

Algorithm

- OP frequency
 - ONNX OP statistics
- Model Compression
 - Model Quantization
 - Neural Arch Search
 - Knowledge Distillation
 - Pruning
 - reference
- Data type
 - data type

OP frequency

OP frequency

ONNX OP statistics

alphabet order

3 Abs
4325 Add
20 And
2 ArgMax
101 AveragePool
3467 BatchNormalization
3152 Cast
8 CategoryMapper
27 Ceil
407 Clip
2 Compress
2706 Concat
7714 Constant
788 ConstantOfShape
6448 Conv
2 ConvInteger
4 ConvTranspose
1 CumSum
3227 DequantizeLinear
503 Div
100 Dropout
82 DynamicQuantizeLinear
5 DynamicQuantizeLSTM
97 Equal
24 Erf
194 Exp
132 Expand
40 Flatten
52 Floor

12 FusedMatMul
7681 Gather
251 Gemm
29 GlobalAveragePool
8 Greater
2 Hardmax
91 Identity
150 InstanceNormalization
228 LeakyRelu
52 Less
1 LessOrEqual
91 Log
30 Loop
90 LRN
5 LSTM
861 MatMul
96 MatMulInteger
2 Max
522 MaxPool
4 Min
3536 Mul
3 Neg
688 NonMaxSuppression
843 NonZero
41 Not
4 OneHot
160 Pad
234 Pow
100 PRelu
1 preprocess
276 QLinearAdd
5 QLinearAveragePool
67 QLinearConcat
1191 QLinearConv
4 QLinearGlobalAveragePool
72 QLinearLeakyRelu
30 QLinearMatMul
62 QLinearMul
10 QLinearSigmoid

1577 QuantizeLinear
7 Range
100 Reciprocal
6 ReduceMax
466 ReduceMean
65 ReduceMin
8 ReduceSum
4794 Relu
2397 Reshape
43 Resize
48 RoiAlign
2 Round
2 Scan
12 Scatter
36 ScatterElements
2229 Shape
90 Sigmoid
2276 Slice
230 Softmax
135 Split
274 Sqrt
1700 Squeeze
879 Sub
249 Sum
146 Tanh
47 Tile
60 TopK
1469 Transpose
7973 Unsqueeze
22 Upsample
15 Where

Frequency order

1 CumSum
1 LessOrEqual
1 preprocess
2 ArgMax
2 Compress

2 ConvInteger
2 Hardmax
2 Max
2 Round
2 Scan
3 Abs
3 Neg
4 ConvTranspose
4 Min
4 OneHot
4 QLinearGlobalAveragePool
5 DynamicQuantizeLSTM
5 LSTM
5 QLinearAveragePool
6 ReduceMax
7 Range
8 CategoryMapper
8 Greater
8 ReduceSum
10 QLinearSigmoid
12 FusedMatMul
12 Scatter
15 Where
20 And
22 Upsample
24 Erf
27 Ceil
29 GlobalAveragePool
30 Loop
30 QLinearMatMul
36 ScatterElements
40 Flatten
41 Not
43 Resize
47 Tile
48 RoiAlign
52 Floor
52 Less
60 TopK
62 QLinearMul
65 ReduceMin
67 QLinearConcat
72 QLinearLeakyRelu
82 DynamicQuantizeLinear
90 LRN
90 Sigmoid
91 Identity
91 Log
96 MatMulInteger
97 Equal
100 Dropout
100 PRelu
100 Reciprocal

101 AveragePool
132 Expand
135 Split
146 Tanh
150 InstanceNormalization
160 Pad
194 Exp
228 LeakyRelu
230 Softmax
234 Pow
249 Sum
251 Gemm
274 Sqrt
276 QLinearAdd
407 Clip
466 ReduceMean
503 Div
522 MaxPool
688 NonMaxSuppression
788 ConstantOfShape
843 NonZero
861 MatMul
879 Sub
1191 QLinearConv
1469 Transpose
1577 QuantizeLinear
1700 Squeeze
2229 Shape
2276 Slice
2397 Reshape
2706 Concat
3152 Cast
3227 DequantizeLinear
3467 BatchNormalization
3536 Mul
4325 Add
4794 Relu
6448 Conv
7681 Gather
7714 Constant
7973 Unsqueeze

Model Compression

Model compression techniques are mainly classified into quantization (Zafir et al., 2019), pruning (Hoefler et al., 2021; Vadera and Ameen, 2022), Neural Architecture Search (NAS) (Sun et al., 2019), and knowledge distillation

Model Quantization

Quantization Granularity

Quantization is a magic spell to reduce the memory footprint of a model. But often quantization leads to a drop in the accuracy of the model. This is where the Granularity of quantization comes into the picture. Selecting the right granularity helps in maximizing the quantization without much drop in the accuracy performance

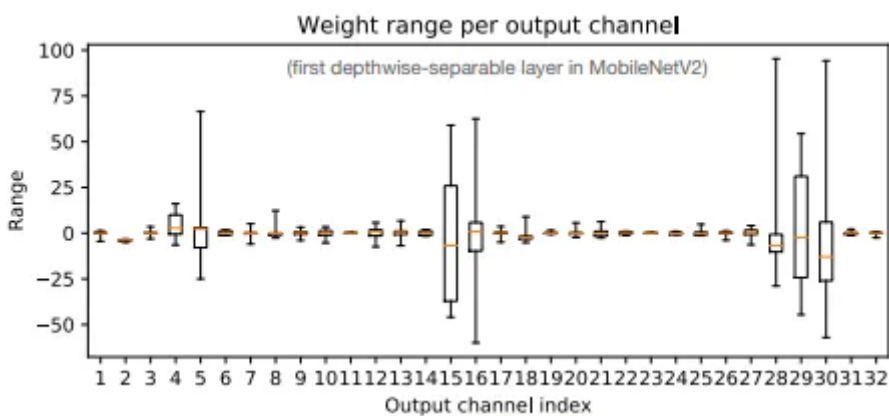
Per-Tensor Quantization

In per-tensor quantization, the same quantization parameters are applied to all elements within a tensor. Applying the same parameters across tensors will cause a drop in accuracy.

Per-Channel Quantization

In per-channel quantization, different quantization parameters are applied to each channel of a tensor independently. This often leads to a lower error while quantizing compared to per tensor quantization.

Per-channel quantization captures variations in different channels more accurately. This usually helps in CNN models where the range of weights varies over different channels.



Mix-precision Quantization

different layer using different precision.

Quantization Method

Quantization Method	Accuracy	Model Size Reduction	Inference Speed	Ease of Implementation
GGUF	High (adaptive quantization of parameter groups)	Significant (grouped quantization)	High (optimized through grouped quantization)	Moderate (requires parameter grouping)
AWQ	High (asymmetric quantization for accuracy)	Moderate	Moderate to High (asymmetric approach)	Moderate (asymmetric scaling)
GPTQ	Moderate (post-training fine-tuning)	Moderate	Moderate (benefits from post-quantization)	Easy (applied post-training)
GGML	High (mixed-precision within groups)	Significant (effective group-wise precision)	High (mixed-precision within groups)	Complex (group-wise mixed-precision)
PTQ	Moderate (quick and easy implementation)	High (straightforward implementation)	High (quick deployment)	Easy (quick application)
QAT	Very High (trained with quantization in mind)	Moderate to High (requires complex training)	High (optimized for quantization during training)	Complex (requires training with quantization)
Dynamic Quantization	Moderate (applied during inference)	High (applied at inference time)	Moderate to High (quick deployment)	Easy (inference-based)
Mixed-Precision	High (selectively applied different precisions)	Moderate to High (selectively applied)	High (selective precision boosts performance)	Complex (selective precision assignment)

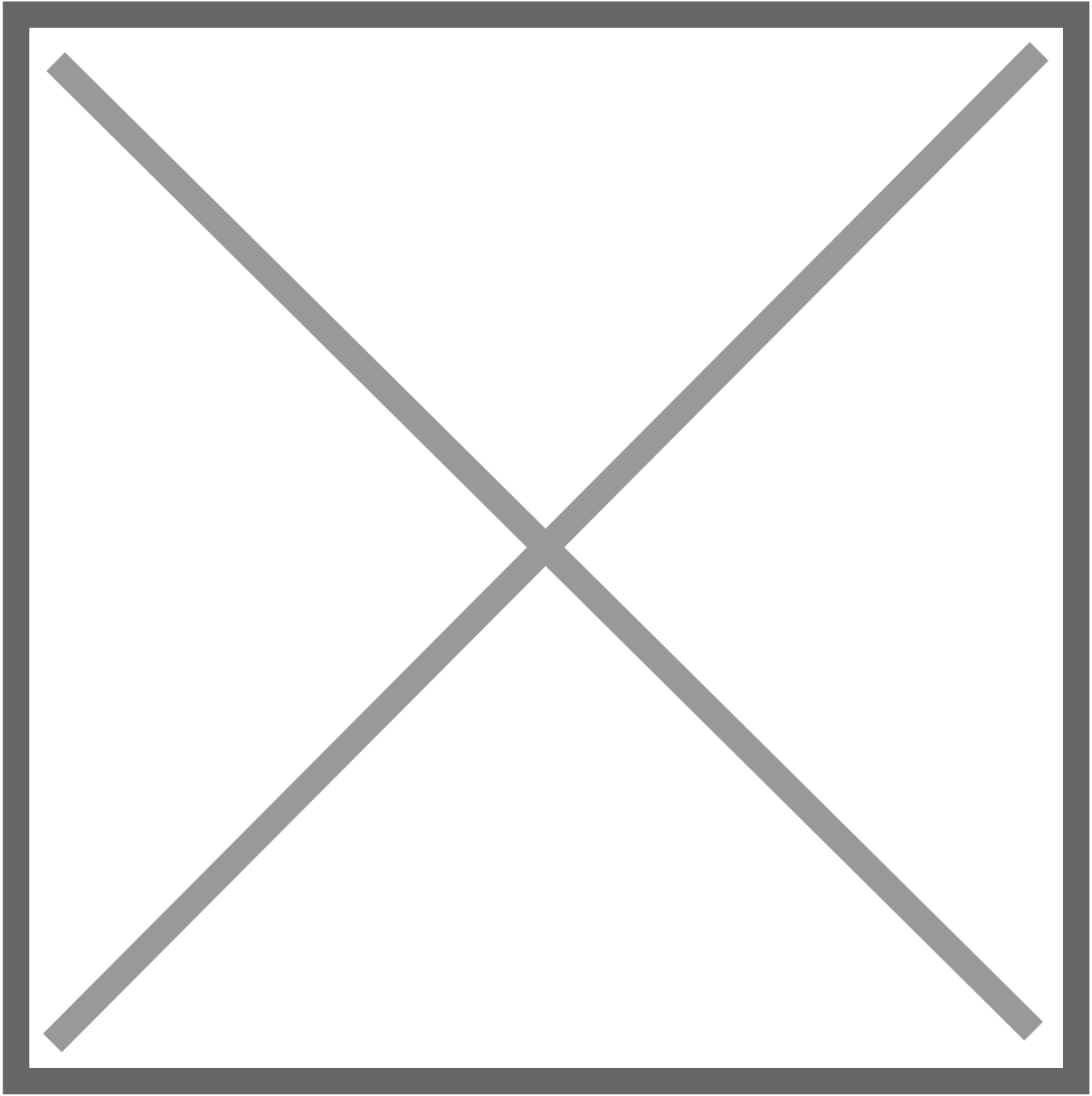
Model Compression

Neural Arch Search

Knowledge Distillation

Model Compression

Pruning



INT8 



[-128, 127]  8bit 



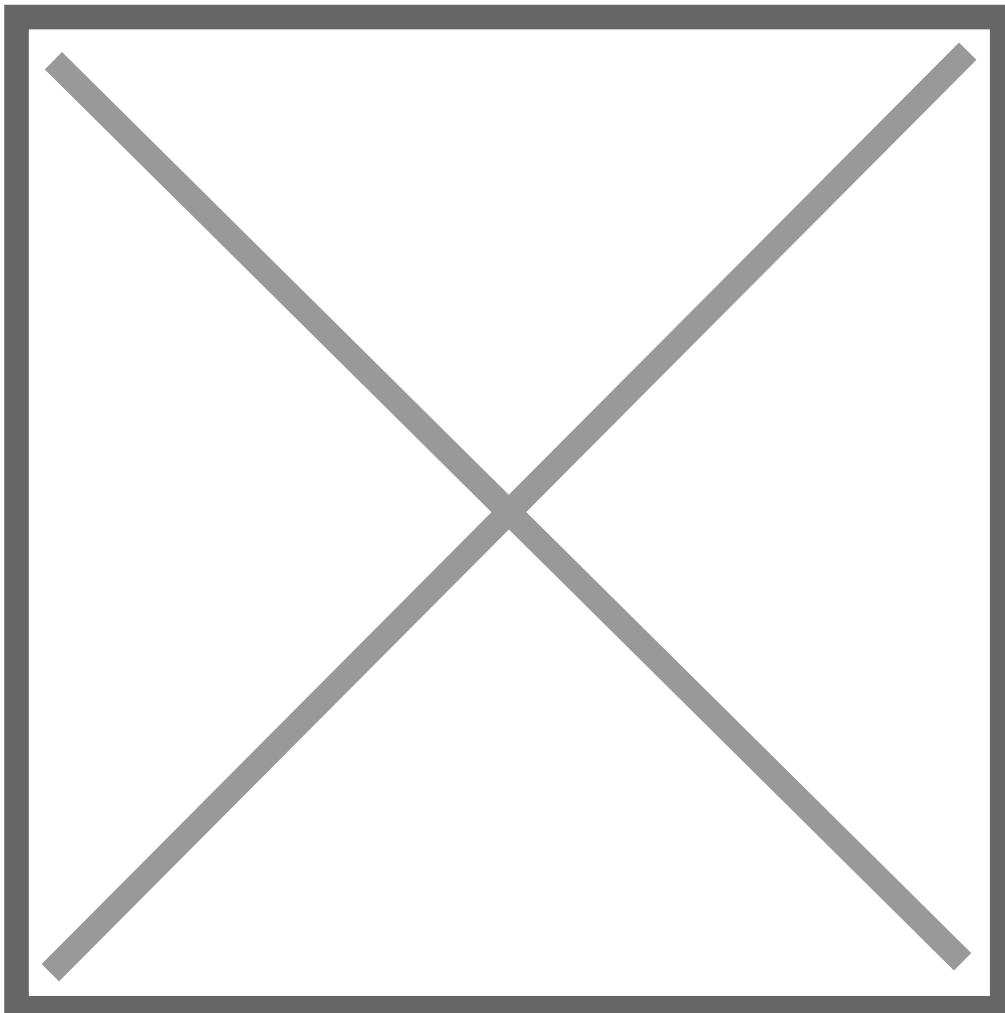
●TensorFloat-32(TF32) 19 BF16 FP16

●Float16 (FP16) 5 10 1 FP16 FP32 (loss scaling)

●Bfloat16 (BF16) 8 (FP32) 7 1 BF16 FP32 FP16 3 BF16 FP16

●(Un)signed Char (INT8) 7 1 [-128, 127] UINT8 [0, 255]

●Unsigned Char (INT4 * 2) 1 [-8, 7] INT4 3



CNN vs. LLM

LLM [] CNN []

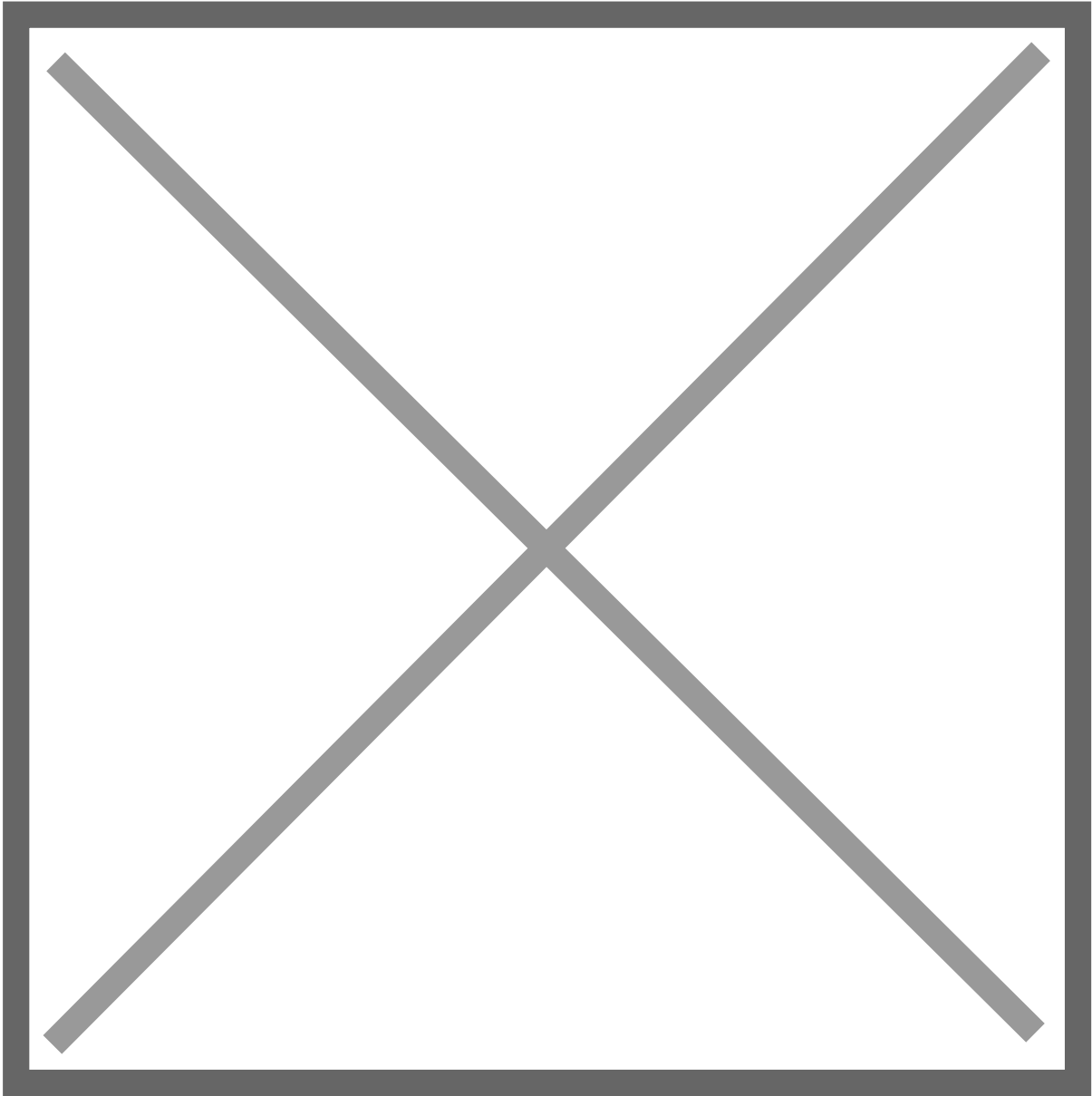
[] [] CNN [] per-channel [] per-layer [] kernel [] LLM
 [] [] [] vector [] (tensor [])
 [] [] [] A*B=C []
 per-tensor [] A [] B [] INT
 [] [] per-vector [] LLM [] per-vector [] LLM
 [] [] X [] W
 [] [] FP16 [] k [] k
 [] 2 [] k [] 256 [] 128 [] per-group []

[]

[] **INT4** [] CNN [] INT8 [] LLM
 [] CNN [] INT4
 [] [] 13B [] scale
 [] scale [] fp16 [] fp32 [] scale [] INT4
 [] []

[] [] CNN [] loss
 [] [] per-channel
 [] [] scale
 [] [] LLM [] head
 [] [] tensor [] hidden dim
 [] [] [50, 100]
 [] [] [0.5, 1.2], [] LLM []
 attention [] softmax [] Attention Is Off By One [] tensor []
 per-tensor [] LLM
 [] []

[] [] CNN [] LLM [] cuda [] INT
 [] [] Core
 [] []
 [] []
 [] [] []

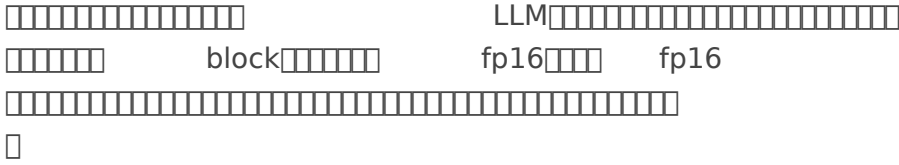


X[] activations[] sequence length[] hidden-size
[]
[]

[] X W [] int8
[] fp16
[] (W []
token[]) []
element-wise[] concat[] []

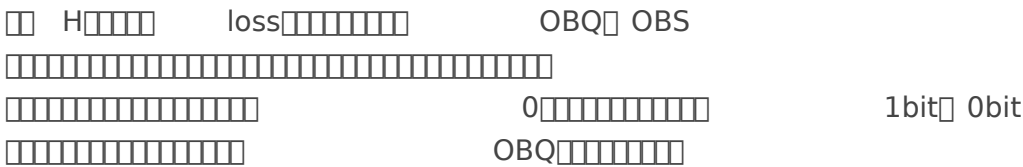
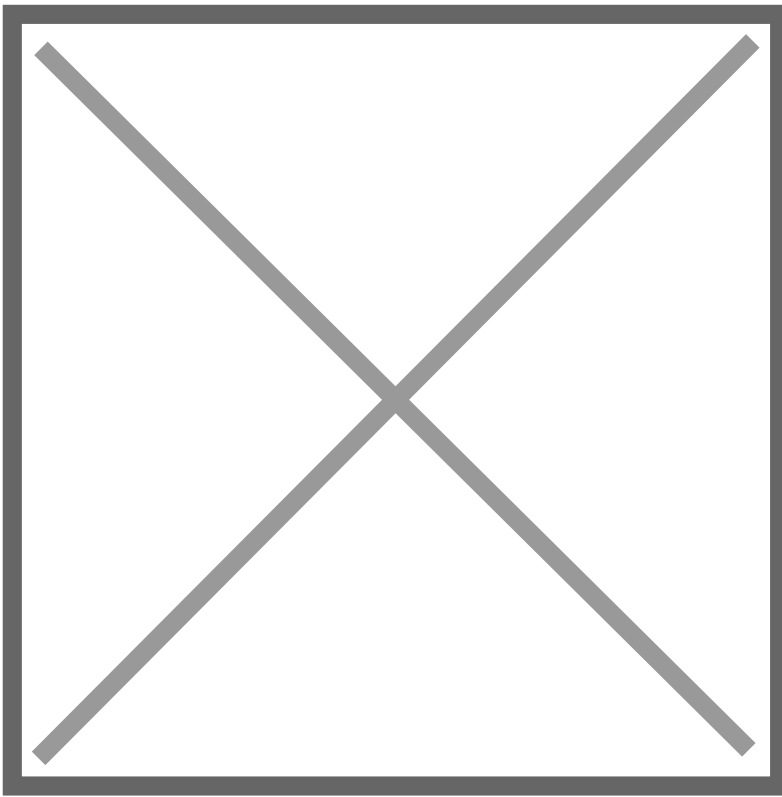
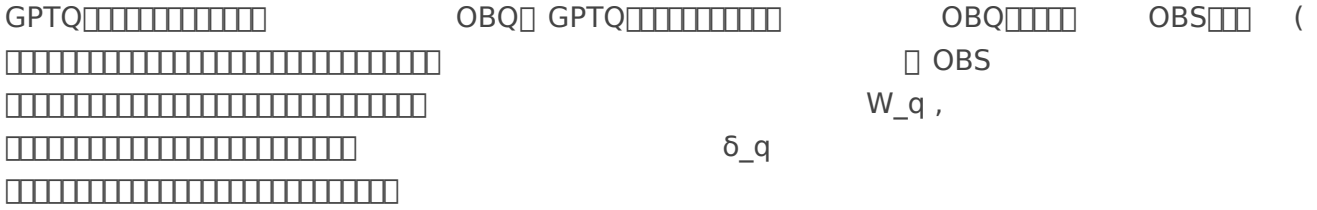
LLM.INT8()[] gemm[] X[] W
[] threshold[] 6
[] 6 []

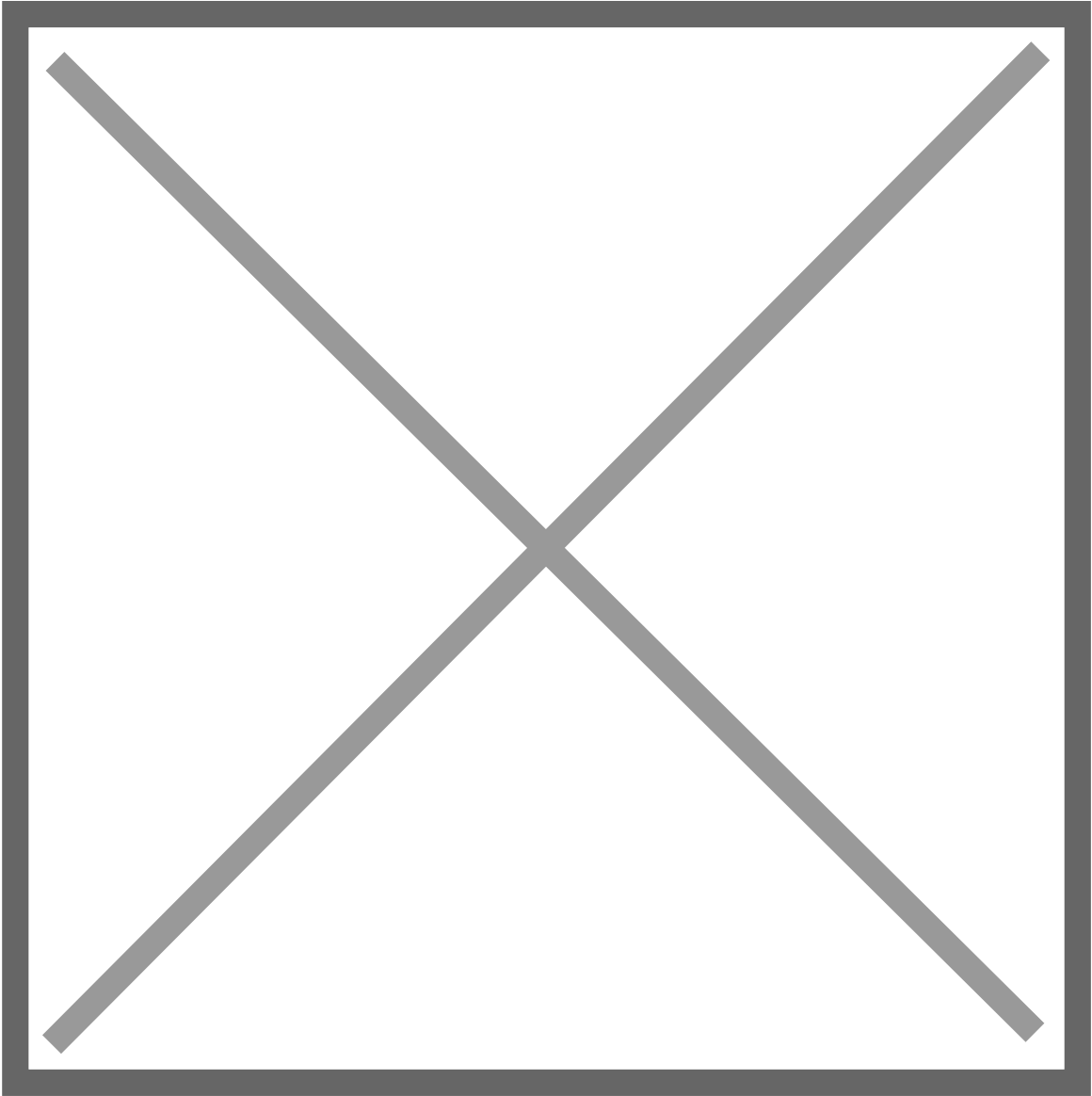
LLM.int8() [] LLM.int8() [] BLOOM-
176B [] FP16 [] 15% [] 23% [] ([] T5-3B [] T5-11B)



block

GPTQ: [][][][][][]



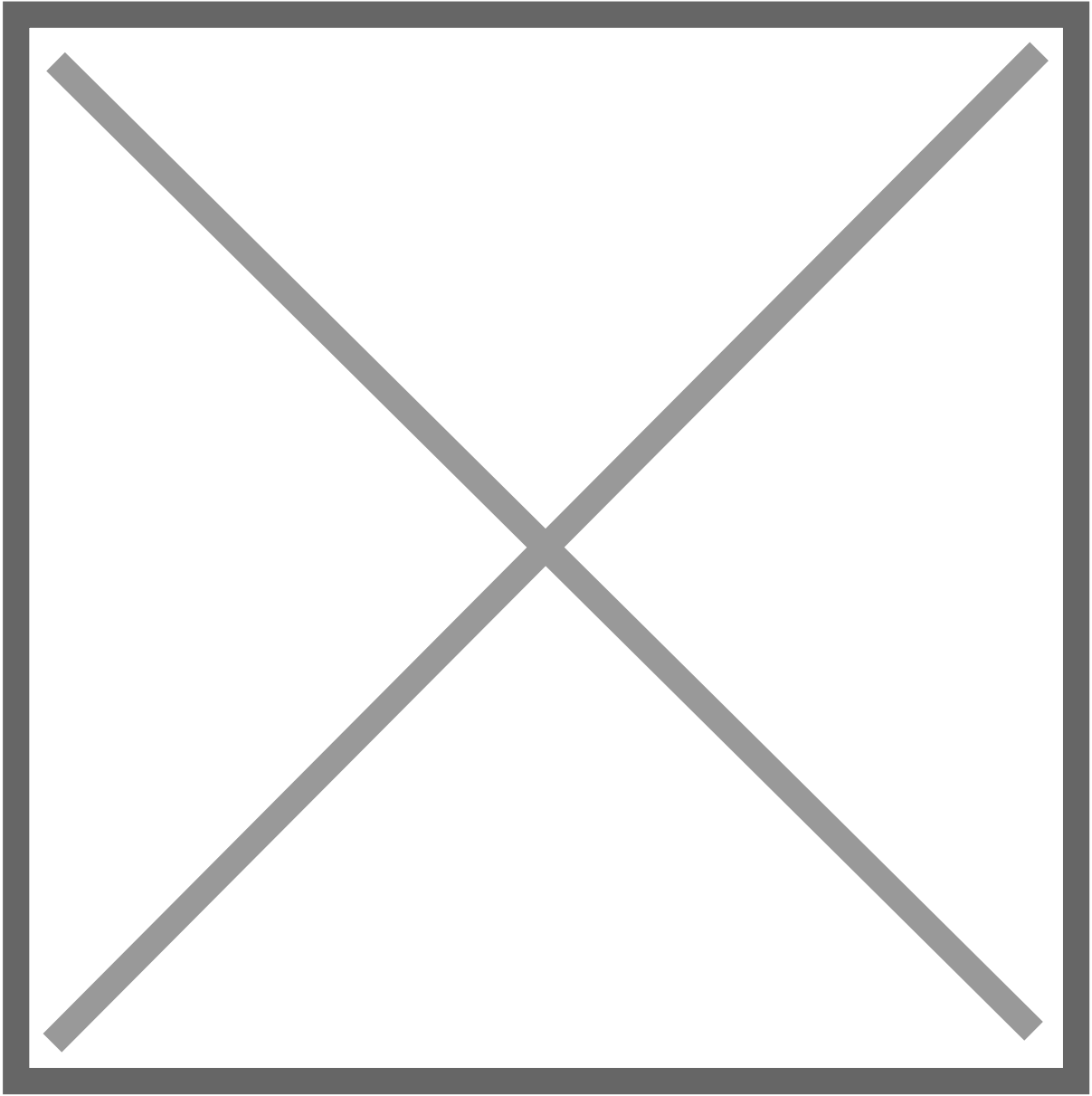


quant(q) q quant grid

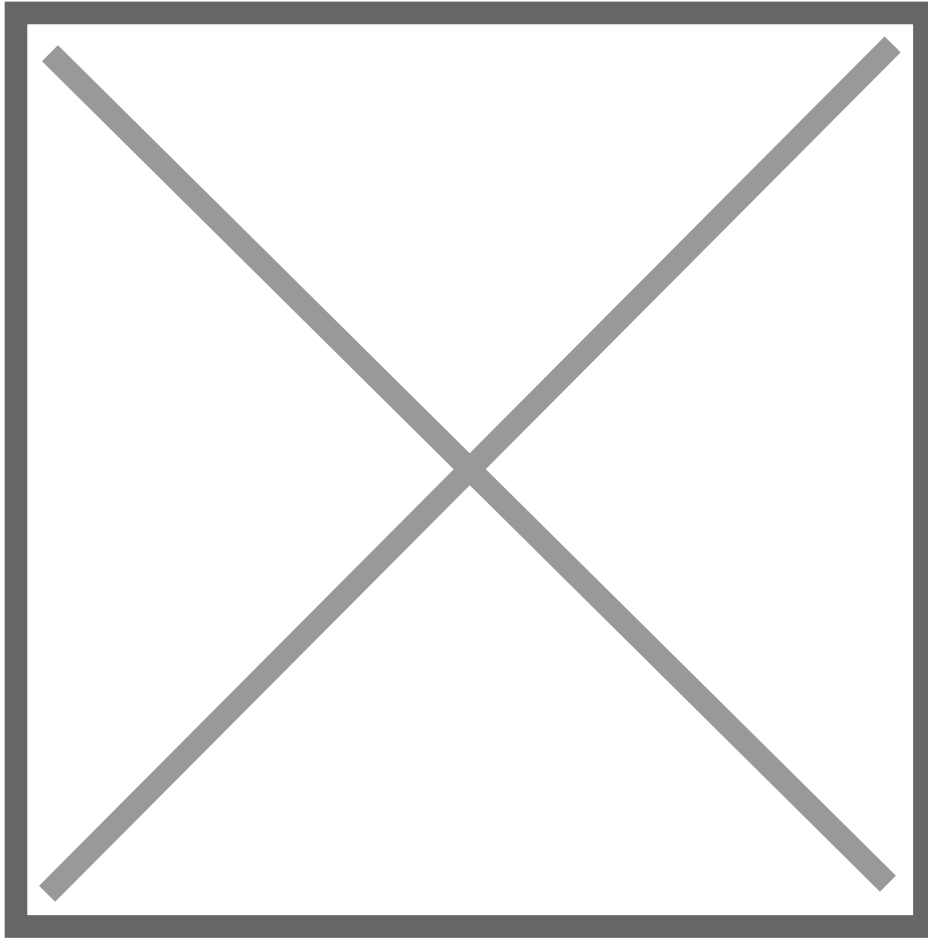
quant(q) OBS OBQ

OBS

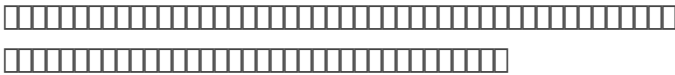
q /



GPTQ OBS



-



GPU



- Scale Zero

Group size



GPTQ LLM

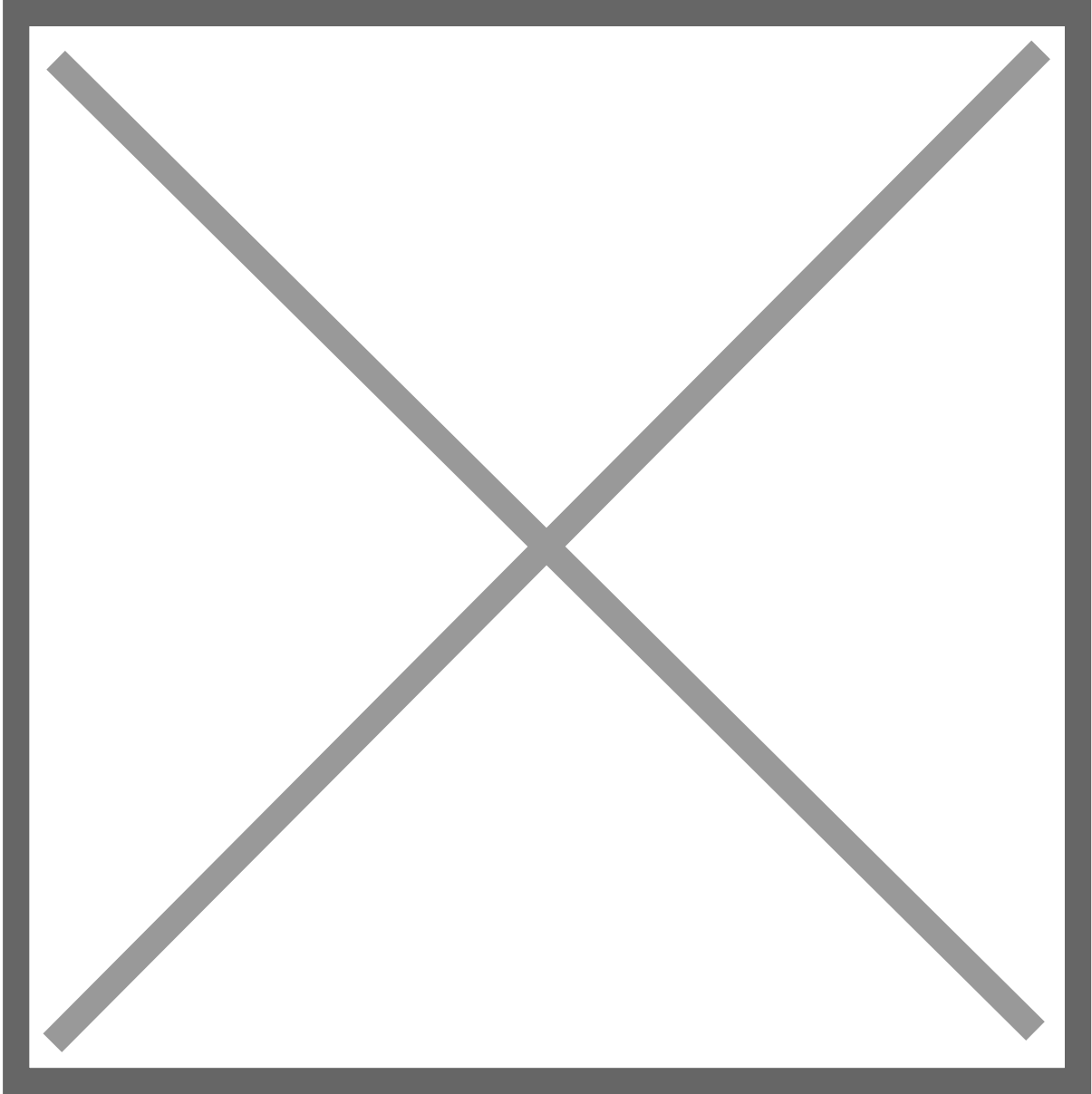


SmoothQuant: W8A8 scale

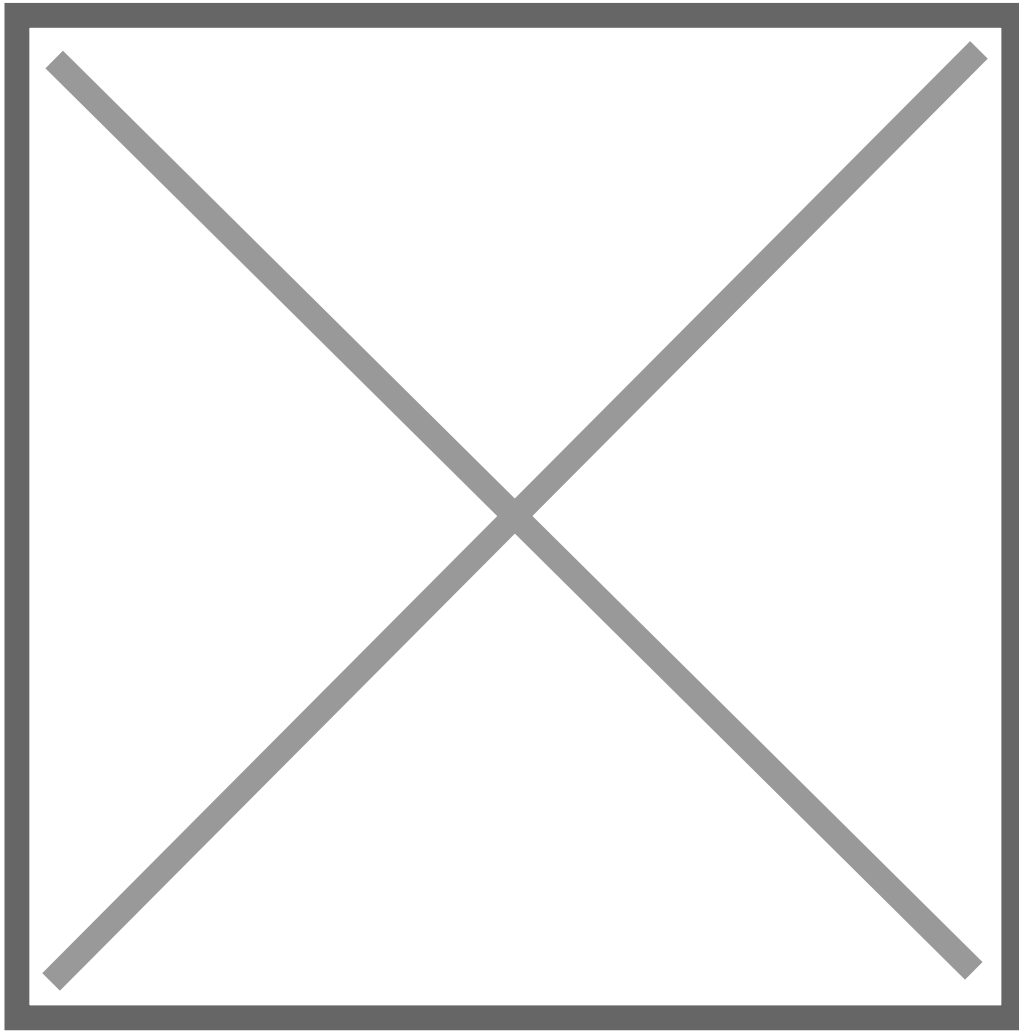
30% LLM channel activation outlier (e.g. GLM-130B)
int8 -128 127 outlier -128 127

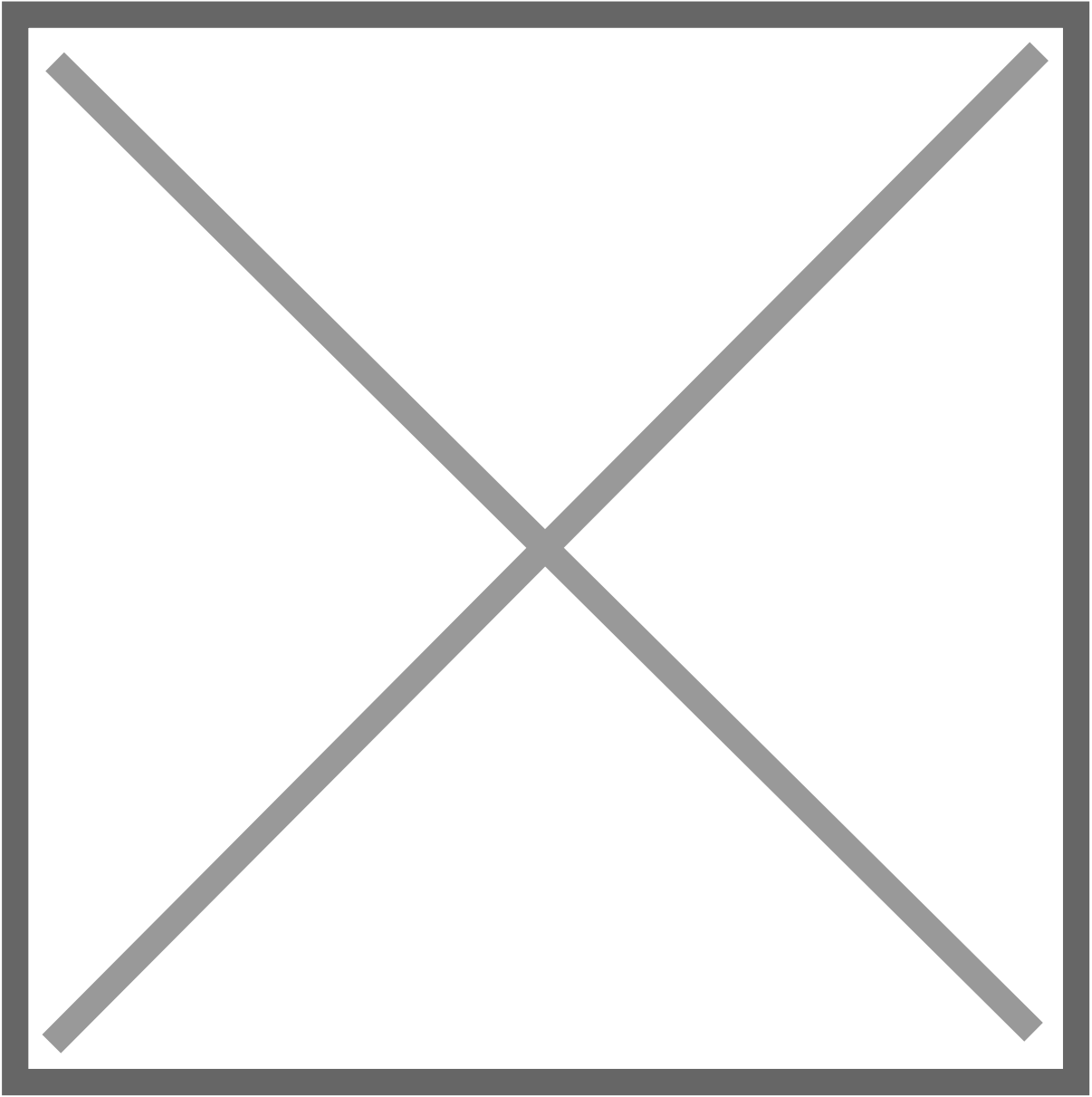
20,20] [-10,10] 8bit 5bit 4bit outlier

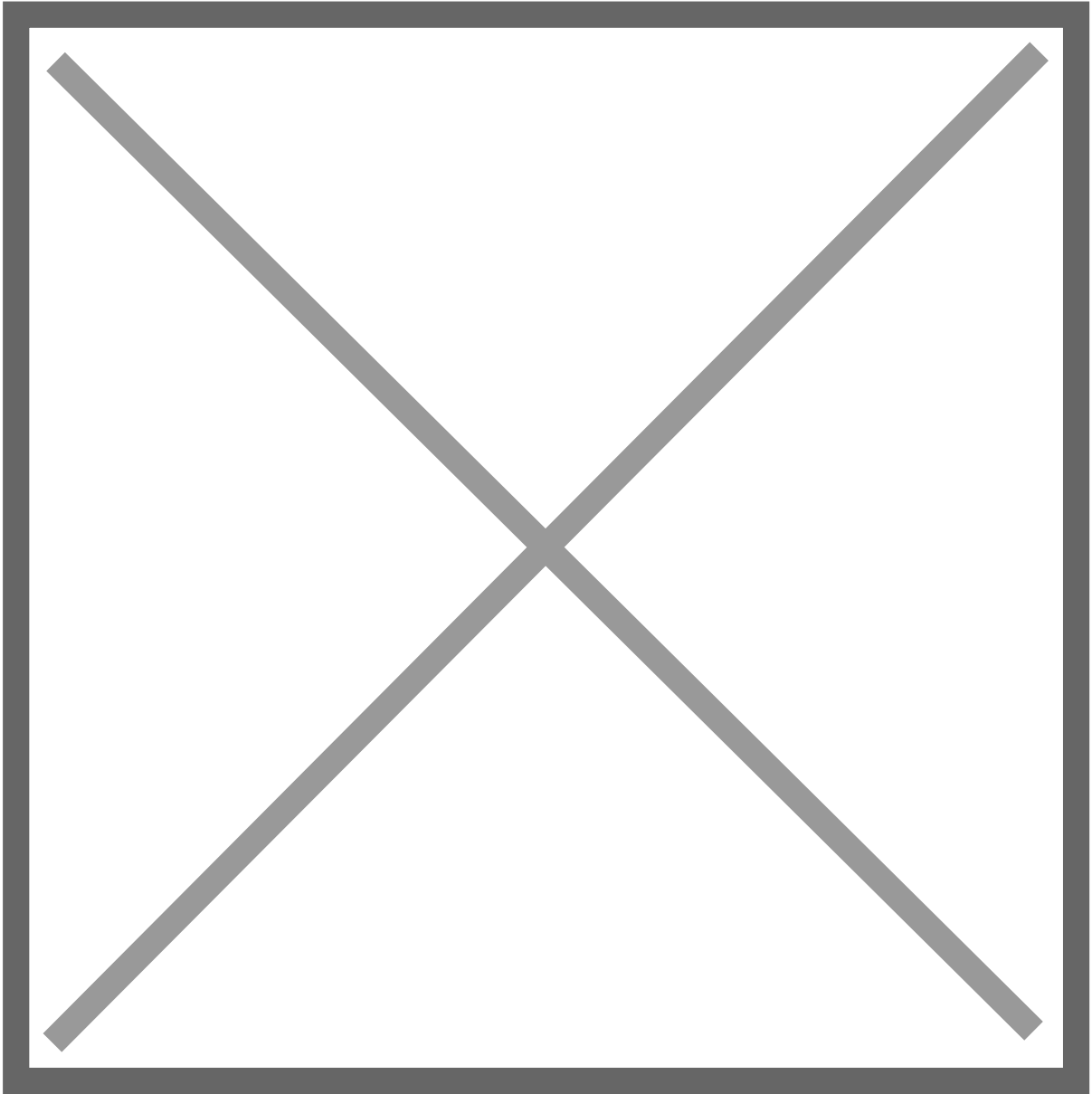
weight activation INT8 INT4 activation weight



SQ activation weight fp32 X act W weight s s







LLM weight-only SmoothQuant LLM

SmoothQuant

SmoothQuant

SmoothQuant

W8A8

SmoothQuant

INT4 bit

SmoothQuant

SmoothQuant

SmoothQuant

SmoothQuant

SmoothQuant

SmoothQuant

SmoothQuant

GPT2 SmoothQuant

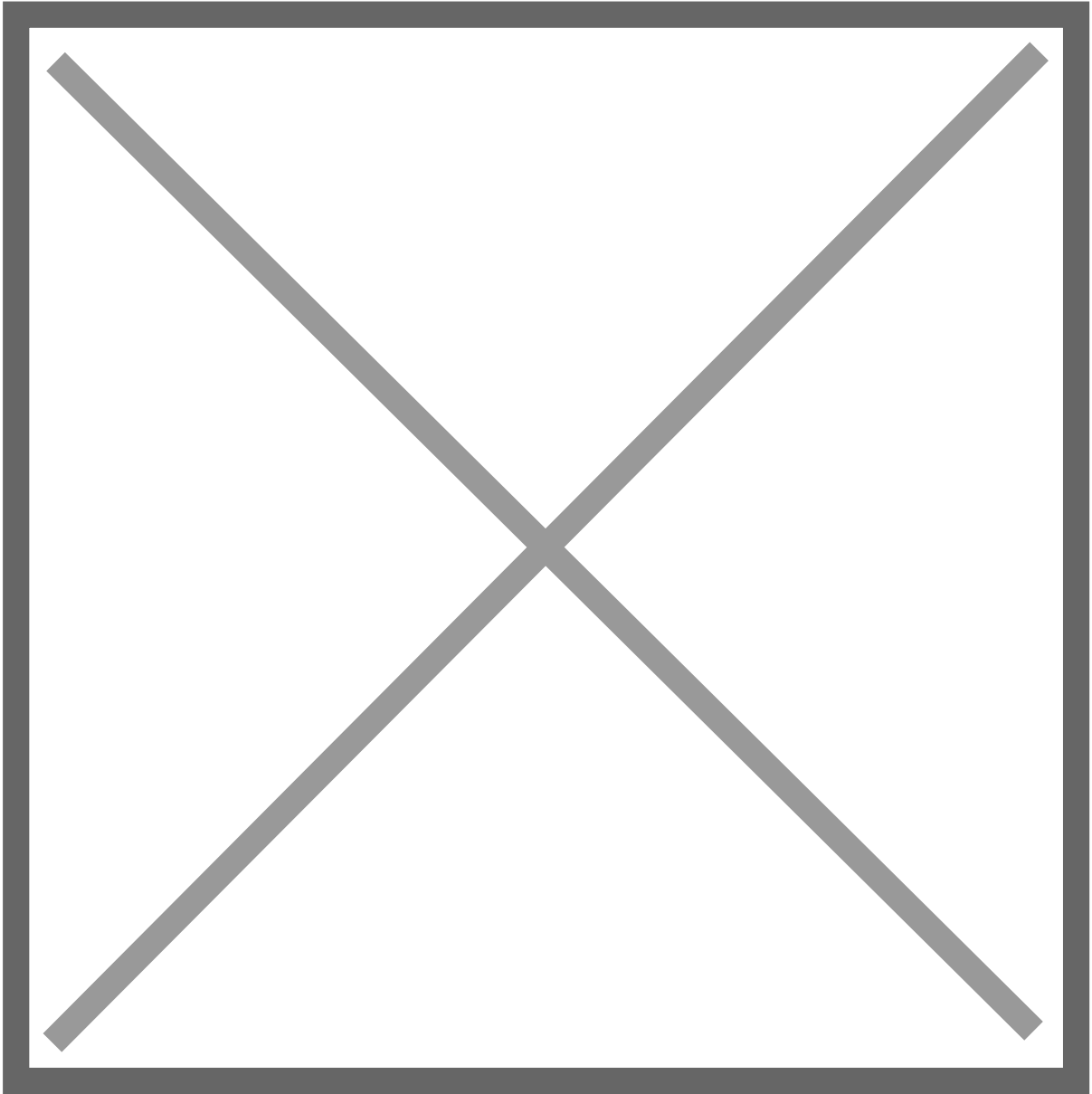
LLM

SmoothQuant

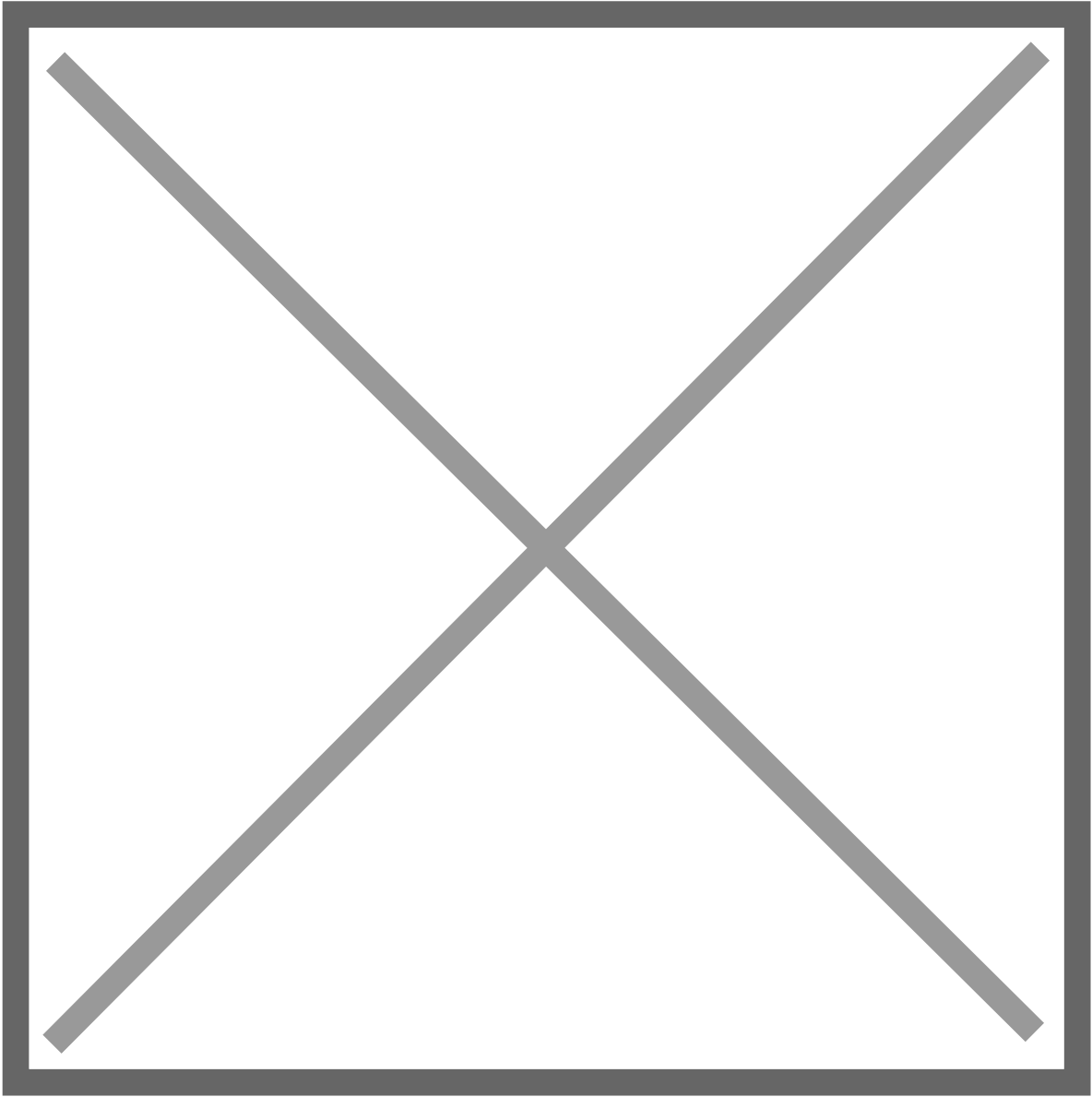
weight-only

TensorRT-LLM TensorRT LLM

AWQ: SmoothQuant

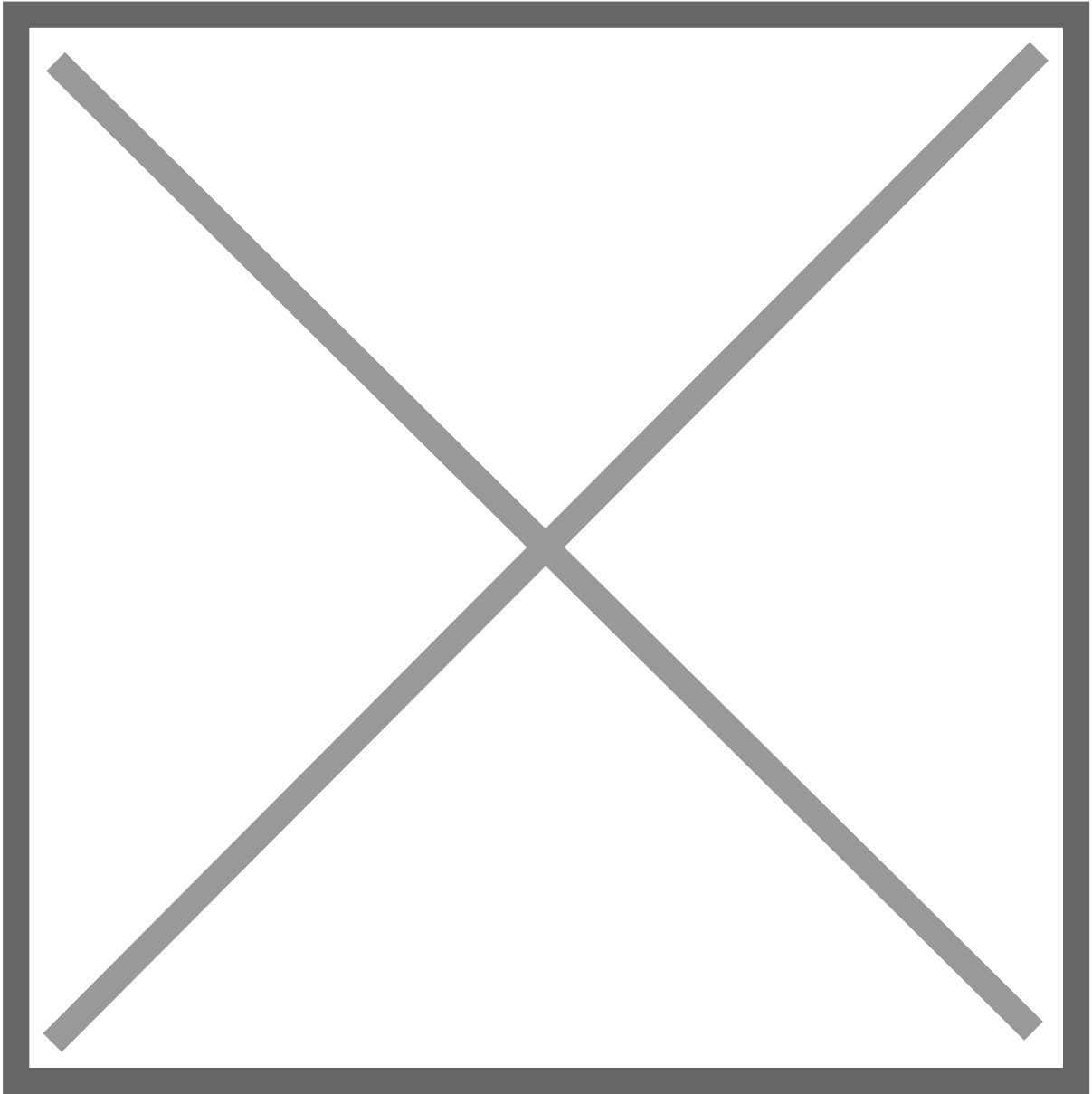


channel (a) 1%
fp16 (b)
(int8+fp16) LLM.int8() tensor
1% group
SmoothQuant
s X



□□□□□□

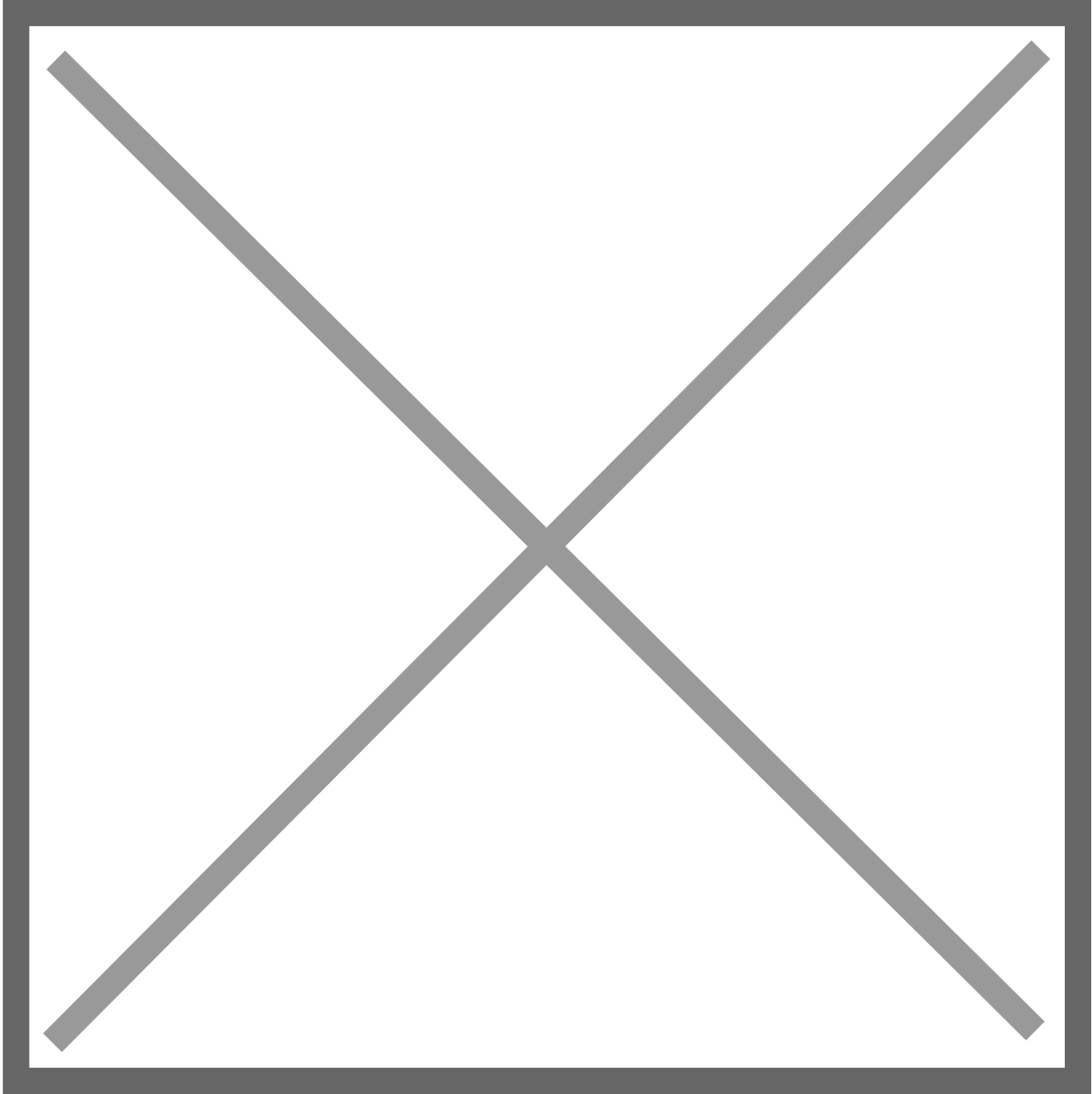
s□□□□□□□□



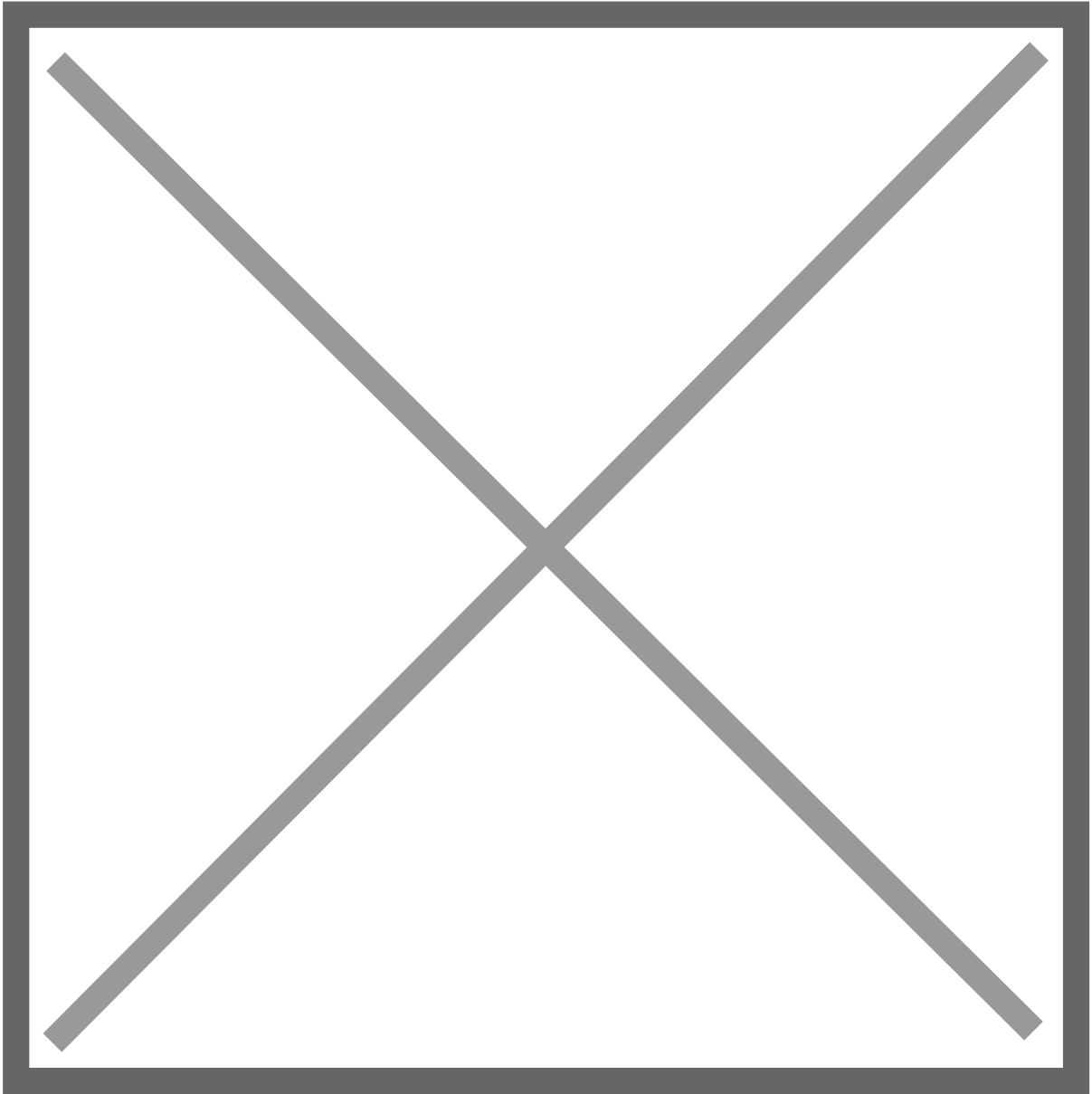
```
int sX = meanc_out|X| int sW = meanc_out|W| int sX_group group
int weight int group alpha beta [0, 1]
int sW
int group weight
int
int AWQ kernel int int int
int smooth quant AWQ weight-only SQ W8A8
int AWQ GPTQ
kernel
```

Llama. CPP K-quant

[LLaMA.cpp](#) [] Georgi Gerganov [] Meta [] LLaMA [] C/C++
 [] tensor [] GGML [] llama.cpp [] GPU
 [] Meta [] LLaMA [] cuBLAS []
 6 [] 15 [] PR [] llama.cpp [] CUDA [] cuda kernel [] llama [] GPU
 []



llama.cpp [] INT4 [] K-quant [] Q
 [] x [] x bit [] y [] 0 [] y [] 1
 [] scale [] zero point [] y [] K [] K-quant [] K
 [] Small, Medium [] Large [] Q4_0 [] Q8_0 [] INT4 []
 INT8 []



$y = K$ K-quant K-quant 16 x 8 " " " " 16
8 scale 16
1 fp16 K_2 16
size 16*8 " "K-
quant K_1 group size = 8, 65B 37.5GB

llama.cpp K-quant fp16 3~4 (token/s)
llama.cpp c++ cuda
+gemm GPT
KV Cache token vector-matrix GEVM
matrix-matrix GEMM llama.cpp

llama.cpp llama +K-quant

llama.cpp llama.cpp

Kernel

QLora

Lora Qlora

QLora

lora

A B W LoRA AB W AB

QLora

Block-wise Quantization block scale block 64, K-quant

Quantile Quantization

Quantile Quantization. Quantile

block-wise quantization block scale float32, block 64 32/64 = 0.5bit 4bit 0.5bit 12.5% scale outliner 256 block 8bit c 8/64 + 32/256 = 0.127 bit

QLora LLM



- LLM weight-only
- LLM CNN per-vector + scale

Data type

Data type

data type

