



Institución
Universitaria
Reacreditada en Alta Calidad

Quantization – Cuantificación

INTEGRACIÓN DE ML EN EMBEBIDOS Y EDGE COMPUTING

Somos Innovación Tecnológica con *Sentido Humano*



Alcaldía de Medellín

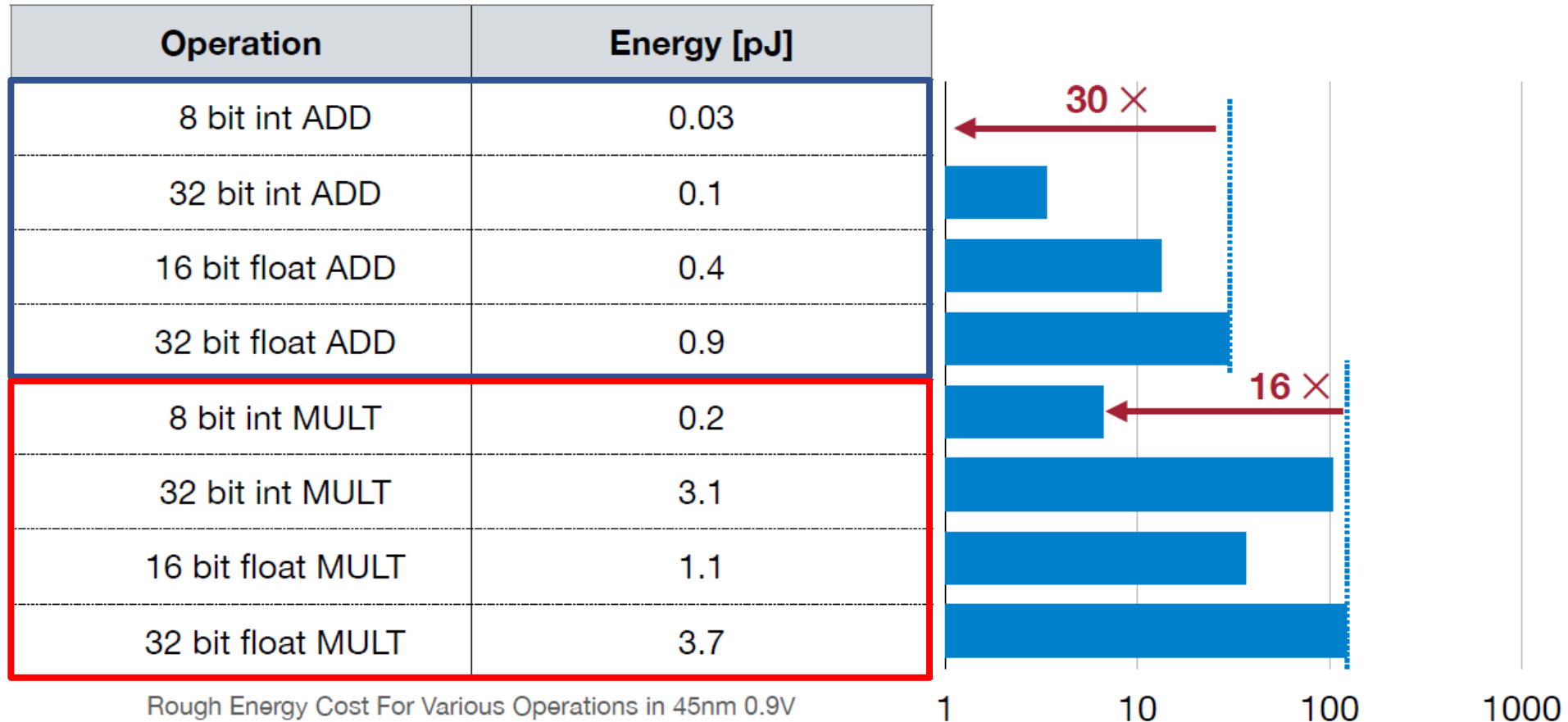


Institución
Universitaria
Reacreditada en Alta Calidad

Contenido

1. Costos de operaciones
2. Cuantificación
3. Cuantificación basada en K-means

Costos de Operaciones con bits



1  = 200 × +

Unsigned Integer

- n -bit Range: $[0, 2^n - 1]$

Signed Integer

- Sign-Magnitude Representation
 - n -bit Range: $[-2^{n-1} - 1, 2^{n-1} - 1]$
 - Both 000...00 and 100...00 represent 0

Two's Complement Representation

- n -bit Range:
- 000...00 represents 0
- 100...00 represents -2^{n-1}

Enteros

0	0	1	1	0	0	0	1
---	---	---	---	---	---	---	---

× × × × × × × ×

$$2^7 + 2^6 + 2^5 + 2^4 + 2^3 + 2^2 + 2^1 + 2^0 = 49$$

Sign Bit

1	0	1	1	0	0	0	1
---	---	---	---	---	---	---	---

× × × × × × × ×

$$- 2^6 + 2^5 + 2^4 + 2^3 + 2^2 + 2^1 + 2^0 = -49$$

1	1	0	0	1	1	1	1
---	---	---	---	---	---	---	---

× × × × × × × ×

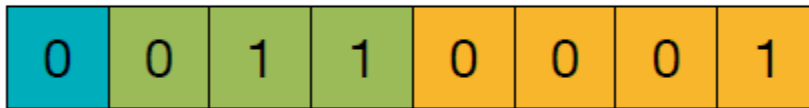
$$-2^7 + 2^6 + 2^5 + 2^4 + 2^3 + 2^2 + 2^1 + 2^0 = -49$$

Fixed-Point Number



Integer . Fraction

“Decimal” Point

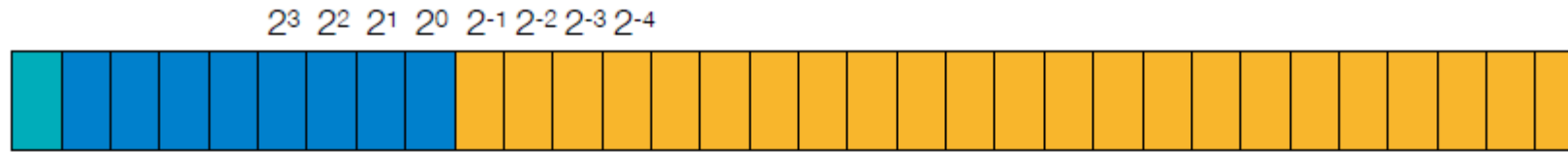


$$\begin{array}{cccccccc} \times & \times & \times & \times & \times & \times & \times & \times \\ -2^3 & +2^2 & +2^1 & +2^0 & +2^{-1} & +2^{-2} & +2^{-3} & +2^{-4} \end{array} = 3.0625$$



$$\begin{array}{cccccccc} \times & \times & \times & \times & \times & \times & \times & \times \\ (-2^7 & +2^6 & +2^5 & +2^4 & +2^3 & +2^2 & +2^1 & +2^0) \end{array} \times 2^{-4} = 49 \times 0.0625 = 3.0625$$

Floating-Point Number



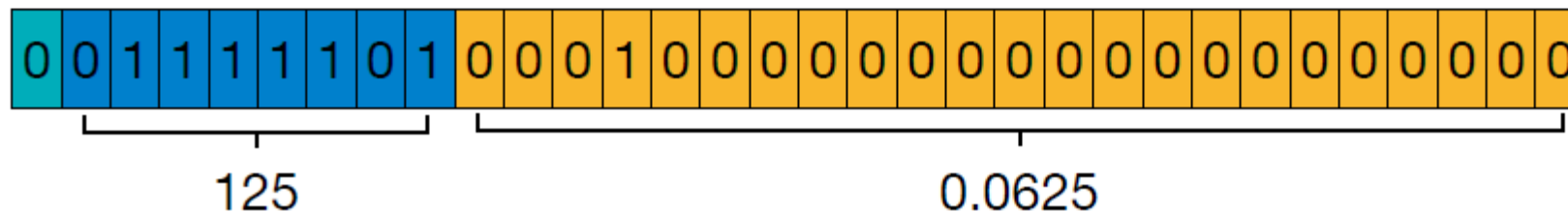
Sign 8 bit Exponent

23 bit Fraction

$$(-1)^{\text{sign}} \times (1 + \text{Fraction}) \times 2^{\text{Exponent}-127} \quad \leftarrow \quad \text{Exponent Bias} = 127 = 2^{8-1}-1$$

(significant / mantissa)

$$0.265625 = 1.0625 \times 2^{-2} = (1 + \underline{0.0625}) \times 2^{\underline{125}-127}$$



Floating-Point Number

Exponent Width → Range; Fraction Width → Precision

IEEE 754 Single Precision 32-bit Float (IEEE FP32)



IEEE 754 Half Precision 16-bit Float (IEEE FP16)



Google Brain Float (BF16)



Exponent (bits)	Fraction (bits)	Total (bits)
8	23	32
5	10	16
8	7	16

Floating-Point Number

Exponent Width → Range; Fraction Width → Precision

IEEE 754 Single Precision 32-bit Float (IEEE FP32)



IEEE 754 Half Precision 16-bit Float (IEEE FP16)



Google Brain Float (BF16)



Exponent (bits)	Fraction (bits)	Total (bits)
8	23	32
5	10	16
8	7	16



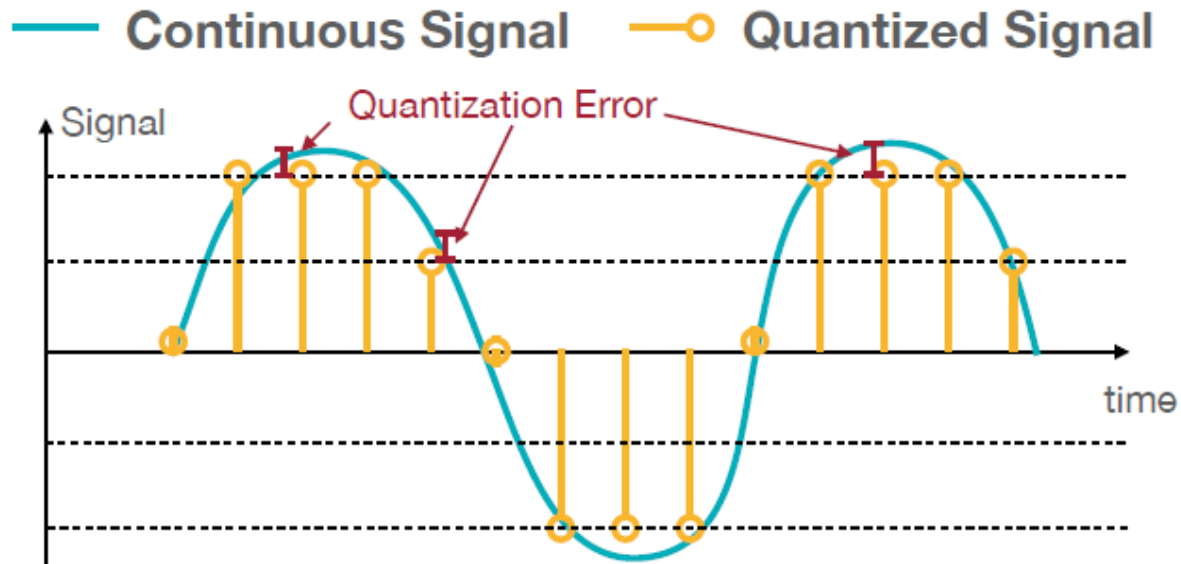
Institución
Universitaria
Reacreditada en Alta Calidad

Contenido

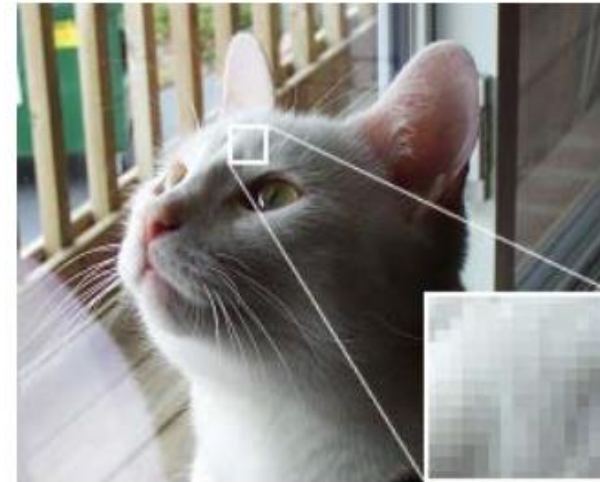
1. Costos de operaciones
2. Cuantificación
3. Cuantificación basada en K-means

Que es cuantificación?

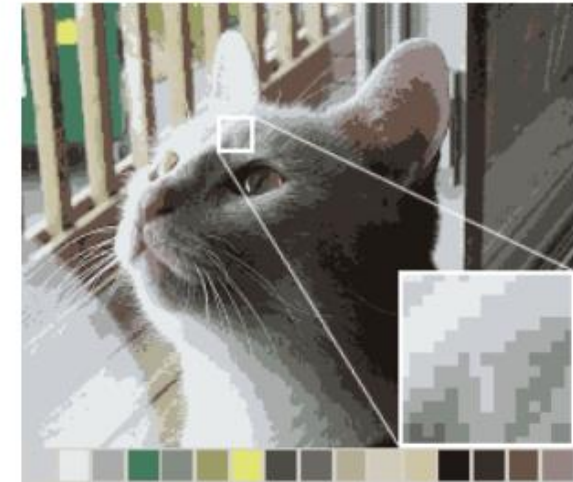
La cuantificación es el proceso de restringir una entrada desde un conjunto continuo o con valores grandes a un conjunto discreto.



Original Image



16-Color Image



Images are in the public domain.

“Palettization”

Neural Network Quantization

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

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

1	-2	0	-1
-1	-1	-2	1
-2	1	-1	-2
1	-1	0	0

$(-1) \times 1.07$

1	0	1	1
1	0	0	1
0	1	1	0
1	1	1	1

**K-Means-based
Quantization**

**Linear
Quantization**

**Binary/Ternary
Quantization**

Almacenamiento

Calculo

Pesos de punto flotante	Pesos Enteros; Libro de código de punto flotante	Pesos Enteros	Pesos binarios/ternarios (-1,0,+1)
Aritmética de punto flotante	Aritmética de punto flotante	Aritmética Entera	Operaciones de bits



Institución
Universitaria
Reacreditada en Alta Calidad

Contenido

1. Costos de operaciones
2. Cuantificación
3. Cuantificación basada en K-means

Weight Quantization

weights
(32-bit float)

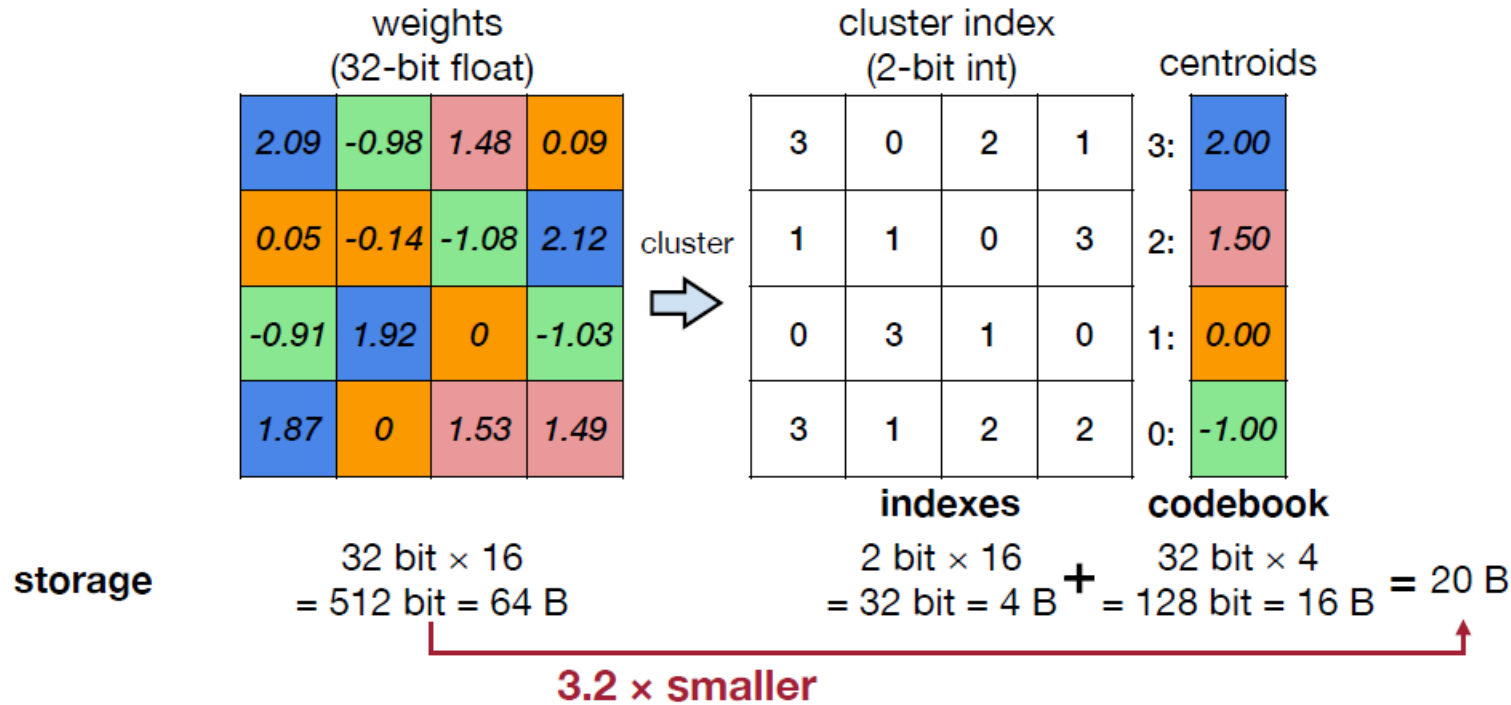
2.09	-0.98	1.48	0.09
0.05	-0.14	-1.08	2.12
-0.91	1.92	0	-1.03
1.87	0	1.53	1.49

~~2.09, 2.12, 1.92, 1.87~~



2.0

K-means Weight Quantization



reconstructed weights
(32-bit float)

2.00	-1.00	1.50	0.00
0.00	0.00	-1.00	2.00
-1.00	2.00	0.00	-1.00
2.00	0.00	1.50	1.50

quantization error

0.09	0.02	-0.02	0.09
0.05	-0.14	-0.08	0.12
0.09	-0.08	0	-0.03
-0.13	0	0.03	-0.01

Assume N -bit quantization, and #parameters = $M \gg 2^N$.

$$32 \text{ bit} \times M \\ = 32M \text{ bit}$$

$$N \text{ bit} \times M \\ = NM \text{ bit}$$

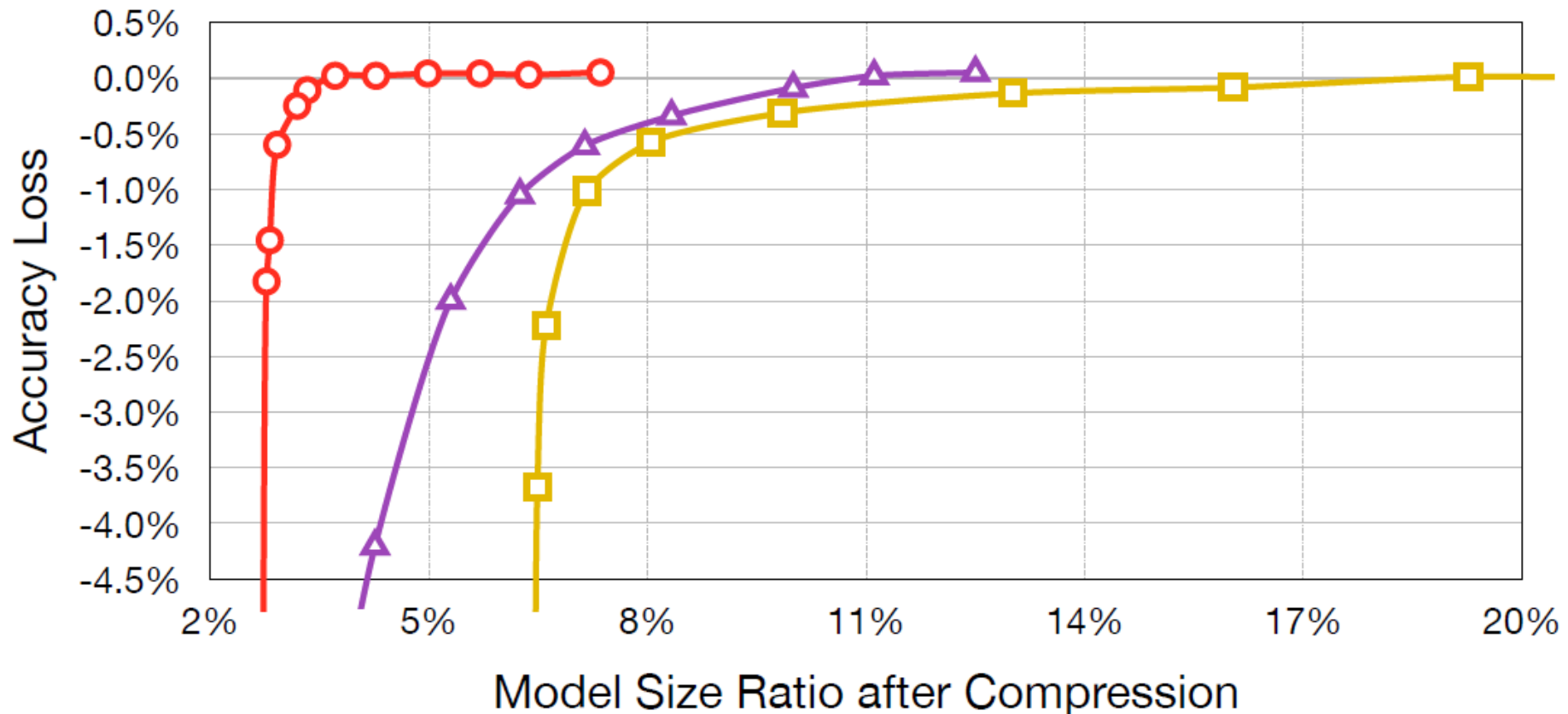
~~$$32 \text{ bit} \times 2^N \\ = 2^{N+5} \text{ bit}$$~~

32/ N \times smaller

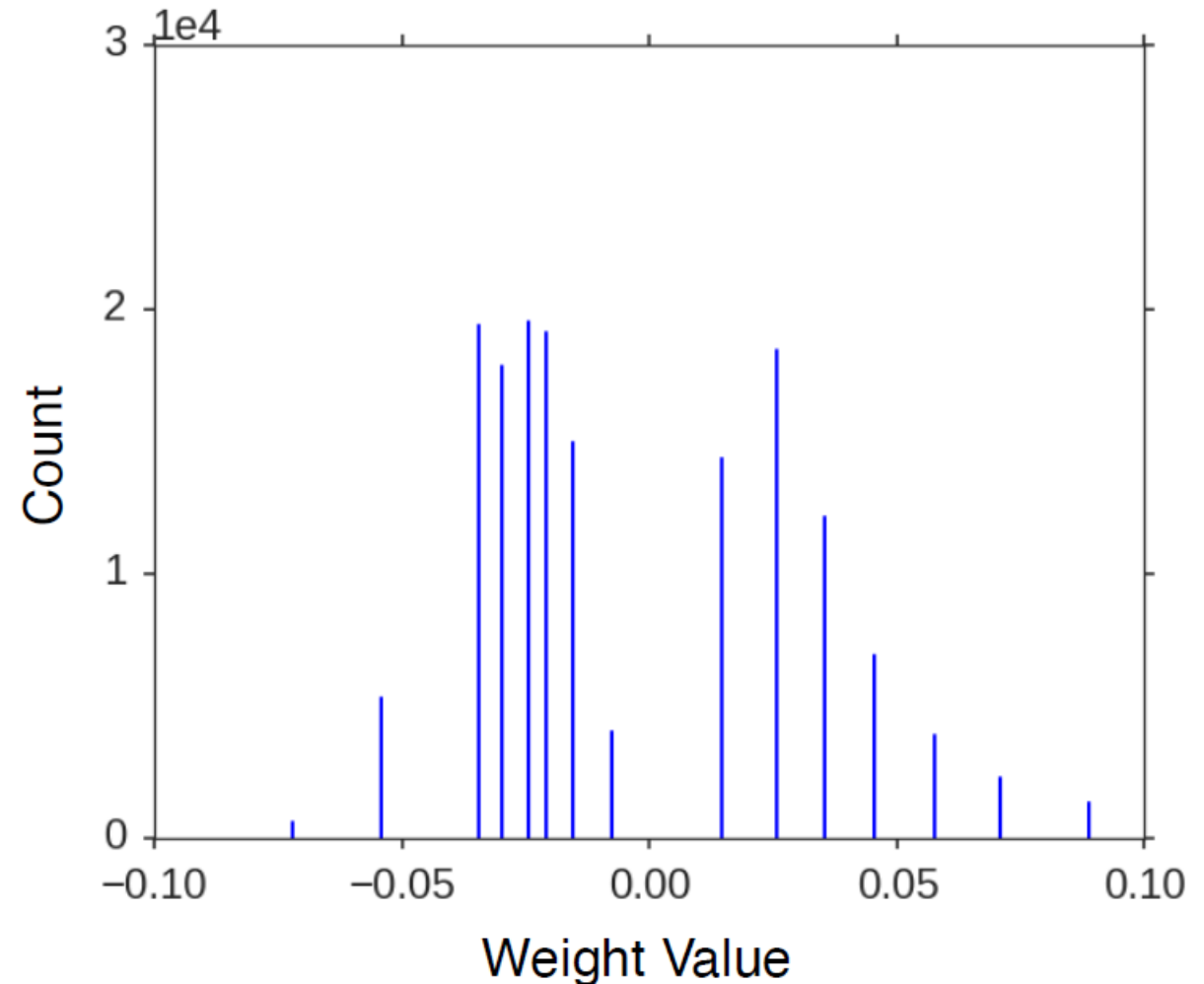
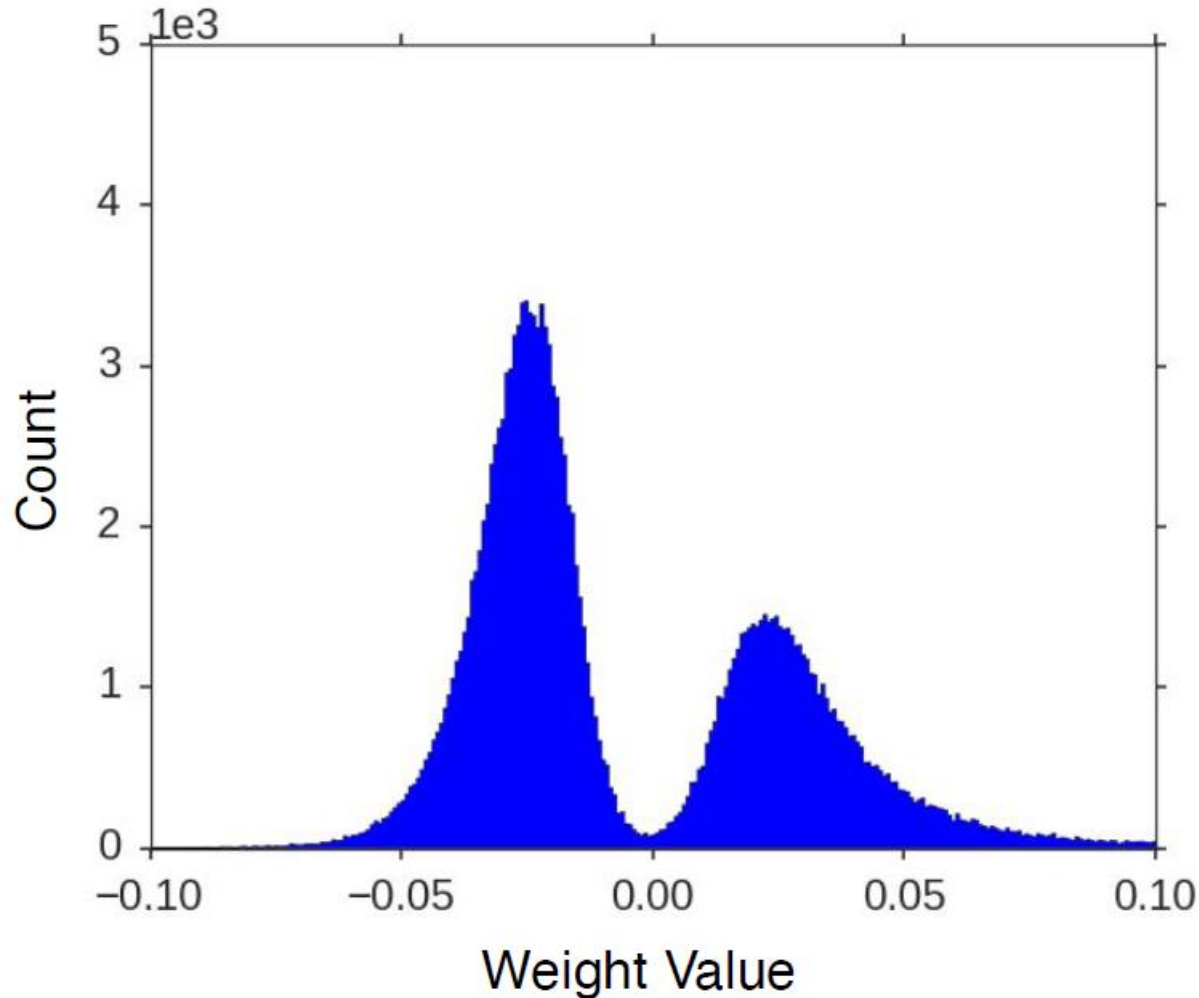
K-means Weight Quantization

Accuracy vs. compression rate for AlexNet on ImageNet dataset

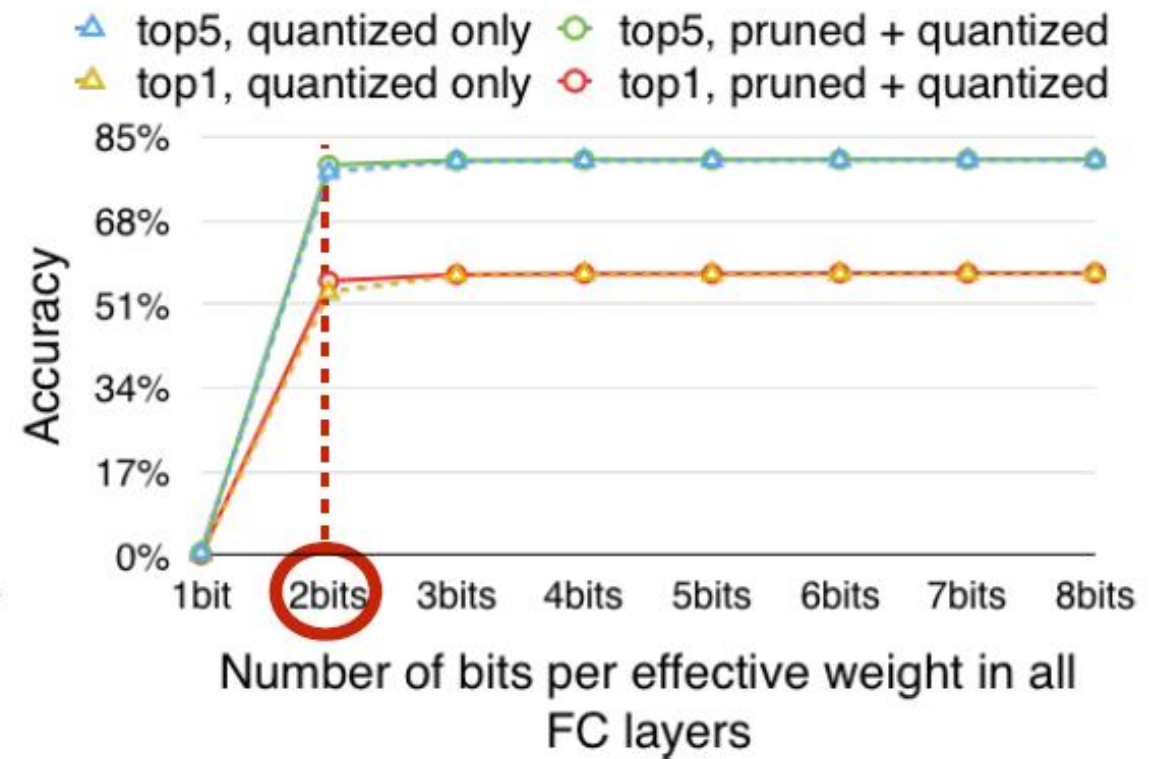
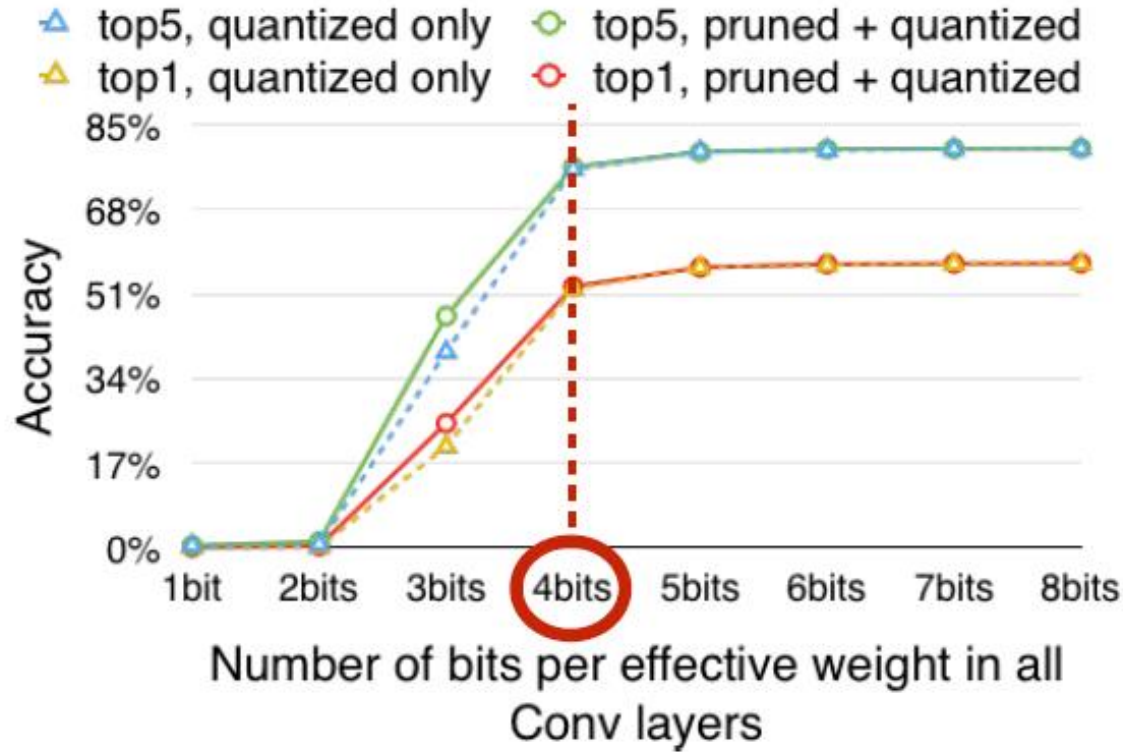
○ Pruning + Quantization △ Pruning Only □ Quantization Only



K-means Weight Quantization

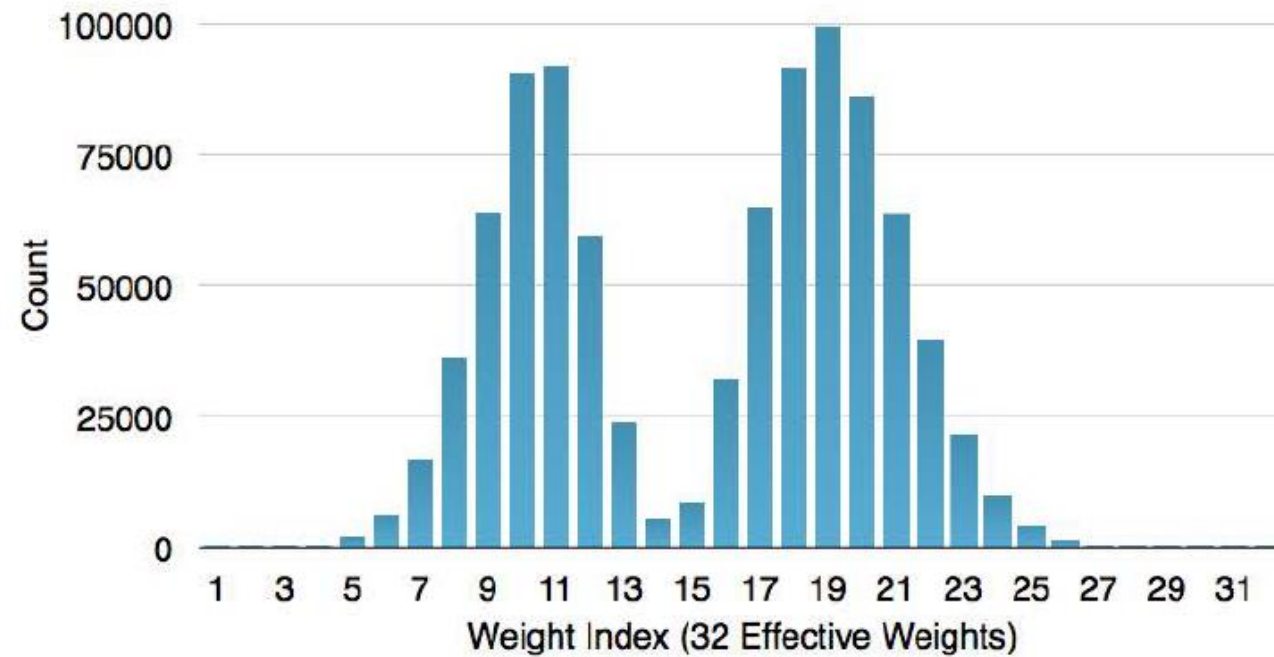
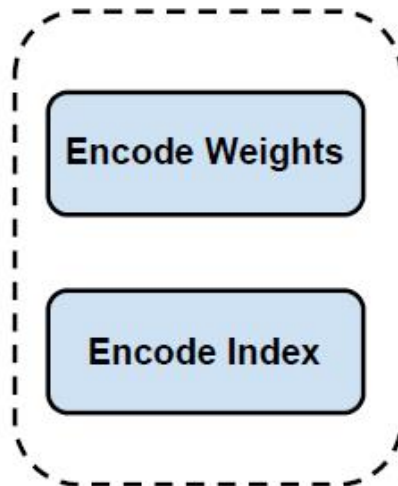


Cuantos bits se necesitan?

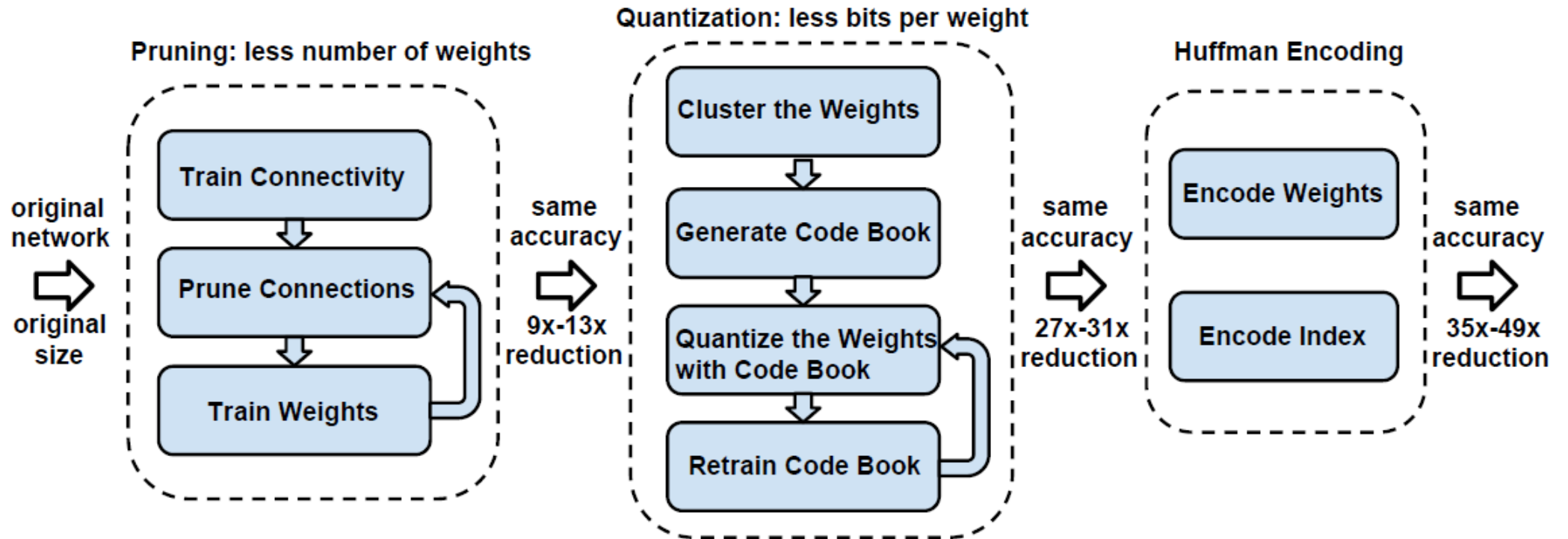


Huffman Coding

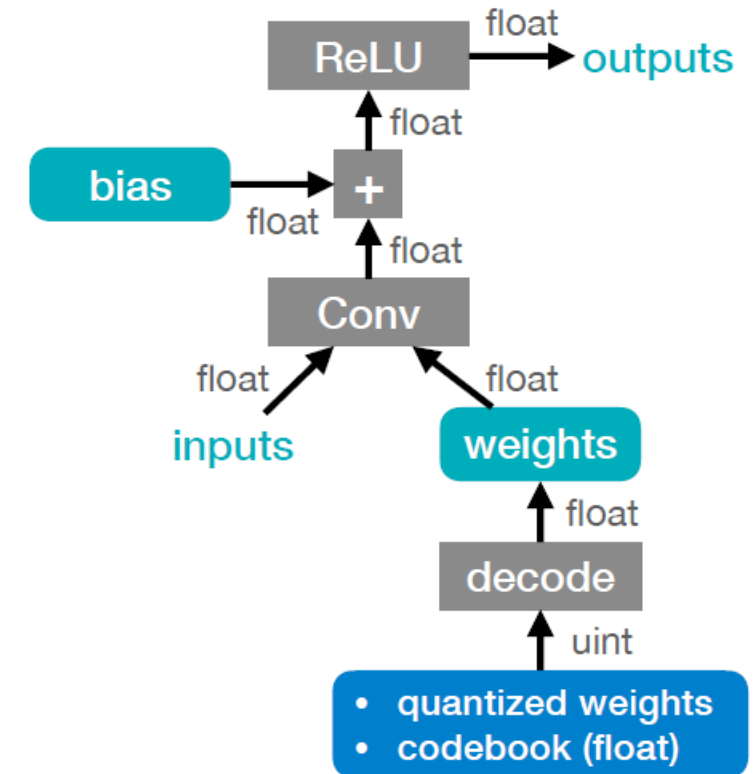
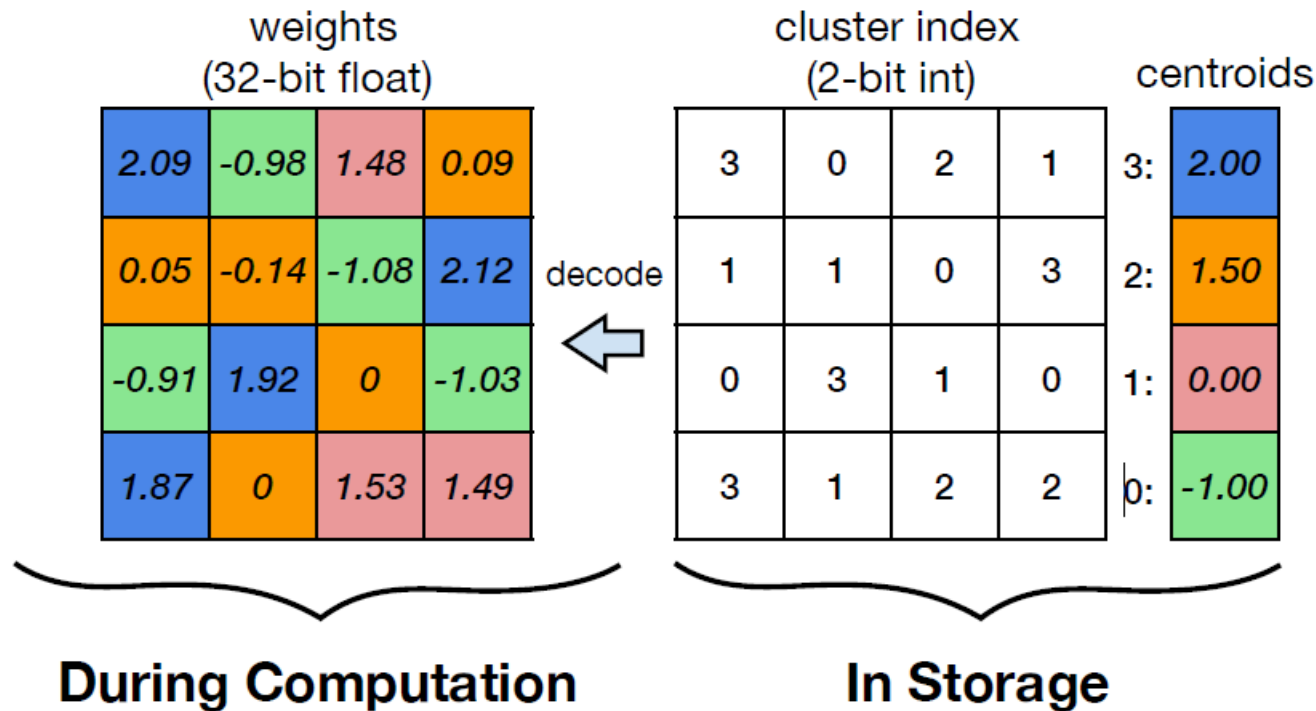
Huffman Encoding



Resumen de Deep Compression



K-means-based Weight Quantization



Los pesos son descomprimidos usando una lookup table (por ejemplo codebook) durante la inferencia. La cuantificación de pesos basada en K-Means sólo ahorra costes de almacenamiento de un modelo de red neuronal.

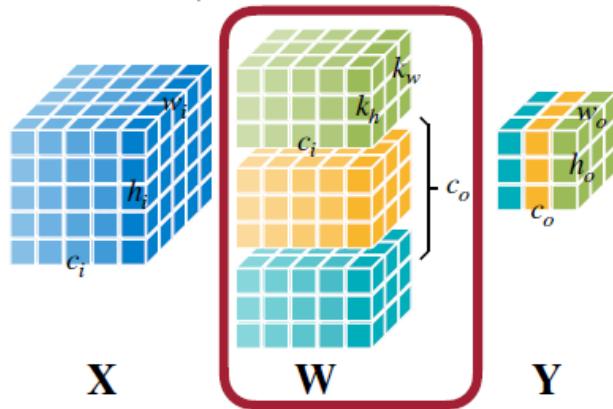
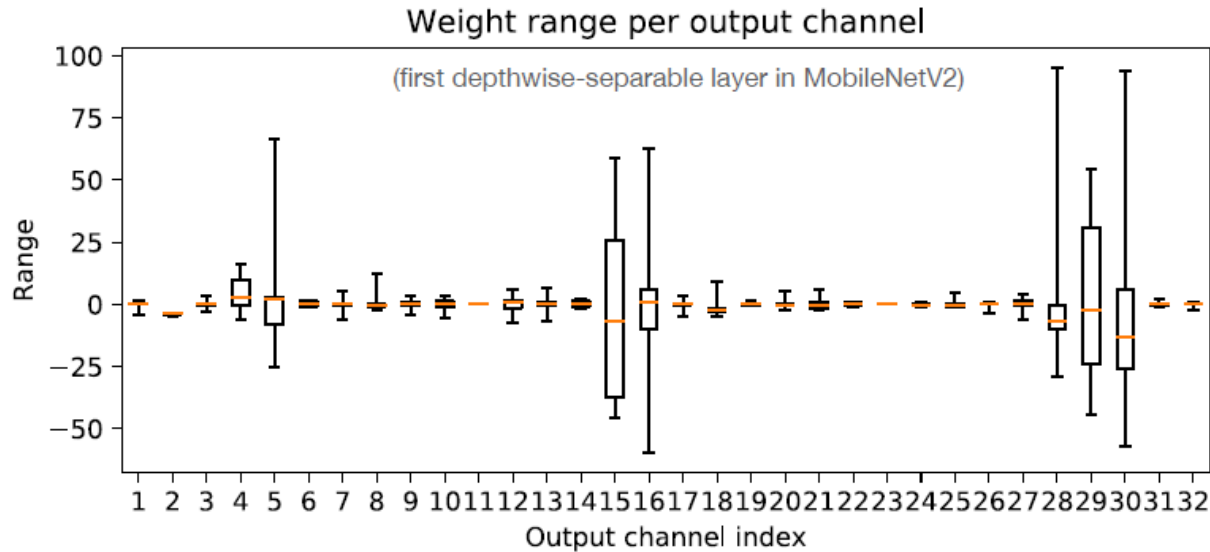
- Todo el computo y accesos de memoria son aun de punto flotante.



Contenido

1. Costos de operaciones
2. Cuantificación
3. Cuantificación basada en K-means
4. **Post-training quantization**

Post-training Quantization



- $|r|_{\max} = |W|_{\max}$
- Usando una sola escala S para todos los pesos del tensor (**Per-Tensor Quantization**)
Funciona bien para modelos grandes.
El acierto cae para modelos pequeños.
- Falla en casos como
Diferencias muy grandes (mas de 100x) en los rangos de los pesos para diferentes canales de salida – pesos atípicos.
- Solución: **Per-channel Quantization**

Post-training Quantization

Per-Channel Quantization

OC	ic					
	2.09	-0.98	1.48	0.09	$ r _{\max} = 2.09$	$S_0 = 2.09$
	0.05	-0.14	-1.08	2.12	$ r _{\max} = 2.12$	$S_1 = 2.12$
	-0.91	1.92	0	-1.03	$ r _{\max} = 1.92$	$S_2 = 1.92$
	1.87	0	1.53	1.49	$ r _{\max} = 1.87$	$S_3 = 1.87$

1	0	1	0	2.09	0	2.09	0
0	0	-1	1	0	0	-2.12	2.12
0	1	0	-1	0	1.92	0	-1.92
1	0	1	1	1.87	0	1.87	1.87

Quantized

Reconstructed

$$\|W - S \odot q_W\|_F = 2.08$$

Per-Tensor Quantization

$$|r|_{\max} = 2.12$$

$$S = \frac{|r|_{\max}}{q_{\max}} = \frac{2.12}{2^{2-1} - 1} = 2.12$$

1	0	1	0	2.12	0	2.12	0
0	0	-1	1	0	0	-2.12	2.12
0	1	0	0	0	2.12	0	0
1	0	1	1	2.12	0	2.12	2.12

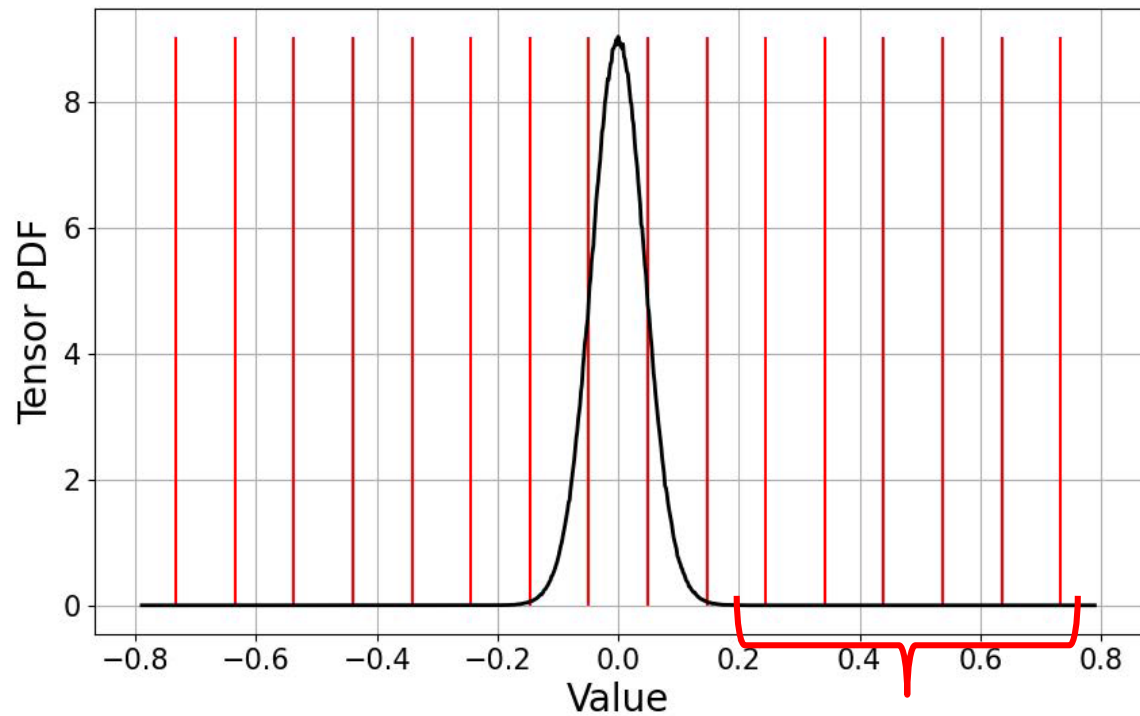
Quantized

Reconstructed

$$\|W - Sq_W\|_F = 2.28$$

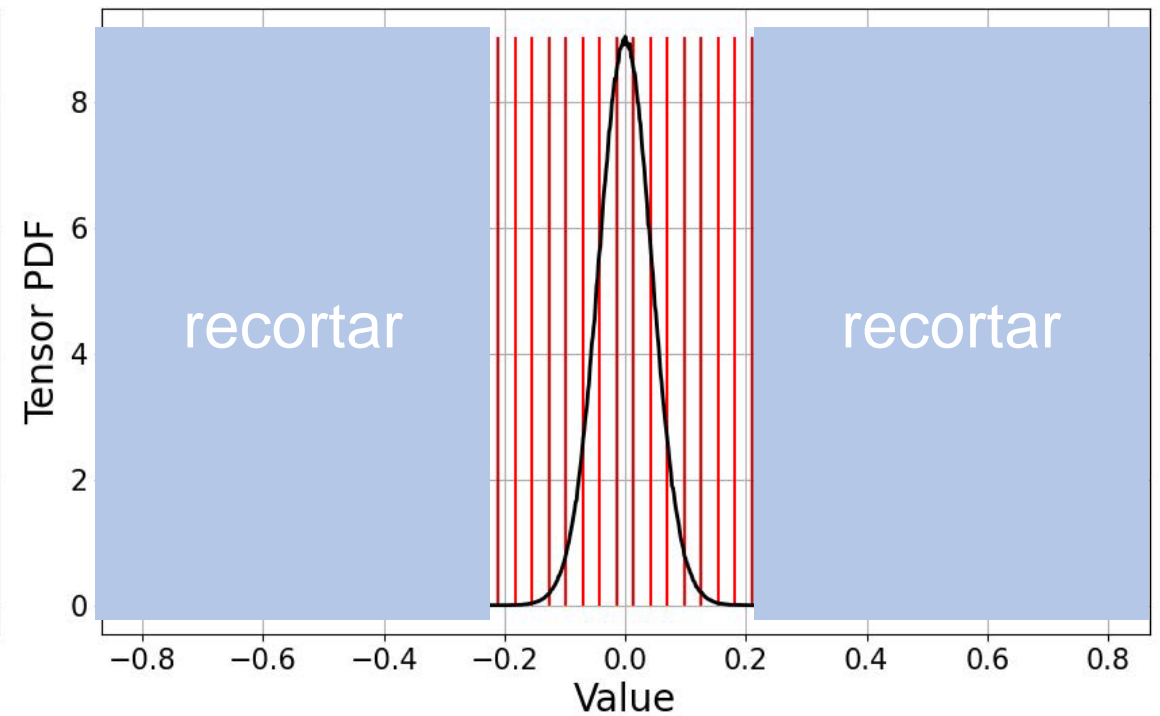
Cuantificación de rango dinámico

Cuantificación de toda la escala



Valores de baja densidad

Cuantificación recortada





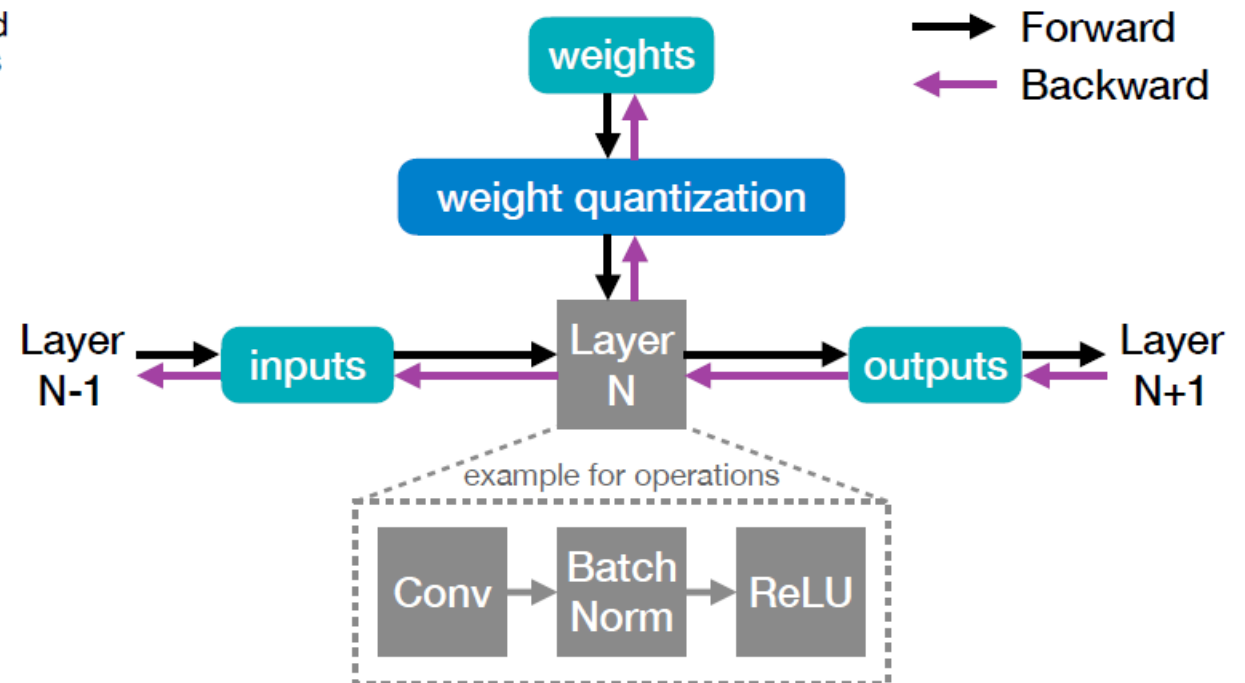
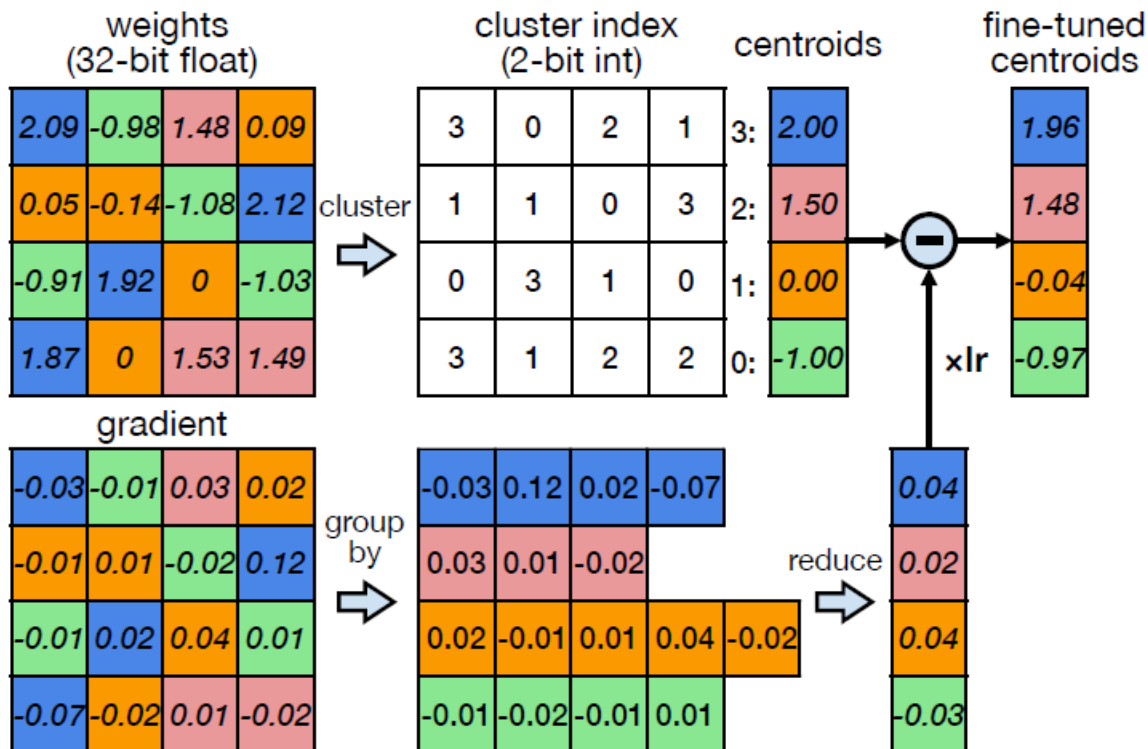
Contenido

1. Costos de operaciones
2. Cuantificación
3. Cuantificación basada en K-means
4. Post-training quantization
5. **Quantization-Aware Training**

Quantization-Aware Training

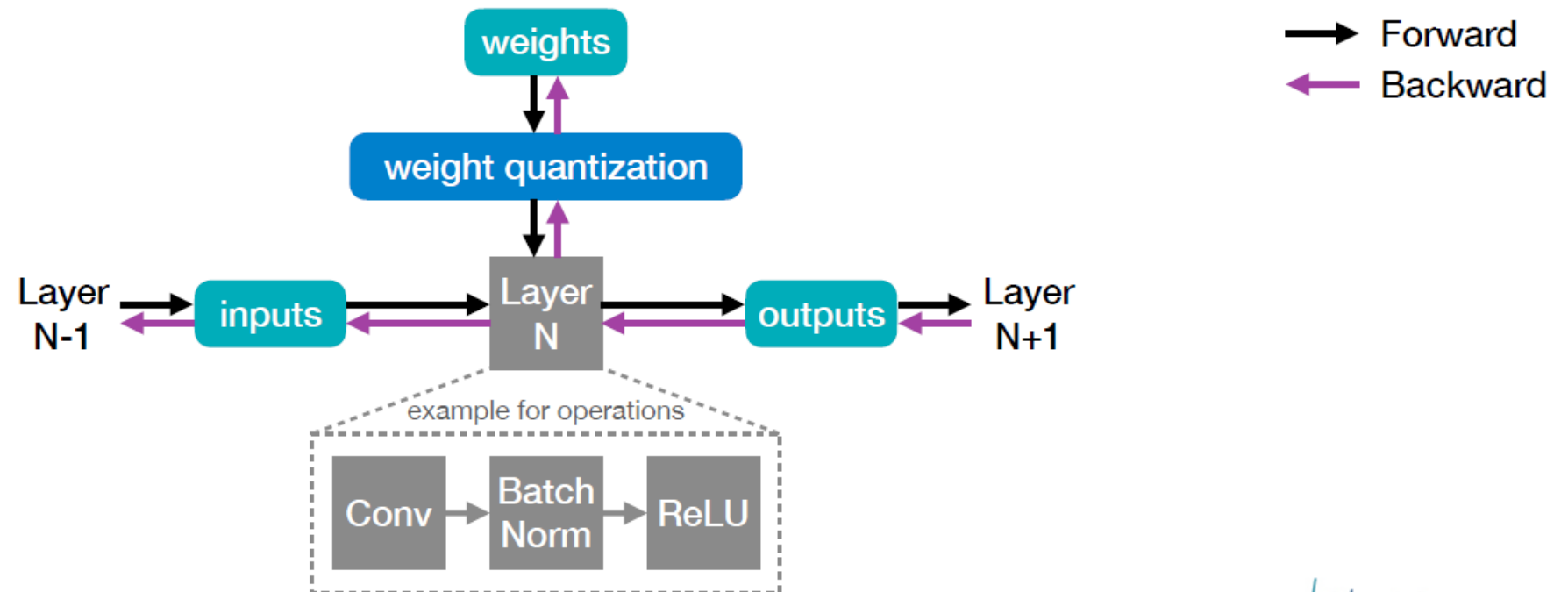
Entrenar el modelo considerando la cuantificación.

- Minimizar la pérdida del acierto, especialmente cuantificaciones agresivas con 4 bits o menos.
- Usualmente, fine-tuning sobre un modelo pre-entrenado con punto flotante provee un mayor acierto que entrenando desde cero.



Quantization-Aware Training

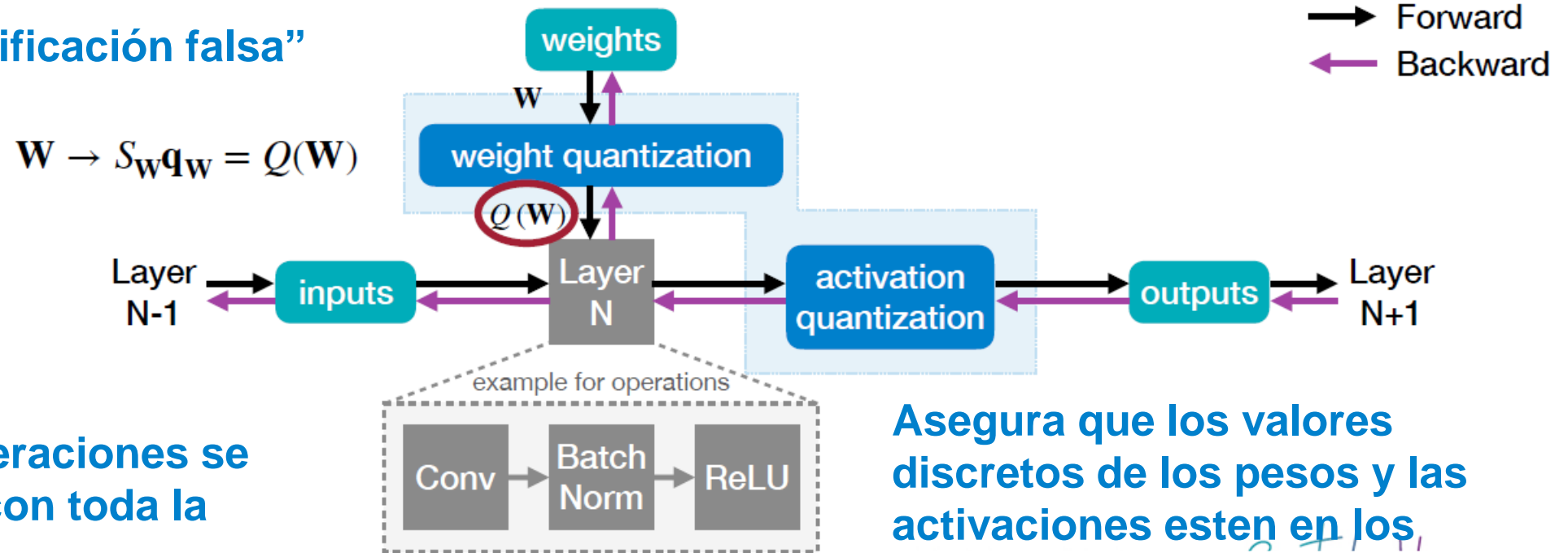
- Una copia precisa de los pesos se mantiene durante el entrenamiento.
- Los pequeños gradientes son acumulados sin pérdida de precisión.
- Una vez el modelo es entrenado, solo los pesos cuantificados son usados para la inferencia.



Quantization-Aware Training

- Una copia precisa de los pesos se mantiene durante el entrenamiento.
- Los pequeños gradientes son acumulados sin pérdida de precisión.
- Una vez el modelo es entrenado, solo los pesos cuantificados son usados para la inferencia.

“Simulada/Cuantificación falsa”



Estas operaciones se realizan con toda la precision

Asegura que los valores discretos de los pesos y las activaciones estén en los límites



Institución
Universitaria
Reacreditada en Alta Calidad

Bibliografía

- Han, Efficient Deep Learning - Lecture 5, Quantization.



Institución
Universitaria
Reacreditada en Alta Calidad

¡Gracias!

Somos Innovación Tecnológica con *Sentido Humano*



Alcaldía de Medellín