

# Getting *free* Bits Back from Rotational Symmetries in LLMs

Jiajun He, Gergely Flamich, José Miguel Hernández-Lobato

University of Cambridge

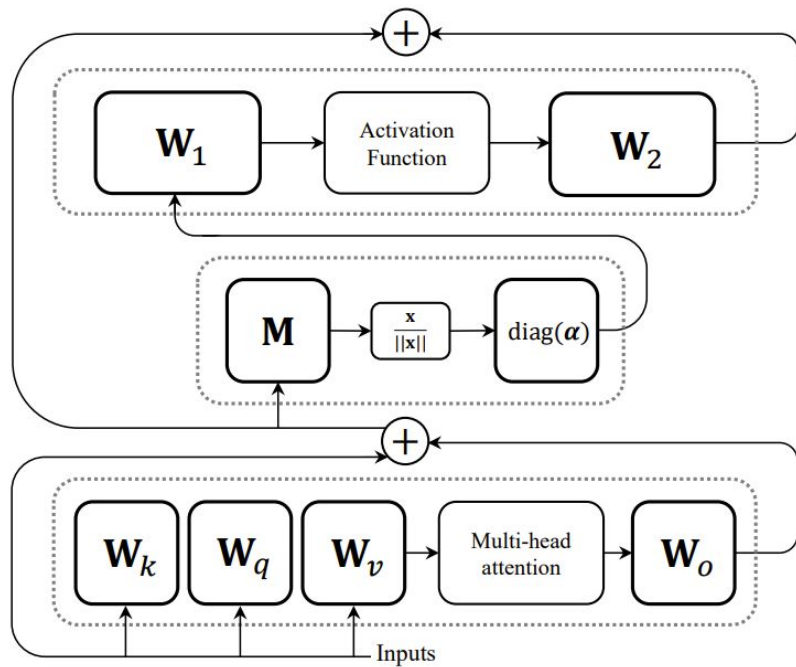
After pruning LLMs,

we can further save 3-5% additional bits *for free*

in storage and transmission

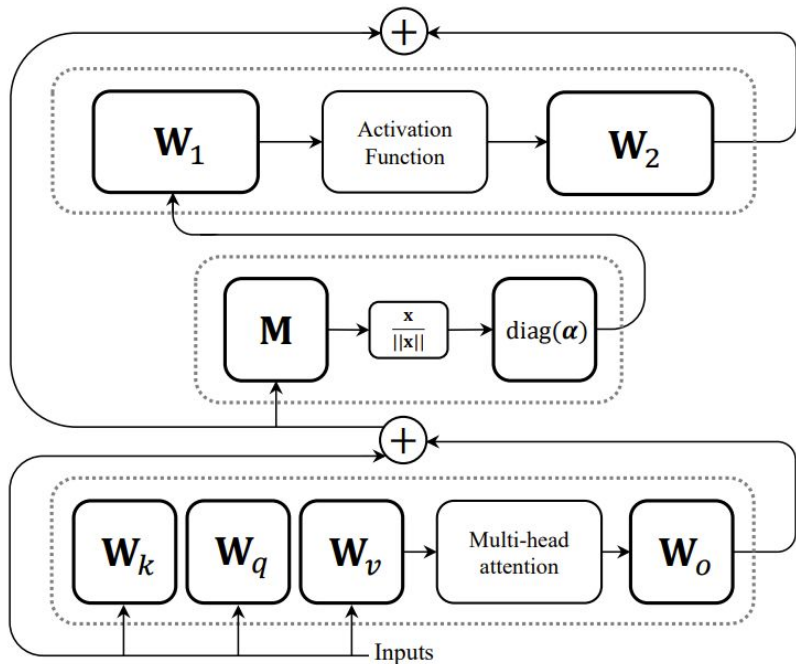
# **SliceGPT and bits-back coding**

# SliceGPT

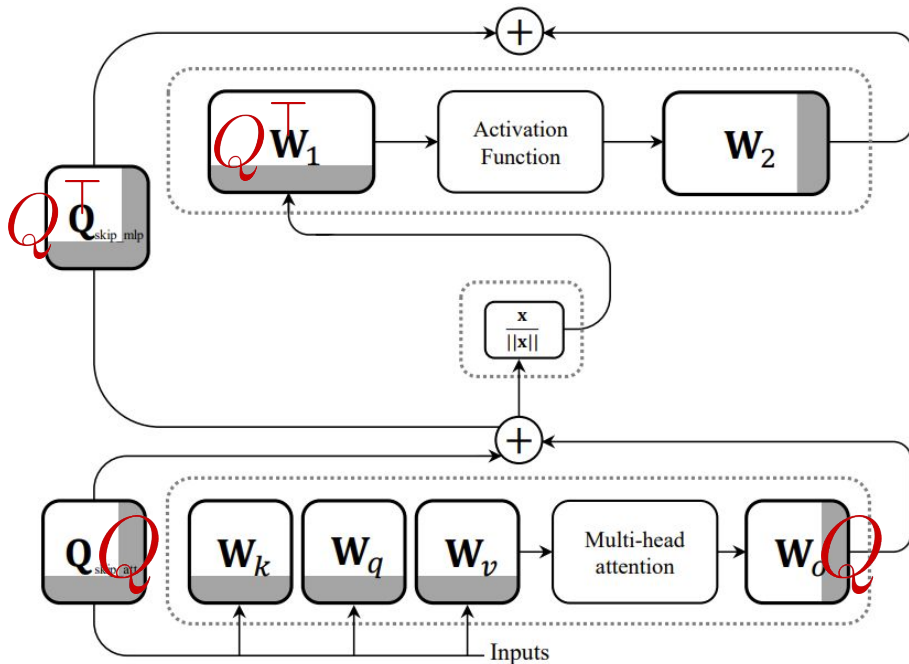


(a) A standard transformer block.

# SliceGPT introduces rotational symmetries:

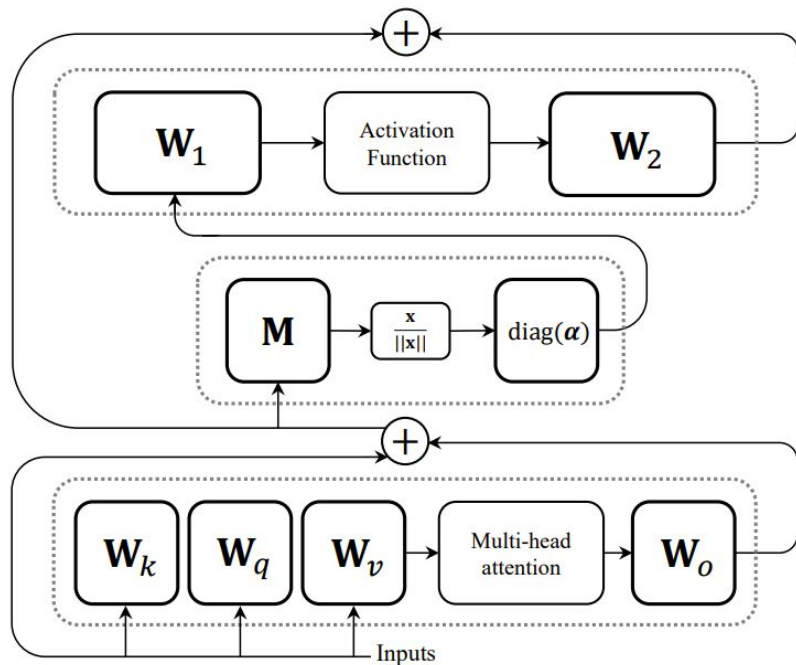


(a) A standard transformer block.

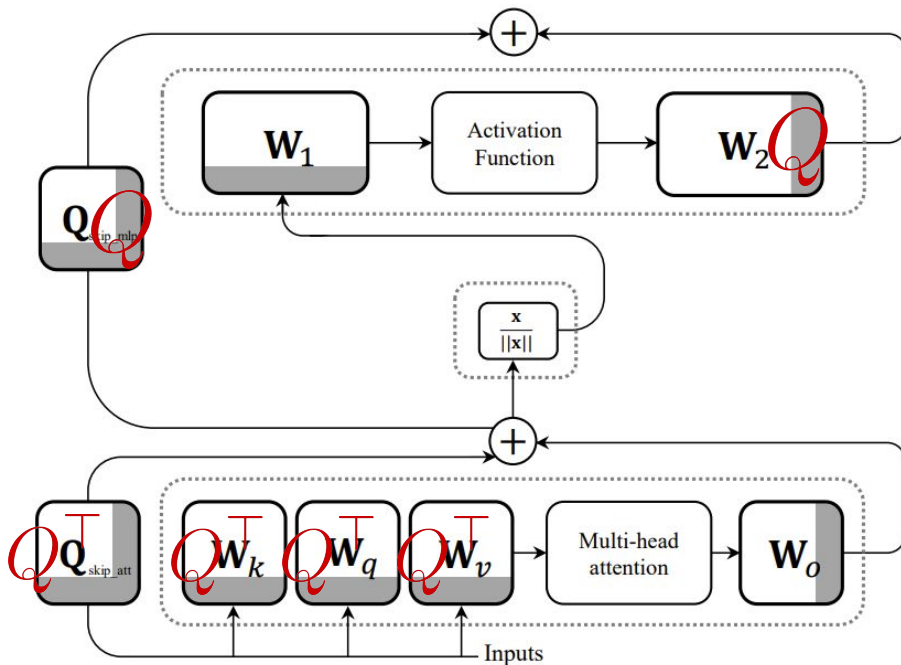


(b) A transformer block with SliceGPT.

# SliceGPT introduces rotational symmetries:



(a) A standard transformer block.







(b) A transformer block with SliceGPT.


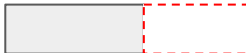

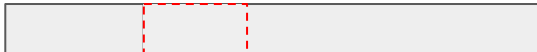
SliceGPT introduces rotational symmetries:

we can write each transformer block as:

$$f(x, W) = f(x, QW)$$

source coding:  + x = 

“AC way”:  + x =  +   
= 


“bits-back way”:  + x =  +   
= 



$x \sim P$ , symmetric  $P$  + good code for  $P$

How to compress  $|x|$ ,  $x \sim P$ ?

1. Remove the last bit of the initial bitstream.
2. If it's 0, encode  $|x|$ ; if 1, encode  $-|x|$ .

“bits-back way”: 

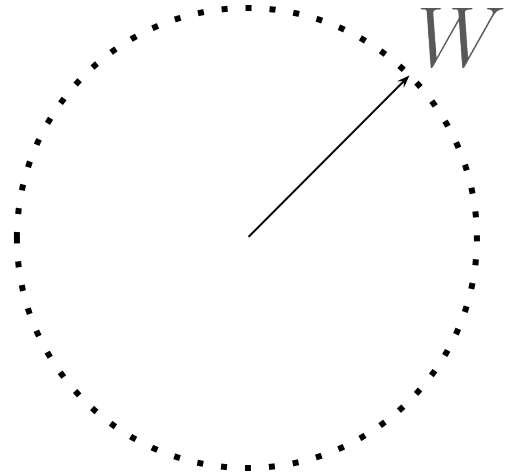
**Getting free bits back from rotational symmetries**

$$f(x, W) = f(x, QW)$$

**encode:**

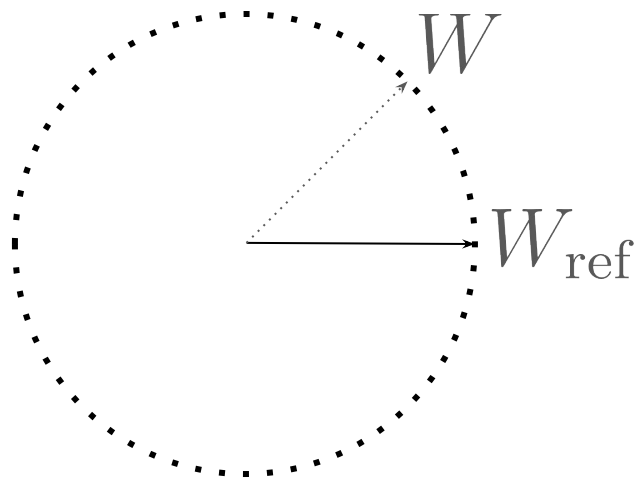
**encode:** start with a weight matrix and some initial bits

Initial bits



**encode step 1:** rotate weight matrix to a “canonical” direction

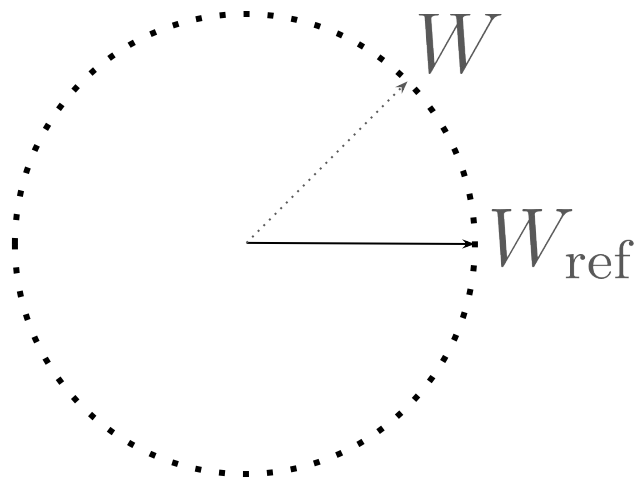
Initial bits



**encode step 1:** rotate weight matrix to a “canonical” direction

$$\text{svd}(W) = U \underbrace{\{\Sigma V^T\}}_{\bar{W}_{\text{ref}}}$$

Initial bits

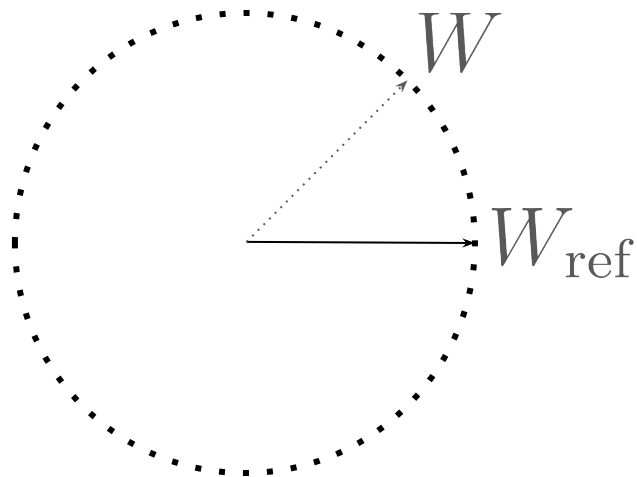


**encode step 1:** rotate weight matrix to a “canonical” direction

$$\text{svd}(W) = U \underbrace{\Sigma V^T}_{\bar{W}_{\text{ref}}}$$

i.e., define  $W_{\text{ref}} W_{\text{ref}}^T$  to be diagonal

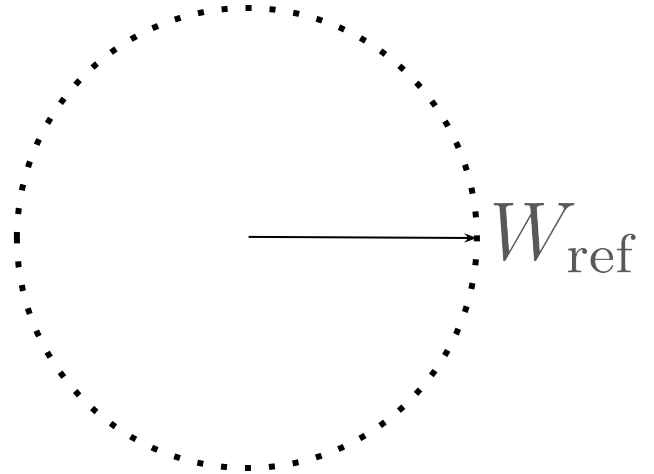

Initial bits





## encode step 2: decode a rotation from the bitstream

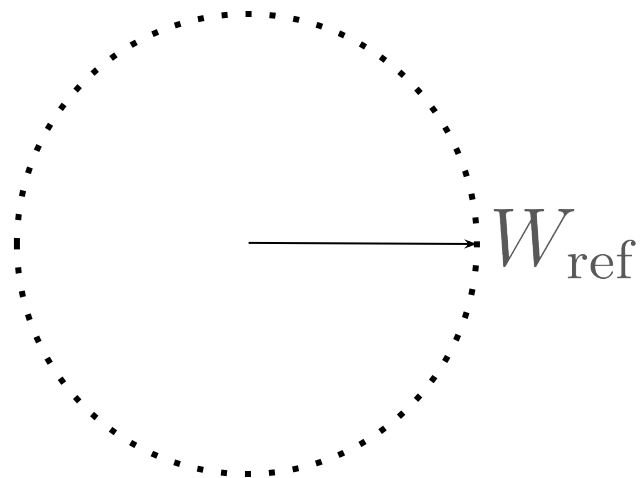
Initial bits



## encode step 2: decode a rotation from the bitstream



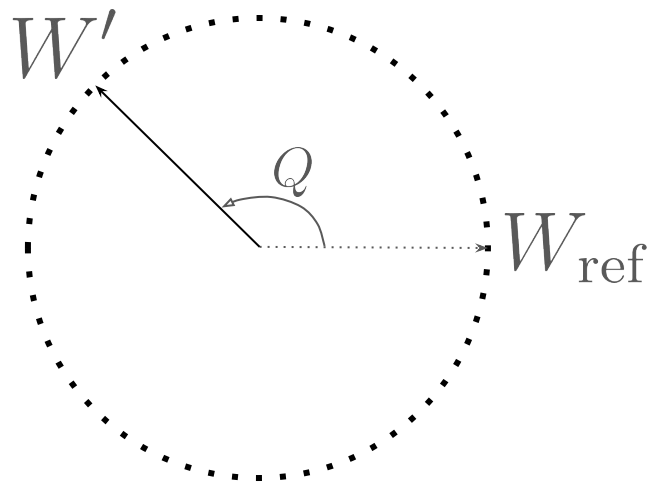
**encode step 3:** rotate weight by the decoded rotation matrix



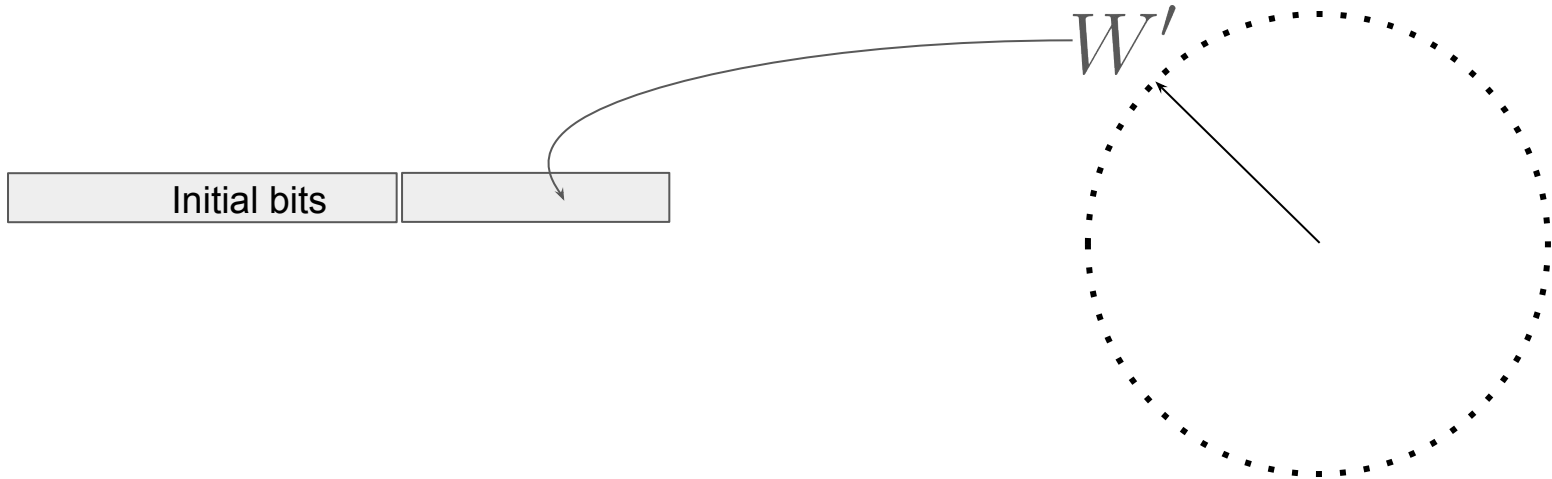
**encode step 3:** rotate weight by the decoded rotation matrix



$$W' = QW_{\text{ref}}$$

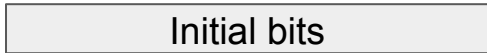


## encode step 4: encode the rotated weight matrix

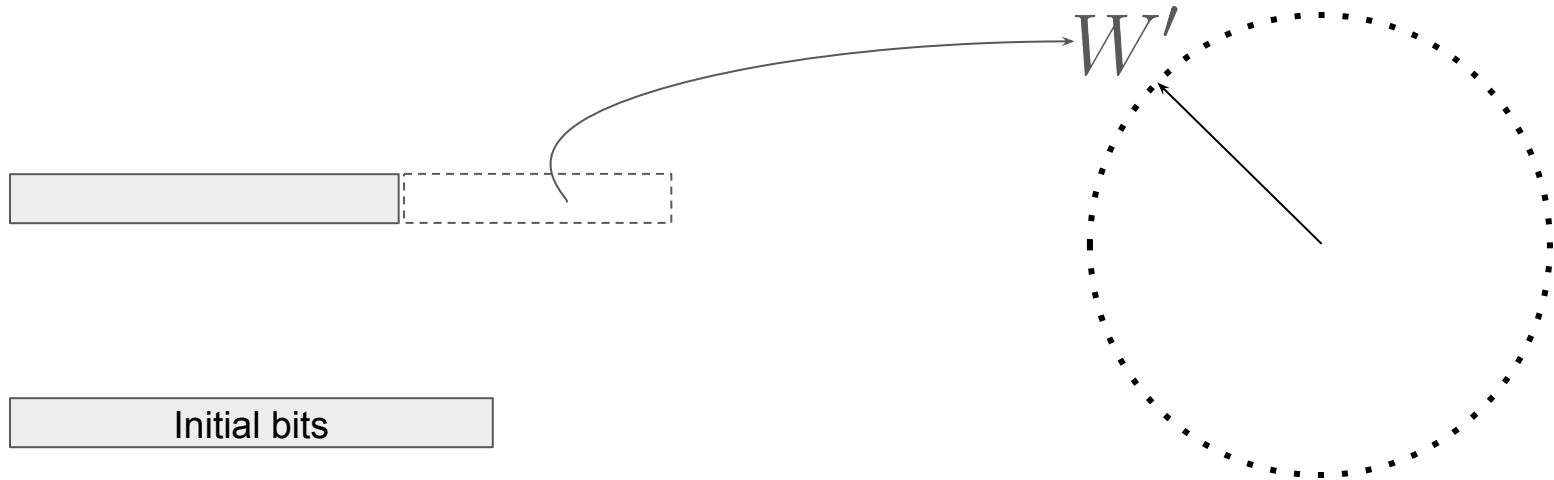


**decode:**

**decode:** start with some bits



# decode step 1: decode the rotated weight matrix

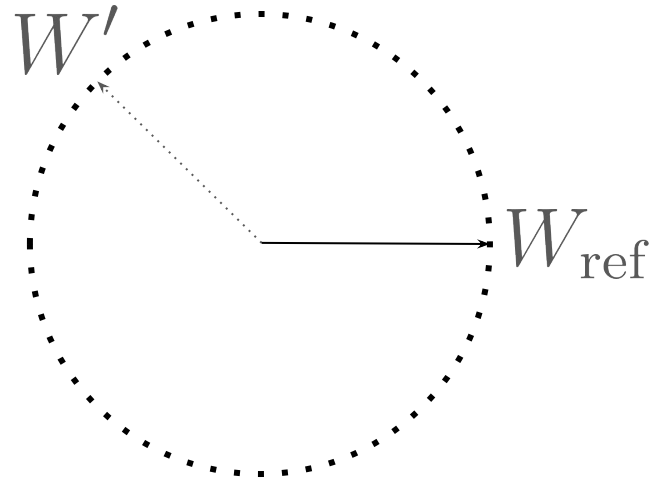




**decode step 2:** rotate it to “canonical” direction



Initial bits

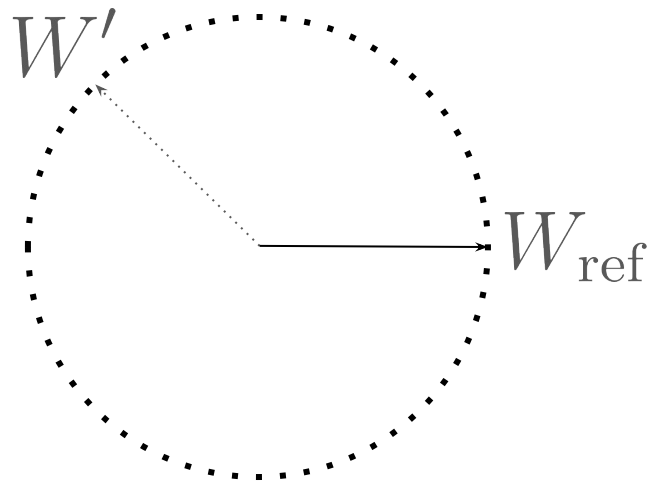


## decode step 2: rotate it to “canonical” direction

recall we define  $\text{svd}(W) = U \begin{bmatrix} \Sigma & V^T \\ \bar{W}_{\text{ref}} \end{bmatrix}$



Initial bits



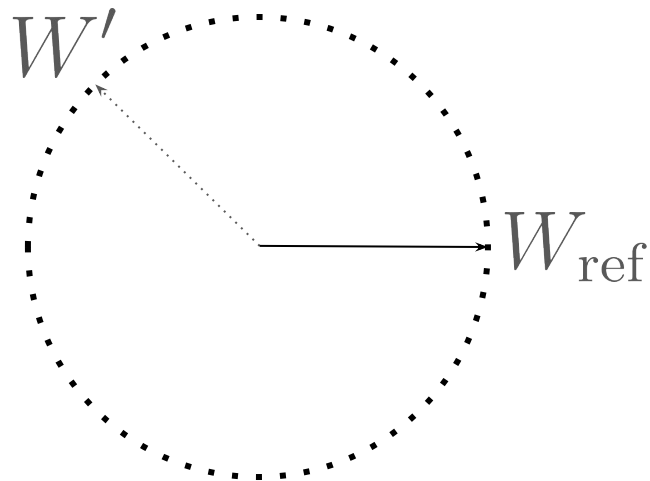
## decode step 2: rotate it to “canonical” direction

recall we define  $\text{svd}(W) = U \begin{bmatrix} \Sigma \\ V^\top \end{bmatrix} \overline{W}_{\text{ref}}$

$$\text{svd}(W') = \begin{bmatrix} \sigma \\ \hat{Q} \end{bmatrix} \Sigma V^\top \sigma, \quad \sigma = \text{diag}(\pm 1, \dots, \pm 1)$$



Initial bits



## decode step 2: rotate it to “canonical” direction

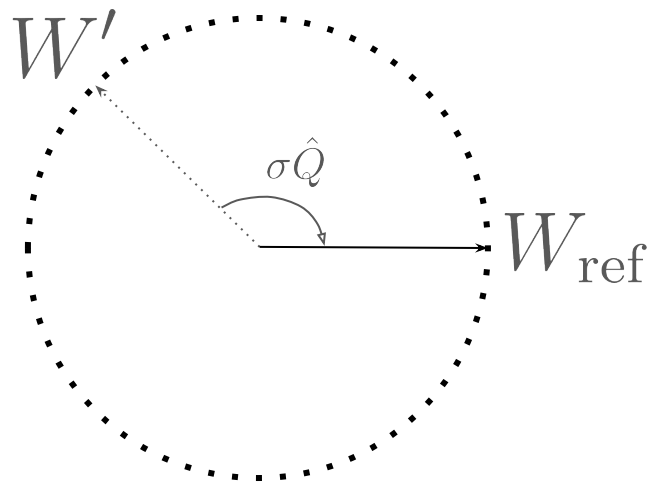
recall we define  $\text{svd}(W) = U \begin{bmatrix} \Sigma \\ V^\top \end{bmatrix} \bar{W}_{\text{ref}}$

$$\text{svd}(W') = \begin{bmatrix} \sigma \hat{Q} \\ \Sigma V^\top \end{bmatrix} \sigma, \sigma = \text{diag}(\pm 1, \dots, \pm 1)$$

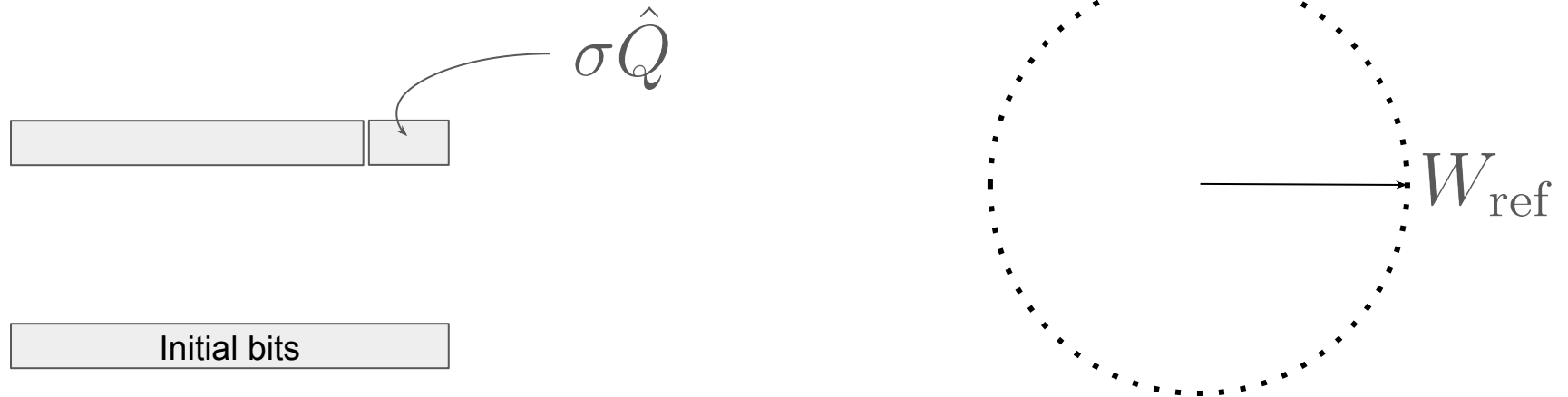
$\hat{Q}$



Initial bits



## decode step 3: encode the recovered rotation matrix



## encoding:

1. rotation to canonical direction
2. decode a rotation
3. rotate weight matrix
4. encode the rotated matrix

## decoding:

3. encode the rotation
2. rotate weight to canonical direction
1. decode the rotated matrix

## encoding:

1. rotation to canonical direction
2. decode a rotation
3. rotate weight matrix
4. encode the rotated matrix

## decoding:

3. encode the rotation
2. rotate weight to canonical direction
1. decode the rotated matrix



**does rotation need infinite precision?**

## encoding:

1. rotation to canonical direction
2. decode a rotation
3. rotate weight matrix
4. encode the rotated matrix

## decoding:

3. encode the rotation
2. rotate weight to canonical direction
1. decode the rotated matrix

 **does rotation need infinite precision?**

 **yes. but just using float16 also works well!**



encoding:

1. rotation to canonical direction
2. decode a rotation
3. rotate weight matrix
4. encode the rotated matrix

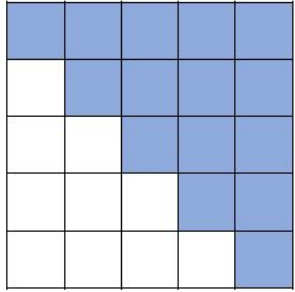
decoding:

3. encode the rotation
2. rotate weight to canonical direction
1. decode the rotated matrix

 **does rotation need infinite precision?**

 **yes. but just using float16 also works well!**

encoding:



canonical direction

rotation

decoding:

3. encode the rotation

2. rotate weight to canonical direction

1. decode the rotated matrix

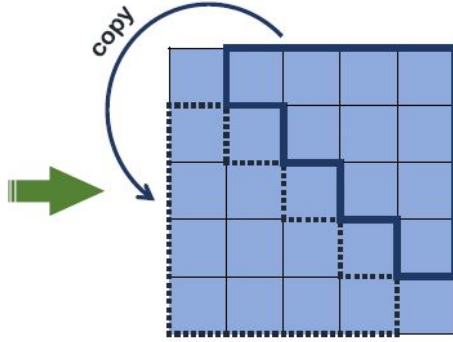
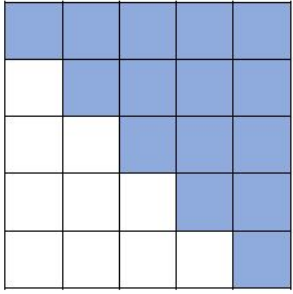
Decode  $D(D + 1)/2$  float weight matrix

4. encode the rotated matrix

🙄 does rotation need infinite precision?

😊 yes. but just using float16 also works well!

encoding:



decoding:

3. encode the rotation

2. rotate weight to canonical direction

1. decode the rotated matrix

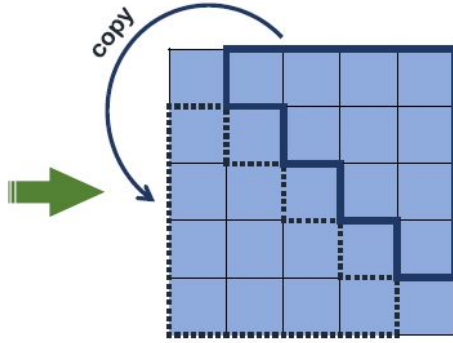
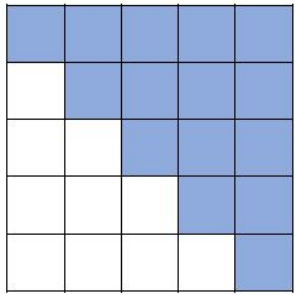
Decode  $D(D + 1)/2$  float

Form a symmetric matrix

😞 does rotation need infinite precision?

😊 yes. but just using float16 also works well!

encoding:

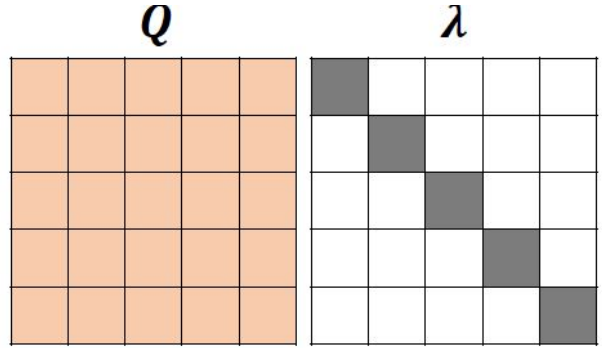


Decode  $D(D + 1)/2$  float

Form a symmetric matrix

decoding:

Eigenvalue  
Decomposition



Return eigenvectors as the rotation matrix  
Encode  $D$  eigenvalues to bitstream

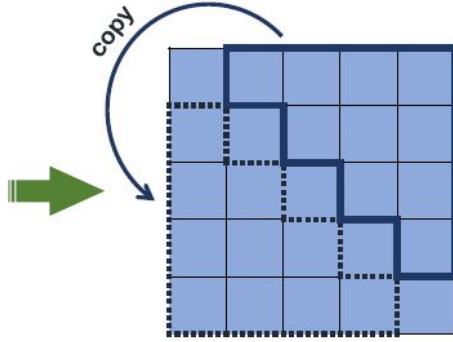
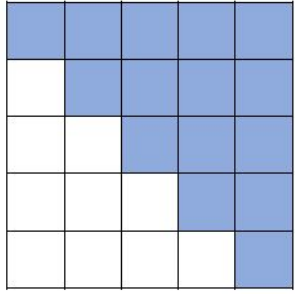
4. encode the rotated matrix

1. decode the rotated matrix

😞 does rotation need infinite precision?

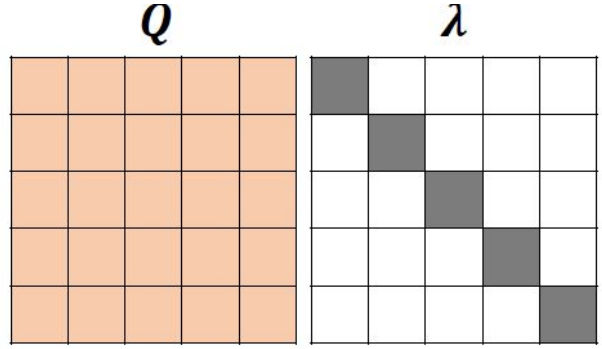
😊 yes. but just using float16 also works well!

encoding:



decoding:

Eigenvalue  
Decomposition



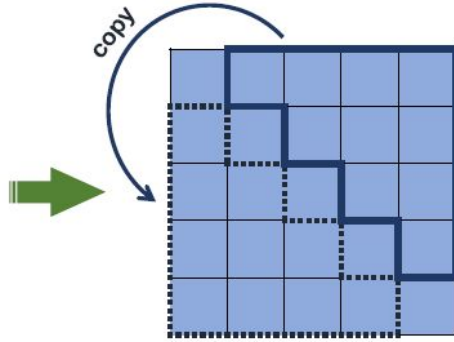
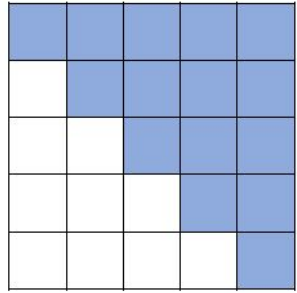
Return eigenvectors as the rotation matrix  
Encode  $D$  eigenvalues to bitstream

Decode  $D$  eigenvalues ( $\lambda$ ) from bitstream

😞 does rotation need infinite precision?

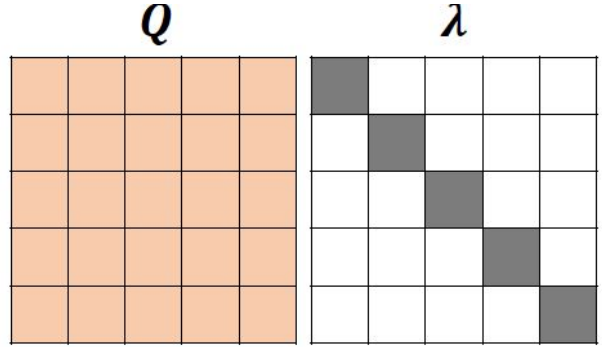
😊 yes. but just using float16 also works well!

encoding:



decoding:

Eigenvalue  
Decomposition



Return eigenvectors as the rotation matrix  
Encode  $D$  eigenvalues to bitstream

Encode the upper triangular and diagonal elements

$$Q \text{diag}(\lambda) Q^T$$



Decode  $D$  eigenvalues ( $\lambda$ ) from bitstream

🙄 does rotation need infinite precision?

😊 yes. but just using float16 also works well!

## encoding:

1. rotation to canonical direction
2. decode a rotation
3. rotate weight matrix
4. encode the rotated matrix

## decoding:

3. encode the rotation
2. rotate weight to canonical direction
1. decode the rotated matrix



**but there may be numerical error...**

## encoding:

1. rotation to canonical direction

2. decode a rotation

3. rotate weight matrix

4. encode the rotated matrix

## decoding:

3. encode the rotation

2. rotate weight to canonical direction

1. decode the rotated matrix



**but there may be numerical error...**



## encoding:

## decoding:

1. rotation to canonical direction

2. decode a rotation



3. encode the rotation

3. rotate weight matrix

2. rotate weight to canonical direction

4. encode the rotated matrix

1. decode the rotated matrix



**but there may be numerical error...**



**but there may be numerical error...**



**We can send correction code if the error is too large!**

 **but there may be numerical error...**

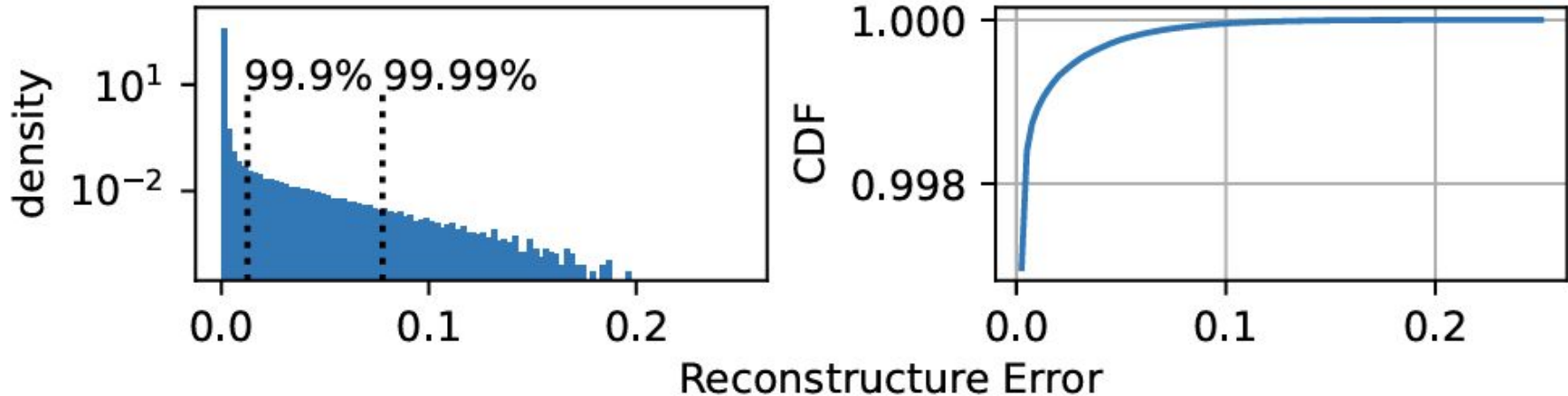
 **We can send correction code if the error is too large!**

**Will this correction code becomes too large? NO!**

🙄 **but there may be numerical error...**

💡 **We can send correction code if the error is too large!**

**Will this correction code becomes too large? NO!**



# Results

# Results

Model	SliceGPT Slicing	Compress Rate after SliceGPT	Compress Rate after bits-back	Performance (before/after bits-back)			
				PPL ( $\downarrow$ )	PIQA ( $\%$ , $\uparrow$ )	WinoGrande ( $\%$ , $\uparrow$ )	HellaSwag ( $\%$ , $\uparrow$ )
OPT-1.3B	20%	-9.53%	-13.77%	<b>16.59</b> /16.60	<b>64.91</b> /64.80	<b>54.78</b> /54.38	45.26/ <b>45.32</b>
	25%	-14.84%	-18.61%	<b>17.78</b> /17.86	<b>63.55</b> /63.33	52.80/ <b>53.28</b>	<b>43.20</b> /43.11
	30%	-20.53%	-23.81%	<b>19.60</b> /19.66	<b>60.88</b> /60.50	52.88/ <b>53.28</b>	<b>40.25</b> /40.06
OPT-2.7B	20%	-9.19%	-13.84%	<b>13.89</b> /13.95	<b>68.44</b> /68.12	<b>58.88</b> /58.72	<b>51.35</b> /51.17
	25%	-15.07%	-19.09%	<b>14.85</b> /14.87	<b>66.70</b> /66.76	57.30/ <b>57.70</b>	<b>48.41</b> /48.38
	30%	-20.88%	-24.43%	<b>16.31</b> /16.33	64.64/ <b>64.69</b>	55.80/ <b>56.04</b>	44.52/ <b>44.57</b>
OPT-6.7B	20%	-9.29%	-14.07%	<b>11.63</b> /11.71	72.91/ <b>73.01</b>	<b>61.33</b> /61.17	60.53/ <b>60.55</b>
	25%	-15.16%	-19.29%	<b>12.12</b> /12.15	71.00/ <b>71.22</b>	60.30/ <b>60.77</b>	<b>57.76</b> /57.55
	30%	-21.18%	-24.84%	<b>12.81</b> /12.91	69.31/ <b>69.42</b>	<b>59.75</b> /59.59	<b>53.64</b> /52.94
OPT-13B	20%	-9.18%	-14.01%	<b>10.75</b> /10.77	74.27/74.27	<b>64.96</b> /64.88	65.74/ <b>65.79</b>
	25%	-15.27%	-19.51%	11.08/ <b>11.07</b>	<b>74.27</b> /73.72	63.46/ <b>63.93</b>	<b>63.48</b> /63.09
	30%	-21.29%	-24.97%	<b>11.55</b> /11.59	72.69/ <b>73.01</b>	61.96/ <b>62.43</b>	<b>60.12</b> /60.05
Llama-2-7B	20%	-9.38%	-14.13%	<b>6.86</b> /6.98	<b>69.53</b> /69.42	64.17/ <b>64.72</b>	<b>58.96</b> /58.89
	25%	-15.34%	-19.53%	<b>7.56</b> /7.59	67.03/ <b>67.57</b>	62.98/ <b>63.38</b>	<b>54.29</b> /53.93
	30%	-21.45%	-25.09%	<b>8.63</b> /8.69	<b>64.69</b> /64.09	<b>62.75</b> /62.12	<b>49.13</b> /49.07

# Results

**3-5% additional bits saving**

Model	SliceGPT Slicing	Compress Rate after SliceGPT	Compress Rate after bits-back	Performance (before/after bits-back)			
				PPL (↓)	PIQA (% , ↑)	WinoGrande (% , ↑)	HellaSwag (% , ↑)
OPT-1.3B	20%	-9.53%	-13.77%	<b>16.59</b> /16.60	<b>64.91</b> /64.80	<b>54.78</b> /54.38	45.26/ <b>45.32</b>
	25%	-14.84%	-18.61%	<b>17.78</b> /17.86	<b>63.55</b> /63.33	52.80/ <b>53.28</b>	<b>43.20</b> /43.11
	30%	-20.53%	-23.81%	<b>19.60</b> /19.66	<b>60.88</b> /60.50	52.88/ <b>53.28</b>	<b>40.25</b> /40.06
OPT-2.7B	20%	-9.19%	-13.84%	<b>13.89</b> /13.95	<b>68.44</b> /68.12	<b>58.88</b> /58.72	<b>51.35</b> /51.17
	25%	-15.07%	-19.09%	<b>14.85</b> /14.87	<b>66.70</b> /66.76	57.30/ <b>57.70</b>	<b>48.41</b> /48.38
	30%	-20.88%	-24.43%	<b>16.31</b> /16.33	64.64/ <b>64.69</b>	55.80/ <b>56.04</b>	44.52/ <b>44.57</b>
OPT-6.7B	20%	-9.29%	-14.07%	<b>11.63</b> /11.71	72.91/ <b>73.01</b>	<b>61.33</b> /61.17	60.53/ <b>60.55</b>
	25%	-15.16%	-19.29%	<b>12.12</b> /12.15	71.00/ <b>71.22</b>	60.30/ <b>60.77</b>	<b>57.76</b> /57.55
	30%	-21.18%	-24.84%	<b>12.81</b> /12.91	69.31/ <b>69.42</b>	<b>59.75</b> /59.59	<b>53.64</b> /52.94
OPT-13B	20%	-9.18%	-14.01%	<b>10.75</b> /10.77	74.27/74.27	<b>64.96</b> /64.88	65.74/ <b>65.79</b>
	25%	-15.27%	-19.51%	11.08/ <b>11.07</b>	<b>74.27</b> /73.72	63.46/ <b>63.93</b>	<b>63.48</b> /63.09
	30%	-21.29%	-24.97%	<b>11.55</b> /11.59	72.69/ <b>73.01</b>	61.96/ <b>62.43</b>	<b>60.12</b> /60.05
Llama-2-7B	20%	-9.38%	-14.13%	<b>6.86</b> /6.98	<b>69.53</b> /69.42	64.17/ <b>64.72</b>	<b>58.96</b> /58.89
	25%	-15.34%	-19.53%	<b>7.56</b> /7.59	67.03/ <b>67.57</b>	62.98/ <b>63.38</b>	<b>54.29</b> /53.93
	30%	-21.45%	-25.09%	<b>8.63</b> /8.69	<b>64.69</b> /64.09	<b>62.75</b> /62.12	<b>49.13</b> /49.07

# Results

Model	SliceGPT Slicing	Compress Rate after SliceGPT	Compress Rate after bits-back	Performance (before/after bits-back)			
				PPL (↓)	PIQA (% , ↑)	WinoGrande (% , ↑)	HellaSwag (% , ↑)
OPT-1.3B	20%	-9.53%	-13.77%	<b>16.59</b> /16.60	<b>64.91</b> /64.80	<b>54.78</b> /54.38	45.26/ <b>45.32</b>
	25%	-14.84%	-18.61%	<b>17.78</b> /17.86	<b>63.55</b> /63.33	52.80/ <b>53.28</b>	<b>43.20</b> /43.11
	30%	-20.53%	-23.81%	<b>19.60</b> /19.66	<b>60.88</b> /60.50	52.88/ <b>53.28</b>	<b>40.25</b> /40.06
OPT-2.7B	20%	-9.19%	-13.84%	<b>13.89</b> /13.95	<b>68.44</b> /68.12	<b>58.88</b> /58.72	<b>51.35</b> /51.17
	25%	-15.07%	-19.09%	<b>14.85</b> /14.87	<b>66.70</b> /66.76	57.30/ <b>57.70</b>	<b>48.41</b> /48.38
	30%	-20.88%	-24.43%	<b>16.31</b> /16.33	64.64/ <b>64.69</b>	55.80/ <b>56.04</b>	44.52/ <b>44.57</b>
OPT-6.7B	20%	-9.29%	-14.07%	<b>11.63</b> /11.71	72.91/ <b>73.01</b>	<b>61.33</b> /61.17	60.53/ <b>60.55</b>
	25%	-15.16%	-19.29%	<b>12.12</b> /12.15	71.00/ <b>71.22</b>	60.30/ <b>60.77</b>	<b>57.76</b> /57.55
	30%	-21.18%	-24.84%	<b>12.81</b> /12.91	69.31/ <b>69.42</b>	<b>59.75</b> /59.59	<b>53.64</b> /52.94
OPT-13B	20%	-9.18%	-14.01%	<b>10.75</b> /10.77	74.27/74.27	<b>64.96</b> /64.88	65.74/ <b>65.79</b>
	25%	-15.27%	-19.51%	11.08/ <b>11.07</b>	<b>74.27</b> /73.72	63.46/ <b>63.93</b>	<b>63.48</b> /63.09
	30%	-21.29%	-24.97%	<b>11.55</b> /11.59	72.69/ <b>73.01</b>	61.96/ <b>62.43</b>	<b>60.12</b> /60.05
Llama-2-7B	20%	-9.38%	-14.13%	<b>6.86</b> /6.98	<b>69.53</b> /69.42	64.17/ <b>64.72</b>	<b>58.96</b> /58.89
	25%	-15.34%	-19.53%	<b>7.56</b> /7.59	67.03/ <b>67.57</b>	62.98/ <b>63.38</b>	<b>54.29</b> /53.93
	30%	-21.45%	-25.09%	<b>8.63</b> /8.69	<b>64.69</b> /64.09	<b>62.75</b> /62.12	<b>49.13</b> /49.07

**negligible influence on performance**



# Encoding and Decoding time

## GPU:

Model Name	OPT-1.3B		OPT-2.7B		OPT-6.7B		OPT-13B	
Slicing	20%	30%	20%	30%	20%	30%	20%	30%
Encoding time	15 s	13 s	30 s	24 s	2.5 min	1.7 min	6.5 min	4.1 min
Decoding time	6 s	5 s	14 s	11 s	1.2 min	45 s	2.5 min	2 min

## CPU:

Model Name	OPT-1.3B		OPT-2.7B		OPT-6.7B		OPT-13B	
Slicing	20%	30%	20%	30%	20%	30%	20%	30%
Encoding time	3.9 min	3.5 min	8 min	6.5 min	30 min	25 min	84 min	68 min
Decoding time	1.5 min	1.5 min	3.5 min	2.5 min	12 min	10 min	30 min	24 min

**In summary:**

save 3-5% additional bits,

## **In summary:**

save 3-5% additional bits,  
no influence on performance,

## In summary:

save 3-5% additional bits,

no influence on performance,

a little overhead in model loading/saving time