

# Kronecker Products and GPT2

Ayoub @ Uni Passau

March 8, 2024

## Abstract

What do we do differently:

- We put distillation to test, and see if does improve improve the learning quality over vanilla supervised learning. **Spoiler alert:** NOPE. (Others, all use distillation)
- We don't compress the attention matrices.
- We try different compression schemes and see how far we can push down the size.
- We use multiple Kronecker Factors.
- We try different compression ratios for different layers.

## Contents

<b>1</b>	<b>Results:</b>	<b>2</b>
<b>2</b>	<b>Introduction and Setup</b>	<b>2</b>
<b>3</b>	<b>Training Set-up</b>	<b>2</b>
<b>4</b>	<b>Kronecker Decomposition</b>	<b>2</b>
4.1	The 95M model . . . . .	3
4.2	Different decomposition schemes: . . . . .	3
<b>5</b>	<b>Impact of down-scaling: How small can we go</b>	<b>4</b>
<b>6</b>	<b>Initialization</b>	<b>5</b>
6.1	The fuck is Van Loan . . . . .	5
6.2	Pruning based Initialization . . . . .	5
6.3	Questions . . . . .	5

## 1 Results:

Early results on wikitext103, wikitext2 [1] and Lambada [2]. **KronyPT** is our model, and the suffixes 1350 and 3950 refer to different checkpoints. Both models are already outperforming distilGPT2 ([3]) on the 3 datasets. And the 3950 checkpoint is outperforming **KnGPT2** ([4]) on Lambada, while slightly still behind **TQCompressedGPT2** [5]. We measure perplexity (duh) using the Hugging Face interface that was used by the other papers.

# Params	Model	Datasets		
		wikitext-103	wikitext-2	Lambada
124M	GPT2	29.16	24.67	45.28
82M	DistilGPT2	44.53	36.48	76.00
81M	<b>KronyPT-81M-1350</b>	41.98	34.99	-
81M	<b>KronyPT-81M-3950</b>	-	-	64.92
81M	TQCompressedGPT2	40.28	32.25	64.72
81M	KnGPT-2 (Huawei)	40.97	32.81	67.62

Table 1: Comparison of different models on various datasets.

## 2 Introduction and Setup

A lot of work has been done on the compression of LLMs using factorization methods, including Kronecker Products. Most notably for encoder based based, and rarely applied to decoder based methods ([4], [6], [5]).

We study how far we can push the compression ratio, and then impact of down-scaling on the quality of the model. We also investigate adding multiple Kronecker factors.

## 3 Training Set-up

Add here the basic set-up that I used to train the latest models, basically Chinchilla recommendations. (for later)

## 4 Kronecker Decomposition

In this section, we study different Kronecker decomposition setups, and the percentage of compression it would lead to. So far we only decompose the weights `c_fc.weight` and `c_proj.weight` (each has  $2.3M$  in the original GPT2-small architecture.). Each transformer layer (there are 12 in total) has two of these weights. They count to  $56.6M$  in total (45% of GPT2  $124M$ ). Hence any significant reduction to these matrices would lead to a remarkable compression ratio of the whole model. We choose not to compress the other weights, namely,

attention weights and embedding matrix for various reasons that we will expose later on.

## 4.1 The 95M model

The most basic strategy is to divide one of the dimensions of each  $W$  by 2, this would lead to a 95M model. The parameter  $p_1 = \text{c\_fc.weight}$  (resp.  $p_2 = \text{c\_proj.weight}$ ) has a shape of (3072, 768) (resp. (768, 3072)). We first try the following decomposition:  $p_1 = \underbrace{W_{11}}_{(3072,384)} \otimes \underbrace{W_{12}}_{(1,2)}$  and  $p_2 = \underbrace{W_{21}}_{(384,3072)} \otimes \underbrace{W_{22}}_{(2,1)}$

This decomposition would lead to to reduction of 28M (50%). The new network would have approx 95M. Our goal is to eventually reach the 82M mark, similar to DistilGPT2, and other Factorized models (inserts other refs here).

## 4.2 Different decomposition schemes:

It is reasonable to aim for a decomposition that guarantees the maximum rank we can get. Since the Rank of Kronecker products is multiplicative, i.e.,  $\text{rank}(A \otimes B) = \text{rank}(A) \cdot \text{rank}(B)$ <sup>1</sup>, we can easily compute the rank of each possible decomposition. In our case, we have  $W \in R^{(m,n)}$  where  $m = 3072$  and  $n = 768$ . Hence, for each layer of GPT2, we aim to find the “optimal”  $A \in R^{(m_1,n_1)}$ , and  $B \in R^{(m_2,n_2)}$ , i.e.,:

$$W \approx A \otimes B, \quad m = m_1 m_2, n = n_1 n_2$$

W.l.o.g, for each decomposition  $(A, B)$  the maximum rank we can reach is  $\min(m_1, n_1) \cdot \min(m_2, n_2)$ . And each of the reduced decompositions would have exactly  $m_1 n_1 + m_2 n_2$  parameters. Hence, theoretically, the maximum rank we can get is 384 of a (3072, 768) matrix. The following table summarizes all possible combinations, alongside the reduction it would lead to per layer, and the total number of parameters in GPT2, for only those decompositions of rank 384.

The following table summarizes some of the decompositions that we can try without degrading the rank. We are particularly interested in 3 models, the 67M, the 81M and the 96M. (Need to add this) Furthermore, we add more factors to the 67M leading to a XXM model and see how loss curve behaves.

<sup>1</sup>Link to proof: <https://math.stackexchange.com/questions/707091/elementary-proof-for-textrank-left-a-otimes-b-right-textrank-a-cdot>

Name	Dimension	params	Model size
67M	(64, 32)	3200	67,893,504
	(64, 48)	3840	67,908,864
	(96, 32)	3840	67,908,864
	(64, 64)	4672	67,928,832
	(128, 32)	4672	67,928,832
	(96, 48)	5120	67,939,584
	(96, 64)	6528	67,973,376
	(128, 48)	6528	67,973,376
	(128, 64)	8480	68,020,224
	(96, 96)	9472	68,044,032
	(192, 48)	9472	68,044,032
	(128, 96)	12480	68,116,224
	(192, 64)	12480	68,116,224
	(128, 128)	16528	68,213,376
	(256, 64)	16528	68,213,376
	...	...	...
	(1024, 256)	262153	74,108,376
	(768, 384)	294920	74,894,784
	(1024, 384)	393222	77,254,032
81M	(768, 768)	589828	81,972,576
	(1536, 384)	589828	81,972,576
	(1024, 768)	786435	86,691,144
96M	(1536, 768)	1179650	96,128,304
GPT2	(3072, 768)	2359297	124,439,832

Table 2: Comparison of different models on various datasets.

## 5 Impact of down-scaling: How small can we go

In this section, we study the impact of down-scaling on the compressed models, we train 4 different architectures, with variant dimensions:

- 67M: diverges // try with higher batch size (the least we can get)
- 81M: Super.
- 81M: with a Variant scheme (67M with multiple factors)
- 95M: duh duh duh. (the highest we can get)

Keep in mind:

- You can probably stabilize training using various tricks. It's just an endless loop. We are not going to play the What-IF game. One single config, that's it.

- One could try to have a mixed strategy:
  - higher levels of compression in the early layers.
  - higher levels of compression in the late layers.
  - Every odd layer
  - In the middle
- 
- 

## 6 Initialization

In our set-up, we blend/merge new parameters into already pre-trained. It seems wastful, to randomly initiliaze the newly introduce Kroncker Products if we can come up with a clever / better way. Here we introduce the Van Loan method, and another pruning based initialization trick, specifically crafted for the 95M model. And eventually both of them.

### 6.1 The fuck is Van Loan

XX

### 6.2 Pruning based Initialization

XX

### 6.3 Questions

- When we plug new KP matrices to the pre-trained model, is it useful to freeze the other weights that are already pre-trained? Or using different learning rates schemes maybe?
- When don't don't freeze the weights. Does the network rely on other pre-trained methods?
  - **\*\*The idea that I'm challenging here\*\***: A lot of work on NNs compression using factorized methods, claim that the network only need to be (post-)trained on small
  - But, what the fail to mention or elaborate on at least, is that not the weights matrices of the network are decomposed.
  - And with all the residual connection that are present in most attention networks, one could suspect, that maybe the weights that were not decomposed are taking over...

- A good remedy for this is to freeze the original weights during the post-training, and only allow the new dropped in (factorized matrices) to be updated with backprop.
- This strategy could also be implemented in distillation. We refer to this method as **forced distillation**.
- When we freeze, we investigate how useful is it to distill matrices one by one, rather than drop in the matrices all at once. And also investigate if the order of dropping has any significance, bottom up or top bottom...

## References

- [1] S. Merity, C. Xiong, J. Bradbury, and R. Socher, “Pointer sentinel mixture models,” *arXiv preprint arXiv:1609.07843*, 2016.
- [2] D. Paperno, G. Kruszewski, A. Lazaridou, Q. N. Pham, R. Bernardi, S. Pezzelle, M. Baroni, G. Boleda, and R. Fernández, “The lambda dataset: Word prediction requiring a broad discourse context,” *arXiv preprint arXiv:1606.06031*, 2016.
- [3] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, “Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter,” in *NeurIPS EMC<sup>2</sup> Workshop*, 2019.
- [4] M. Tahaei, E. Charlaix, V. Nia, A. Ghodsi, and M. Rezagholizadeh, “Kroneckerbert: Significant compression of pre-trained language models through kronecker decomposition and knowledge distillation,” in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2116–2127, 2022.
- [5] V. Abronin, A. Naumov, D. Mazur, D. Bystrov, K. Tsarova, A. Melnikov, I. Oseledets, S. Dolgov, R. Brasher, and M. Perelshtein, “Tqcompressor: improving tensor decomposition methods in neural networks via permutations,” *arXiv preprint arXiv:2401.16367*, 2024.
- [6] A. Edalati, M. Tahaei, A. Rashid, V. P. Nia, J. J. Clark, and M. Rezagholizadeh, “Kronecker decomposition for gpt compression,” *arXiv preprint arXiv:2110.08152*, 2021.