

How Do Transformers Fill in the Blanks?

A Case Study on Matrix Completion

Pulkit Gopalani¹ Ekdeep Singh Lubana^{1,2} Wei Hu¹

Abstract

Completing masked sequences is an important problem in language modeling, and analyzing how Transformer models perform this task is crucial for understanding their mechanisms. In this direction, we formulate the low-rank matrix completion problem as a masked language modeling (MLM) task, and train a BERT model to solve this task. We find that BERT succeeds in matrix completion and outperforms the classical nuclear norm minimization method. Moreover, the mean-squared-error (MSE) loss curve displays an early plateau followed by a sudden drop to near-optimal values, despite no changes in the training procedure or hyper-parameters. To gain interpretability insights, we examine the model’s predictions, attention heads, and hidden states before and after this transition. Concretely, we observe that (i) the model transitions from simply copying the masked input to accurately predicting the masked entries; (ii) the attention heads transition to interpretable patterns relevant to the task; and (iii) the embeddings and hidden states encode information relevant to the problem.

1. Introduction

This paper investigates the behavior of Transformers trained on the classical mathematical task of low-rank matrix completion [5] to gain insights into the mechanisms of Transformers and their training process. In this setup, we assume access to a matrix with some fraction of its entries missing, and would like to complete the missing entries assuming the ground truth matrix is low-rank. By treating a matrix as a sequence of tokens, we find that training a BERT model [11] in an online manner can successfully solve this problem to a small error. Moreover, BERT can outperform the classical nuclear norm minimization algorithm for matrix completion, suggesting that BERT does not simply recover this classical

¹University of Michigan, Ann Arbor ²CBS, Harvard University. Correspondence to: Pulkit Gopalani <gopalani@umich.edu>.

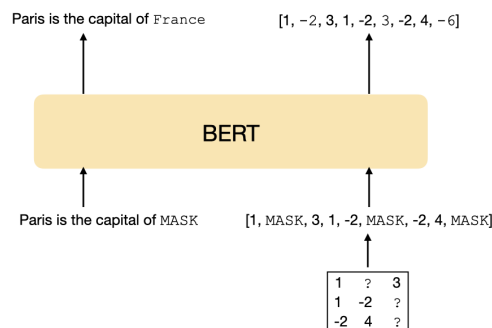


Figure 1: Matrix completion using BERT; missing entries in the matrix are encoded as masked positions in MLM.

algorithm (Please see Appendix A for details on low-rank matrix completion and nuclear norm minimization).

Further, MSE loss curve during training undergoes a *sudden decrease* (Fig. 2), marking the transition to a model that generalizes well. Such a sudden decrease was also observed in [7] for BERT trained in natural language setups. We find that this decrease in loss marks an *algorithmic shift* from the pre-transition model simply copying the input (predicting 0 at masked positions) to the post-transition model accurately predicting missing values at masked positions.

2. BERT Solves Matrix Completion

We train a BERT model (with parameters θ) to predict missing entries given a masked matrix \tilde{X} as input. For the model output $\hat{X} := \hat{X}(\tilde{X}; \theta) \in \mathbb{R}^{n \times n}$, the training objective $L(\theta)$ is the mean-squared-error (MSE) loss over all entries

$$L(\theta) = \frac{1}{n^2} \sum_{i,j=1}^n (X_{ij} - \hat{X}_{ij})^2.$$

In experiments, data for matrix completion is generated as

$$X = UV^T; \quad U, V \in \mathbb{R}^{n \times r}, \quad U_{ij}, V_{ij} \stackrel{\text{iid}}{\sim} \text{Unif}[-1, 1]$$

$\forall i, j \in [n] \times [r]$ so that X has rank at most r . To mask entries at random, we sample binary matrices $M \in \{0, 1\}^{n \times n}$ such that $M_{ij} = 0$ with probability p_{mask} , indicating that the element at position (i, j) is masked in the input matrix. Define $\Omega := \{(i, j) \mid M_{ij} = 1\}$.

Apart from the MSE (L) over all entries, we separately track MSE over observed and masked entries,

$$L_{obs} = \frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} (X_{ij} - \hat{X}_{ij})^2$$

$$L_{mask} = \frac{1}{|\Omega^C|} \sum_{(i,j) \in \Omega^C} (X_{ij} - \hat{X}_{ij})^2.$$

In the subsequent sections, we analyse a 4-layer, 8-head BERT model trained upto MSE $\sim 4e-3$ on 7×7 , rank-2 inputs for training dynamics and interpretability. Please see Appendix B for details about experiments (including those with larger model and input matrices), and Appendix D for comparison of our approach to nuclear norm minimization.

3. Matrix Completion Capability Emerges during Training

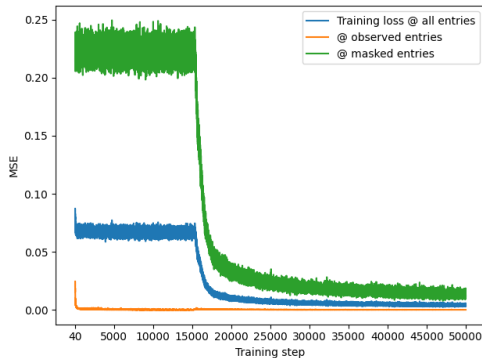


Figure 2: Sharp drop in training loss shortly after step 15000

We observe a sharp decrease in training loss at approximately step 15000 (Fig. 2). Interestingly, this decrease in total loss is driven almost exclusively by the corresponding decrease in L_{mask} , since L_{obs} is very close to 0 both before and after the drop. We hypothesize that this sudden drop is caused by an algorithmic shift; that is, the model switches to a different, more accurate algorithm for prediction at missing entries and hence L_{mask} rapidly decreases following that shift. Moreover, since L_{obs} barely changes during this algorithmic shift, we hypothesize that (1) the model has 2 distinct mechanisms for prediction at masked and observed entries after the algorithmic shift; and (2) that the mechanism for prediction at observed positions is not significantly affected by this algorithmic shift.

We note that this sudden drop is in MSE loss, i.e., not in a discontinuous metric like accuracy. Hence, the idea of sudden emergence being an artifact of discontinuous metrics and poorly defined evals [31] is unlikely to explain the full story in our setting.

3.1. Before the Algorithmic Shift – Copying

In this section, we demonstrate that before the algorithmic shift, the model simply copies the input at all positions in the matrix, and *actual computation* for missing entries does not occur at this stage. In the following, Pre-shift denotes the model at step 4000, and post-shift denotes model at the end of training (step 50000).

3.1.1. VERIFYING COPYING VIA TOKEN INTERVENTION

To rigorously verify the copying hypothesis for the pre-shift model, we use the following method – replace the masked elements in the 7×7 , rank-2 input by the token corresponding to a real value m . For such input, we would like to see whether the model implements copying and outputs m at the masked positions. For model output \hat{X} on this input, MSE at observed positions is L_{obs} , and for masked positions the MSE is defined as

$$L'_{mask} = \frac{1}{|\Omega^C|} \sum_{(i,j) \in \Omega^C} (\hat{X}_{ij} - m)^2.$$

L_{obs} and L'_{mask} for this experiment averaged over 512 samples are compiled in Row 1, Table 1. The small loss values verify the copying hypothesis – model output matches the ground truth at observed positions, while at masked positions it outputs a value nearly equal to m , the replacement mask value. Note that when the mask token is MASK (i.e., no replacement), we set $m = 0$, indicating that the model is outputting 0 at the masked locations.

To further confirm this hypothesis, we sample random 7×7 matrices for input; i.e., all entries in the matrix are i.i.d. uniformly in $[-1, 1]$. Observe that these matrices do not necessarily have a low-rank structure. Using these input matrices on the same pre-shift model as before, we find that model still copies the input (Row 2, Table 1) similar to the rank-2 case.

3.1.2. ROLE OF ATTENTION HEADS

Attention heads at this stage (Fig. 9a) do not appear to attend to any specific tokens in an interpretable manner. Since the model is simply copying the masked input, we hypothesize that attention heads (that combine different tokens) are inconsequential to the model output. To quantitatively verify this hypothesis we use *uniform ablation*: simply replace the softmax probabilities in the attention head by $1/n^2$ for all elements i.e. equally attend to all tokens (Sec. 4.6, [18]). With these ablations, there is negligible change in model performance at both observed and masked positions. Averaged over 256 samples, $L_{obs} = 3.4e-4$ and $L_{mask} = 0.2236$ when using all attention heads; whereas, on ablating all heads, these values are $3.2e-4$ and 0.2236 respectively. Clearly in this case, attention heads do not substantially affect the model prediction.

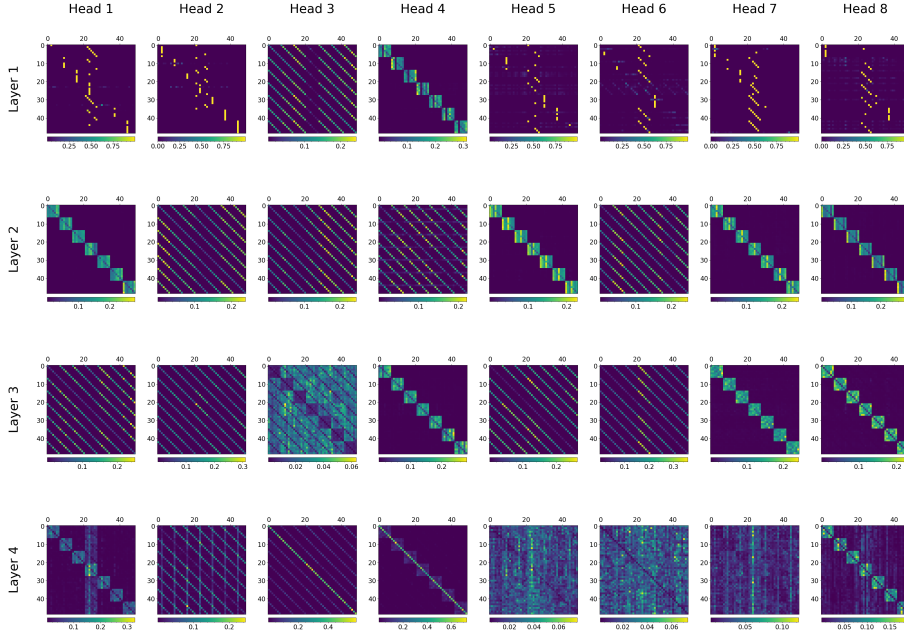


Figure 3: Attention heads in post-shift model show distinct regions in the input they attend to.

As further confirmation, we replace the key, query and value weights in the pre-shift model by those from the post-shift model. Averaged over 256 samples, L_{obs} is $5e-3$, that is similar to the optimal total MSE obtained at the end of training, while $L_{mask} = 0.2246$, similar to that obtained without replacing the weights.

3.2. After the Algorithmic Shift – Matrix Completion

We analyze the post-shift model separately for missing and observed entries, with a focus on the role of attention heads given the apparent interpretable patterns in Fig. 3. We find that for observed entries, the model output is still not substantially affected by the attention heads. However, for masked entries, the effect is substantial. Moreover, at this stage, the attention heads exhibit interpretable structure, that indicates the model is *actually* solving matrix completion instead of trivial copying.

3.2.1. OBSERVED ENTRIES

To check the effect of attention heads, we uniformly ablate *all* attention heads in the post-shift model. Averaged over 256 samples, this leads to $L_{obs} = 9.2e-5$ when using all attention heads, compared to $3.7e-3$ with ablation (close to the total MSE at model convergence). However, L_{mask} increases from 0.0128 to 0.2183, essentially the value in the pre-shift model. Further, we replace attention (key, query, value) weights in the post-shift model by weights from a pre-shift model. Indeed, averaged over 256 samples,

$L_{obs} = 9.5e-4$ in this case, supporting our claim.

Finally, since attention crucially depends on the position of elements, we randomly permute the positional embeddings in the post-shift model. That is, the embedding originally encoding position i in the input now represents position $\pi(i)$ for some random permutation $\pi : [n^2] \rightarrow [n^2]$. Averaged over 256 samples, $L_{obs} = 2.4e-4$, whereas $L_{mask} = 0.5687$, implying that the observed positions are negligibly affected compared to masked positions due to this intervention. This supports our ‘sub-algorithm’ hypothesis; from an intuitive viewpoint, positional information is not required for copying.

3.2.2. MISSING ENTRIES

To confirm that attention heads causally affect the model output for missing entries, in addition to uniform ablations, we perform *causal interventions* (activation patching) [37] on the hidden states just after the attention heads. This involves replacing the hidden state after an attention head for input A with the hidden state obtained at the same attention head, but for a different input A' . Ideally, if that head is causally relevant to the output, then such an intervention should steer the model towards the output for A' , instead of A . We find in our case that for $A = X$ and $A' = -X$, such an intervention simultaneously on all attention heads clearly steers the model output at missing entries towards $-X$ (more details in Appendix J).

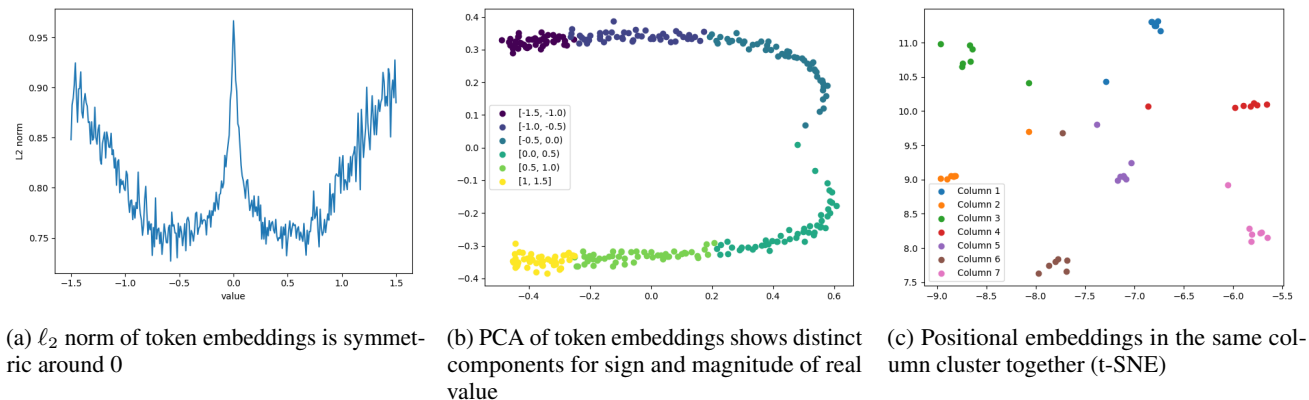


Figure 4: Embeddings in the post-shift model display interpretable behavior.

Denote attention head H in layer L by the tuple (L, H) . We can group the attention heads depending on the specific regions of the input matrix they attend to: (a) the same row as the query element – the ‘block-diagonal’ patterns, e.g. $(2, 1)$; (b) the same column as the query element – the ‘parallel-off-diagonal’ patterns, e.g. $(2, 2)$; (c) the query element itself – the ‘diagonal’ patterns, specifically in the last layer, e.g. $(4, 3)$. There are also some other attention heads that do not neatly fit into either of these 3 categories – for example, all heads in layer 1 except $(1, 2)$, $(1, 3)$; $(3, 3)$; $(4, 2)$, $(4, 5-7)$. In this context, we note that uniformly ablating heads $(3, 3)$, $(4, 2)$, $(4, 5-7)$ gives $L_{obs} = 9.36e-5$, $L_{mask} = 0.01575$ compared to $L_{obs} = 9.44e-5$, $L_{mask} = 0.01428$ without ablation, i.e. these uninterpretable heads do not significantly affect the output.

Additional Experiments We investigate attention heads when the observed entries in the input are arranged in a structured manner to analyse the function of individual attention heads (Appendix E). Further, we also probe hidden states of intermediate layers to understand the working of the model (Appendix I). Moreover, positional and token embeddings of the model also exhibit structure relevant to the problem, in some cases *before* the algorithmic shift occurs (Appendix F).

4. The Curious Case of Embeddings

Interpretable Embeddings In the post-shift model, positional and token embeddings also exhibit interesting properties related to the input elements and structure. For instance, the ℓ_2 norm of token embeddings corresponding to values from -1.5 to 1.5 is symmetric w.r.t. 0 as seen in Fig. 4a. Further, the PCA of token embeddings in Fig. 4b shows that the embeddings have a separable structure based on the sign of the real-valued input (y -axis), and continuous variation w.r.t. the absolute value of the real-valued input (x -axis).

The t-SNE projection of positional embeddings also show an interesting clustering pattern; positions in the same column tend to cluster together as seen in Fig. 4c. This is especially important because we have not used any marker tokens to mark the end of a row or column. Additionally, the ℓ_2 norm of positional embeddings (Fig. 8) is nearly constant across positions, except for a drop at positions around 21 – 26; that is, most of the middle row of the 7×7 input. This can be understood as the model marking the ‘origin’ of the position range from 1 to 49, and use it in subsequent computation.

Do embeddings change abruptly? Unlike attention heads (Fig. 9), embeddings might not abruptly change with the algorithmic shift. Similar to [23], we compute the top-2 principal components of the token embeddings at the final step (50000), and project the token embeddings at intermediate training steps on these components. The results (Fig. 7) show that the embeddings align very closely to the final arrangement before the actual drop in loss. This hints towards a possibility that even though the model undergoes a sudden algorithmic shift, some components evolve beforehand and possibly are a driving force behind the shift.

5. Discussion

To conclude, BERT can be trained to solve low rank matrix completion as an instance of masked language modeling. We also showed that a sudden drop in training loss marks an algorithmic shift in the model from merely copying the input to actually solving the task. The question of *why* this shift occurs suddenly rather than gradually, as well as a concrete characterization of the algorithm used by the model for missing entries are interesting directions for future work.

Acknowledgements We thank Yu Bai, Andrew Lee, and anonymous reviewers for their helpful suggestions. WH acknowledges support from the Google Research Scholar Program.

References

- [1] Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [2] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. In *The Eleventh International Conference on Learning Representations*, 2023.
- [3] Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes, 2017.
- [4] Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [5] Emmanuel J. Candès and Benjamin Recht. Exact matrix completion via convex optimization. *Found. Comput. Math.*, 9(6):717–772, 2009.
- [6] Francois Charton. Linear algebra with transformers. *Transactions on Machine Learning Research*, 2022.
- [7] Angelica Chen, Ravid Shwartz-Ziv, Kyunghyun Cho, Matthew L Leavitt, and Naomi Saphra. Sudden drops in the loss: Syntax acquisition, phase transitions, and simplicity bias in MLMs. In *The Twelfth International Conference on Learning Representations*, 2024.
- [8] Xiang Cheng, Yuxin Chen, and Suvrit Sra. Transformers implement functional gradient descent to learn non-linear functions in context, 2024.
- [9] Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What does BERT look at? an analysis of BERT’s attention. In Tal Linzen, Grzegorz Chrupała, Yonatan Belinkov, and Dieuwke Hupkes, editors, *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, Florence, Italy, August 2019. Association for Computational Linguistics.
- [10] Arthur Conmy, Augustine N. Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [12] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- [13] Deqing Fu, Tian-Qi Chen, Robin Jia, and Vatsal Sharan. Transformers learn higher-order optimization methods for in-context learning: A study with linear models, 2023.
- [14] Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [15] Michael Hanna, Ollie Liu, and Alexandre Variengien. How does GPT-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [16] Yifan Hou, Jiaoda Li, Yu Fei, Alessandro Stolfo, Wangchunshu Zhou, Guangtao Zeng, Antoine Bosselut, and Mrinmaya Sachan. Towards a mechanistic interpretation of multi-step reasoning capabilities of language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Singapore, December 2023. Association for Computational Linguistics.

- [17] Adam Jermyn and Buck Shlegeris. Multi-component learning and s-curves - ai alignment forum. *Eleventh International Conference on Learning Representations*, 2023.
- [18] Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. Revealing the dark secrets of BERT. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4365–4374, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [19] János Kramár, Tom Lieberum, Rohin Shah, and Neel Nanda. Atp*: An efficient and scalable method for localizing llm behaviour to components, 2024.
- [20] Nayoung Lee, Kartik Sreenivasan, Jason D. Lee, Kangwook Lee, and Dimitris Papailiopoulos. Teaching arithmetic to small transformers, 2023.
- [21] Michael A. Lepori, Thomas Serre, and Ellie Pavlick. Uncovering intermediate variables in transformers using circuit probing, 2023.
- [22] Tom Lieberum, Matthew Rahtz, János Kramár, Neel Nanda, Geoffrey Irving, Rohin Shah, and Vladimir Mikulik. Does circuit analysis interpretability scale? evidence from multiple choice capabilities in chinchilla, 2023.
- [23] Ziming Liu, Ouail Kitouni, Niklas S Nolte, Eric Michaud, Max Tegmark, and Mike Williams. Towards understanding grokking: An effective theory of representation learning. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 34651–34663. Curran Associates, Inc., 2022.
- [24] Ziming Liu, Eric J Michaud, and Max Tegmark. Omnigrok: Grokking beyond algorithmic data. In *The Eleventh International Conference on Learning Representations*, 2023.
- [25] Arvind V. Mahankali, Tatsunori Hashimoto, and Tengyu Ma. One step of gradient descent is provably the optimal in-context learner with one layer of linear self-attention. In *The Twelfth International Conference on Learning Representations*, 2024.
- [26] Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. Circuit component reuse across tasks in transformer language models, 2024.
- [27] Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. In *The Eleventh International Conference on Learning Representations*, 2023.
- [28] Neel Nanda, Andrew Lee, and Martin Wattenberg. Emergent linear representations in world models of self-supervised sequence models. In Yonatan Belinkov, Sophie Hao, Jaap Jumelet, Najoung Kim, Arya McCarthy, and Hosein Mohebbi, editors, *Proceedings of the 6th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, Singapore, December 2023. Association for Computational Linguistics.
- [29] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. In-context learning and induction heads. *Transformer Circuits Thread*, 2022. <https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html>.
- [30] Anna Rogers, Olga Kovaleva, and Anna Rumshisky. A primer in BERTology: What we know about how BERT works. *Transactions of the Association for Computational Linguistics*, 8, 2020.
- [31] Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language models a mirage? In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [32] Elena Voita, David Talbot, Fedor Moiseev, Rico Senrich, and Ivan Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In Anna Korhonen, David Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Florence, Italy, July 2019. Association for Computational Linguistics.
- [33] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, Joao Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 35151–35174. PMLR, 23–29 Jul 2023.
- [34] Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability

in the wild: a circuit for indirect object identification in GPT-2 small. In *The Eleventh International Conference on Learning Representations*, 2023.

- [35] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Huggingface’s transformers: State-of-the-art natural language processing, 2020.
- [36] Zhiwei Xu, Yutong Wang, Spencer Frei, Gal Vardi, and Wei Hu. Benign overfitting and grokking in reLU networks for XOR cluster data. In *The Twelfth International Conference on Learning Representations*, 2024.
- [37] Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models: Metrics and methods. In *The Twelfth International Conference on Learning Representations*, 2024.

A. Low Rank Matrix Completion

Low-rank matrix completion is a well-studied problem in machine learning and statistics. This problem finds applications in recommender systems, where given an incomplete matrix of user ratings on some items, the goal is to recover the missing entries *assuming* the ground truth matrix is low-rank. For a matrix $X \in \mathbb{R}^{n \times n}$, denote its observed (visible) entries by the set $\Omega \subset [n] \times [n]$, and the set of missing entries by $\Omega^C = [n] \times [n] \setminus \Omega$. Formally, the problem is

$$\min_U \text{rank}(U) \quad \text{s.t. } U_{ij} = X_{ij} \quad \forall (i, j) \in \Omega.$$

Nuclear norm minimization Since rank is not a convex function of the matrix entries, nuclear norm minimization [5] is a widely used convex optimization approach to low-rank matrix completion. The modified optimization problem is,

$$\min_U \|U\|_* \quad \text{s.t. } U_{ij} = X_{ij} \quad \forall (i, j) \in \Omega \quad (1)$$

where $\|U\|_*$ denotes the nuclear norm (sum of singular values) of matrix U . A regularized version of this problem for $\lambda > 0$ is

$$\min_U \left[\frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} (U_{ij} - X_{ij})^2 + \lambda \|U\|_* \right]. \quad (2)$$

B. Experimental details

Data Preprocessing We tokenize real values as follows: discretize the range $[-10, 10]$ (all matrix entries in our experiments are in this range) in steps of size $\epsilon = 0.01$, and assign token IDs to these values with IDs starting from 1; the mask token (MASK) is assigned token ID 0. Input to the transformer is the tokenized masked sequence $X_{mask} = \text{TOK}(\text{Vec}(X \odot M))$, where TOK denotes tokenization, Vec denotes vectorizing the $n \times n$ matrix to a n^2 -dimensional vector and \odot denotes the element-wise product. Due to this discretization, in experiments, MSE will be with the rounded-off version of X to 2 decimals.

Training We use the BERT model implementation from the HuggingFace library [35], with ‘absolute’ positional embeddings and no dropout. Since the model maps a sequence of discrete token IDs to a sequence of real values, we compute the MSE loss between the real valued model output, and the discretized real values in the ground truth matrix. For example, for input sequence $[0.12, 0.45, 0.87] \in \mathbb{R}^3$, corresponding discretized and masked token sequence [“0.12”, “MASK”, “0.87”], and output $[x_1, x_2, x_3] \in \mathbb{R}^3$, the MSE loss is $\frac{1}{3} [(x_1 - 0.12)^2 + (x_2 - 0.45)^2 + (x_3 - 0.87)^2]$.

We use the MSE loss on *all* elements of the input and output matrices for training. We additionally fix the masking probability $p_{\text{mask}} = 0.3$ in all cases. Using a 12-layer, 12-heads per layer BERT model with a linear read-out layer, the train loss is optimized using Adam with constant step size $5e-5$ for 50000 epochs (without weight decay or warmup). Data at each step is obtained by sampling 256 train and 64 test matrices. Since the training is ‘online’, train and test losses are nearly identical at all points in training, and thus we will not separately analyze them.

We additionally train a smaller BERT model with 4 layers and 8 heads per layer, with step-size $1e-4$ on square matrices of order 7 and rank-2. We use this smaller model primarily to keep our interpretability analyses tractable; in any case the attention heads are similar to those in larger models (Appendix L).

The model converges to a total MSE of the order $1e-3$ (i.e. solves matrix completion well) for all runs – square matrices of order 7, 10, 12, 15 and rank 2, 3, 3, 4 respectively. For the 4-layer 8-heads case, we obtain comparable performance (final total MSE $\sim 4e-3$) to the 12-layer, 12-head model. The plot in Figure 2 demonstrates the loss evolution over the course of training the model in the smaller model setup. All discussion in the subsequent sections is specifically for the 4-layer, 8-heads model on a 7×7 rank-2 input. In addition, ‘pre-shift’ denotes this model at step 4000, while ‘post-shift’ denotes model at the end of training i.e. step 50000.

C. Pre-shift copying

Input Samples ↓	Mask = "MASK"		Mask = "0.44"		Mask = "-0.24"	
	L'_{mask}	L_{obs}	L'_{mask}	L_{obs}	L'_{mask}	L_{obs}
Rank-2 matrices	3e-4	3.3e-4	4e-4	2.8e-4	3.7e-4	2.7e-4
Random matrices	2.8e-4	3.5e-4	3.7e-4	3e-4	3.6e-4	2.8e-4

Table 1: Pre-shift model implements copying, predicting the value for mask token at missing entries.

D. Nuclear Norm Minimization

Is it possible that the model is implicitly implementing nuclear-norm minimization to predict missing entries in the input?

We use CVXPY to solve matrix completion using nuclear-norm minimization at various levels of p_{mask} comparing it to the output of a BERT model trained on $p_{mask} = 0.3$. We find that BERT performs better than nuclear norm minimization w.r.t. MSE; at the same time, the nuclear norm of BERT solution is larger (Fig. 5)

Further, we also solve the regularized version of the above problem (Eq. 2) to attempt to match the performance of BERT at some $\lambda > 0$, and verify if the model indeed implicitly optimizes such an objective. We find that this is not the case, as for various values of λ , BERT still outperforms regularized MSE minimization; please see Appendix D for details.

We use the regularized version of the nuclear norm minimization problem as detailed in Sec. D, and obtain the following L, L_{obs}, L_{mask} for various values of λ . We average our results over 256 samples generated in the same way as the training data for BERT (including rounding off to 2 decimal places) for the sake of comparison.

λ	L_{obs}	L_{mask}	L
0.0005	1.015e-5	0.040728	0.012173
0.001	3.686e-5	0.040456	0.01211
0.0015	7.959e-5	0.040505	0.012155
0.002	0.00013769	0.040734	0.012264
0.005	0.00078591	0.043402	0.013516

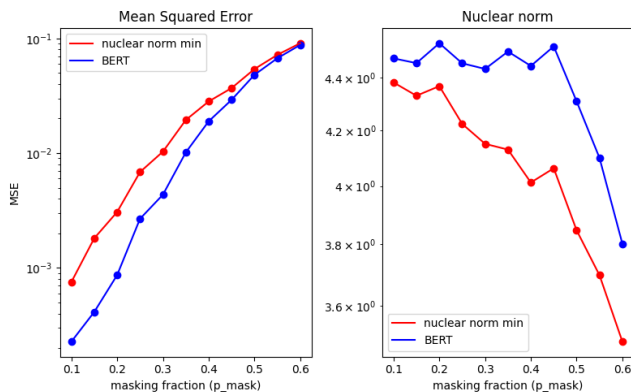


Figure 5: BERT outperforms nuclear norm minimization.

E. Attention Heads with Structured Mask

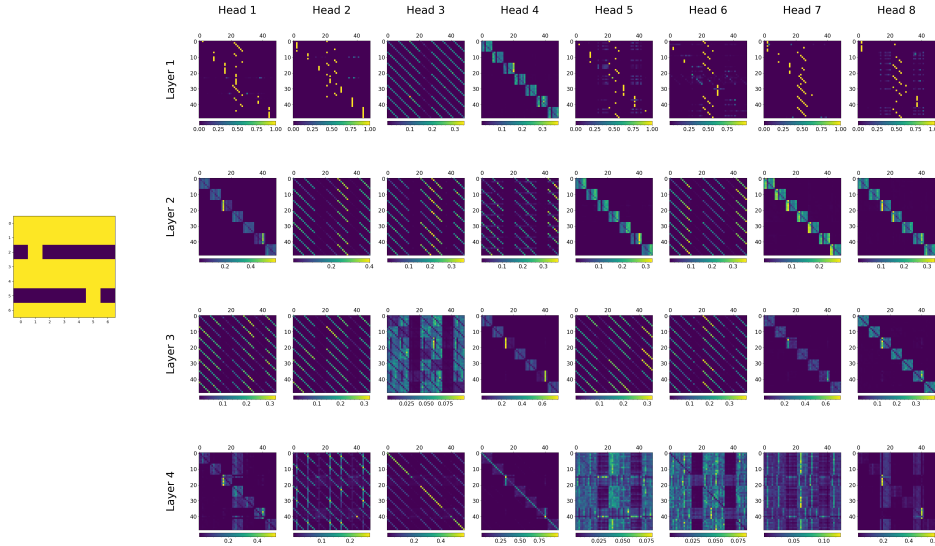


Figure 6: Attention heads for a structured mask attend to specific entries in the input. Left: structured mask (blue denotes missing entries)

Since the maps in Fig. 3 are averaged over multiple random masks and input matrices, it is difficult to extract more specific details about the algorithm, apart from the coarse-grained insights as above. To remedy this, we generate inputs with specific mask structure, see for example Fig. 6. This implies that for different input matrices, the mask i.e. Ω^C remains the same. This step helps us highlight how an attention head attends to input elements based on the element being masked or observed. From the results in Fig. 6, it is evident that different attention heads focus on specific parts of the input. For instance,

1. (2, 1), (3,4) and (4,8) are significantly active only at the masked rows, and in those cases has maximal attention at the only observed positions in those rows. In other words, this head acts as a ‘masked-only’ head.
2. (4,3) and (4,4) correspond roughly to an identity map, slightly deviating in the masked rows. In these cases, again the maximal attention score corresponds to the only observed position in these rows. That is, this head acts as an ‘identity-map’ head.
3. Further, there are multiple ‘parallel off-diagonal’ heads that completely ignore the masked rows for their computation. These heads include (2,2–4), (2,6); (3,2), (3,3), (3,5). Additionally, there are also attention heads like (3,1), (3,6) that attend to only the observed element of each masked row. Collectively these heads act as ‘observed-only’ heads, attending to only observed entries, and using this information to compute missing entries.
4. There also exist attention heads that respond systematically to changes in the mask. For example, consider attention heads (2, 5), (2, 7), (2, 8) in Fig. 10. For each row, these heads attend to the element in the 6th and 2nd column respectively for part (a) and (b). On a closer look, the connecting link between these two mask patterns is that, the longest contiguous unmasked column is exactly the column that these heads attend to. We hypothesize that this information is somehow used by the model in its inner computation for masked entries.
5. Finally, Heads (1,1–2), (1, 5–8) do not fall in any of the categories above. These heads are mostly static across different mask / input variations (for example, comparing Fig 3 and 6), and the patterns suggest that these heads almost exclusively focus on the middle row of the input matrix and some other elements. A possible function of these heads is to process positional and token embeddings (input to the first layer) so that this information can be used appropriately in the subsequent layers.

F. Embeddings

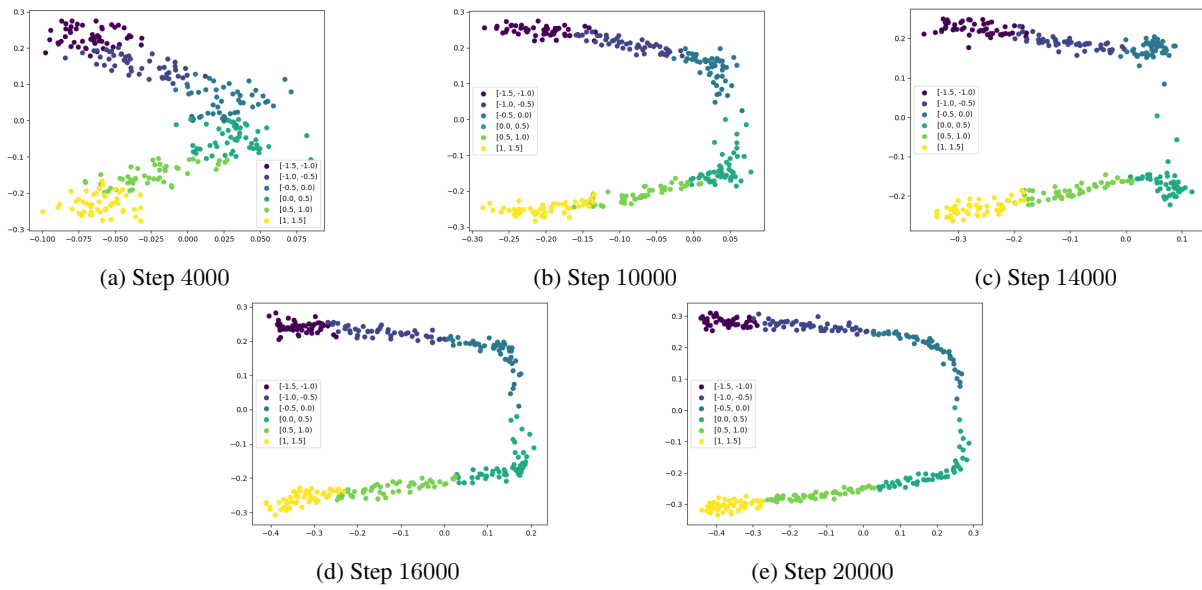


Figure 7: Projection of token embeddings along principal components of embeddings at step 50000.

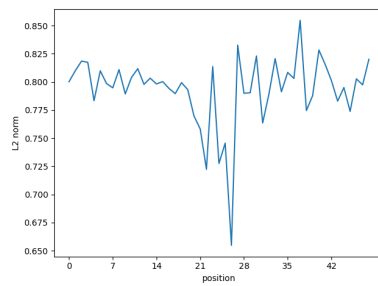
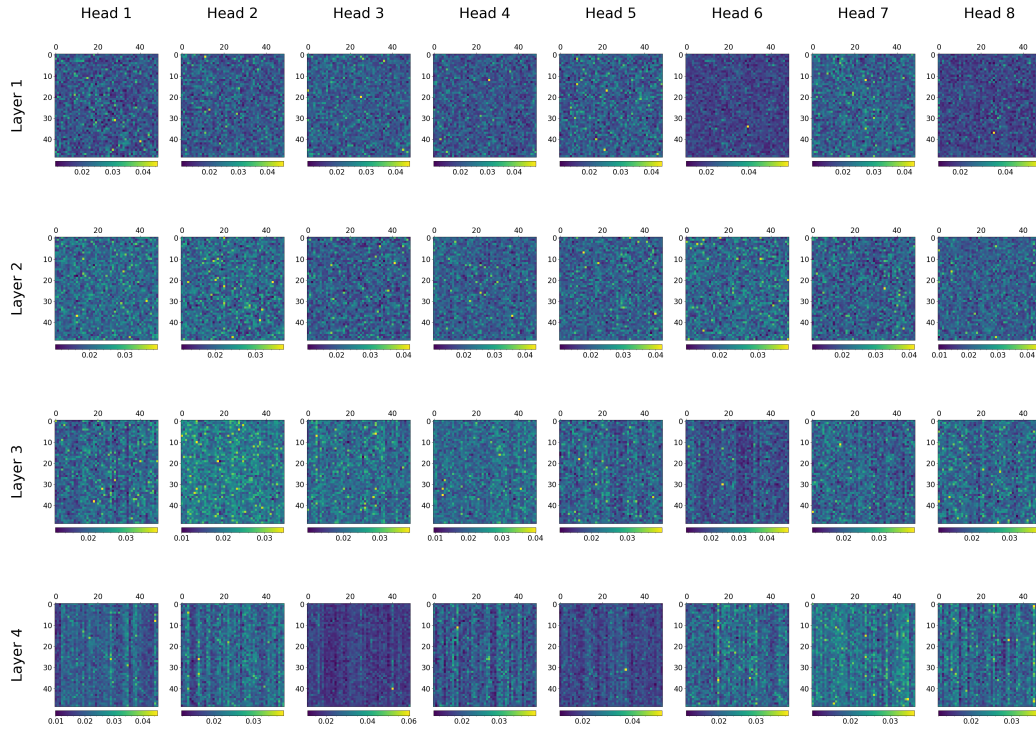
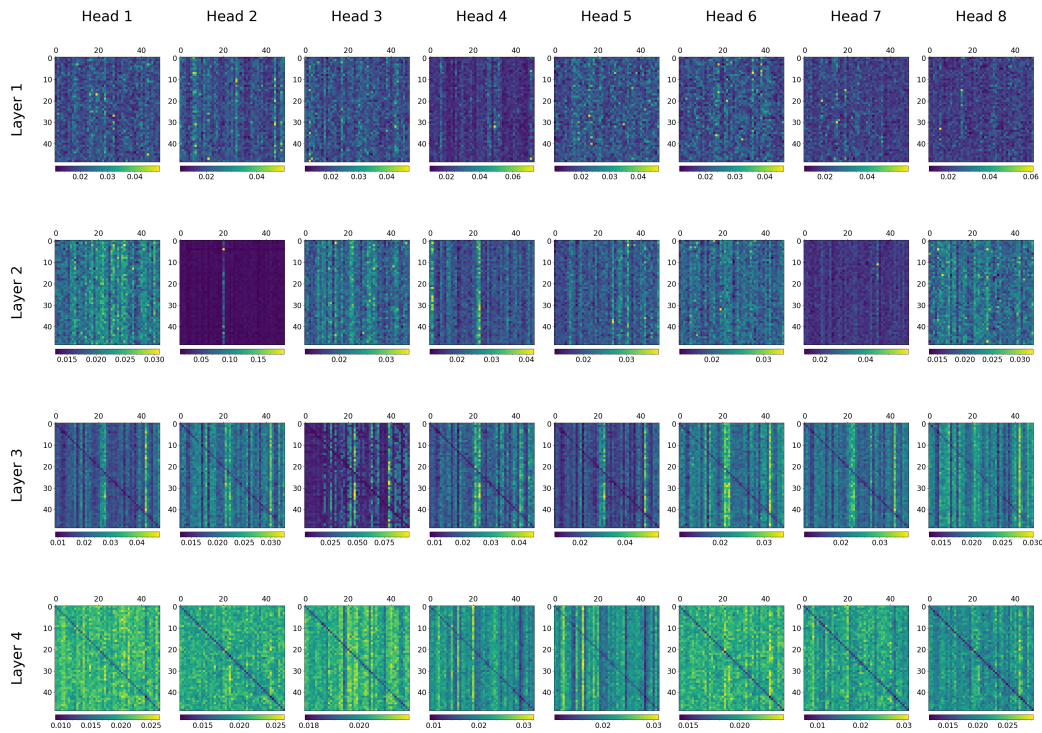


Figure 8: ℓ_2 norm of positional embeddings in post-shift model.

G. Evolution of Attention Heads during training

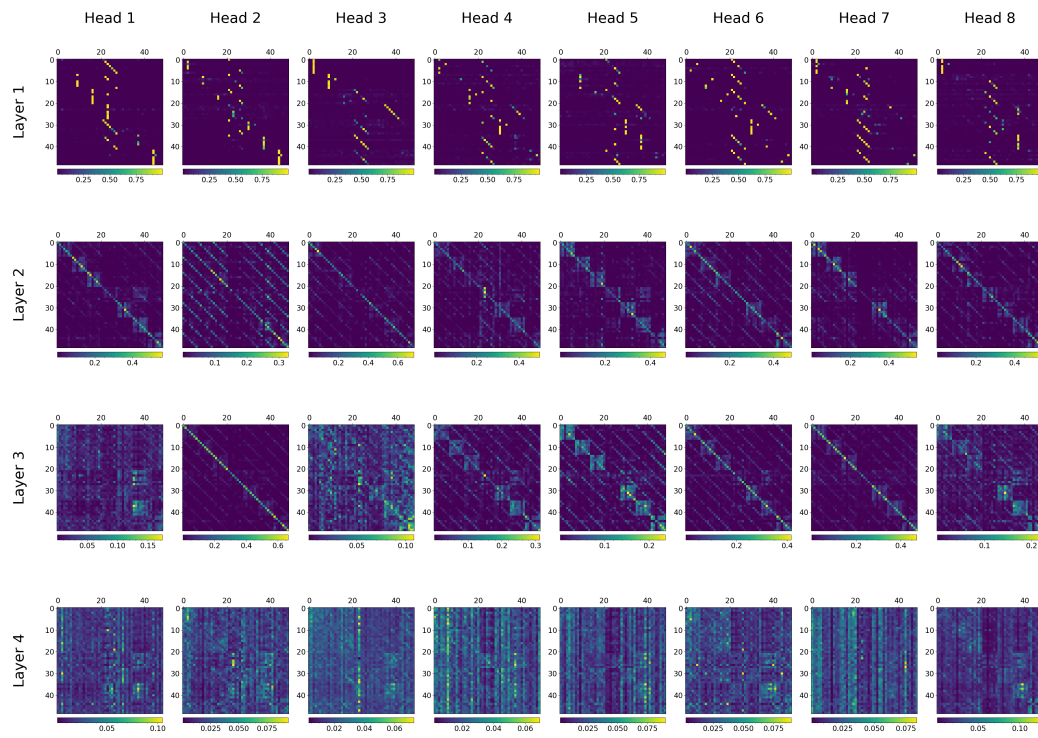


(a) Step 4000

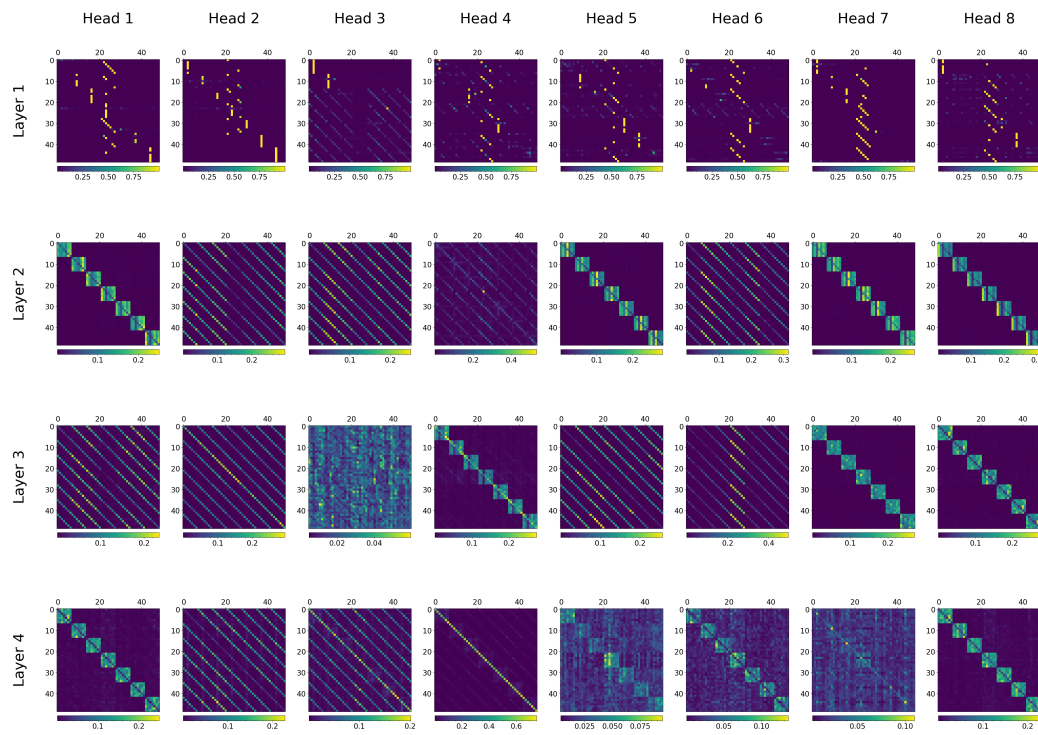


(b) Step 14000

How Do Transformers Fill in the Blanks? A Case Study on Matrix Completion



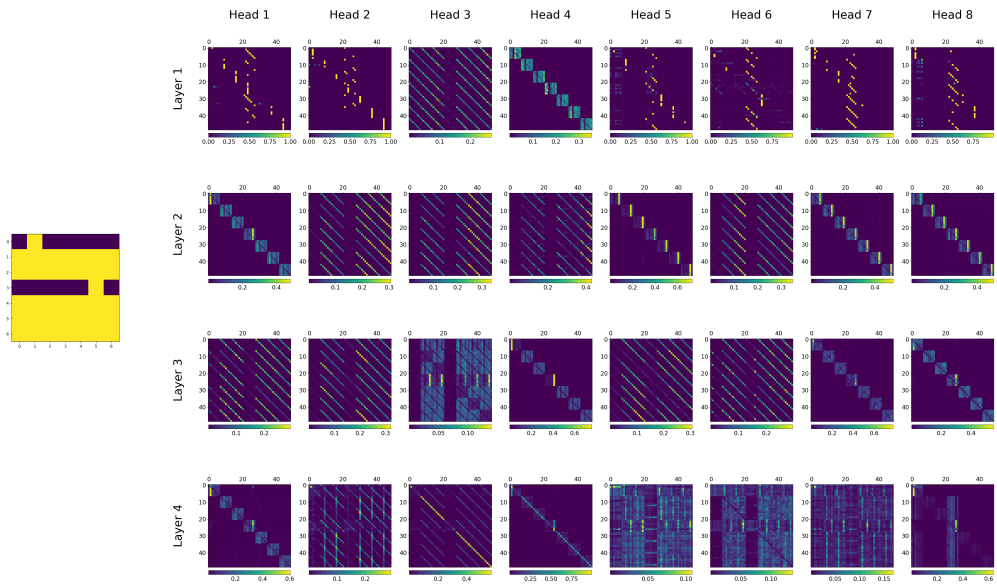
(c) Step 16000



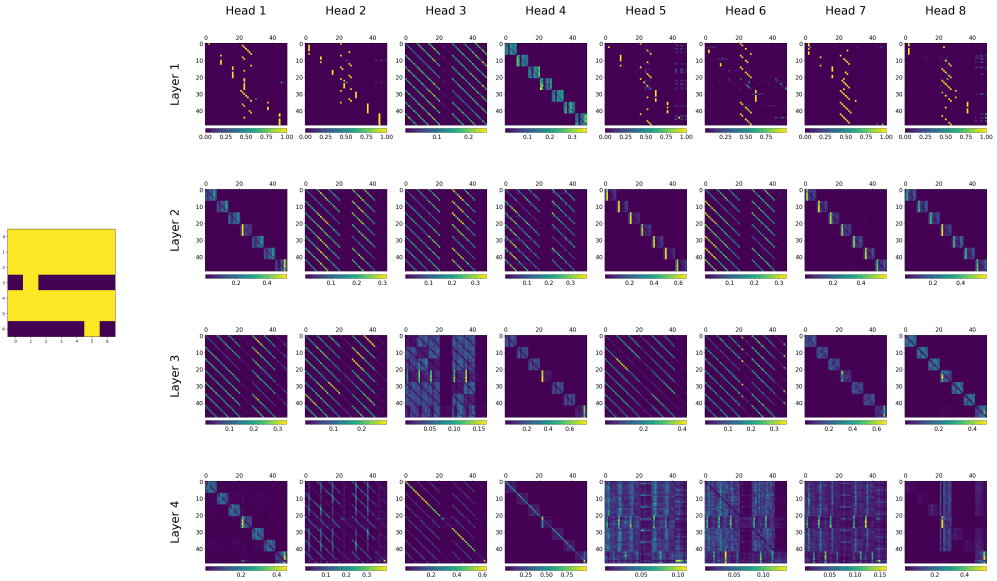
(d) Step 20000

Figure 9: Attention heads across various training steps.

H. Attention Heads with structured mask



(a)



(b)

Figure 10: Attention heads and corresponding masks; blue denotes masked position in the input matrix.

I. Probing

We probe for properties of the input matrix in the hidden states of the model, to concretely determine how the model computes the output. We use our 12-layer model in this case, for enhancing contrast between probing in different layers. Specifically, for every element in the input, we use a linear probe [3] on its hidden state after a given layer, mapping the hidden state to the n -dimensional masked row that this element belongs to. Missing entries are replaced by 0, and the linear probe is fit using least squares. The results for this experiment in Fig. 11 demonstrate that, layer 3 and 4 in the model correspond quite strongly to the probe target, compared to other layers. This suggests that the model tracks input information in its intermediate layers and uses it for computation.

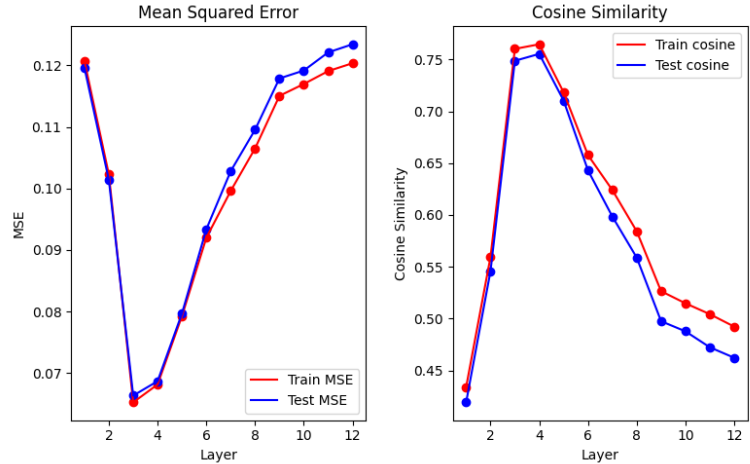


Figure 11: Layer 3 and 4 store information about the rows of the masked input matrix.

J. Causal effect of Attention heads

To verify whether attention heads actually contribute towards the model output, or are simply a side-effect of some other latent factor in the model, we employ 2 methods used earlier to quantify the contribution of attention heads in transformers.

- Uniform Ablation** Following the methodology in (Sec 4.6, [18]), for a square matrix input of order n , we set each element of the $n^2 \times n^2$ softmax attention matrix to $1/n^2$. That is, attend equally to all tokens in the input sequence, and remove any learned information about attending to specific positions in the input.
- Causal Interventions** In the uniform ablation setup, it is possible that setting the softmax probabilities to a given value might change the distribution of resultant hidden states, and consequently degrade model performance. A more principled technique to analyse the effect of a specific component is to replace the hidden state just after that component by hidden states on a different input, and analyse how this affects the final output [37]. In our case, we intervene on attention heads by replacing the hidden state after an attention head for input matrix X by the hidden state for input $(-X)$. Importantly, this change does not affect properties like rank of the input, and hence the hidden states obtained are from the same distribution as those for input X .

Pre-shift In the pre-shift model, we want to demonstrate that removing attention heads does not affect the model predictions significantly. For this, we uniformly ablate all attention heads in the pre-shift model, and measure the effect averaging over 256 samples. We get that $L_{obs} = 3.4e-4$ and $L_{mask} = 0.2236$ when using all attention heads; whereas, on ablating all heads, these values are $3.2e-4$ and 0.2236 respectively. Clearly, in the pre-shift model, attention heads do not substantially affect the model prediction.

Post-shift In the post-shift model, we want to demonstrate that the attention heads causally affect the output. Using uniform ablation, we get that $L_{obs} = 9.2e-5$ and $L_{mask} = 0.0128$ when using all attention heads; whereas, on ablating all heads, these values are $3.7e-3$ and 0.2183 respectively.

From these observations, we could claim causal effect of attention heads for prediction at missing entries. A stronger test however is through causal interventions,

- Step 1 Extract the hidden states for all attention layers from the model on some input matrix X ; call these h_+ . Concretely, these hidden states are obtained just after the matrix product of the softmax attention probabilities and the value matrix and hence before the output matrix product.

Step 2 Change the input to the model to $-X$, however, also replace the hidden states *just after* the attention layers with h_+ obtained in Step 1. Call the output of the model in this setup as $f_p(-X, X)$.

We observe that, the MSE between $f_p(-X, X)$ and X , averaged over 256 samples at masked positions is approximately 0.014 (this is comparable to optimal L_{mask}), compared to the MSE between $f_p(-X, X)$ and $-X$ being 0.8066. This demonstrates that the attention heads are causally relevant to the model output for missing entries.

K. Related Work

Mathematical problem solving capabilities of Transformers have been a topic of interest lately [20; 6; 4]. In fact, [20] show that learning addition from samples is equivalent to low-rank matrix completion. Further, [6] show that it is possible to train a transformer based model to solve various linear algebraic tasks e.g. eigendecomposition, matrix inversion, etc.; however, to the best of our knowledge, interpretability studies for such tasks have not been conducted before. For interpretability in simpler math tasks, [15] mechanistically analyse GPT-2 small on predicting whether a number is ‘greater-than’ a given number, by formulating the problem as a natural language task. [30; 32; 9] analyse BERT from an interpretability perspective. More recently, there has been a line of research works analysing decoder based models to reverse-engineer the mechanisms employed by these models, termed as ‘mechanistic interpretability’ [12; 26; 27; 34; 10; 28; 21; 22; 19; 29; 16]. We note that our setting is distinct from the recent work on solving mathematical tasks like linear regression through ‘in-context’ learning in transformers [4; 1; 8; 13; 14; 2; 25; 33]. Whether our model learns to implicitly ‘implement’ an optimization procedure as shown in some of these works is an open question.

Further, [23; 24; 27; 36; 17] analyse ‘grokking’, the sudden emergence of generalization during model training. In the context of training dynamics of MLM, [7] analyses ‘breakthroughs’ (sudden drop in loss and associated improvement in generalization capabilities of the model), specifically for BERT. They show that the breakthrough marks the transition of the model to a generalizing one. Their work however is focused on language tasks, distinct from our setting which is more mathematical in nature. We also note that their work is not in the online training setting; our setup is online in the sense of sampling new data at every step of training.

L. Attention Heads for larger inputs

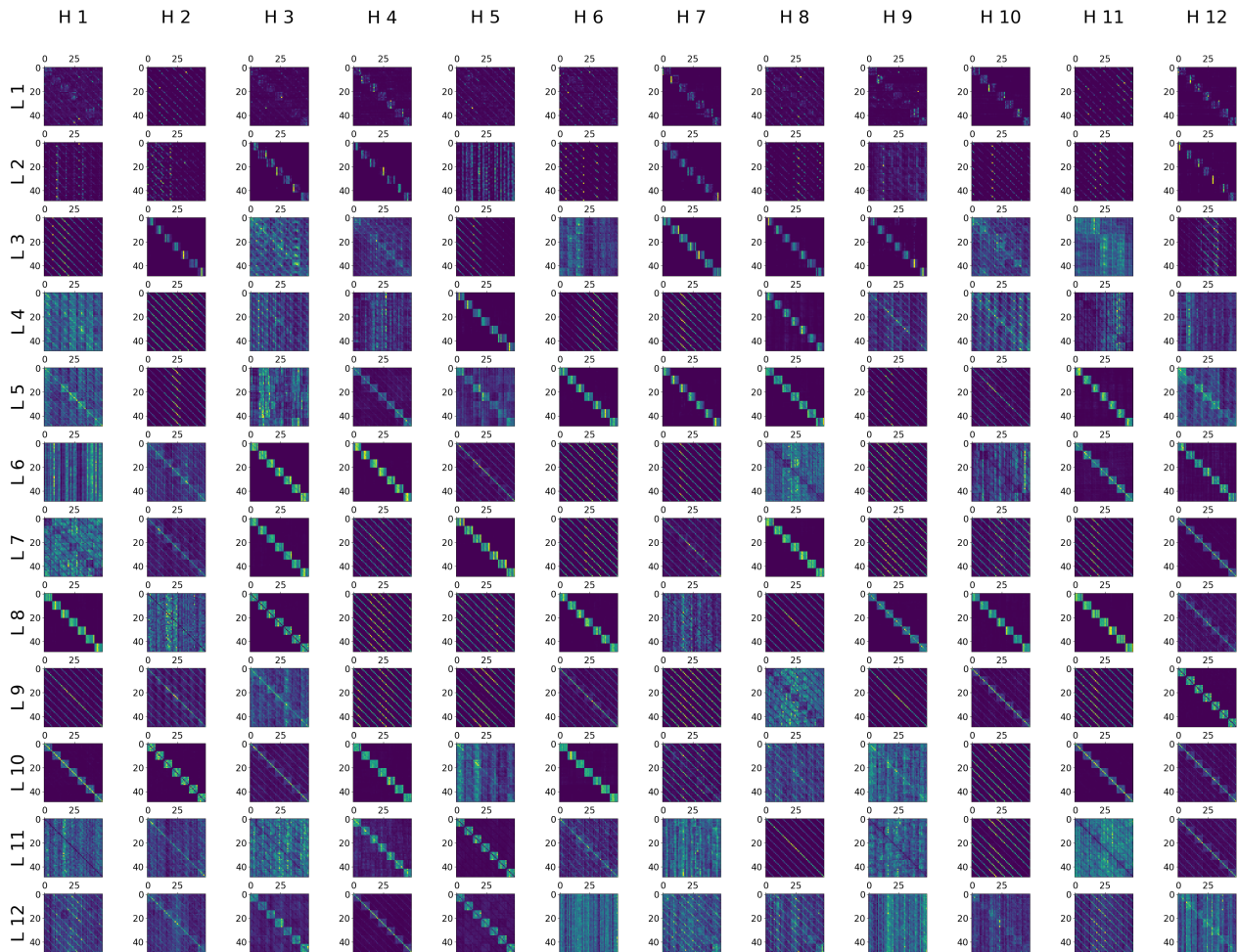


Figure 12: Attention heads in 12 layers, 12-heads model on 7×7 rank-2 input

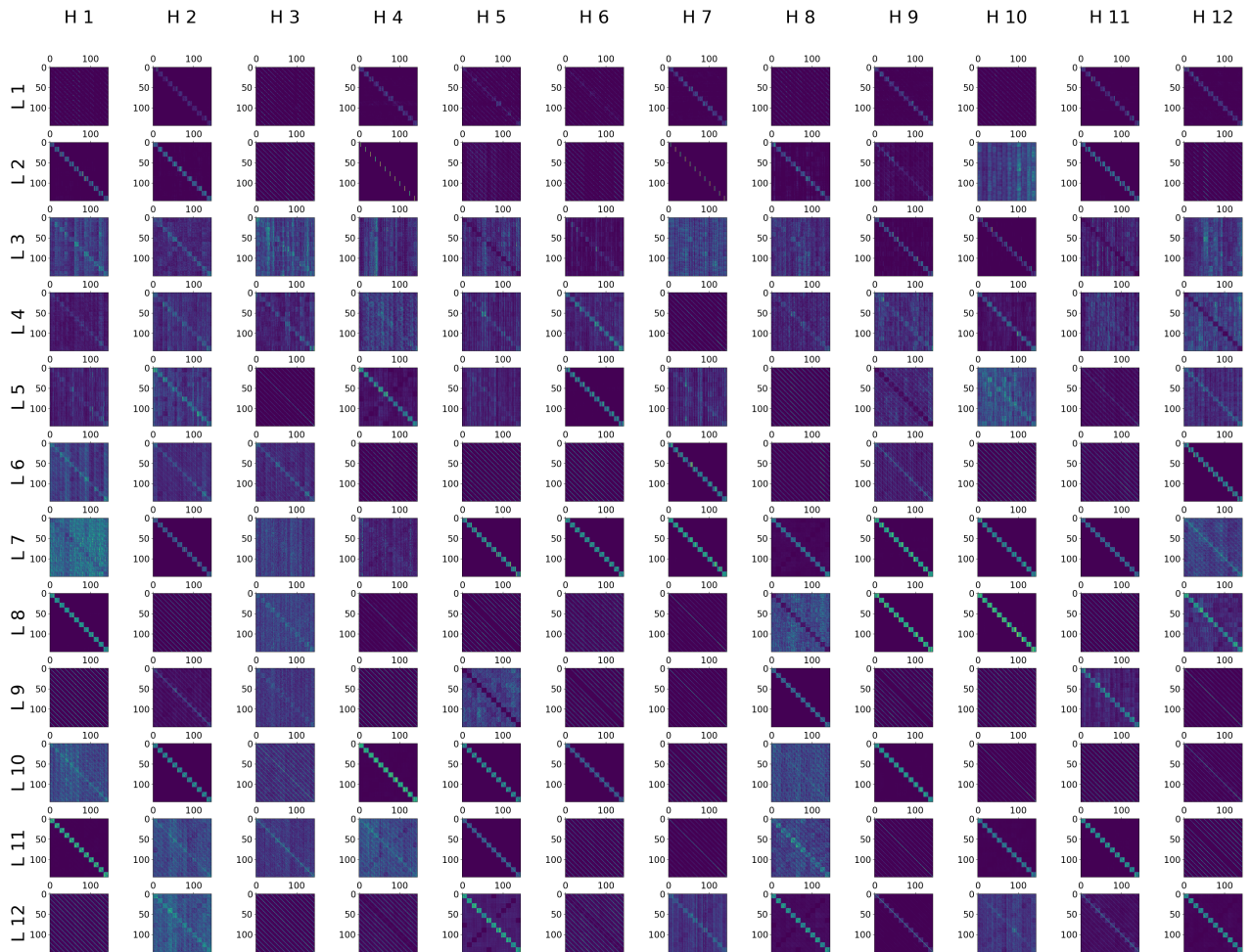


Figure 13: Attention heads in 12 layers, 12-heads model on 12×12 rank-3 input

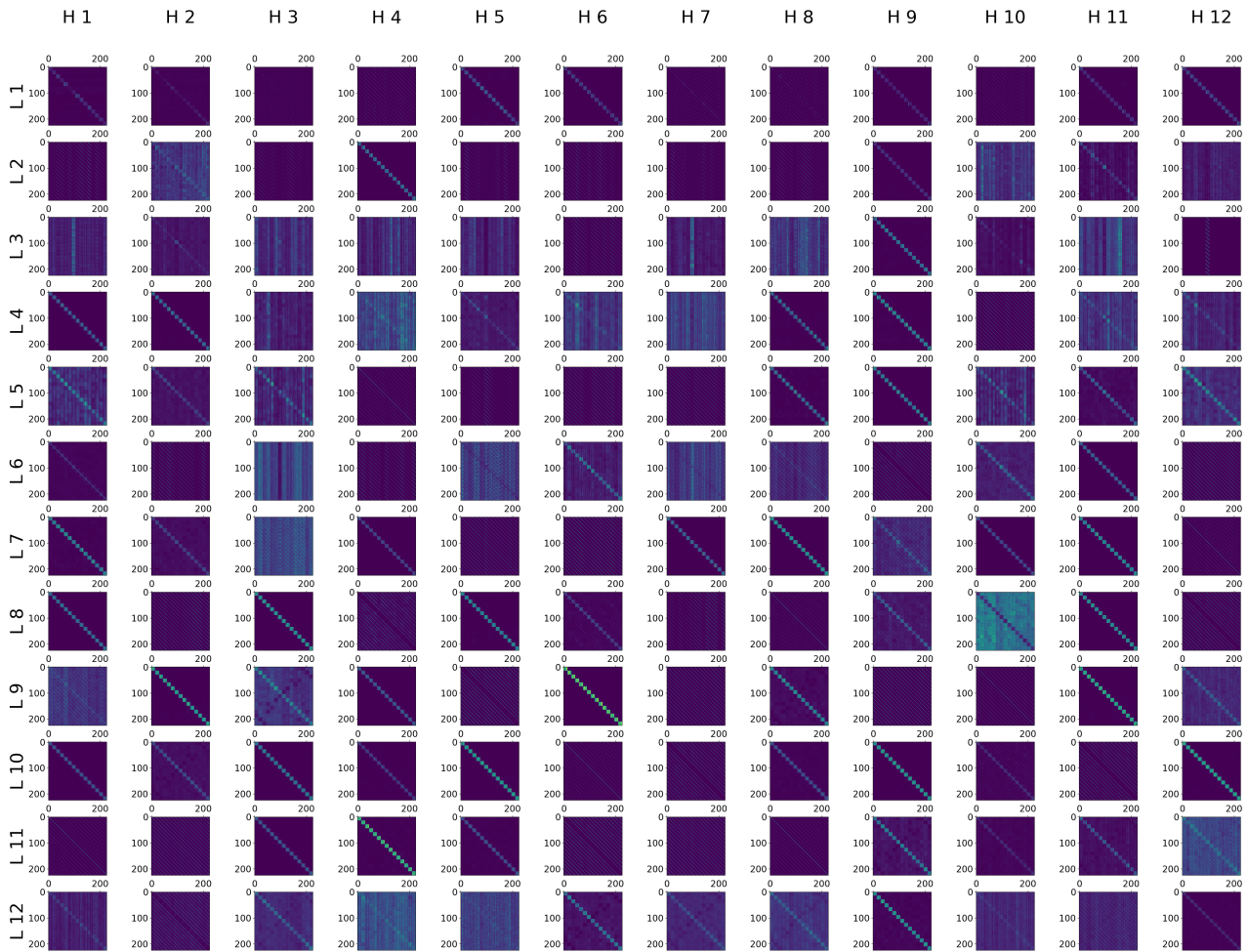


Figure 14: Attention heads in 12 layers, 12-heads model on 15×15 rank-4 input