

# Today's Agenda:

- ① Quantization methods in training / inference
- ② A little bit on distillation
- ③ A little more on MoEs (content that Trevor didn't get to)

## Recall from last time

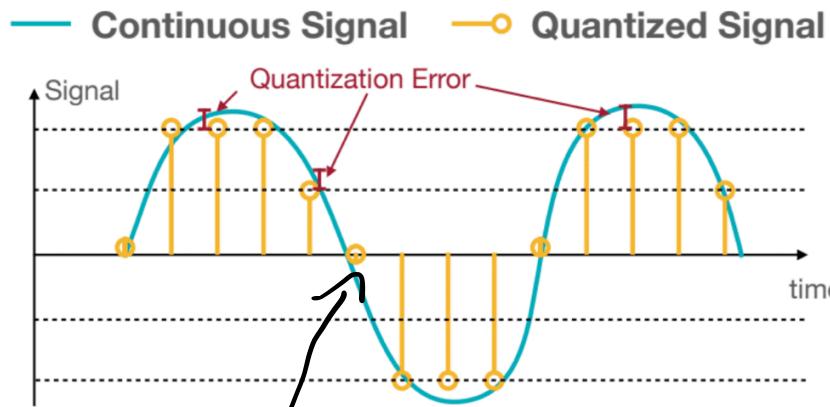
	Exponent bits	Fraction bits	Total
IEEE 754 FP32	8	23	32
IEEE 754 FP16	5	10	16
Google Brain Float 16	8	7	16
Nvidia FP8 EFM4	4	3	8
Nvidia FP8 E5M2	5	2	8

↳ used for gradient in BW pass.

- #exponent bits defines range / dif. windows that can exist
- fraction bits define precision (within any bucket)
  - \* 2<sup>nd</sup> fp8 type can represent larger normal values

## What is quantization?

"process of constraining an input from a continuous or otherwise large set of values to a discrete set"



diff. btwn input and quantized value =

Credits: MIT 6.5940: TinyML and Efficient Deep Learning Computing

quantization error

## K-means Quantization for ML

- use INTEGER WEIGHTS, floating-point

computation

### K-means quantization (1<sup>st</sup> For inference)

Original FP32 Weights			
2.09	-0.98	1.48	0.09
0.05	-0.14	-1.08	2.12
-0.91	1.92	0	-1.03
1.87	0	1.53	1.49

→ Centroids

3:	2.00
2:	1.50
1:	0.00
0:	-1.00

2 is  
cluster  
center  
for all  
blue  
values

→ idea: take original fp32 weight matrix, cluster values

→ instead of storing raw weights, store index to cluster centroid

→ cluster centroids themselves are stored w/ fp32

\*use K-means to find cluster centroids

↳ so what we actually store is:

3	0	2	1
1	1	0	3
0	3	1	0
3	1	2	2

+ cluster  
centroids  
themselves

WHAT ARE STORAGE SAVINGS?

(for inference case)

→ original:  $D \times D$  matrix in fp32

$32D^2$  bits (or  $4D^2$  bytes)

→ compressed:  $D \times D$  matrix of indices to  
cluster centroids ; suppose  $S = \text{size of index}$

\* # of bits needed to represent a centroid index =

$\log_2(S)$  bits  $\rightarrow \log_2(S) \cdot D^2$  bits

→ also storing centroids themselves

in fp32  $\rightarrow 32S$  bits or  $8S$  bytes

Original:

$32D^2$  bits

or  $4D^2$  bytes

compressed:

$(\log_2(S) \cdot D^2 + 32S)$  bits or

$\log_2 S \cdot D^2 / 8 + 4S$  bytes

→ consider  $D=4$ ,  $S=4$ :

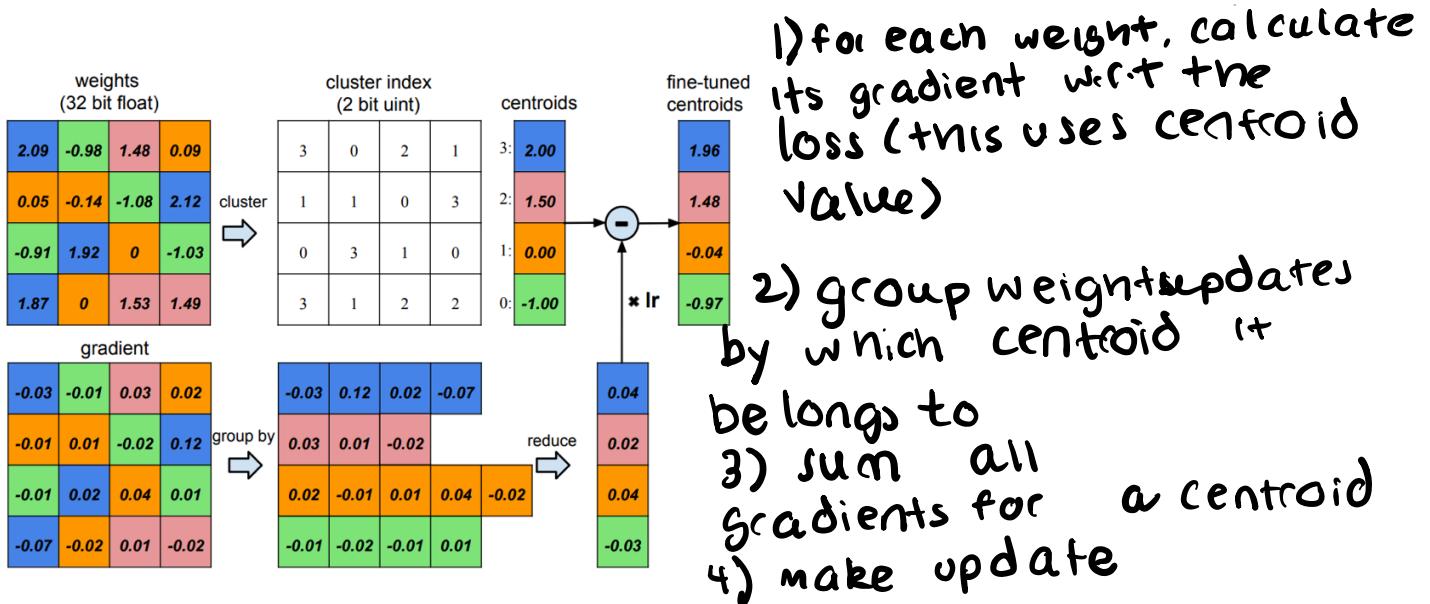
→ original = 64 Bytes

→ compressed =  $4 + 16 = 20$  Bytes

]  
3.2x  
smaller!

# -WHAT HAPPENS FOR TRAINING, when we have to update weights?

## Idea: update cluster centroids (in fp32)



Song Han et al., Deep Compression, 2016 ICLR (Best Conference Paper)

→ steps 1 → 3:

$$\frac{\partial L}{\partial C_k} = \sum_{i,j} \frac{\partial L}{\partial w_{i,j}} \text{ for } w_{i,j} \text{ in centroid } c$$

$$\rightarrow \text{step 4: } C^{(s+1)} \rightarrow C^{(s)} - l.r. \cdot \frac{\partial L}{\partial C^{(s)}}$$

+ general compression

rate for a set of weights  $w$ , represented w/  $s$  centroids, each using  $b$  bits:

$$r = \frac{w \cdot b}{w \cdot \log_2(s) + s \cdot b}$$

in prev example,  
 $\rightarrow b=32$ ,  
 $s=4$

→ Can we do even better? what if  
weights are distributed non-uniformly?

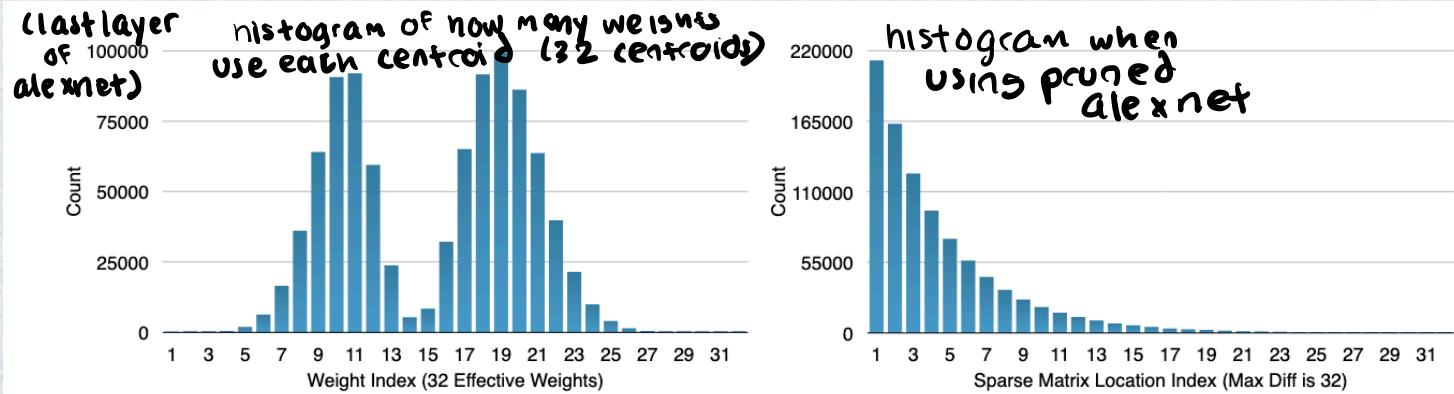


Figure 5: Distribution for weight (Left) and index (Right). The distribution is biased.

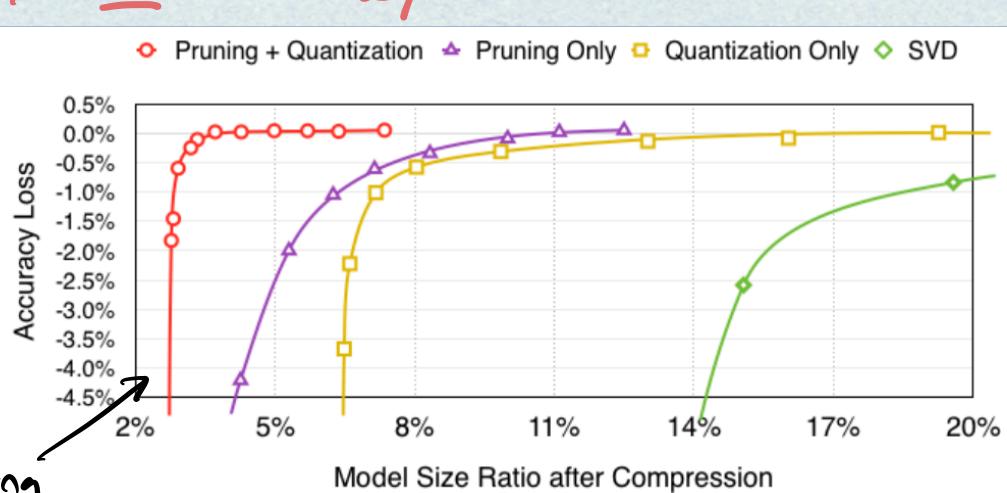
Song Han et al., Deep Compression, 2016 ICLR (Best Conference Paper)

**Observation: use HUFFMAN codes**

- Instead of using a constant # of bits for each centroid(codeword)-use variable-length codewords
- intuition: can use FEWER bits for more frequent-values ; results in 20 % storage savings w/ NO accuracy loss

### Results

- within a method, as we increase compression, accuracy ↓



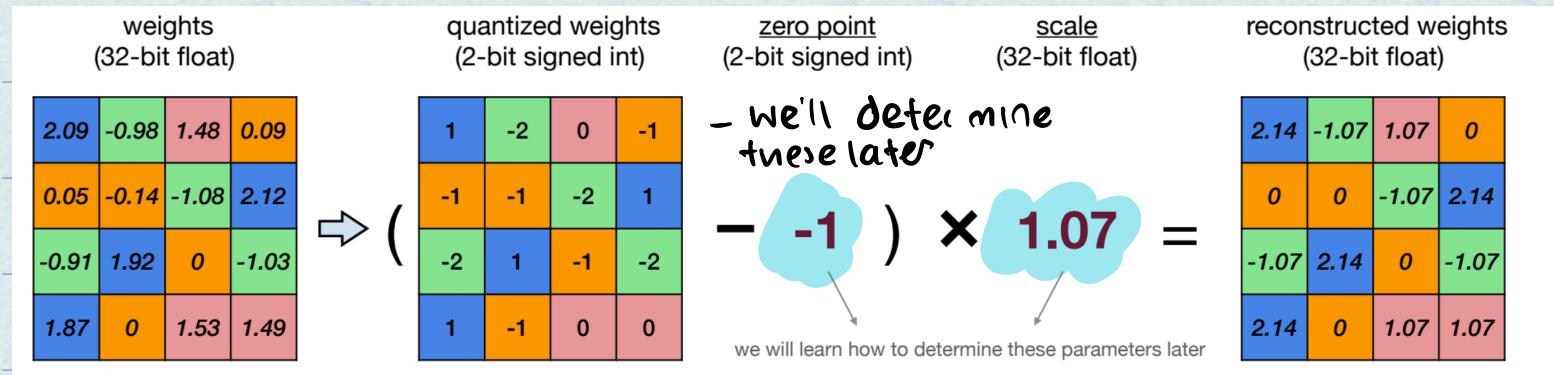
Pruning AND quantization is best

## LINEAR QUANTIZATION

→ integer weights + integer arithmetic

→ "affine" transformation of FPs to

Integers - involves a scaling + offset

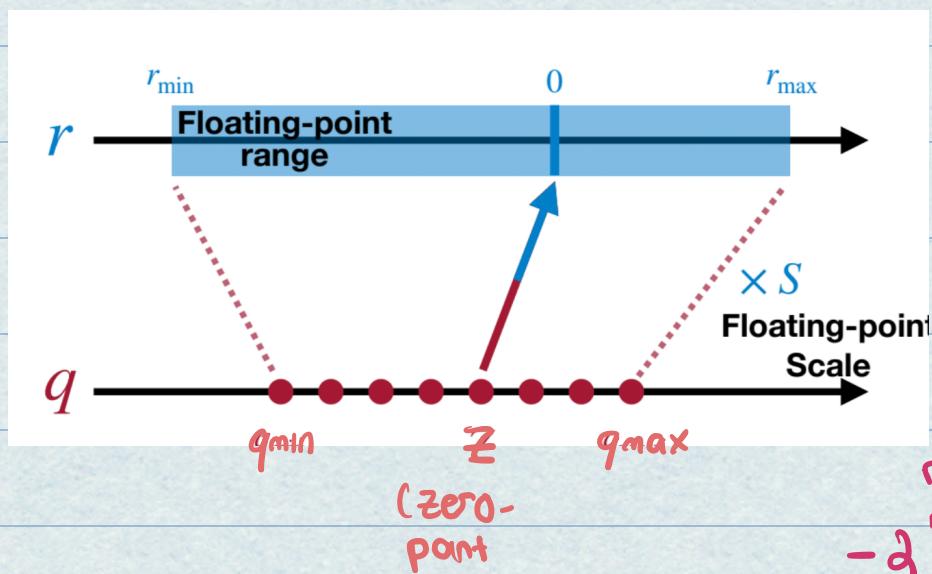


MIT 6.5940 Lecture: TinyML and Efficient Deep Learning Computing (Song Han)

idea:  $r \approx s(q - z)$

→ so only need to store  $q$  (quantized integer version),  $z$  (offset int) and  $s$  (scaling factor)

How do we determine offset and scaling factor?



→ we know  $r_{\min}$

and  $r_{\max}$  from original weights

→ we know  $q_{\min}$

and  $q_{\max}$ :

if we are quantizing to N bits →

range is

$-2^{N-1} \dots 2^{N-1}$

→ we can just solve for  $s$

$$r_{\max} = S(q_{\max} - z) \rightarrow r_{\max} - r_{\min} = S(q_{\max} - q_{\min})$$

$$r_{\min} = S(q_{\min} - z)$$

\* now calc  $z$ :

$$S = \frac{r_{\max} - r_{\min}}{q_{\max} - q_{\min}}$$

$$z = q_{\min} - \frac{r_{\min}}{S} \approx \text{round}\left(q_{\min} - \frac{r_{\min}}{S}\right)$$

CONSIDER a MATRIX MULTIPLICATION

$$y = w x$$

$$s_y (q_y - z_y) = s_w (q_w - z_w) \cdot s_x (q_x - w x)$$

$$q_y = \frac{s_w s_x}{s_y} \underbrace{(q_w - z_w)(q_x - z_x)}_{\text{all integer math!}} + z_y$$

↳ this can be  
rescaled to  
an n-bit integer

QUICK OVERVIEW OF KNOWLEDGE DISTILLATION

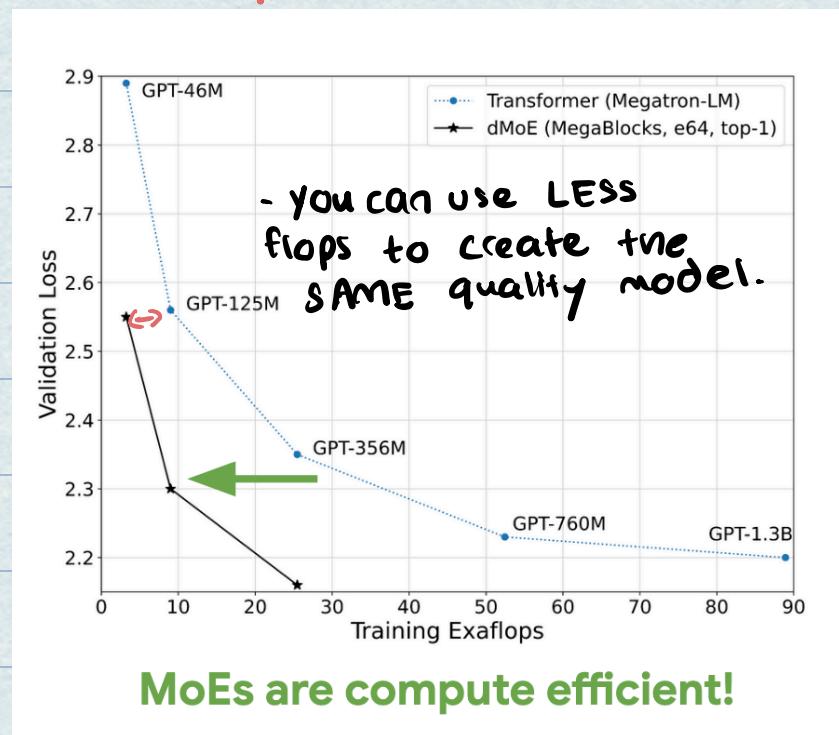
- general idea: guide large "teacher-model" to  
teach training of "student" model

- can do this at many dif. layers, including:

- output logits
- intermediate weights
- intermediate features
- gradients
- sparsity patterns

More Mixture of Experts (what Trevor didn't get to )

## \* Efficiency of MOEs in Inference vs. Training



\* MOEs will be compute efficient - we less compute for some quality

- you can actually translate theoretical compute gain into K runtime wins - large batches

RECAP on AI of Dense model vs. AI of Expert:

$$AI_{dense} = \frac{ops}{bytes} = \frac{N_{tokens} \cdot d_{model} \cdot d_{ff} \cdot 2}{N_{tokens} \cdot d_{model} + d_{model} \cdot d_{ff} + N_{tokens} \cdot d_{ff}}$$

assumes you write  $\text{fout}$  intermediate result

2 projection:  $( (N_t \times d_{model}) (d_{model} \times d_{ff}) ) - (d_{ff} \cdot d_{model})$

$$AI_{moe} \approx \frac{topk}{\text{num experts}}$$

$AI_{dense} \rightarrow$  we do FFN compute topic times, but load num experts times as many weights

\* so this means moe, have low AI compared to dense (will use HW ineffectively)

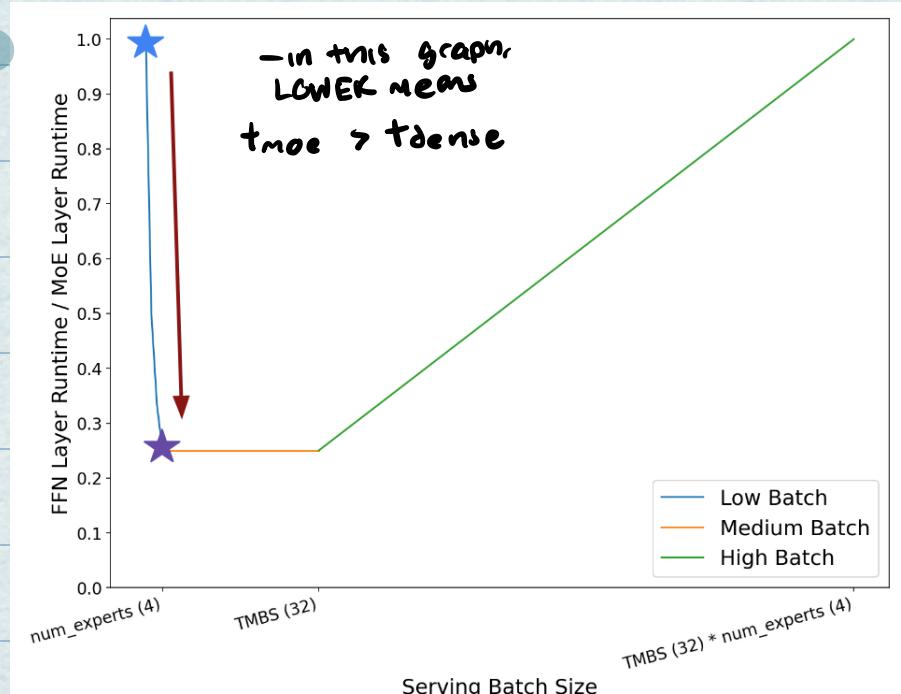
↳ as trevor said, expert-level parallelism can help w/ this!

\* At serving time - batch size can be small which CHANGES speedup calculus

Region 1 = low batch serving

↳ perf limiter will be  
MEMORY to load  
weights

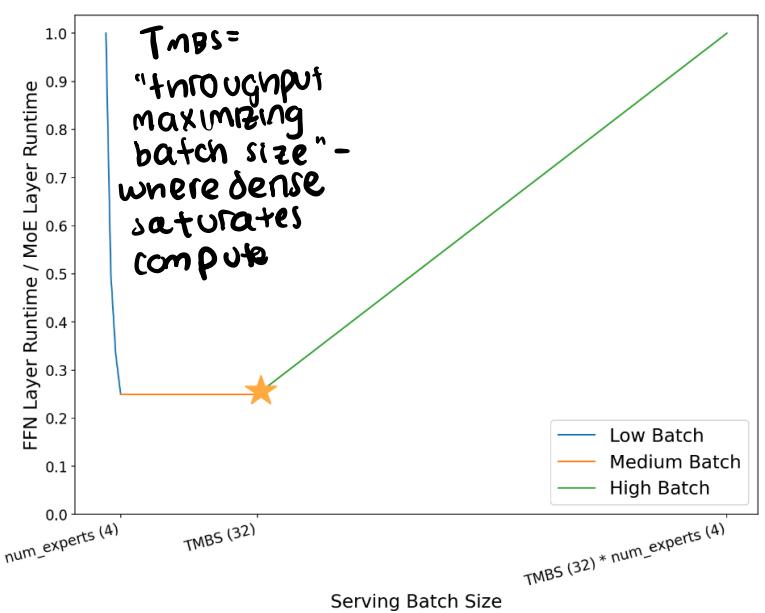
$B=1 \rightarrow \text{MoE / dense}$   
touch some # of weights  
(1 expert's weights are loaded)



$B = \text{num\_experts} \rightarrow$  all experts are touched,  
so MoE loads topic as many weight

Region 2 - medium batch serving

- here  $\frac{t_{\text{dense}}}{t_{\text{moe}}}$  is constant
- as all experts are loaded, but compute utilization for dense ramps up
- recall - moe has LOWER AI (flops/byte) - so not yet saturating compute



## Region 3: High batch

- when  $\text{batch} = \text{TMBS}$ .

$\text{num\_experts}$  -

$\text{MoE}$  will saturate

compute and  $\text{tdense} =$

$\text{t}_{\text{MoE}}$

\* remember this  $\text{MoE}$

has overall many params - so

i) more accurate - so we may be happy

w/ tradeoff even if we're in "inefficient"

region in the middle

