^n$ ,

$$w_{min} \leq w_i \leq w_{max},$$

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{P} \mathbf{w} - \mathbf{q}^T \mathbf{w},$$

s.t.  $\mathbf{w} \in \mathbb{Z}^n$ ,

$$w_{min} \leq w_i \leq w_{max},$$

$$L(\boldsymbol{\lambda}) = \max -\gamma,$$

$$\text{s.t. } \begin{bmatrix} \mathbf{P} - \text{diag}(\boldsymbol{\lambda}) & q + \frac{1}{2}\boldsymbol{\lambda} \\ (q + \frac{1}{2}\boldsymbol{\lambda})^T & \gamma \end{bmatrix} \succeq 0,$$

$$\boldsymbol{\lambda} \geq 0,$$

$$E = \min_{\hat{\mathbf{w}}} \|\mathbf{X}\hat{\mathbf{w}} - \mathbf{X}\mathbf{w}_{orig}\|_2^2$$

s.t.  $\|\hat{\mathbf{w}}\|_0 \leq s$ ,

$$E(\mathbf{w}) = \min_{\mathbf{w}, \mathbf{m}} \|\mathbf{X}(\mathbf{m} \odot \mathbf{w}) - \mathbf{X}\mathbf{w}_{orig}\|_2^2,$$

s.t.  $\|\mathbf{m}\|_1 = s$ ,

$$-\mathbf{m} \odot l \leq \hat{\mathbf{w}} \leq \mathbf{m} \odot u$$

$$l, u > 0, m_i \in \{0, 1\},$$

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Quantization vs Pruning

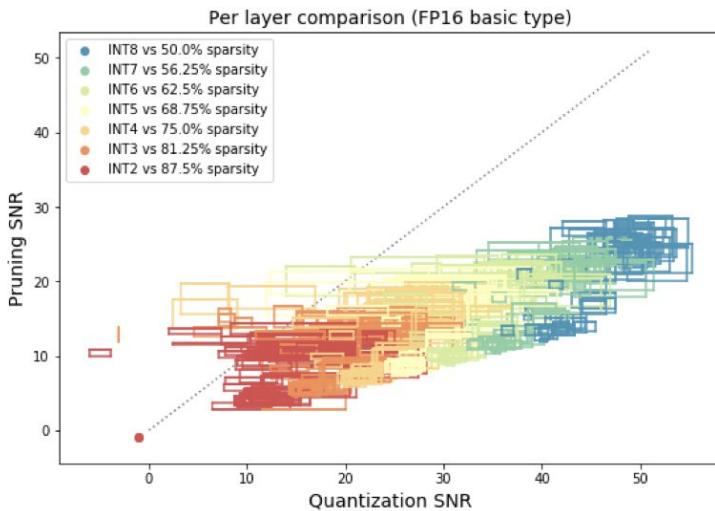


Figure 4: Comparison in the post-training scenario. Each box corresponds to a subset of one of 10 layers from the 4 different models that were used, with 7 different bit-width comparison points. The ranges of the box indicate the lower and higher-bounds found by the algorithms.

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Quantization vs Pruning

Model	Orig.	Metric	Method	8b	7b	6b	5b	4b	3b	2b
Resnet-18	69.7	acc.	quant. pruning	<b>70.5</b> 70.3	<b>70.5</b> 70.1	<b>70.6</b> 69.9	<b>70.3</b> 69.5	<b>70.0</b> 69.3	<b>68.9</b> 68.3	<b>67.3</b> 66.8
Resnet-50	76.1	acc.	quant. pruning	76.4 <b>76.6</b>	76.4 <b>76.4</b>	76.4 76.2	76.3 76.1	76.2 75.9	75.5 75.4	72.3 <b>74.3</b>
MobileNet-V2	71.7	acc.	quant. pruning	<b>71.9</b> 68.1	<b>72.0</b> 65.6	<b>71.7</b> 61.9	<b>71.6</b> 56.3	<b>70.9</b> 48.0	<b>68.6</b> 34.0	<b>59.1</b> 21.2
EfficientNet	75.4	acc.	quant. pruning	<b>75.2</b> 72.5	<b>75.3</b> 70.9	<b>75.0</b> 68.1	<b>74.6</b> 63.6	<b>74.0</b> 56.4	<b>71.5</b> 44.5	<b>60.9</b> 27.1
MobileNet-V3	67.4	acc.	quant. pruning	<b>67.7</b> 65.6	<b>67.6</b> 64.4	<b>67.1</b> 62.4	<b>66.3</b> 60.2	<b>64.7</b> 56.1	<b>60.8</b> 31.7	<b>50.5</b> 0.0
ViT	81.3	acc.	quant. pruning	<b>81.5</b> 76.6	<b>81.4</b> 76.6	<b>81.4</b> 76.2	<b>81.0</b> 73.1	<b>80.4</b> 72.4	<b>78.4</b> 71.5	<b>72.2</b> 69.4
DeepLab-V3	72.9	mIoU	quant. pruning	<b>72.3</b> 65.2	<b>72.3</b> 62.8	<b>72.4</b> 56.8	<b>71.9</b> 47.7	<b>70.8</b> 32.9	<b>63.2</b> 18.6	<b>17.6</b> 10.0
EfficientDet	40.2	mAP	quant. pruning	<b>39.6</b> 34.5	<b>39.6</b> 33.0	<b>39.6</b> 30.9	<b>39.2</b> 27.9	<b>37.8</b> 24.2	<b>33.5</b> 17.9	<b>15.5</b> 8.0
OPT-350m	14.8	perpl.	quant. pruning	<b>14.8</b> 18.0	<b>14.8</b> 19.7	<b>14.9</b> 22.6	<b>15.0</b> 27.2	<b>15.3</b> 35.4	<b>15.9</b> 53.5	<b>19.9</b> 101.4

Table 1: Comparison of QAT and magnitude pruning with fine-tuning given equal model size and equal number of epochs for fine-tuning.

# Эффективные системы машинного обучения, ВМК МГУ

Занятие 1. Квантизация и спарсификация  
Феоктистов Дмитрий  
04 октября 2024

# Quantization vs Pruning

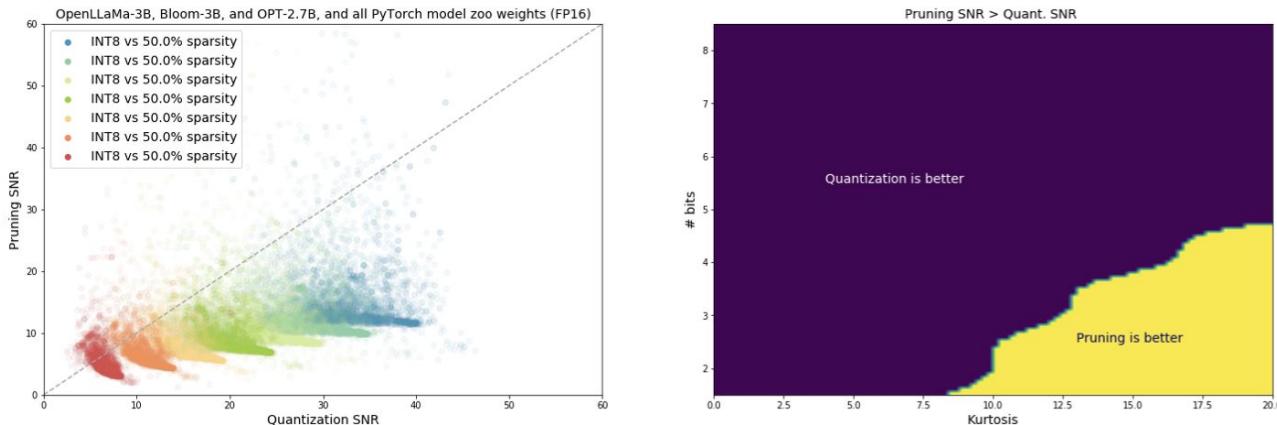


Figure 3: (left) Comparison on all the weights from PyTorch model zoo (46 models) combined with 3 large language models (Bloom-3b, Llama-3b, OPT-2.7b). (left) Pruning SNR versus quantization SNR for every tensor. (right) Pruning is preferable at high compression ratios for tensors with high sample kurtosis values.

# Спасибо за внимание!

Если остались вопросы, то задавайте их в чат или пишите преподавателю в tg:  
[@tradelik](https://t.me/trandelik)