

Compression | Models are getting larger

<u>Alexnet</u>	2012	→ 8 layers	1.4 GFlop	} <u>16x</u>
		~16% Error.		
	<u>2015</u>	→ 152	22.6 GFlop	
		<u>3.5% Error</u>		

The increase in size → increased accuracy
→ need for acceleration

- run a network faster (performance, runtime, inf/sec)
- run a network more efficiently (energy, monetary cost)

Objectives of acceleration inf/s inf/\$

For inference (previous lecture) → cascades
→ replication
(just focus on the fast forward passes forward pass)

Training: → fast back-prop. (gradients)
→ update.

Today's lecture: How do we reduce cost for all three operations?

1) Using reduced precision arithmetic
→ increase the amount of data we can fit into a given processing unit and hence optimize util.

2) Compression \rightarrow pruning. (reduce the ops we are performing)

3) Better algorithm. (fastest alg possible)
Low-rank approx. of fundamental ops in DNNs.

Reduced Precision

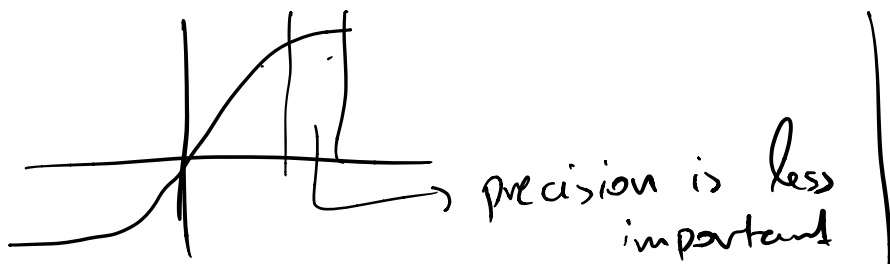
\rightarrow We get reduced storage

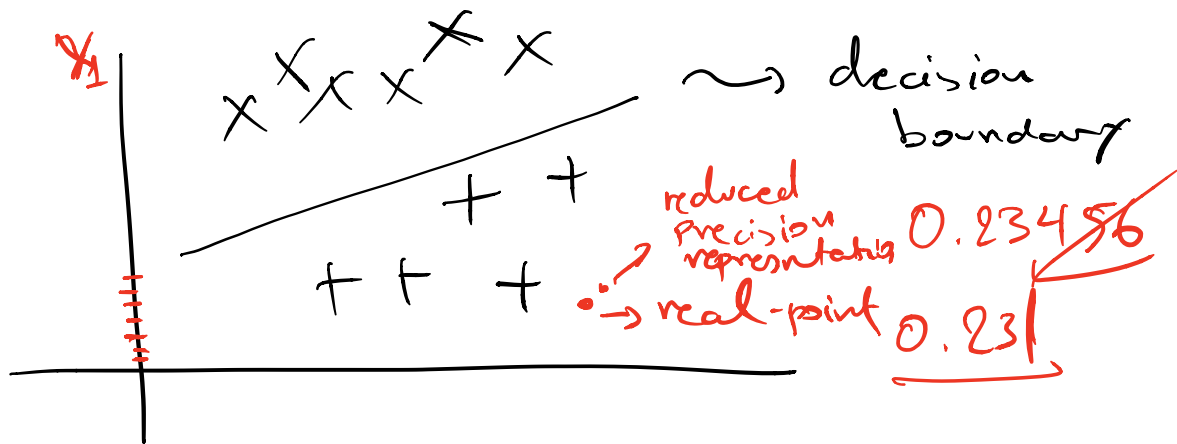
\rightarrow We get reduced energy (because we need less complex hardware "fewer transistors" when we use reduced precision.)

\rightarrow Improved performance (speed) \nearrow

\rightarrow Sometimes: has little effect on the accuracy.

Examples $\max(0, x)$, $\underline{x < 0}$

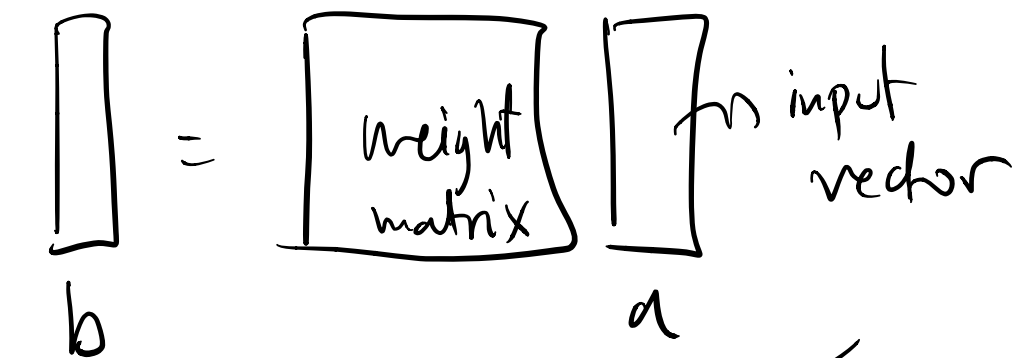




DNN: \rightarrow relu, or other non-linear.

$$b_i = f\left(\sum_j w_{ij} a_j\right)$$

\uparrow output
 \uparrow Params
 \uparrow input variables



\rightarrow Multiplication + Addition.

$$b = f_i(Wa) \quad \text{forward.}$$

\downarrow \downarrow
 Matrix vector

$$W_{ij} = W_{ij} + \eta a_i g_j \quad \text{updates}$$

\downarrow \downarrow
 learning rate Multiplication
 Addition

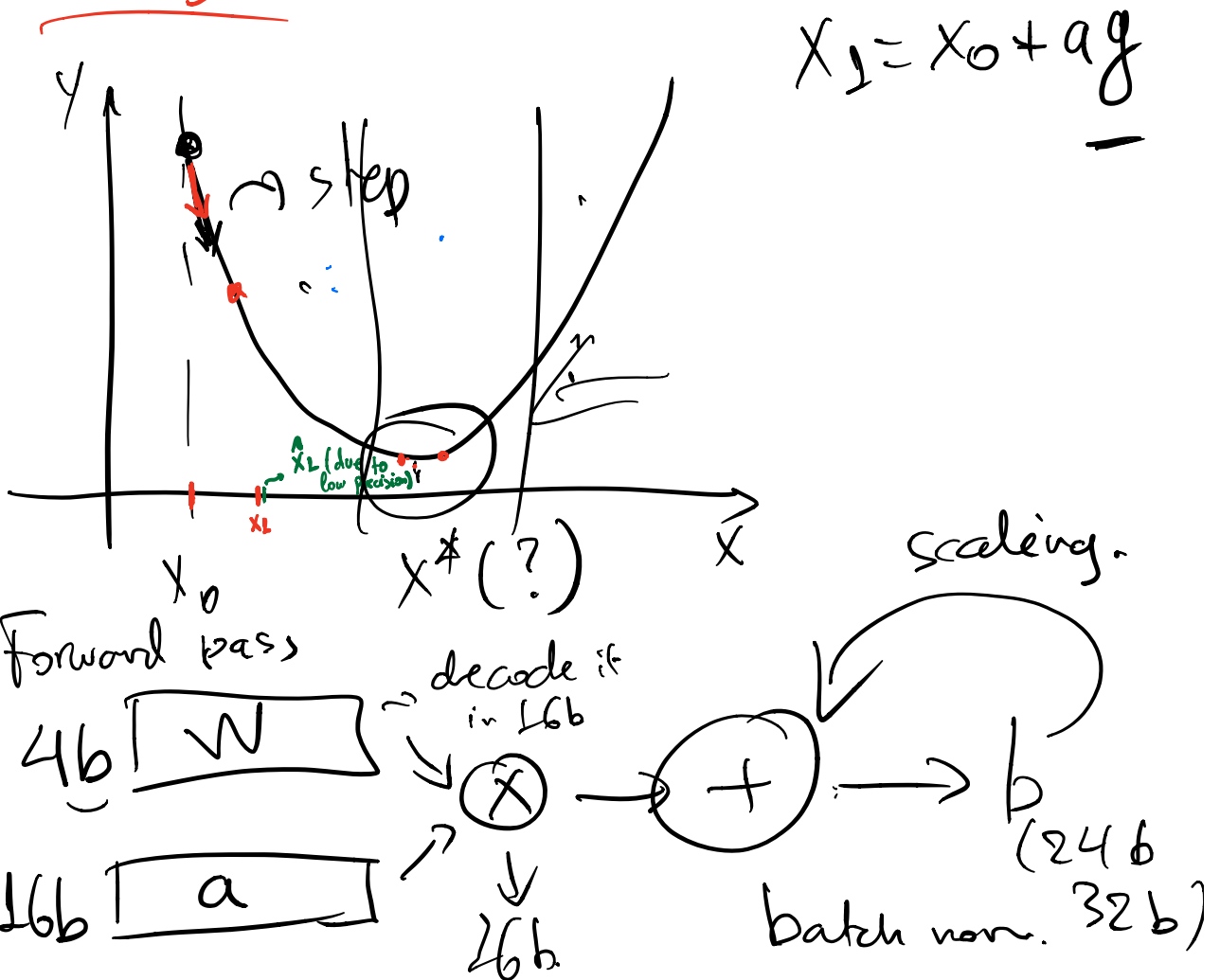
Add, Mult units Range

FP32 (full precision)	10^{-38} .. 10^{38}	6×10^{-6} %
FP16	6×10^{-5} .. 6×10^4	0.05%
INT32	2×10^9	0.5
INT16	6×10^4	0.5
INT8	0-127	0.5

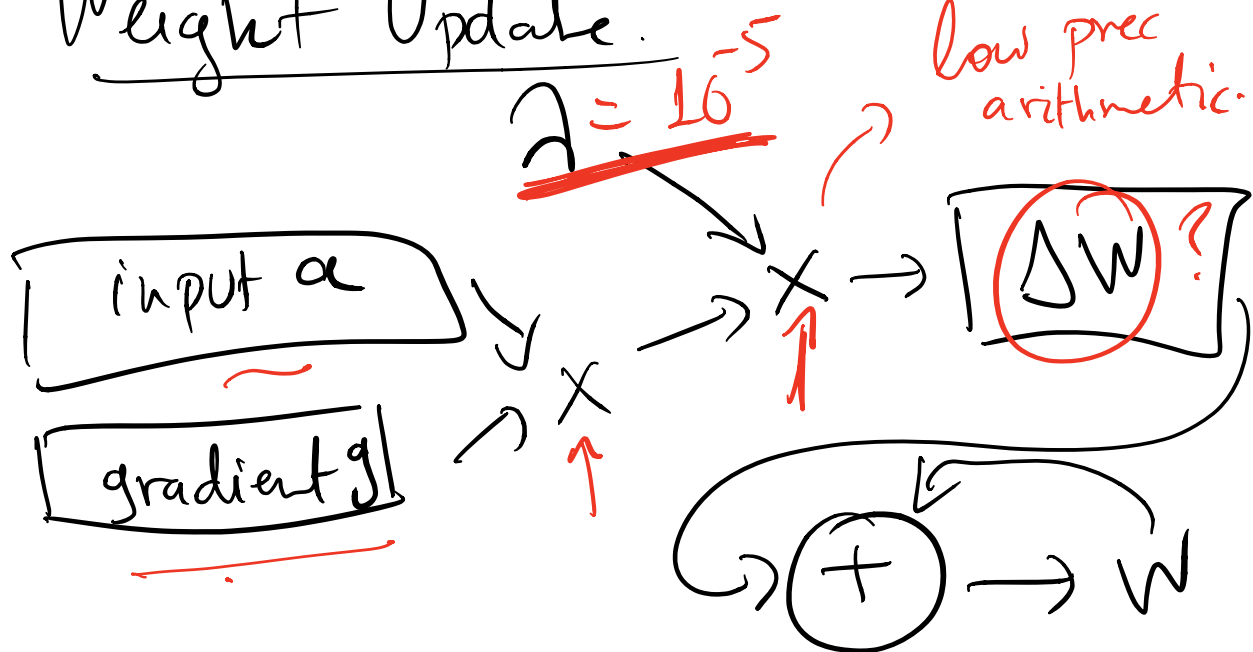
SEM

design.

Mixed precision says:
use different arithmetic for
you inputs / outputs / gradients
weight et.



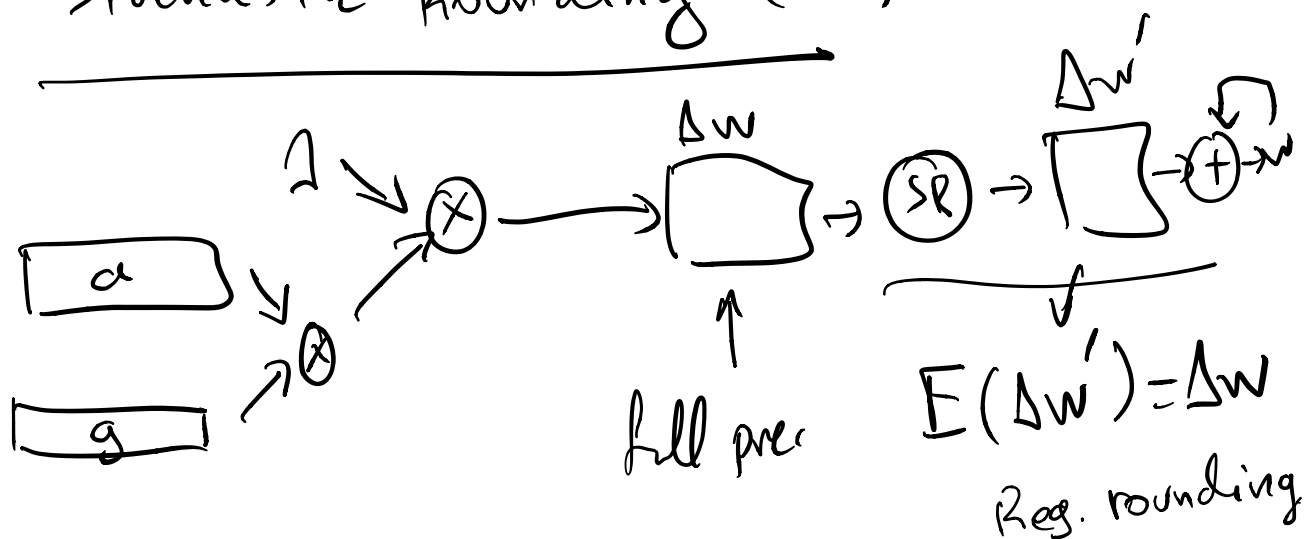
Weight Update.



What is the effect of low precision here (assuming quantization / rounding)

10^{-5} \leftarrow very low learning rate.

To avoid $\Delta W = 0$ due to low-prec.
Stochastic Rounding (SR)



I want to add 0.3 to 0
 100 times.

$0.3 \rightarrow 0$
 \downarrow
 $0 \neq 30$

$$\text{Round}(x) = \begin{cases} \lfloor x \rfloor, & \text{w. prob. } 1 - (x - \lfloor x \rfloor) \\ \lfloor x \rfloor + 1, & \text{w. } x - \lfloor x \rfloor \end{cases}$$

$$\text{Round}(0.3) = \begin{cases} 0, & 70\% \\ 1, & 30\% \end{cases} \quad \underline{\underline{E[\text{Sum}] = 30}}$$

Chris De Sa ISCA '17 discusses
efficient stochastic rounding schemes

Summary

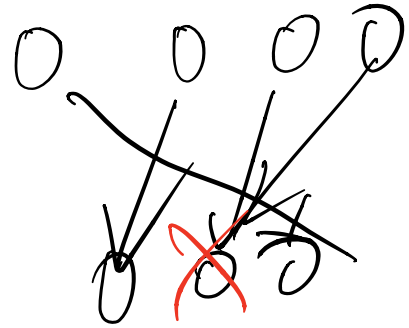
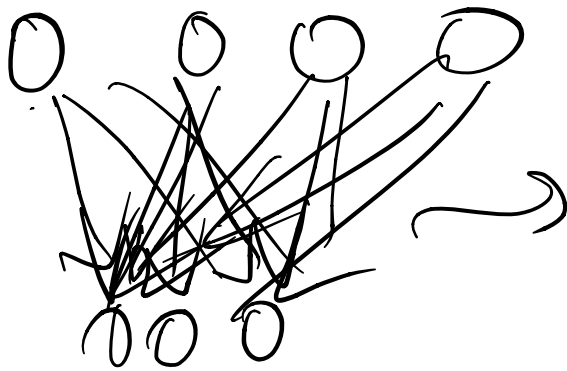
Reduced prec. \rightarrow save memory
space
bandwidth

\rightarrow Important points

High prec \rightarrow Batch norm \rightarrow cert.

Stochastic Rounding during use.

2. Pruning



after training.

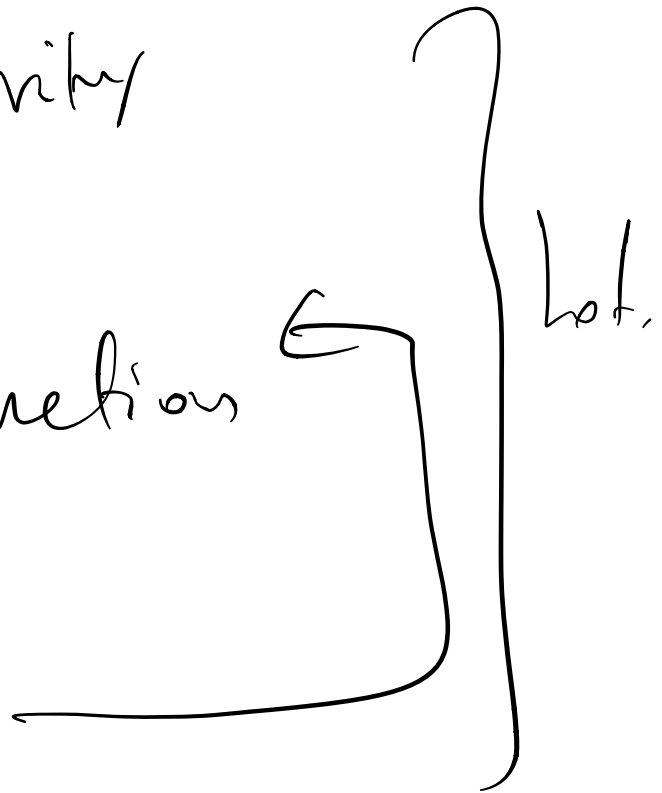
Train ^{w. Full} Connectivity



Prune Connections



Retrain

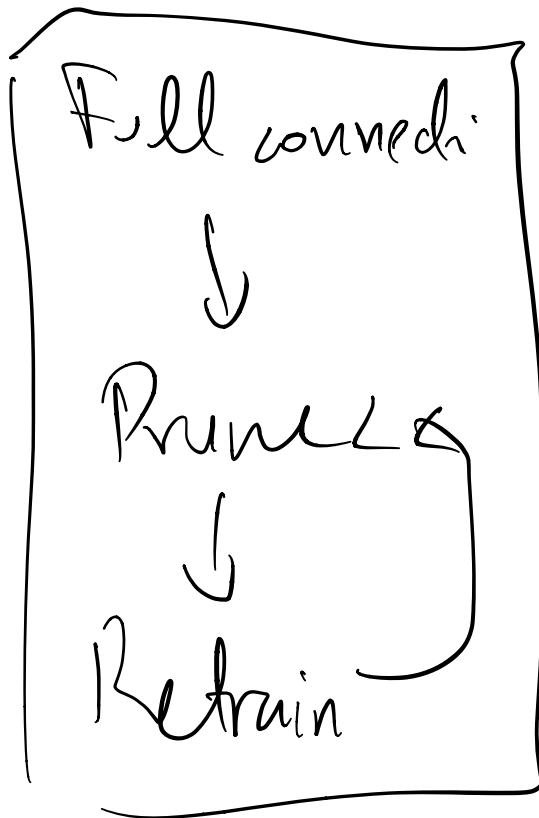


After pruning (whatever params are left)

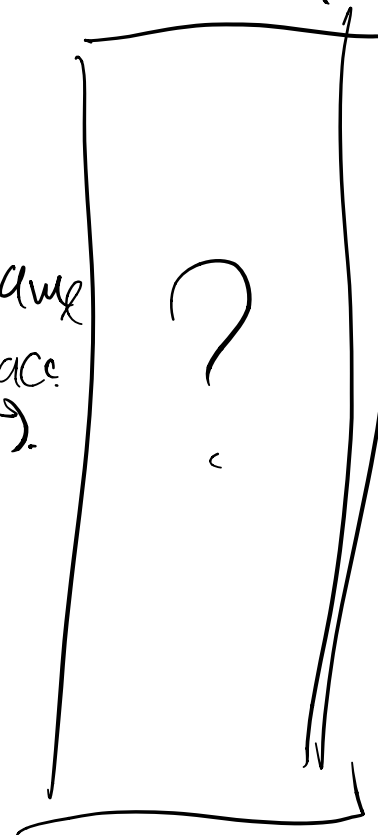
→ Quantize

Quantize

original network →



same case →



Pruning

reduces the size of my net.

Want to group weights w.
similar values into one
weight.

2.09, 2.12, 1.92

1.87 -

assigns these params into
a cluster $k \rightarrow$ 2.00

int 0 \downarrow

\downarrow FP26

4x FP16 vs. 4x int8 + FP16

This requires changing how one performs updates and learning.

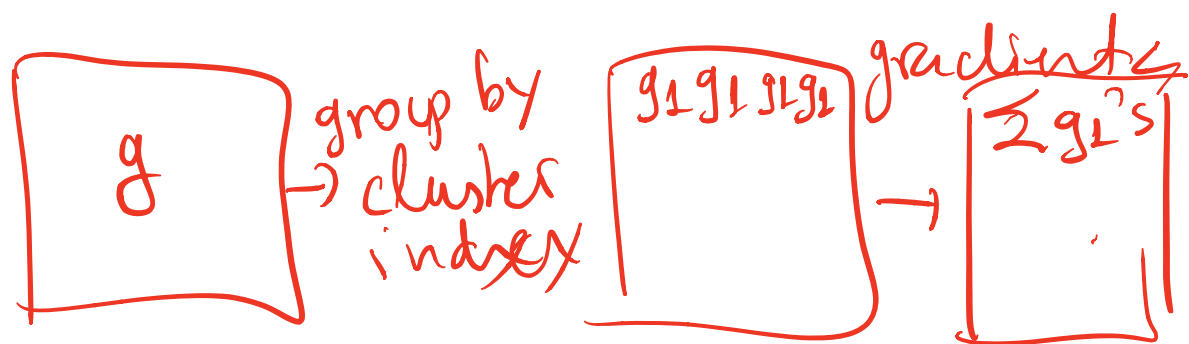
$$\boxed{W \text{ Matrix}} = \boxed{W \text{ Matrix}} + 2 \begin{matrix} \boxed{} \\ \boxed{} \\ \boxed{} \end{matrix} \begin{matrix} \boxed{g} \\ \boxed{g} \\ \boxed{g} \end{matrix}$$

cluster index a centroids.

$$\boxed{W \text{ Matrix}} = \boxed{\frac{\text{int } 8}{\text{int } 1}} \begin{matrix} \boxed{} \\ \boxed{} \\ \boxed{} \end{matrix} \rightarrow k \text{ FP16}$$

Now I only update the cluster
centroids and the indexes

they don't care about



Some W 's are changing frequently

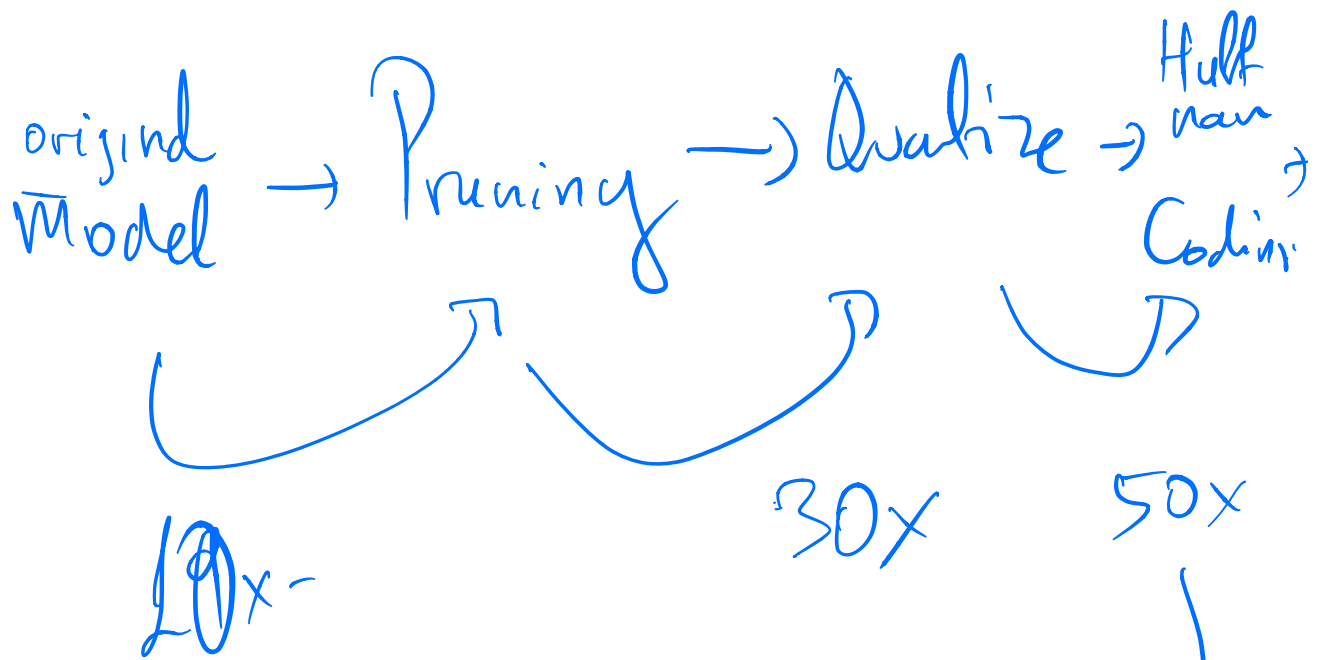
What does it mean that W is changing frequently?

from an opt perspective. (we have not converged)

freq. weights \rightarrow less bits

in-freq weights \rightarrow increased bits

(Huffman Coding)



SqueezeNet (great paper)