

---

# Counting in Small Transformers: The Delicate Interplay between Attention and Feed-Forward Layers

---

Freya Behrens<sup>1</sup> Luca Biggio<sup>2</sup> Lenka Zdeborová<sup>1</sup>

## Abstract

Next to scaling considerations, architectural design choices profoundly shape the solution space of transformers. In this work, we analyze the solutions simple transformer blocks implement when tackling the histogram task: counting items in sequences. Despite its simplicity, this task reveals a complex interplay between predictive performance, vocabulary and embedding sizes, token-mixing mechanisms, and feed-forward layer capacity. We identify two theoretical counting strategies transformers adopt, relation-based and inventory-based counting, each defining distinct learning regimes for the task. These strategies dictate how functionality is distributed between attention and feed-forward layers. We further show that adding softmax and beginning-of-sequence tokens allow for more robustness when embedding dimensions are comparatively small. Empirical introspection of trained models closely confirms both the learning regimes of the various architectures and the formation of these strategies during training. We demonstrate how a basic task that requires only aggregation and selection is significantly impacted by minor design changes.

## 1. Introduction

Transformers are the key neural network behind many recent deep learning advances, most notably large language models (LLMs). Their success is partly due to their versatility in processing diverse data types, including text, images, and video, represented as sequences of tokens (Liu et al., 2021; Girdhar et al., 2019; Brown et al., 2020). While scale has

been a key factor in unleashing the potential of these models, it is remarkable that their architecture still largely follows the same simple template of the original transformer model proposed by Vaswani et al. (2017). At its core, a single transformer block primarily alternates two basic components: the token-mixing attention mechanism and a standard fully connected multi-layer perceptron. At a high level, the attention mechanism mixes the tokens, while the multi-layer perceptron applies a nonlinear feature transformation identically to each token. Despite the widespread use of transformers, there is no clear consensus on the distinct roles of their components, how they interact, or if they can be substituted with alternative modules (Tolstikhin et al., 2021; Dordevic et al., 2024; Gu & Dao, 2024). In particular, the specific contribution of each architectural element to the model’s hypothesis space –the range of algorithms it can learn and implement in practice– remains opaque (Weiss et al., 2021; Delétang et al., 2023; Abbe et al., 2023; Ouellette et al., 2023).

In this work, we investigate this question for algorithms that compare and aggregate information from a mechanistic perspective (Cammarata et al., 2020; Olah et al., 2020; Elhage et al., 2021; Michaud et al., 2023; Ouellette et al., 2023) and focus on the histogram task (Weiss et al., 2021) as a prototypical problem. It consists of predicting the number of appearances of each token in the input sequence processed by the model – counting - and solving it requires both comparison and aggregation. Even modern language models with up to 8B parameters currently fail to solve this task robustly and efficiently in-weights, rather than via chain-of-thought, (Appendix A). This failure, despite the task’s apparent simplicity, motivates a study of the relative role of different architectural components and their impact on the final solutions implemented by the model in a *controlled* setting. To this end, we focus on models following the architectural template of primitive transformer blocks, i.e. alternating a token-mixing attention mechanism and a multi-layer perceptron.

In our analysis, we provide explicit theoretical constructions (parameter configurations) for a range of such architectures reaching perfect accuracy in a model-dependent hyperparameter regime. In a subsequent step, we compare these

---

<sup>1</sup>Statistical Physics of Computation Laboratory, École polytechnique fédérale de Lausanne (EPFL), Lausanne, Switzerland <sup>2</sup>Department of Computing Sciences, Università Bocconi, Milan, Italy. Correspondence to: Freya Behrens <freya.behrens@epfl.ch>.

algorithms with the performance and mechanistic behavior of models trained from data. Our findings reveal that this class of models is capable of implementing strikingly different solutions for the histogram task, with a strong dependence on the scale of the model’s hyperparameters and the type of token-mixing mechanism utilized. Our main contributions are as follows:

- We identify two algorithmic strategies for solving the histogram task: *relation-based counting*, which uses local pairwise token comparisons, and *inventory-based counting*, which counts all tokens in the alphabet and extracts the desired count at a given position.
- We clarify how the emergence of these strategies depends on the architecture-specific inductive biases. Relation-based counting is memory and compute-efficient, leveraging attention for comparisons. Inventory-based counting relies on token-mixing for aggregation, but uses feed-forward layers to implement a comparison to the complete alphabet with higher computational demands.
- We show how the embedding size required for perfect solutions relates to alphabet size and the maximal sequence length. Due to the discrete nature of the task, near-orthogonality is sufficient in many cases. Gains in prediction accuracy compared to models of the same capacity can arise from the softmax operator and dot-product attention reducing noise from linear dependence.

Our code is publicly available at <https://github.com/SPOC-group/counting-attention>.

## 2. Background and Notation

**Architecture.** As inputs, we consider sequences of tokens  $\mathbf{x} = (x_1, x_2, \dots, x_L) \in \mathcal{T}^L$  from the alphabet  $\mathcal{T} = \{1, \dots, T\}$ . A sequence of outputs is  $\mathbf{y} = (y_1, \dots, y_L)$  has the same length as the input sequence, where each output token  $y_\ell \in \{1, \dots, C\}$ , with  $C \leq L$ . We analyze several 1-layer model architectures where a token-mixing mechanism is followed by a per-token feature transformation. Formally, we consider a models  $F : \mathcal{T}^L \rightarrow \mathcal{C}^L$  defined for the positions  $\ell = 1, \dots, L$  as

$$F(\bar{\mathbf{x}})_\ell = \arg \max_{c \in \{1, \dots, C\}} f(\bar{x}'_\ell)_c; \quad \bar{x}'_\ell = \bar{x}_\ell + [\mathbf{A}(\bar{\mathbf{x}})\bar{\mathbf{x}}]_\ell \quad (1)$$

with the token mixing matrix  $\mathbf{A} : \mathbb{R}^{L \times d} \rightarrow \mathbb{R}^{L \times L}$  and the token-wise feature transformation  $f : \mathbb{R}^d \rightarrow \mathbb{R}^C$ . The embedding is  $\bar{\mathbf{x}} \in \mathbb{R}^{L \times d}$ , where  $\bar{x}_\ell$  denotes its  $\ell$ -th row, is obtained by passing the input sequence  $\mathbf{x}$  into a standard embedding layer (learnable lookup-table) of dimension  $d$ . We refer to the embedding associated with token  $t \in \mathcal{T}$  as  $e_t \in \mathbb{R}^d$  or  $e_{x_\ell} \in \mathbb{R}^d$  for the embedding of the token  $x_\ell$  at position  $\ell$ . We do not include positional embeddings due to

the inherent permutation equivariance of the histogram task. We refer to the vector  $\bar{x}'_\ell$ , for each position  $\ell = 1, \dots, L$ , as the mixed token. Note that we assume that all operations in the network are executed with infinite precision. We comment on this assumption when it becomes problematic.

**Token Mixing.** We consider two types of mixing mechanisms  $\mathbf{A}$  with different activation functions and a single head. We refer to the case where the function  $\mathbf{A}$  is constant in  $\bar{\mathbf{x}}$  as *linear mixing* (`lin`), e.g.  $\mathbf{A}_{\text{lin}}(\bar{\mathbf{x}}) = A$  and  $\mathbf{A}_{\text{lin+sftm}}(\bar{\mathbf{x}}) = \text{softmax}(A)$ , where  $A \in \mathbb{R}^{L \times L}$  is a learnable matrix and the softmax operator is applied row-wise. The number of learnable parameters is therefore  $L^2$ .

As an alternative mixing structure, which we refer to as *dot-product mixing* (`dot`), we consider the popular attention mechanism which constructs the matrix  $\mathbf{A}$  to be explicitly dependent on the inputs, i.e.

$$\mathbf{A}_{\text{dot}}(\bar{\mathbf{x}}) = \frac{1}{\sqrt{d}} \bar{\mathbf{x}} W_Q W_K^T \bar{\mathbf{x}}^T \quad (2)$$

and  $\mathbf{A}_{\text{dot+sftm}}(\bar{\mathbf{x}}) = \text{softmax}(\mathbf{A}_{\text{dot}}(\bar{\mathbf{x}}))$  where  $W_Q$  and  $W_K$  are learnable  $d \times d$  matrices. Note that, without loss of generality, we assume the value matrix to be the identity. Then the number of parameters for dot-product mixing is  $2d^2$ . In line with previous work (Weiss et al., 2021), for architectures employing the dot-product mixing, we also analyze models utilizing the so-called beginning-of-sequence (BOS) token. This special token, indicated with the symbol  $\$,$  is appended to the original input  $\mathbf{x}$  resulting in a new sequence  $\tilde{\mathbf{x}} = (\$, x_1, x_2, \dots, x_L)$  of length  $L + 1$ . We will refer to the architecture that includes the BOS token as `bos`.

**Feature Transformation.** The feature transformation is a single hidden layer perceptron with ReLU activations. The hidden layer is of dimension  $p$ . The function  $f$  is applied identically to every mixed token  $\bar{x}'_\ell$  for  $\ell = 1, \dots, L$ , as:

$$f(\bar{x}'_\ell) = \text{ReLU}(\bar{x}'_\ell W_1 + b_1) W_2 + b_2 \quad (3)$$

where  $f(\bar{x}'_\ell) : \mathbb{R}^d \rightarrow \mathbb{R}^C$  and where the weights have the appropriate dimensions to accommodate a hidden layer of size  $p$ , i.e.  $W_1 \in \mathbb{R}^{d \times p}$ ,  $b_1 \in \mathbb{R}^p$ ,  $W_2 \in \mathbb{R}^{p \times C}$  and  $b_2 \in \mathbb{R}^C$ .

## 3. Experimental Setup

**Task and Dataset.** We consider a simple algorithmic task that is referred to as *histogram*: given a sequence of tokens, the goal is to return a sequence of the same length where each entry represents the number of times the corresponding input token appears in the entire sequence. For example, given  $\mathbf{x} = [A, B, D, D, B, B]$ , the output will be  $\mathbf{y} = [1, 3, 2, 2, 3, 3]$ . We define the count of a token  $t$  in the sequence  $\mathbf{x}$  at position  $\ell$  as  $\text{hist}_{\mathbf{x}}(\ell)$ . In our experiments,

we consider i.i.d. distributions of sequences of length  $L$  from an input alphabet of size  $T$ , where  $L \leq T$ . Our sampling strategy relies on first sampling a set of partitions, and then assigning a token to each partition (see App. D for details). This allows for a close to uniform distribution over the values of  $\mathbf{y}$ .

**Models and Training.** We measure the performance on the histogram task of the four different variants of the token mixing models described in Sec. 2, i.e. `lin` and `dot`, with or without the softmax (`+sftm`). The token embeddings are jointly learned with the model parameters. We consider the dimension of the embedded tokens  $d$ , and the hidden layer size  $p$  of the feature transformation. We also consider the model `bos(+sftm)` where every input sequence is prefixed with the BOS token prior to entering a dot-product mixing layer (with softmax).

Models are trained with Adam with a learning rate of  $10^{-3}$  with cross-entropy loss for 500 epochs with a batch size of 32. We consider the online learning setting where for each new epoch we generate a dataset of fresh 10,000 samples. The accuracy is computed from 3,000 independent samples.

## 4. Learning Regimes in Counting

In order to understand the contributions of the different architectural components, we analyze the performance of the above-stated models with varying mixing mechanisms in different learning regimes characterized by the embedding dimension  $d$  and the number of hidden neurons  $p$  of the feed-forward module.

Fig. 1 shows the accuracy attained by learned models for sequences of length  $L = 10$  with  $T = 32$  different input tokens. We observe that the models exhibit both high and low accuracy across various parameter regimes, with a strong dependence on the architecture. Fig. 2 further clarifies that the parameter efficiency under different architectures varies substantially. To investigate the underlying mechanisms we devise theoretical constructions and mechanistic interpretations of the learned solutions. We delineate two regimes in each of the parameters: for the embedding dimension  $d$  we distinguish the regime of non-orthogonal embeddings ( $d < T$ ) and of possibly orthogonal embeddings ( $d \geq T$ ). For hidden layer size  $p$  we distinguish the regime where models can sense only a constant number of directions/features ( $p = 1$ ) or one scaling as the alphabet size ( $p = T$ ). In App. E.8 we show that a similar phenomenology persists for 2-layer transformers.

### 4.1. $d \geq T$ : Orthogonal token embeddings are separable

When the model dimension  $d$  is equal or greater than the alphabet size  $T$ , tokens can be represented by embeddings that are mutually orthogonal. Assuming for  $t \in \mathcal{T}$  there are

such mutually orthogonal embeddings  $e_t \in \mathbb{R}^d$  with a norm of 1, the overlap  $\langle e_s, e_t \rangle = 0$  for distinct tokens  $t \neq s$  and it is 1 when  $t = s$ . In such a scenario, a linear combination of token embeddings preserves magnitude (count information) about elements from the alphabet. A weighted sum of tokens, denoted as  $e' = \sum_{t \in \mathcal{T}} \alpha_t e_t$ , can be broken down into the original tokens using projections on the original token embeddings, where  $\alpha_t = \langle e_t, e' \rangle / \|e_t\|_2^2$ .

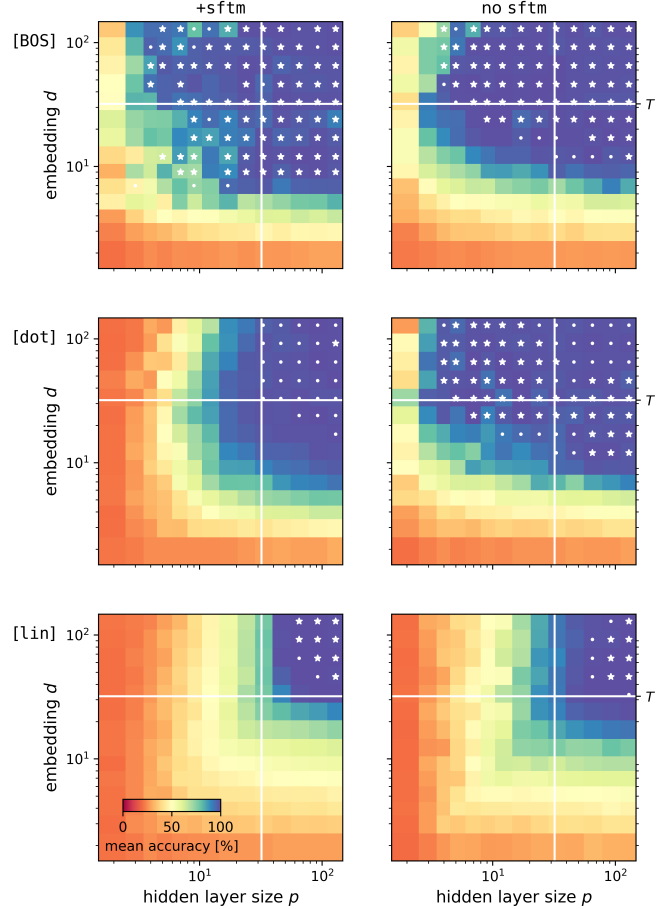


Figure 1. Accuracy on the histogram task for different 1-layer architectures. Mean accuracy for varying embedding size  $d$ , hidden layer size  $p$ , for fixed  $T = 32$  and  $L = 10$  for the different token mixing mechanisms `dot`, `bos` and `lin`. (Left) Models with softmax; (Right) Models without softmax. Average over 5 runs for every  $d, p \in \{1, 2, 3, 4, 6, 8, 12, 16, 23, 32, 45, 64, 91, 128\}$ . Vertical and horizontal white lines indicate  $p = T$  and  $d = T$  respectively. White stars (dots) mark a 100% (> 99%) accuracy configuration was found for at least one of the five runs.

In the following, we use this property to *theoretically* construct the weights for all models that solves the task when  $d \geq T$ . Remarkably, the constructions require different numbers of hidden neurons  $p$  depending on the mixing

mechanism. This demonstrates the interplay of the mixing layer and the feature transform: for some mixing mechanisms, the latter needs to implement *inventory-based counting* (IC) (requiring  $p \geq T$ ), and for others, *relation-based counting* (RC) (where  $p \geq 1$  is sufficient).

#### 4.1.1. RELATION-BASED COUNTING: LEVERAGING DOT-PRODUCT MIXING

When an extra beginning-of-sequence token  $t_{\text{BOS}}$  is available in `bos`, it can be used to extract information about a token’s count  $\text{hist}_{\mathbf{x}}(\ell)$  in the attention layer of the network through its attention score (Kazemnejad et al., 2023). In the literature, the beginning (or end) of sequence tokens have been linked to model-internal computations, such as counting. In (Weiss et al., 2021), it is shown that the RASP language can solve the histogram task with one layer and one attention head. We confirm empirically that `bos` and `bos+sftm` reach (close to) 100% accuracy whenever  $d > T$ , and we verify that a relation-based counting algorithm can be theoretically implemented in these two architectures by construction.

**Proposition 4.1** (RC with BOS token). *For `bos` and `bos+sftm` and a given  $L \geq 2$ , there each exists a configuration of weights that solves the histogram task at 100% accuracy, given that  $d \geq T > 2$  and  $p = 1$ .*

We prove this by construction in App. B.2.3-B.2.2 and we provide the intuition of the proof in the following. For `bos` we set the  $t_{\text{BOS}}$  embedding to  $e_{\text{BOS}} = \sum_{t \in \mathcal{T}} e_t$  and take the mutually orthogonal token embeddings  $e_t$  to have norm 1. Assuming that  $t_{\text{BOS}}$  is at the first position of the sequence of now length  $L + 1$ , a simple dot-product operation in the attention mechanism (with  $Q, K = d^{\frac{1}{4}} \mathbb{I}_d$ ) will lead to an attention matrix with entries:

$$a_{\ell m} = \begin{cases} T & \text{if } \ell = m = 1 \\ 1 & \text{if } (\ell > 1, m = 1) \text{ or } (\ell, m > 1, x_\ell = x_m) \\ 0 & \text{if } \ell, m > 1, x_\ell \neq x_m \end{cases}.$$

Projecting the mixed token  $\bar{x}'_\ell$  onto the  $t_{\text{BOS}}$  we obtain  $\langle \bar{x}'_\ell, e_{\text{BOS}} \rangle = T + \text{hist}_{\mathbf{x}}(\ell) + 1$ , i.e.  $e_{\text{BOS}}$  is the single relevant direction for the prediction. Its magnitude relates linearly to  $\text{hist}_{\mathbf{x}}(\ell)$ . A single hidden neuron  $p = 1$  suffices and the output layer can transfer the count into a categorical representation. For `bos+sftm` one needs to further account for the non-linearity of the softmax as described in App. B.2.2.

In the learned models, some instances in the given regime indeed achieve 100% accuracy. While their weights do not correspond exactly to the relation-based counting algorithm described previously, they exhibit similar properties. In Fig. 3, we show for `bos+sftm`, that  $t_{\text{BOS}}$  indeed plays a special role in the *learned* model: in the attention matrix

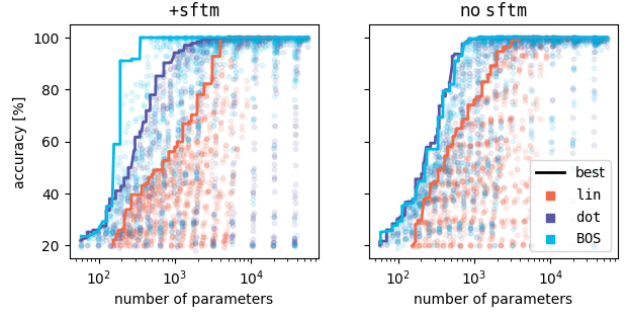


Figure 2. Test accuracy vs. total number of learned parameters. The data is the same as generated for Fig. 1, every data point is the a single experiment and we show the convex hull in solid lines.

its activation can be interpreted as a proxy for the number of occurrences of  $x_\ell$ , as it has different values for tokens that occur a different amount of times. Other entries of the attention matrix are comparatively low when the compared tokens are the same and high when they are different. The comparison operation naturally provided by the dot-product allows the model to extract the count of the same tokens, for each token in the sequence. We also show in Fig. 3 how the presence of the  $t_{\text{BOS}}$  determines the final prediction through the application of  $f$ . Surprisingly, the `dot` model (without the softmax) reaches a an empirical performance comparable to `bos` in the regime  $d \geq T$  and  $p = 1$ , even though it does not have an extra token available.

**Proposition 4.2** (RC with tagged embeddings). *For `dot` and a given  $L, T > 2$ , there exists a configuration of weights that solves the histogram task at 100% accuracy, given that  $d \geq T > 2$  and  $p = 1$ .*

We prove this in App. B.2.1. Intuitively, the construction uses a single common direction  $e_{\text{cnt}}$  that is added to the otherwise mutually orthogonal token embeddings. A dot-product mixing then leads to  $a_{\ell m} = a_{\neq} > 0$  when  $x_\ell$  is different from  $x_m$ , and  $a_{\ell m} = a_{=} > 0$  when tokens are the same. Then, the number of counts can be easily extracted from the dot-product  $\langle e_{\text{cnt}}, \bar{x}'_\ell \rangle$  of the counting token with the mixed token  $\bar{x}'_\ell$ , i.e.  $\langle e_{\text{cnt}}, \bar{x}'_\ell \rangle \propto 1 + \text{hist}_{\mathbf{x}}(\ell)a_{=} + (L - \text{hist}_{\mathbf{x}}(\ell))a_{\neq}$ . We can, therefore, obtain a perfect accuracy implementation in the regime where  $d \geq T$  with only a single hidden neuron. This is in line with the observed empirical performance by `dot` even without access to a BOS token.

**Dot-product attention with softmax fails to implement relation-based counting.** Since the dot-product mechanism can naturally be used in relation-based counting, one might expect the `dot+sftm` model to implement the same mechanism. However, and maybe surprisingly so, we em-



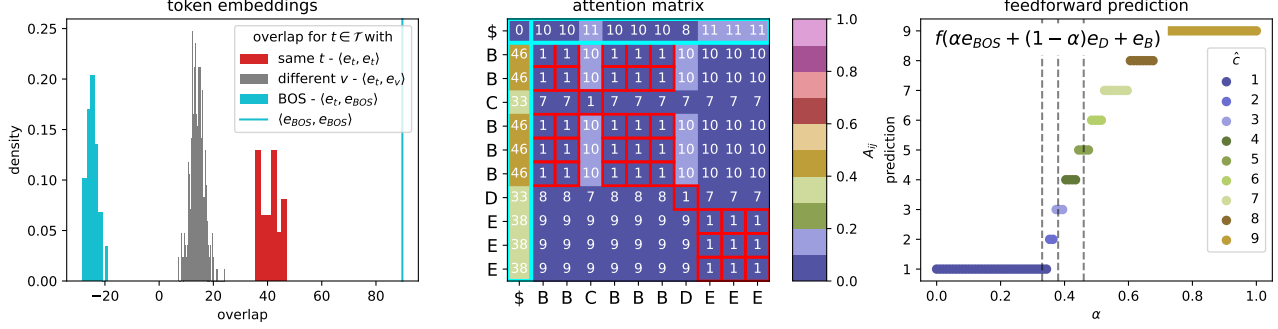


Figure 3. Relation-based counting with *bos+sftm* ( $T = 32, L = 10, p = 2, d = 45$ ). This model achieves 99.9% accuracy. It was selected as the best model from all our experiments with  $p = 2$ . (Left) The tokens overlap (cosine similarity) with the **same** tokens, different tokens and the **BOS** all concentrate around different values. (Middle) This is reflected in the attention matrix after the application of the row-wise softmax. The  $t_{BOS}$  ('\$') in the first column  $a_{\ell,0}$  becomes a proxy for the count of  $x_\ell$ . (Right) To demonstrate that the feedforward network is only sensitive to this direction, we show its count predictions for a mix of tokens and fix  $\bar{x}' = \alpha e_{BOS} + (1 - \alpha)e_D + e_B$ . The contribution  $\alpha$  of the BOS token to the intermediate  $\bar{x}'$  is varied and  $D, B$  are two specific elements of the alphabet  $\mathcal{T}$ , represented by their embeddings  $e_D$  and  $e_B$ . The  $y$ -axis shows the predicted class of the feedforward layer  $f(\bar{x}')$  for a given  $\alpha$ . We measure  $a_{\ell,0}$ , the actual contribution of the BOS token's for the sequence in the middle plot for the letters that occur 1, 3 and 5 times. The same experiment is repeated for different elements of the alphabet in App. E.4, showing independence of the count prediction on the other tokens present in the sequence.

pirically observe a marked difference between *dot* and *dot+sftm* in Fig. 1. *dot+sftm* only starts performing close to 100% accuracy when both the model dimension  $d$  and the number of hidden neurons  $p$  are larger than the number of tokens  $T$ . To understand why it fails to learn for  $p = 1$ , we show the attention matrix of *dot+sftm* in Fig. 4. Notably, it is based on the semantics, as  $(A_{\text{dot+sftm}})_{\ell m}$  is higher when  $x_\ell = x_m$  than otherwise. However, the normalization effect of the softmax activation prevents the development of a meaningful counter subspace that is needed in the relation-based algorithm. As a result of normalization, the attention scores are  $\sum_m a_{\ell m} = 1$ , so any direction present in all tokens (and by the symmetry of the task, it would need to be present in all tokens) would be uninformative after the token mixing – its weight would be one regardless of the input sequence and would therefore not carry information about the count. Before, the model *bos+sftm* circumvented this problem by adding the extra token with a special functionality that does not need to be counted. Because this is not possible for *dot+sftm*, the architecture fails to perform well for  $p = 1$  – it now needs to measure more than one direction in the feed-forward module.

In the following, we show that a solution of the histogram task can still be achieved through an inventory-based counting algorithm with  $p \geq T$ . We detail this in the following section, for the example of *lin*. The statement for *dot+sftm* is given in App. B.3.

#### 4.1.2. INVENTORY-BASED COUNTING: MEMORIZATION IN THE FEED-FORWARD LAYER

When the feed-forward hidden layer has one neuron for each distinct token available in the alphabet, it can detect as many directions. This allows the feed-forward layer to extract the information of any token direction separately and thereby implement a custom comparison operation that works for all of the tokens in the alphabet. While this is less parameter efficient and requires memorizing the complete alphabet, it enables the model to solve the task.

**Proposition 4.3** (IC with memorization in the feed-forward layer). *For  $\text{lin}$  and  $\text{lin+sftm}$  and a given  $L, T > 2$  there exists a configuration of weights which solves the histogram task for  $p \geq T$  and  $d \geq T$ .*

We describe examples of such constructions in App. B.3.1 and B.3.2. Again, several solutions exist due to symmetries, and in the following we give an intuition for one of them.

In the linear mixing layer  $A_{\text{lin}}$  we set a constant value  $a = 1/L$  so that the result of the mixing is simply a position-independent linear combination of the input. The count  $\text{hist}_x(\ell)$  can be extracted after the residual connection where we add  $\bar{x}'_\ell = \bar{x}_\ell + e_{x_\ell}$ . By setting the columns of the matrix  $(W_1)_t = e_t$  we can extract the count information up to the factor  $1/a$

$$\begin{aligned} \text{hist}_x(\ell) &= \frac{1}{a} \sum_{t \in \mathcal{T}} \text{ReLU}(\langle \bar{x}'_\ell, (W_1)_t \rangle - 1) \\ &= \frac{1}{a} \sum_{t \in \mathcal{T}} \text{ReLU}(\langle \bar{x}'_\ell, e_t \rangle - 1) = \frac{1}{a} \langle \bar{x}'_\ell, e_{x_\ell} \rangle \end{aligned}$$

Note that, due to the  $-1$  bias term, only the hidden neu-

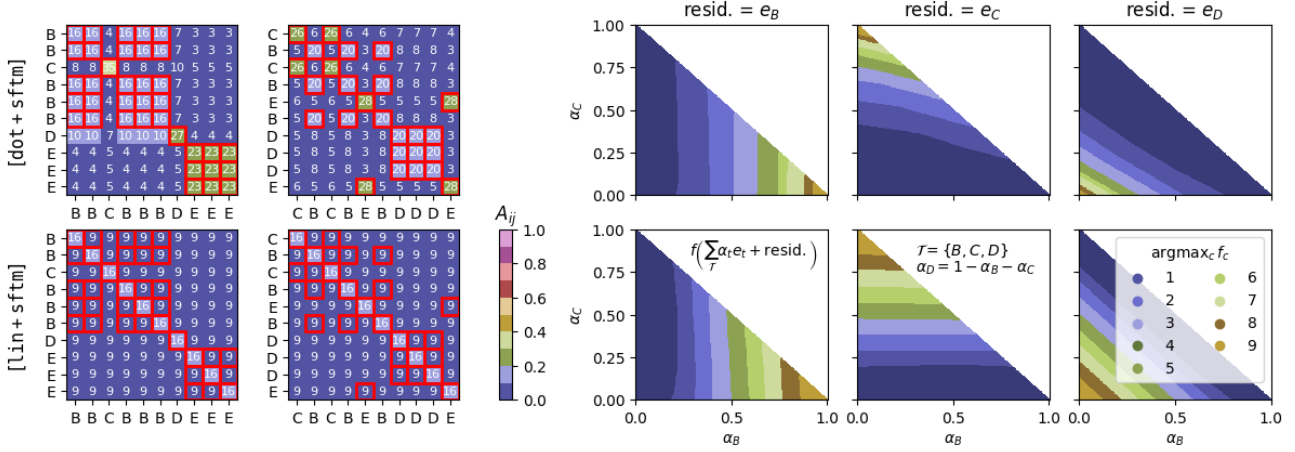


Figure 4. Inventory-based counting with *dot+sftm* ( $p = 32, d = 32$ ) and *lin+sftm* ( $p = 64, d = 64$ ). We have  $T = 32, L = 10$ . The models achieves 99.47% and 100% accuracy respectively. (Left columns) The attention matrix for two different sequences. It differentiates between same (red squares) and different tokens for *dot+sftm* but is invariant to the semantics for *lin*. For *dot+sftm*, any counting direction that could emerge in token space is not informative due to the softmax normalization, so  $p \geq T$  is required (see Fig. 1). (Right columns) We test the feed-forward layer  $f$  in isolation, by feeding it with an artificial mix of learned embeddings for the three tokens  $B, C$  and  $D$ . We plot the class predicted by  $f$  for inputs constructed as  $\tilde{x}' = \alpha_B e_B + \alpha_C e_C + \alpha_D e_D + R$ , where  $1 = \alpha_B + \alpha_C + \alpha_D$  and we change the residual  $R \in \{e_B, e_C, e_D\}$  from left to right, for each of the models *dot+sftm* and *lin+sftm*. The prediction strongly depends on the coefficient  $\alpha_t$  associated with the token  $t$  present in the residual connection and only weakly on the others. The non-linear scaling of the decision boundaries is due to the softmax activation function.

ron for token  $(W_1)_t = e_t = x_\ell$  that occurs in the residual connection has a non-zero activation. The output layer  $W_2$  can then be designed to activate the correct output vector corresponding to the count  $hist_x(\ell)$  (see App. B.4). Since  $a \in [0, 1]$  and  $\sum_{m=1}^L a_{\ell m} = 1$  the same procedure can be implemented by a matrix which is passed through the softmax operator for *lin+sftm*. In practice, in this construction the feed-forward module is correlated with the complete alphabet, acting as an inventory, or look-up table.

In Fig. 4, we inspect the attention matrix  $A_{lin}$  and the feature transformation  $f$  which is *learned* for *lin+sftm* in the regime where  $p \sim T \sim d$ . The mixing has an off-diagonal of  $\sim 0.09$  and a diagonal of  $\sim 0.16$ . Feeding the feature transformation  $f$  with a weighted combination of 3 tokens,  $B, C, D$ , we observe that the final prediction of the network depends mainly on the coefficient  $\alpha_t$  corresponding to the token embedding fed through the residual connection. This behavior is reflected in Fig. 4 (right) and suggests that the feature transformation must have encoded the information of the token embedding in its weights, hence requiring at least  $p = T$  hidden neurons.

**Superpositioned and selective implementations.** Some of the models capabilities include one another. For example, the models that can implement relation-based counting for  $p = 1$  can also implement the solutions for inventory-based counting for  $p \geq T$ . It is unclear, whether the memory-intensive solution is preferred when the memory is available,

or if the efficient solution is learned nonetheless. Curiously, in Fig. 1, we observe that the model *dot* (which is capable of RC) witnesses a very slight decrease in maximal learned performance from 100% accuracy to 99% despite its capacity being *increased* to  $p = T$  when inventory-based counting can in principle be implemented. In App. E.5 we investigate the singular value decomposition of  $W_1$ , for learned models with  $\geq 99\%$  accuracy and  $p, d \geq T$ . We find that the largest  $T = 32$  singular values are larger than the surplus singular values when  $p > T$  for models that can implement only IC. This behavior is less pronounced for models that can implement RC, where the largest singular value is often relatively much larger than the following  $T = 32$ , but still show a small dip after the  $T = 32$  singular values. Understanding which algorithm is implemented in this regime, or if it is a superposition of the two, thus requires further investigation.

#### 4.2. $d < T$ : Non-orthogonal embeddings and the discrete nature of counting

The scenario where  $d < T$  fundamentally differs from the one explored in Section 4.1 because the embeddings for different tokens can no longer be mutually orthogonal. Some token pairs then have a non-zero overlap due to their linear dependence, causing the mixing of tokens to entangle count information across different directions in the embedding space.

This phenomenon is illustrated for *dot* in Fig. 5, where learned models with smaller  $d$  tend to overcount items in

the input, and observe a less spread distribution of overlaps. Nevertheless in Fig. 1 we observe a number of results that empirically show almost perfect accuracy solutions with  $d < T$  both for models with RC or IC. Indeed, the discrete nature of the histogram task, i.e. the fact that every token can only be mapped to  $L$  distinct counts, makes the prediction inherently more robust to the effect of noise stemming from entangled embeddings. This concept is illustrated in Fig. 8 in App. B.4 for the `dot+sftm` model. As long as the value of the logits in the final output layer falls within the margin between two counts the model still solves the task with perfect accuracy. The relative size of this margin decreases when  $L$  is increased, making the task harder when more classes need to be distinguished.

In the following, we link concepts on optimally placing decision boundaries for noise robustness to a characterization of this entanglement noise, measured by the mutual coherence of the token embedding set (i.e., the maximum absolute overlap between pairs of distinct embeddings). The mutual coherence of a set of  $T$  vectors of dimension  $d$  is lower bounded by the Welch bound (Welch, 1974). This gives a means to understand the size of  $d$  a given task with  $T, L$  requires at least.

**Proposition 4.4** (Robustness via bounded mutual coherence). *Given  $L \geq 5, T \geq 2$  and assuming that the Welch bound is attained for a given  $T, d$ , there exists a construction that solves the histogram task with*

$$(lin, lin+sftm; p = T): \left\lceil \frac{T(2L-3)^2}{T-1+(2L-3)^2} \right\rceil \leq d,$$

$$(dot, bos; p = 1): \left\lceil \frac{T(2L-3)^2}{T-1+(2L-3)^2} \right\rceil + 1 \leq d,$$

$$(dot, bos; p = T): \left\lceil \frac{T(L-1)}{T-1+(L-1)} \right\rceil \leq d.$$

We provide additional background and the proofs in App. C.2. The idea is to use constructions analogous to the RC and IC with orthogonal embeddings, while keeping track on how the errors of non-zero overlaps between pairs of different embeddings propagate through the model. For a given  $L$  and  $T$  this provides an upper bound on the maximal mutual coherence that is tolerated for a perfect solution. This can be connected to the dimensionality  $d$  via the Welch bound. Evaluating the bounds for the setting in Fig. 1, we obtain, in the order of the above list,  $d \geq 29, 30, 7$ . Generally it is hard to generate matrices that attain the Welch bound and manually we did not succeed to find them for  $d = 29, 30$ . However we can indeed create an explicit construction a for `dot` and  $p = T$  which attains  $d = 12$ , as provided in the supplementary code and in correspondence with Fig. 1. While this bound does not reach the  $d$  as indicated by the Welch bound, the mutual coherence of the embedding matrix we use is close to the maximally allowed value of  $\mathcal{M} = 0.299 < 1/3$ .

The previous results apply specifically to models without the softmax operator in the token mixing step – models with this non-linearity can be more robust and attain even smaller  $d$ , as clearly visible in Fig. 1. The idea is that a softmax function with a high enough inverse temperature can non-linearly scale down the attention scores for different token pairs relative to those of the same tokens. Thereby, the noise introduced in the dot-product layer through pairs of different embeddings becomes *arbitrarily* close to zero after applying the softmax.

**Proposition 4.5** (Robustness via softmax error-reduction). *Given  $T, L > 2$ , there exist weight configurations that solve the histogram task for the parameter combinations (`bos+sftm`;  $p = 1$ ) and (`dot+sftm`;  $p = T$ ) with  $\lceil \log_2(T+1) \rceil + 2 \leq d$ .*

Put simply, this construction requires that there are token embeddings for  $t, s = 1, \dots, T$  and  $s \neq t$  with  $\epsilon > 0$  such that  $\langle e_t, e_t \rangle = 1$  and  $\langle e_t, e_s \rangle < 1 - \epsilon$ . This is fulfilled when every token is the binary encoding of its value, modulo minor modifications due to the RC mechanism for `bos+sftm`. Setting the softmax temperature high enough as a function of  $L$  allows for the contributions from non-equal tokens to be decreased relative to the ones of same tokens. Evaluating this function for Fig. 1, we obtain  $d = 7$ , which closely corresponds to the most parameter efficient solutions of the histogram task that we observe. As  $L$  grows, we require stronger concentration from the softmax by adjusting its temperature. Since real-world networks execute finite computations, computational instabilities or collapses might occur. It is therefore not clear that this correspondence will hold for all values of  $L$ .

In App. C.3.1 we show that this bound can be even further improved for `bos+sftm` to a *constant*  $d = 4$ , but at the cost of increasing the temperature further as a function of  $T$ , in addition to  $L$ , requiring yet more accurate computations. While in practice we are able to hand construct and code a solution for the  $L, T$  under consideration (available in the supplementary code), these models are not learned. However the very large numerical precision required might lead to learning difficult learning dynamics and be the reason why we do not observe any learned solutions of the histogram task in this regime.

### 4.3. Extensive $T, L$ ablations and two layer architectures

Appendix E provides a number of ablations for  $T = 15, 64$  and  $L = 5, 15, 30$ , for different architectures and hidden layer size  $p$  and embedding size  $d$ . Our predictions from the theory are largely consistent with the performance for different alphabet sizes  $T$  and small  $L$ , similarly to  $T = 32$  and  $L = 10$ . However, for very large  $L = 30$  it seems that the keeping the training budget constant lead to a weaker performance in the feasible regimes.

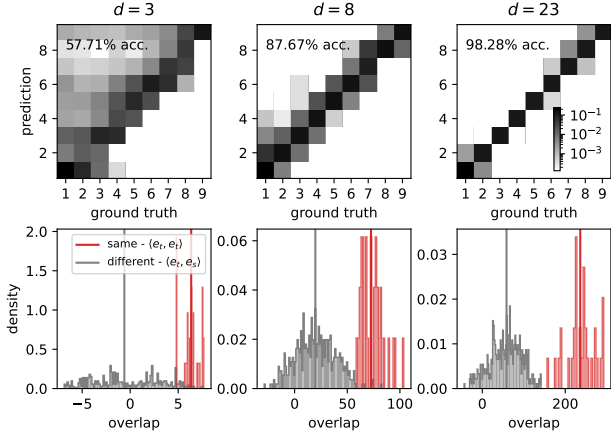


Figure 5. Introspecting the regime with Entangled Embeddings with  $\text{dot}$  ( $T = 32, L = 10, p = 32$ ). We show trained instances of  $\text{dot}$  with varying the model dimension  $d$ . (Top) The confusion matrix of ground truth and predicted counts. (Bottom) The overlap distribution between same and different token embeddings.

In the same section, we show that for learned 2-layer attention blocks the general phenomenology for the single layer models from Fig. 1 is very coherent. Only in the critical regimes along the transition from not possible to possible transitions, there are some larger differences in average performance. This indicates that 2-layer networks do not necessarily operate differently in what they learned.

## 5. Related Work

**Counting and Interpretability.** The emergence of algorithmic capabilities in transformers (Olsson et al., 2022; Power et al., 2022) has led to numerous investigations aimed at reverse-engineering trained models into human-understandable mechanisms (Zhong et al., 2023; Nanda et al., 2023; Quirke & Barez, 2024), including a variety of counting scenarios (Gould et al., 2023; Chollet et al., 2020; Ouellette et al., 2023; Cui et al., 2024; Golkar et al., 2024). In our work, we consider the histogram task introduced via the RASP(-L) programming language (Weiss et al., 2021; Abbe et al., 2023). While RASP predicts that the histogram single layer transformers with one head require a BOS token (Nye et al., 2021), we find that this is not necessary by providing constructions. While many works in interpretability focus on causal interventions (Vig et al., 2020; Meng et al., 2022) to understand the computational mechanisms of models or assign relevance scores to their components (nostalgabraist, 2020; Elhage et al., 2021), we focus on the hyperparameter scaling of several distinct models in relation to their performance and explicit constructions of different algorithms, similar to e.g. (Zhong et al., 2023; Quirke & Barez, 2024; Nichani et al., 2025). We give precise theoretical conditions on the model configurations that lead to

perfect explicit constructions and link them to introspection on learned models. In concurrent but independent work to ours, Yehudai et al. (2024) explore a very similar version of the histogram task, and find, as we do, that the numerical precision limits the counting size. While they do not examine different architectures, they find rigorous upper bounds.

**Memorization and Feed-forward Layers.** The role of feed-forward layers as memorization modules is investigated in the context of factual recall for language models (Geva et al., 2021; Meng et al., 2022; Chughtai et al., 2024). Henighan et al. (2023) study a double decent phenomenon where the purpose of the feed-forward layer transitions from storing data points to discovering generalizing features as a function of increasing training data diversity (Raventos et al., 2023). In the histogram task, we observe a similar phenomenon as a function of the architecture: the feed-forward layer acts either as a look-up table or a feature detector for a counting subspace.

**Aligning Algorithm and Architecture.** While theoretical work has defined the computational capacity of several (autoregressive) neural networks (Weiss et al., 2021; Yun et al., 2019; Delétang et al., 2023; Liu et al., 2023), hallucinations and failure modes on seemingly trivial tasks in real-world transformers are the rule rather than an exception. Dziri et al. (2023) postulate that this may be due to a misalignment between the computational graph of a model and the task itself. In this work, we show that subtle differences in components such as the mixing type and layer width play a crucial role in terms of algorithmic alignment. Previous work found evidence of superimposed computational graphs in a single model (Elhage et al., 2022); our models shed light on how such entanglement functions in detail for non-orthogonal embeddings.

## 6. Discussion & Conclusion

**Limitations.** Similar to other works in mechanistic interpretability (Zhong et al., 2023), we focus on 1-layer transformers as a simplified model for modern transformers. Our models are not autoregressive and do not account for the impact of causal masks or positional encodings. While more complex models could lead to more intricate interdependencies between the components, potentially limiting the applicability of our findings to such architectures, it seems plausible that similar vector arithmetic could emerge in subspaces of large transformers (Gould et al., 2023; Engels et al., 2025). Given its specificity, it is unclear if and how similar memory-architecture phenomena would emerge for different simple tasks (e.g. sorting or lookup). However, the learning regimes seem to transfer to two layer models, and even large auto-regressive language models fail to reliably perform well on this task, hence motivating the study of



algorithmic emergence in smaller models.

**Summary.** We study how different components of simple transformer models contribute to the emergence of different solutions to the histogram task. The feasibility of solving the histogram task including the comparison and aggregation operations depends on the mixing mechanism, its interaction with the feed-forward transformation, and the softmax function in the attention. Relation-based counting, uses a dot product mixing mechanism and low-capacity feed-forward transformation, and inventory-based counting, which memorizes token embeddings in feed-forward weights, requires more parameters. We characterize feasibility regimes in the phase space of embedding dimension  $d$  and feed-forward dimension  $p$ , confirming that models converge to these mechanisms. Some regimes allow both strategies, and experiments show features of superimposed mechanisms. When  $d < T$ , tokens cannot form an orthogonal basis for linear projection. Models exhibit varying robustness to noise from non-orthogonality. The softmax activation minimizes token similarity in attention, reducing non-orthogonality effects—relevant in real-world models where  $T$  often exceeds model dimensions.

**Future Directions.** Examples of hallucinations and failures in LLMs are as numerous as their successes. Though limited to one task, our analysis highlights how subtle architectural changes can drastically impact predictive power. We expect that similar mechanistic investigations at or close to the regimes where models start failing will be extremely useful to understand how and why models fail in sometimes puzzling manners.

## Acknowledgements

Luca Biggio acknowledges support from the AI for Science postdoctoral fellowship at EPFL.

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## References

Abbe, E., Bengio, S., Lotfi, A., and Rizk, K. Generalization on the unseen, logic reasoning and degree curriculum. In *ICML*, 2023. URL <https://arxiv.org/abs/2301.13105>.

Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G.,

Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. *CoRR*, abs/2005.14165, 2020. URL <https://arxiv.org/abs/2005.14165>.

Cammarata, N., Carter, S., Goh, G., Olah, C., Petrov, M., Schubert, L., Voss, C., Egan, B., and Lim, S. K. Thread: circuits. *Distill*, 5(3):e24, 2020.

Chollet, F., Tong, K., Reade, W., and Elliott, J. Abstraction and reasoning challenge, 2020. URL <https://kaggle.com/competitions/abstraction-and-reasoning-challenge>.

Chughtai, B., Cooney, A., and Nanda, N. Summing up the facts: Additive mechanisms behind factual recall in LLMs, 2024. URL <https://openreview.net/forum?id=P2gnDEHGu3>.

Cui, H., Behrens, F., Krzakala, F., and Zdeborova, L. A phase transition between positional and semantic learning in a solvable model of dot-product attention. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=BFWDIPPLgZ>.

Delétang, G., Ruoss, A., Grau-Moya, J., Genewein, T., Wenliang, L. K., Catt, E., Cundy, C., Hutter, M., Legg, S., Veness, J., and Ortega, P. A. Neural networks and the chomsky hierarchy. In *11th International Conference on Learning Representations*, 2023.

Donoho, D. L. and Elad, M. Optimally sparse representation in general (nonorthogonal) dictionaries via sup minimization. *Proceedings of the National Academy of Sciences*, 100(5):2197–2202, 2003. doi: 10.1073/pnas.0437847100. URL <https://www.pnas.org/doi/abs/10.1073/pnas.0437847100>.

Dordevic, D., Bozic, V., Thommes, J., Coppola, D., and Pal Singh, S. Rethinking attention: Exploring shallow feed-forward neural networks as an alternative to attention layers in transformers (student abstract). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(21):23477–23479, Mar. 2024. doi: 10.1609/aaai.v38i21.30436. URL <https://ojs.aaai.org/index.php/AAAI/article/view/30436>.

Dziri, N., Lu, X., Sclar, M., Li, X. L., Jiang, L., Lin, B. Y., Welleck, S., West, P., Bhagavatula, C., Le Bras, R., Hwang, J., Sanyal, S., Ren, X., Ettinger, A., Harchaoui, Z., and Choi, Y. Faith and fate: Limits of transformers on compositionality. In Oh, A., Naumann, T., Globerson, A.,

- Saenko, K., Hardt, M., and Levine, S. (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 70293–70332. Curran Associates, Inc., 2023.
- Elhage, N., Nanda, N., Olsson, C., Henighan, T., Joseph, N., Mann, B., Askell, A., Bai, Y., Chen, A., Conerly, T., DasSarma, N., Drain, D., Ganguli, D., Hatfield-Dodds, Z., Hernandez, D., Jones, A., Kernion, J., Lovitt, L., Ndousse, K., Amodei, D., Brown, T., Clark, J., Kaplan, J., McCandlish, S., and Olah, C. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- Elhage, N., Hume, T., Olsson, C., Schiefer, N., Henighan, T., Kravec, S., Hatfield-Dodds, Z., Lasenby, R., Drain, D., Chen, C., Grosse, R., McCandlish, S., Kaplan, J., Amodei, D., Wattenberg, M., and Olah, C. Toy models of superposition, 2022.
- Engels, J., Michaud, E. J., Liao, I., Gurnee, W., and Tegmark, M. Not all language model features are one-dimensionally linear. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=d63a4AM4hb>.
- Fickus, M. and Mixon, D. G. Tables of the existence of equiangular tight frames, 2016. URL <https://arxiv.org/abs/1504.00253>.
- Geva, M., Schuster, R., Berant, J., and Levy, O. Transformer feed-forward layers are key-value memories. In Moens, M.-F., Huang, X., Specia, L., and Yih, S. W.-t. (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 5484–5495, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.446. URL <https://aclanthology.org/2021.emnlp-main.446>.
- Girdhar, R., Carreira, J., Doersch, C., and Zisserman, A. Video action transformer network. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 244–253, 2019.
- Golkar, S., Bietti, A., Pettee, M., Eickenberg, M., Cranmer, M., Hirashima, K., Krawezik, G., Lourie, N., McCabe, M., Morel, R., Ohana, R., Parker, L. H., Blancard, B. R.-S., Cho, K., and Ho, S. Contextual counting: A mechanistic study of transformers on a quantitative task, 2024. URL <https://arxiv.org/abs/2406.02585>.
- Gould, R., Ong, E., Ogden, G., and Conmy, A. Successor heads: Recurring, interpretable attention heads in the wild, 2023.
- Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., Yang, A., Fan, A., and et al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Gu, A. and Dao, T. Mamba: Linear-time sequence modeling with selective state spaces. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=tEYskw1VY2>.
- Henighan, T., Carter, S., Hume, T., Elhage, N., Lasenby, R., Fort, S., Schiefer, N., and Olah, C. Superposition, memorization, and double descent. *Transformer Circuits Thread*, 2023.
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.
- Jiang, Q., Li, S., Bai, H., de Lamare, R. C., and He, X. Gradient-based algorithm for designing sensing matrix considering real mutual coherence for compressed sensing systems. *IET Signal Processing*, 11(4):356–363, 2017.
- Jyothi, R. and Babu, P. Telet: A monotonic algorithm to design large dimensional equiangular tight frames for applications in compressed sensing. *Signal Processing*, 195:108503, 2022.
- Kazemnejad, A., Padhi, I., Natesan, K., Das, P., and Reddy, S. The impact of positional encoding on length generalization in transformers. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=Dr1l2gcjzl>.
- Liu, B., Ash, J. T., Goel, S., Krishnamurthy, A., and Zhang, C. Transformers learn shortcuts to automata. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=De4FYqjFueZ>.
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. Swin transformer: Hierarchical vision transformer using shifted windows. *CoRR*, abs/2103.14030, 2021. URL <https://arxiv.org/abs/2103.14030>.
- Meng, K., Bau, D., Andonian, A. J., and Belinkov, Y. Locating and editing factual associations in GPT. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*,

2022. URL <https://openreview.net/forum?id=h6WAS6eE4>.
- Michaud, E. J., Liu, Z., Girit, U., and Tegmark, M. The quantization model of neural scaling. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=3tbTw2ga8K>.
- Nanda, N., Chan, L., Lieberum, T., Smith, J., and Steinhardt, J. Progress measures for grokking via mechanistic interpretability. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=9XFSbDPmdW>.
- Nichani, E., Lee, J. D., and Bietti, A. Understanding factual recall in transformers via associative memories. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=hwSmPOAmhk>.
- nostalgebraist. interpreting GPT: the logit lens — LessWrong, January 2020. URL <https://www.lesswrong.com/posts/AcKRB8wDpdAN6v6ru/interpreting-gpt-the-logit-lens>.
- Nye, M. I., Andreassen, A. J., Gur-Ari, G., Michalewski, H., Austin, J., Bieber, D., Dohan, D., Lewkowycz, A., Bosma, M., Luan, D., Sutton, C., and Odena, A. Show your work: Scratchpads for intermediate computation with language models. *CoRR*, abs/2112.00114, 2021. URL <https://arxiv.org/abs/2112.00114>.
- Olah, C., Cammarata, N., Schubert, L., Goh, G., Petrov, M., and Carter, S. Zoom in: An introduction to circuits. *Distill*, 2020. doi: 10.23915/distill.00024.001. <https://distill.pub/2020/circuits/zoom-in>.
- Olsson, C., Elhage, N., Nanda, N., Joseph, N., DasSarma, N., Henighan, T., Mann, B., Askell, A., Bai, Y., Chen, A., Conerly, T., Drain, D., Ganguli, D., Hatfield-Dodds, Z., Hernandez, D., Johnston, S., Jones, A., Kernion, J., Lovitt, L., Ndousse, K., Amodei, D., Brown, T., Clark, J., Kaplan, J., McCandlish, S., and Olah, C. In-context learning and induction heads, 2022.
- Ouellette, S., Pfister, R., and Jud, H. Counting and algorithmic generalization with transformers. *arXiv preprint arXiv:2310.08661*, 2023.
- Petzka, H., Trimmel, M., and Sminchisescu, C. Notes on the symmetries of 2-layer relu-networks. In *Proceedings of the northern lights deep learning workshop*, volume 1, pp. 6–6, 2020.
- Power, A., Burda, Y., Edwards, H., Babuschkin, I., and Misra, V. Grokking: Generalization beyond overfitting on small algorithmic datasets, 2022.
- Quirke, P. and Barez, F. Understanding addition in transformers. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=rIx1YXVWZb>.
- Raventos, A., Paul, M., Chen, F., and Ganguli, S. Pretraining task diversity and the emergence of non-bayesian in-context learning for regression. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=BtAz4a5xDg>.
- Strohmer, T. and Heath, R. W. Grassmannian frames with applications to coding and communication. *Applied and Computational Harmonic Analysis*, 14(3):257–275, 2003. ISSN 1063-5203. doi: [https://doi.org/10.1016/S1063-5203\(03\)00023-X](https://doi.org/10.1016/S1063-5203(03)00023-X). URL <https://www.sciencedirect.com/science/article/pii/S106352030300023X>.
- Tolstikhin, I., Houlsby, N., Kolesnikov, A., Beyer, L., Zhai, X., Unterthiner, T., Yung, J., Steiner, A. P., Keysers, D., Uszkoreit, J., Lucic, M., and Dosovitskiy, A. MLP-mixer: An all-MLP architecture for vision. In Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=EI2K0XKdnP>.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, pp. 6000–6010, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- Vig, J., Gehrmann, S., Belinkov, Y., Qian, S., Nevo, D., Singer, Y., and Shieber, S. M. Causal mediation analysis for interpreting neural NLP: the case of gender bias. *CoRR*, abs/2004.12265, 2020. URL <https://arxiv.org/abs/2004.12265>.
- Weiss, G., Goldberg, Y., and Yahav, E. Thinking like transformers. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 11080–11090. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/weiss21a.html>.
- Welch, L. R. Lower bounds on the maximum cross correlation of signals (corresp.). *IEEE Trans. Inf. Theory*, 20:397–399, 1974. URL <https://api.semanticscholar.org/CorpusID:20783885>.
- Yehudai, G., Kaplan, H., Ghandeharioun, A., Geva, M., and Globerson, A. When can transformers count to n?, 2024.

Yun, C., Bhojanapalli, S., Rawat, A. S., Reddi, S. J., and Kumar, S. Are transformers universal approximators of sequence-to-sequence functions? *CoRR*, abs/1912.10077, 2019. URL <http://arxiv.org/abs/1912.10077>.

Zhong, Z., Liu, Z., Tegmark, M., and Andreas, J. The clock and the pizza: Two stories in mechanistic explanation of neural networks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=S5wmbQc1We>.



## Appendices

- A Counting with large language models
- B Explicit Constructions for Orthogonal Embeddings  $d = T$ 
  - B.1 Overview
  - B.2 Relation-based counting
  - B.3 Inventory-based counting
  - B.4 Mapping a scalar to a categorical one-hot encoding
- C Explicit Constructions for Linearly Dependent Embeddings  $d < T$ 
  - C.1 Overview
  - C.2 Explicit construction for bounded mutual coherence
  - C.3 Explicit Construction with binary representations and softmax
- D Data Generation
- E Additional Experiments
  - E.1 Best Accuracy
  - E.2 Variability
  - E.3 Model with Random but Fixed Embeddings
  - E.4 BOS mixing token
  - E.5 Singular Value Decomposition of  $W_1$
  - E.6 Varying the number of tokens in the alphabet
  - E.7 Varying the length of the sequence
  - E.8 Models with two layers

## A. Counting with large language models

We prompt language models with 7 and 8 billion parameters to count the inputs of a list. We use the torch implementations of `mistral-7B-v0.1` (Jiang et al., 2023) and `Llama-3.1-8B` (Grattafiori et al., 2024) with `float16` respectively and prompt in the english language. We prompt the model to count items in different lists of ten items which stem from groups of 12 items (letters, months, animals). The groups of items are

- letters: `['A', 'B', 'C', 'D', 'E', 'F', 'G', 'H', 'I', 'J', 'K', 'L']`
- animals: `['dog', 'cat', 'donkey', 'horse', 'elephant', 'giraffe', 'lion', 'tiger', 'bear', 'wolf', 'zebra', 'monkey']`
- months: `['january', 'february', 'march', 'april', 'may', 'june', 'july', 'august', 'september', 'october', 'november', 'december']`

And the four prompt types are

1. (compare, position)  
`[list] By comparing all ten items in the previous list to one another, we find that the [i'th] item of the list appears [prompt],`
2. (compare, semantic)  
`[list] By comparing all ten items in the previous list to one another, we find that [item] appears [prompt],`
3. (summarize, position)  
`[list] By summarizing the occurrences of [category_name] in the previous list of ten items, we find that the [i'th] item of the list appears [prompt],`
4. (summarize, semantic)  
`[list] By summarizing the occurrences of [category_name] in the previous list of ten items, we find that [item] appears [prompt].`

When we generate data we sample the ten items from the list either uniformly at random with replacement, or as we do for the toy model experiments as specified in appendix D.

We then sample 100 lists for each model and generate all possible queries for different items and positions. We also change the category type and sampling strategy and report the confusion matrices of the ground truth over the predictions in Fig. 6.

We observe that the performance is not very good. Querying for a given item semantically rather than its position consistently performs better. For the positional queries the models mainly predict the count 3. Depending on the type of object that is counted, counting more than 5 is correct only rarely for both models. The type of item that is counted also impacts performance. Generally Llama-3.1 seems to perform better than Mistral-7B.

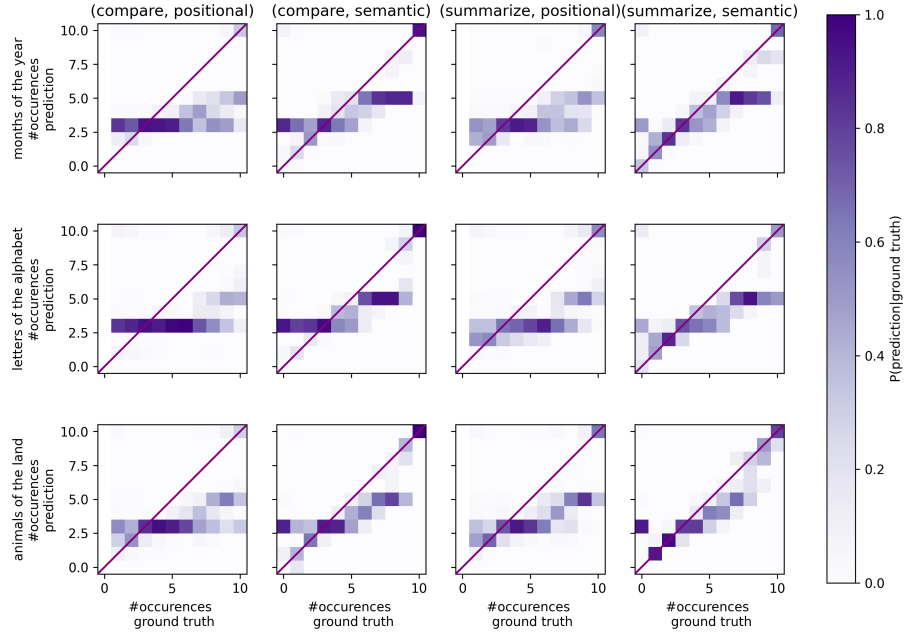


Figure 6. **Mistral 7B-v0.1** predictive performance on counting tasks with temperature zero. Samples of 100 lists of 10 items with each sampling strategy, where we generate all possible queries for different items and positions for each prompt type.

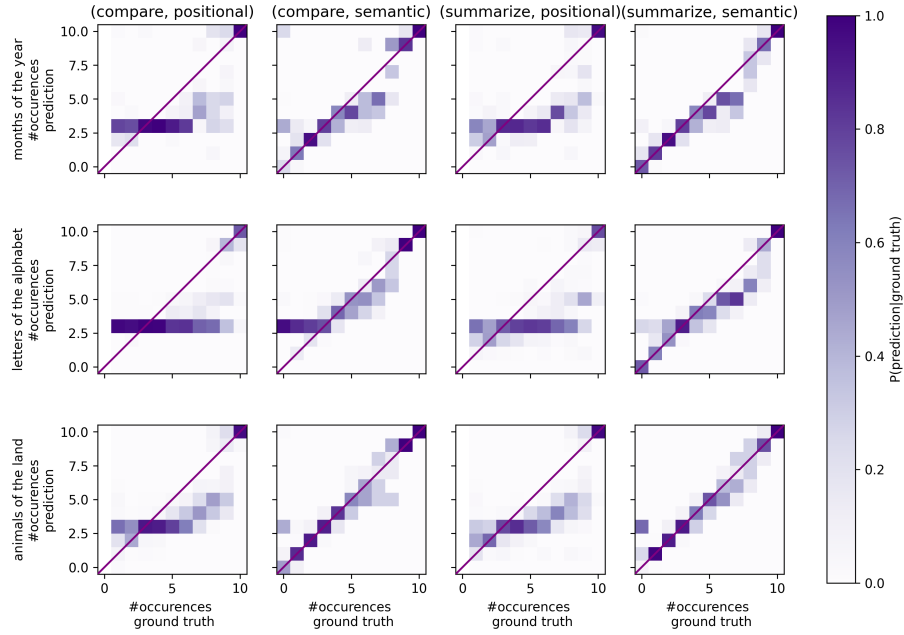


Figure 7. **Llama 3.1** predictive performance on counting tasks with temperature zero. Samples of 100 lists of 10 items with each sampling strategy, where we generate all possible queries for different items and positions for each prompt type.

## B. Explicit Constructions for Orthogonal Embeddings $d = T$

### B.1. Overview

In the parameter regime where  $d \geq T$  there is always an orthonormal basis of size  $T$  in  $\mathbb{R}^d$ , these explicit constructions give the correct prediction for all input token sequences. For all the models we describe below, we define the sum of the hidden layer neurons as:

$$\gamma_\ell = \sum_{i=1}^p \text{ReLU}(z_{\ell,i}) = \sum_{i=1}^p \text{ReLU}(W_1 \tilde{x}'_\ell + b_1)_i \quad (4)$$

In many cases, a simple linear regression can map the scalar  $\gamma_\ell$  to the correct count of tokens  $x_\ell$ , and we describe how to achieve this mapping to the classification problem in Section B.4.

In the following, we characterize which parameters  $W_1, b_1$  in equation 4 allow for a correct mapping in each mechanism. Importantly, the architecture exhibits numerous symmetries due to the feed-forward ReLU network (Petzka et al., 2020). To demonstrate feasibility, we select one specific implementation. In the main text we observe that there is no one-to-one correspondence between our explicit constructions and the learned weights, even though both functions achieve the same perfect accuracy. Throughout, unless otherwise specified, we assume that  $E \in \mathbb{R}^{d \times T}$  is an orthonormal basis of  $\mathbb{R}^d$ , which we will use to create different forms of token embeddings.

The supplementary code at <https://github.com/SPOC-group/counting-attention> contains executable `pytorch` models that have the weight configurations that are used to prove Propositions 4.2-4.3 and B.1, which allows one to test the devised weight configurations for fixed  $T, L, d$  in practice.

### B.2. Relation-based counting

#### B.2.1. (DOT; $p = 1$ )

*Proof of Proposition 4.2.* We set  $T = d > 2$  with  $L \geq 2$  and  $p = 1$ . We choose the embeddings of the tokens of the dot model as

$$e_t = \tilde{e}_t + \tilde{e}_{\text{cnt}} \quad \forall t = 1, \dots, T \quad (5)$$

where the set  $E = \{\tilde{e}_t\}_{t=1}^T$  is an orthonormal basis of an arbitrary but fixed  $T$ -dimensional subspace of  $\mathbb{R}^d$ , and  $\tilde{e}_{\text{cnt}} = \sum_{t=1}^T \tilde{e}_t$ . The key and query matrix are set to the scaled identity  $W_K = W_Q = d^{1/4} I_d$  and hence the mixing layer  $\mathbf{A}_{\text{dot}}$  can be viewed as carrying out the unmodified dot-product operation between all pairs of tokens. The first layer weights  $W_1, b_1 \in \mathbb{R}^d$  can be fixed as

$$W_1 = \tilde{e}_{\text{cnt}} / (T + 1); \quad b_1 = -(1 + L(T + 2)), \quad (6)$$

and the second layer weights  $W_2, b_2 \in \mathbb{R}^L$  follow the recursion

$$(W_2)_1 = -1 + \frac{1}{L + 1}, \quad (b_2)_1 = 0; \quad (7)$$

$$(W_2)_\ell = -1 + \frac{\ell}{L + 1}, \quad (b_2)_\ell = ((W_2)_{\ell-1} - (W_2)_\ell)(\ell - 0.5) + b_{\ell-1}, \quad \forall \ell = 2, \dots, L. \quad (8)$$

Given these parameters, it holds that for tokens  $1 \leq t, s \leq T$  their dot-product is

$$\langle e_t, e_s \rangle = \begin{cases} 2 + T & \text{if } t \neq s, \\ 3 + T & \text{if } t = s. \end{cases} \quad (9)$$

Because of our choice of the query and key matrices, it directly follows that for tokens  $x_\ell, x_m$  at positions  $\ell$  and  $m$  from a given sequence  $\mathbf{x}$ , their attention score is

$$a_{\ell m} = \begin{cases} 2 + T & \text{if } x_\ell \neq x_m, \\ 3 + T & \text{if } x_\ell = x_m. \end{cases} \quad (10)$$



Hence, the mixed token after applying the residual connection is

$$\bar{x}'_\ell = \bar{x}_\ell + \sum_{m:x_\ell=x_m} (T+3)\bar{x}_m + \sum_{m:x_\ell \neq x_m} (T+2)\bar{x}_m \quad (11)$$

so that computing

$$\bar{x}'_\ell W_1 = \left\langle \bar{x}'_\ell, \frac{\tilde{e}_{\text{cnt}}}{1+T} \right\rangle \quad (12)$$

$$= \left\langle \bar{x}_\ell, \frac{\tilde{e}_{\text{cnt}}}{1+T} \right\rangle + \sum_{m:x_\ell=x_m} (T+3) \left\langle \bar{x}_m, \frac{\tilde{e}_{\text{cnt}}}{1+T} \right\rangle + \sum_{m:x_\ell \neq x_m} (T+2) \left\langle \bar{x}_m, \frac{\tilde{e}_{\text{cnt}}}{1+T} \right\rangle \quad (13)$$

$$= 1 + \text{hist}_{\mathbf{x}}(\ell)(T+3) + (L - \text{hist}_{\mathbf{x}}(\ell))(T+2) \quad (14)$$

$$= \text{hist}_{\mathbf{x}}(\ell) + 1 + L(T+2). \quad (15)$$

Then the single hidden unit has the value  $\gamma_\ell = \text{ReLU}(\bar{x}'_\ell W_1 + b_1) = \text{hist}_{\mathbf{x}}(\ell)$ . It is easy to show (analogous to Fig. 8) that the output logits  $c = \gamma_\ell W_2 + b_2$  with  $c \in \mathbb{R}^L$ , correctly identify the count for integer values  $x \in [1, \dots, L]$ . This is because we constructed our recursion such that at a given input  $x = \ell$  we have that  $(W_2)_\ell(\ell - 0.5) + (b_2)_\ell = (W_2)_{\ell-1}(\ell - 0.5) + (b_2)_{\ell-1}$  and  $(W_2)_\ell > (W_2)_{\ell-1}$ , so it holds that

$$\arg \max_{i=1 \dots L} c_i(y) = \begin{cases} 1 & y = 1 \\ 2 & y = 2 \\ \dots & \\ L & y = L, \end{cases} \quad (16)$$

which gives the correct classification output for all possible inputs, and hence solves the histogram task at 100% accuracy.  $\square$

Note, however, that this weight configuration is only one example, and some symmetries in the model can lead to different but also 100% correct algorithms. This is especially important as we compare the regime outlined in the Theorem with the weight configurations learned.

### B.2.2. (BOS+SFTM; $p = 1$ )

*Proof of Proposition 4.1 for bos+sftm.* We set  $T = d > 2$  with  $L \geq 2$  and  $p = 1$  and consider the model dot+sftm. Note that in this model every sequence  $\mathbf{x}$  is prefixed with  $t_{\text{BOS}}$  before it is fed into the embedding and then the mixing layer. Again we use mutually orthogonal embeddings.  $E = \{\tilde{e}_t\}_{t=1}^T$  is an orthonormal basis of an arbitrary but fixed  $T$ -dimensional subspace of  $\mathbb{R}^d$ , and  $\tilde{e}_{\text{cnt}} = \sum_{t=1}^T \tilde{e}_t$ . We set  $e_{\text{BOS}} = \sum_{t=1}^T E_t$ , where  $e_t = E_t$  and the latter is a column of  $E$ . Analogous to the background token from Proposition 4.1 there is only one direction  $p = 1$  to detect in the feedforward model, so we set

$$W_1 = e_{\text{BOS}}; \quad b_1 = -1. \quad (17)$$

For a given token  $x_\ell$  we have that in the dot-product mechanism  $\langle e_{\text{BOS}}, e_{x_\ell} \rangle = 1$ ,  $\langle e_{x_\ell}, e_{x_m} \rangle = 1$  if  $x_m = x_\ell$  and 0 otherwise. Due to the softmax, the mixing coefficient is  $a = e/((k_{x_\ell} + 1)e + (L - k_{x_\ell}))$  (where  $e$  is Euler's number) for comparing  $x_\ell$  to  $t_{\text{BOS}}$  and to all the tokens where  $x_\ell = x_m$ , and  $b = 1/((k_{x_\ell} + 1)e + (L - k_{x_\ell}))$  otherwise, where,  $k_{x_\ell} = \text{hist}_{\mathbf{x}}(\ell)$ . Hence, the mixed token is:

$$\bar{x}'_\ell = ae_{\text{BOS}} + ak_{x_\ell}\bar{x}_\ell + \sum_{x_m \neq x_\ell} b\bar{x}_m + \bar{x}_\ell. \quad (18)$$

Applying  $W_1$  and  $b_1$ , we obtain:

$$\begin{aligned} \gamma_\ell &= aT + ak_{x_\ell} + b(L - k_{x_\ell}) \\ &= aT + ak_{x_\ell} + 1 - a(k_{x_\ell} + 1) \\ &= a(T - 1) + 1 \end{aligned} \quad (19)$$

since  $(k_{x_\ell} + 1)a + (L - k_{x_\ell})b = 1$  by normalization via the softmax function. The value of  $\gamma_\ell$  has a dependence on  $k_\ell$  through  $a$  and can be readout into the correct classification as shown Fig. 8.  $\square$

### B.2.3. (BOS; $p = 1$ )

*Proof of Proposition 4.1 - bos.* We set  $T = d > 2$  with  $L \geq 2$  and  $p = 1$ . The construction of the embeddings and  $e_{\text{BOS}}$  is analogous to the construction from Section B.2.2 for  $\text{bos+sftm}$  in the same setting. However, since no softmax is applied, the mixing coefficients as outputs of  $\mathbf{A}_{\text{dot+sftm}}$  for comparing  $(x_\ell, t_{\text{BOS}})$  or  $(x_\ell, x_m)$  where  $x_\ell = x_m$  is  $a = 1$ . For  $x_\ell \neq x_m$  it is  $b = 0$ . Then from inserting these values in equation 18 and applying  $W_1 = e_{\text{BOS}}$  and  $b_1 = -T$  we obtain

$$\gamma_\ell = k_{x_\ell}. \quad (20)$$

This clearly allows again the single neuron to be read off to the correct result similar to the construction from equation 6.  $\square$

Note that there is a simple alternative construction that uses the tagged embeddings from the constructive proof of Prop. 4.2.

*Alternative Proof of Proposition 4.1 - bos.* We set  $T = d > 2$  with  $L \geq 2$  and  $p = 1$ . We note that by setting  $t_{\text{BOS}}$  to zero we can achieve equivalence to the model  $\text{dot}$ . Since according to Prop. 4.2 there exists a weight configuration for  $\text{dot}$  which solves the histogram task, this configuration will also solve the histogram task for  $\text{bos}$  with  $t_{\text{BOS}} = 0$ .  $\square$

## B.3. Inventory-based counting

### B.3.1. (LIN; $p = T$ ).

*Proof of Proposition 4.3 - lin.* Assume that  $T = d > 2$  with  $L > 2$  and  $p = T$  and the goal is to find a weight configuration for the model  $\text{lin}$ . As embeddings we directly use the orthonormal basis with  $T$  vectors  $e_t$  in  $\mathbb{R}^d$ , where vectors are the embeddings are for the  $T$  tokens. We set

$$\mathbf{A}_{\text{lin}} = \begin{bmatrix} a & a & \cdots & a \\ a & a & & \\ \vdots & & \ddots & \\ a & & & a \end{bmatrix}; \quad W_1 = E; \quad b_1 = -1, \quad (21)$$

where  $a = 1/L$ . We start by writing  $z_{\ell,t}$  for  $t \in \{1, \dots, p = T\}$ , the  $t$ -th activation of the first hidden layer of the feed-forward module

$$z_{\ell,t} = \sum_{m=1}^L a_{\ell m} \langle e_{x_m}, e_t \rangle + \langle e_{x_\ell}, e_t \rangle - 1. \quad (22)$$

If  $e_t = e_{x_\ell}$ , we have

$$z_{\ell,t} = k_{x_\ell} a + 1 - 1 = a k_{x_\ell}, \quad (23)$$

where,  $k_{x_\ell} = \text{hist}_x(\ell)$ , applying the ReLU to this scalar keeps its value unchanged. If  $e_t \neq e_{x_\ell}$ , we have

$$z_{\ell,t} = a k_{e_t} + 0 - 1 = a k_{e_t} - 1 \leq 0. \quad (24)$$

The right hand side of the above equation is negative given our choice of  $a$ , hence applying the ReLU returns 0. This means that, for each token in the input sequence, the contributions of orthogonal tokens cancel, leaving us with a single hidden neuron activated. Hence the count can be read off from  $\gamma_\ell$ . Since only one neuron is activated at a time, the readout from the same procedure as in  $\text{bos+sftm}$  can be applied to all hidden neurons  $z_{\ell,t}$  simultaneously, instead of only one. This allows the model to solve the histogram task.  $\square$

### B.3.2. (LIN+SFTM; $p = T$ )

*Proof of Proposition 4.3 - lin+sftm.* Assume that  $T = d > 2$  with  $L > 2$  and  $p = T$ . With the statement already proven for  $\text{lin}$ , we note that we can construct  $\mathbf{A}_{\text{lin+sftm}}$  such that it is equivalent to  $\mathbf{A}_{\text{lin}}$  from equation 21 via

$$\mathbf{A}_{\text{lin+sftm}} = \begin{bmatrix} a & a & \cdots & a \\ a & a & & \\ \vdots & & \ddots & \\ a & & & a \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} \alpha & \alpha & \cdots & \alpha \\ \alpha & \alpha & & \\ \vdots & & \ddots & \\ \alpha & & & \alpha \end{bmatrix} \right), \quad (25)$$

where  $a = 1/L$  which implicitly defines a choice of  $\alpha$ . This means that the construction is equivalent to  $\text{lin}$  and it follows automatically that also  $\text{lin+sftm}$  can solve the histogram task.  $\square$

B.3.3. (DOT+SFTM:  $p = d = T$ )

**Proposition B.1** (IC for dot+sftm). *For dot+sftm and given  $L, T > 2$  there exists a configuration of weights which solves the histogram task for  $p \geq T$  and  $d \geq T$ .*

*Proof for Proposition B.1.* We assume  $L, T > 2$  and  $p = T$  and  $d = T$  and we consider dot+sftm. As previously for dot in Prop. 4.2, we set the key and query matrix to the scaled identity  $W_K = W_Q = d^{1/4}I_d$ . We use an orthonormal basis of  $\mathbb{R}^d$  to define the parameters  $e_t \in \mathbb{R}^d$  for the token embeddings. In  $\mathbf{A}_{\text{dot+sftm}}$  the pre-softmax mixing weights will be 1 for equal and 0 for different tokens due to the unit-norm token embeddings. Defining  $k_{x_\ell} = \text{hist}_x(\ell)$  for brevity, after the softmax we have that

$$a_{lm} = \begin{cases} \frac{e}{(L-k_{x_\ell})+ek_{x_\ell}} & x_m = x_\ell, \\ \frac{1}{(L-k_{x_\ell})+ek_{x_\ell}} & \text{else.} \end{cases} \quad (26)$$

Hence, for  $e_t \neq e_{x_\ell}$

$$\langle \bar{x}'_\ell, e_t \rangle = \frac{k_{e_t}}{(L-k_{x_\ell})+ek_{x_\ell}} < 1, \quad (27)$$

while for  $e_t = e_{x_\ell}$

$$\langle \bar{x}'_\ell, e_{x_\ell} \rangle = \frac{k_{x_\ell}e}{(k_{x_\ell}e + (L-k_{x_\ell}))} + 1, \quad (28)$$

where the extra summand comes from the residual connection. Hence, by setting

$$W_1 = E; \quad b_1 = -1, \quad (29)$$

and applying the ReLU activation, equation 27 will be 0, while equation 28 will implicitly give us the counts as:

$$k_{x_\ell} = (L\gamma_\ell)/(-e\gamma_\ell + \gamma_\ell + e) \quad (30)$$

While the final layer cannot immediately implement non-linear functions in  $\gamma_\ell$ , it can take advantage of the fact that  $\gamma_\ell$  can take only  $L$  different values, similar to how we constructed  $W_2$  and  $b_2$  in Section B.2.1. Since eventually we need to map the  $L$  values of  $\gamma_\ell$  to the counts  $[1, \dots, L]$  the linear output layer is sufficient to implement this non-linear discrete map. Fig. 8 shows an example for this map for a given example. This allows the model to solve the histogram task.

The statement for  $p > T$  and  $d > T$  follows as we can simply set the surplus of parameters in the hidden layer/embeddings to zero.  $\square$

## B.4. Mapping a scalar to a categorical one-hot encoding

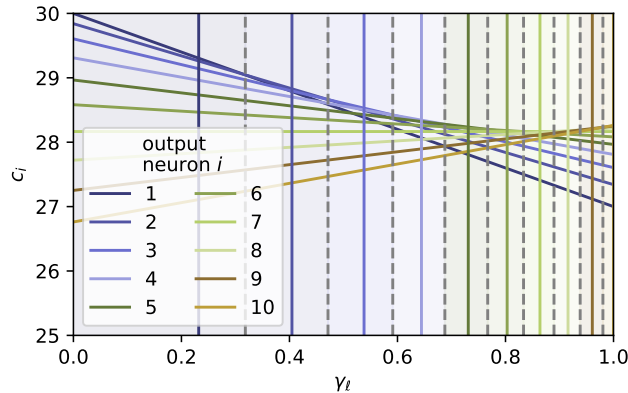


Figure 8. Demonstration of a hidden neuron output  $\gamma_\ell$  which is mapped to  $L = 10$  different neurons  $c_i$  using a single output layer, according to the decision boundaries shown by the dotted lines. After applying the argmax function to the 10 output neurons, the highest value gives the discrete output. The solid lines mark the values a single hidden neuron would achieve for different counts in the explicit construction of the dot+sftm model.

It is straightforward to map a single scalar  $\gamma_\ell$  to a series of neurons which activate one after another. This is needed as the second part of the feed-forward parameters to transform the count measured by the sum of the hidden neurons  $\gamma_\ell$  to the discrete categorical representation of the output vector. Every output logit is a linear function of the hidden neuron’s value. Since in our constructions we only map functions, where the ground truth output logit corresponds to an interval  $[a, b] \in R$ , the superposition of linear functions with increasing slope allows us to realize such a mapping. A visual sample is given in Fig. 8 for `dot+sftm`. In Fig. 9 we show the outputs for the `lin+sftm` model with the best accuracy for  $T = 32$  for every  $p, d$  ran in Fig. 1. While it is possible to learn the count from one hidden neuron only using inventory-based counting for each neuron, for some examples the count information seems to be spread out over several hidden neurons: The output logits are non-linear in the count and can hence not rely on a single hidden neuron only.

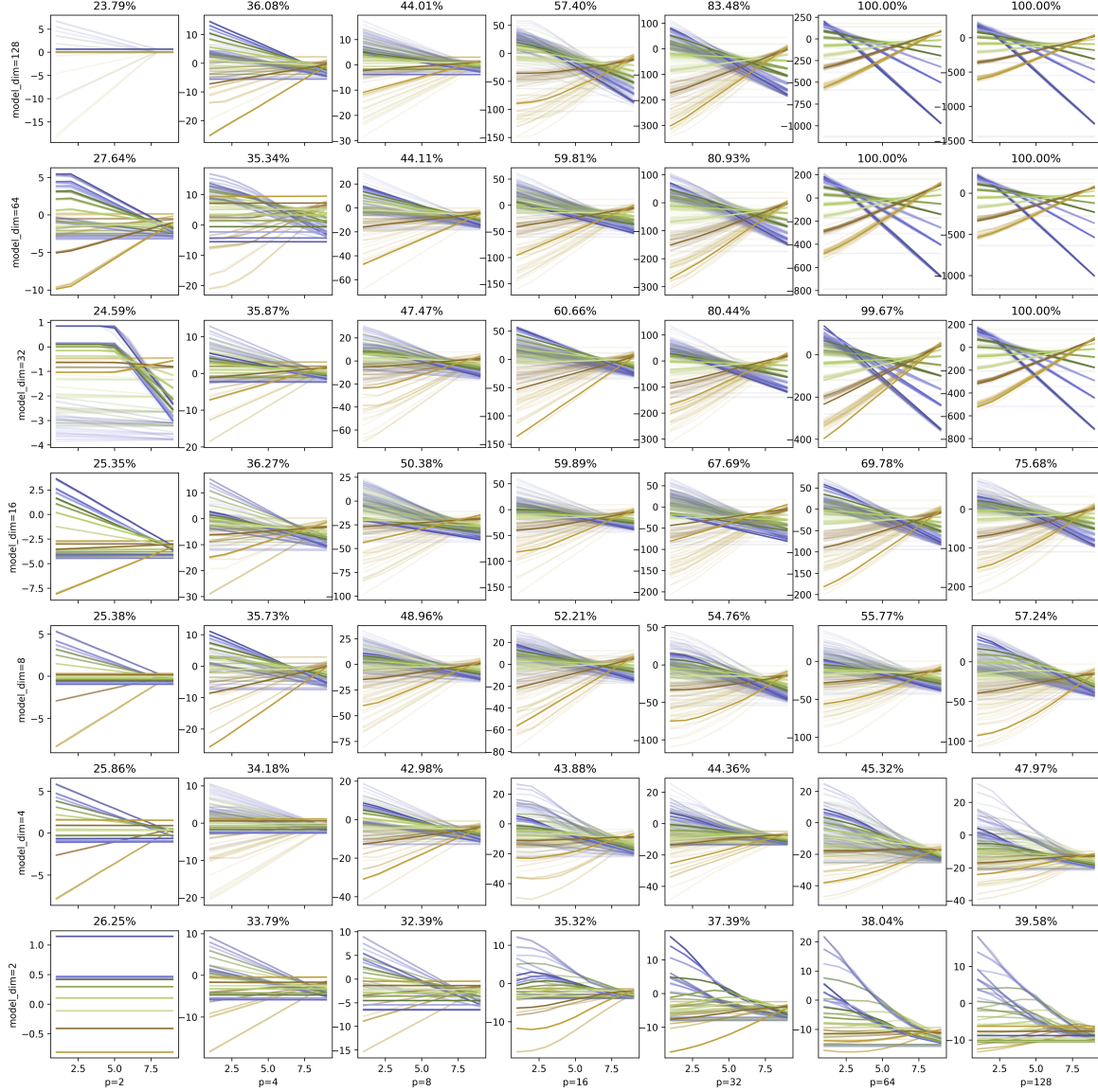


Figure 9. The output neurons  $c_i(\mathbf{x}_\ell)$  visualized for examples of a learned version of `lin+sftm` for several model dimensions  $d$  and hidden layer sizes  $p$ . Note that differently from Fig. 8, in this case the x-axis shows the number of occurrences  $hist_{\mathbf{x}}(\ell = 0) = 1, \dots, L - 1$  of the token  $t$  in an input sequence  $\mathbf{x} = [t, \dots, t, v, \dots, v]$  that contains otherwise only a token  $v \neq t$  (and *not* the activation of a hidden neuron). We show the activations  $c_i$  of the final layer output neurons activations (logits) in terms of the number of occurrences of a given token in the input. The colors represent the different output predictions and are as in the explicit construction from Fig. 8. We show several activations for different tokens  $t \in \mathcal{T}$ , where  $T = 32$ , and we highlight one of the example tokens  $t$  with a wider line. While similar to the explicit construction from Fig. 8, the models with 100% accuracy are not necessarily linear in the count.



## C. Explicit Constructions for Linearly Dependent Embeddings $d < T$

### C.1. Overview

In this section, we discuss the scenario when  $d < T$ , i.e. when the embeddings are necessarily linearly dependent. In that case, we can no longer assume that there exist embeddings with  $\langle e_t, e_s \rangle = 0$  for all  $t \neq s$ . Nonetheless, also in this regime for some models it is possible to provide explicit constructions of the weights that have 100% accuracy. This relies on the fact that the prediction problem is inherently discrete, i.e. it chooses exactly one among  $L$  classes. When we examine  $\gamma_\ell \in \mathbb{R}$  from equation 4 which is mapped to the discrete class through the readout layer (see for example Fig. 8), we notice that the class boundary (the gray dashed class borders) can be placed variably in the margin between the values that  $\gamma_\ell$  assumes for different counts  $k_{x_\ell}$  (solid lines). In the following explicit constructions, our goal is to design embeddings with  $d < T$  in such a way that we maximize the aforementioned margin: there will be pairs of token embeddings in the alphabet that have non-zero similarity  $\langle e_t, e_s \rangle$ , and in equation 4 this will create non-zero terms that will alter the value of  $\gamma_\ell$ . This means that for every possible sequence with  $k$  occurrences of token  $x_\ell$ , the hidden activation  $\gamma_\ell$  will assume values in a certain range  $[\gamma_\ell^{\text{lower}}(k), \gamma_\ell^{\text{upper}}(k)]$ . If these ranges overlap for different  $k$ , the count cannot be identified. However, we construct embeddings such that every for every  $k = 1, \dots, L-1$  it holds that

$$\gamma_\ell^{\text{upper}}(k) < \gamma_\ell^{\text{lower}}(k+1), \quad (31)$$

so we can still use a construction as in Fig. 8 to correctly compute the final count. In the remainder of this section, we introduce explicit constructions with  $d < T$  for a given  $L$ , both for the cases where we have relation-based counting and inventory-based counting (the same argument as above transfers to  $z_{\ell,t}$  from equation 22). Notably, for the explicit constructions we propose, the function of the lowest achievable  $d(p, T, L)$  differs across different mixing types. To summarize:

- For models with **A** constant in the inputs or models *without* softmax activation, our explicit construction relies on an embedding matrix with a small mutual coherence. The mutual coherence is a concept from compressed sensing and coding theory that ensures that the maximal similarity between pairs of vectors is small (Donoho & Elad, 2003). We can upper bound the mutual coherence that the margins of the construction can tolerate to still achieve perfect accuracy in terms of a given  $L$ . At the same time, the mutual coherence of a set of vectors is naturally lower bounded in terms of the number of vectors  $T$  and their respective dimension  $d$ , known as the Welch bound (Welch, 1974). When this bound can be attained and  $T, L$  are given, this leads to the following bounds on  $d$  for the different models, as outlined in Prop. 4.4, as

$$\begin{aligned} (\text{lin}, \text{lin}+\text{sftm}; p=T): \left\lceil \frac{T(2L-3)^2}{T-1+(2L-3)^2} \right\rceil &\leq d, \\ (\text{dot}, \text{bos}; p=1): \left\lceil \frac{T(2L-3)^2}{T-1+(2L-3)^2} \right\rceil + 1 &\leq d, \\ (\text{dot}, \text{bos}; p=T): \left\lceil \frac{T(L-1)}{T-1+(L-1)} \right\rceil &\leq d. \end{aligned}$$

- For **bos+sftm** we rely on the fact that the softmax function accentuates the largest value and thereby can drive attention scores for equal tokens  $a_{ii}$  higher relative to attention scores of non-equal tokens  $a_{ij}$ . This distinguishes it from the previous case, and allows us to state Prop. 4.5 for which we describe an explicit construction that solves the histogram task with

$$\begin{aligned} (\text{bos}+\text{sftm}; p=1): d &\geq \lceil \log_2(T+1) \rceil + 2. \\ (\text{dot}+\text{sftm}; p=T): d &\geq \lceil \log_2(T+1) \rceil + 2. \end{aligned}$$

Notably there is no explicit dependence on  $L$  for the dimension. However, the smaller the dimension  $d$  the more accurate computations and softmax numerical stability are required, as the softmax temperature depends on  $L$ . With infinitely precise computations we show it is even possible to achieve perfect accuracy with  $d = 4$ , but for finite computations this might pose a problem when  $L$  becomes too large.

### C.2. Explicit construction for bounded mutual coherence

We define the *mutual coherence*  $\mathcal{M}$  of a set of  $T$  unit norm vectors  $\{v_1, \dots, v_T\} \subset \mathbb{R}^d$  as

$$\mathcal{M} = \max_{i \neq j} |\langle v_i, v_j \rangle|. \quad (32)$$

This value is lower bounded for a given matrix by the Welch bound (Welch, 1974)

$$\mathcal{M} \geq \sqrt{\frac{T-d}{d(T-1)}} = \mathcal{W}(T, d), \quad (33)$$

and equality can only be attained if  $T < d^2$  (Strohmer & Heath, 2003). There is a large body of work in coding theory and compressed sensing concerning the existence and construction of a set of vectors that attains  $\mathcal{M}$  at or close to  $\mathcal{W}(T, d)$ . Explicit constructions exist but are not known for every combination of  $T$  and  $d$ . A list with existing constructions for the real space for small  $T, d$  can be found in (Fickus & Mixon, 2016), but otherwise gradient-based optimization has been used to find good candidate matrices (Jiang et al., 2017; Jyothi & Babu, 2022).

In order to prove Prop. 4.4, we use the following idea: For a given  $T, L$  and  $p$ , we can derive an upper bound on the mutual information of the embeddings in terms of  $L$ , which is required to obtain perfect accuracy. The form of this upper bound depends on the precise mixing strategy and the choice of  $p$ . Through the Welch lower bound on  $\mathcal{M}$  we can in turn obtain a lower bound on  $d$  in terms of  $L$  and  $T$ . Note that the Welch bound cannot be attained for  $T < d^2$  and in this case the bound on  $d$  is strict.

#### C.2.1. (LIN, LIN+SFTM; $p = T$ )

*Proof of Proposition 4.4 - lin.* To show the bound on  $d$ , we analyze the inventory-based construction for `lin` in equation 21. Given that  $p = T$ , and  $L > 2$  is given, let us assume that there exists set of  $T$  unit norm vectors  $\{e_1, \dots, e_T\} \subset \mathbb{R}^d$  with mutual coherence  $\mathcal{M}$ . We use these vectors as our embeddings.

The value  $z_{\ell, t}$  for  $t = x_\ell$ , with  $W_1 = [e_1, \dots, e_T]$  and  $b_1 = -1$  is

$$z_{\ell, t} = ak_{x_\ell} + a \sum_{m: x_m \neq t} \langle e_{x_m}, e_t \rangle, \quad (34)$$

and using that fact that the mutual coherence bounds the absolute value of the inner product

$$ak_{x_\ell} - a\mathcal{M}(L - k_{x_\ell}) \leq z_{\ell, t} \leq ak_{x_\ell} + a\mathcal{M}(L - k_{x_\ell}). \quad (35)$$

Similarly, for  $t \neq x_\ell$  and  $a = 1/L$  it still holds that

$$z_{\ell, t} \leq ak_t + a(L - k_t)\mathcal{M} - 1 + \mathcal{M} \leq 0, \quad (36)$$

provided that  $\mathcal{M} < 1/(L + 1)$ , for the worst case where  $k_t = L - 1$ . This means that the ReLU sets all hidden neurons  $z_{\ell, t}$  to zero when  $t \neq x_\ell$ , and are therefore no contribution to the final result. Then, defining

$$\gamma_\ell^{\text{lower}}(k) = ak - a\mathcal{M}(L - k), \quad (37)$$

$$\gamma_\ell^{\text{upper}}(k) = ak + a\mathcal{M}(L - k), \quad (38)$$

we have that indeed for a sequence where  $x_\ell$  occurs  $k = 1, \dots, L - 1$  times it holds that

$$0 \leq \gamma_\ell^{\text{lower}}(k) \leq \gamma_\ell(k) \leq \gamma_\ell^{\text{upper}}(k). \quad (39)$$

The first inequality is required due to the ReLU and holds when  $\mathcal{M} < 1/(L - 1)$ . From equation 31 we have the condition that for all  $k = 1, \dots, L$  it holds that

$$\gamma_\ell^{\text{upper}}(k) < \gamma_\ell^{\text{lower}}(k + 1), \quad (40)$$

$$k + (L - k)\mathcal{M} < (k + 1) + (L - k - 1)(-\mathcal{M}), \quad (41)$$

$$\mathcal{M} < \frac{1}{2(L - k) - 1}, \quad (42)$$

and since we assume that there exist at least two different tokens in the sequence, minimizing the bound over  $k$  leaves for  $k = 1$

$$\mathcal{M} < \frac{1}{2L - 3}. \quad (43)$$

which is valid provided that  $L \geq 2$ . Collecting all previous bounds on  $\mathcal{M}$ , we conclude that when  $L \geq 4$  the above construction achieves the correct counts with  $\mathcal{M} < \frac{1}{2L-3}$ .

The Welch bound equation 33 gives an upper bound on  $\mathcal{M}$  in terms of  $T, d$  and therefore yields the final condition

$$d \geq \left\lceil \frac{T(2L-3)^2}{T-1+(2L-3)^2} \right\rceil \quad (44)$$

under which the given weight configuration is able to solve the histogram task with perfect accuracy.  $\square$

For `lin+sftm` the construction and conditions transfer directly, when the constant  $\mathbf{A}_{\text{lin+sftm}}$  is constructed to match  $\mathbf{A}_{\text{lin}}$  exactly.

### C.2.2. (DOT, BOS; $p = 1$ )

*Proof of Proposition 4.4 - dot,  $p = 1$ .* We assume that  $L > 2$  and  $T$  given and we use a similar idea as the relation-based weight configuration from the proof of Prop. 4.2 for `dot` with  $p = 1$ . For the token embeddings, we assume that we have a set of  $T$  unit norm vectors with  $\{v_1, \dots, v_T\} \subset \mathbb{R}^{d-1}$  with mutual coherence  $\mathcal{M}$ , where  $d > 2$ . We set the entries of the  $T$  embedding vectors  $e_t \in \mathbb{R}^d$  to be

$$e_t = \begin{bmatrix} v_t \\ \alpha \end{bmatrix}. \quad (45)$$

The shared counting subspace is defined on the last coordinate of the vectors via  $e_{cnt} = [0, 0, \dots, 1/\alpha]$ . Then

$$\langle e_t, e_t \rangle = 1 + \alpha^2 \quad (46)$$

$$|\langle e_t, e_s \rangle| \leq \mathcal{M} + \alpha^2 \quad (47)$$

The mixed token with the residual connection at position  $\ell$  for a given input sequence is

$$\bar{x}'_\ell = \sum_{m=1}^L \langle e_{x_m}, e_{x_\ell} \rangle e_{x_m} + e_{x_\ell} \quad (48)$$

and the single hidden neuron  $\gamma_\ell$  for a bias term  $b_1 = 0$  and  $W_1 = e_{cnt}$

$$\gamma_\ell = \langle e_{cnt}, \bar{x}'_\ell \rangle = \sum_{m=1}^L \langle e_{x_m}, e_{x_\ell} \rangle \langle e_{x_m}, e_{cnt} \rangle + \langle e_{x_\ell}, e_{cnt} \rangle \quad (49)$$

$$= k_{x_\ell}(1 + \alpha^2) + \sum_{m: x_m \neq x_\ell} \langle e_{x_m}, e_{x_\ell} \rangle + 1 \quad (50)$$

So that we can achieve for a given count  $k$  the  $\gamma_\ell^{\text{lower}}(k) \leq \gamma_\ell(k) \leq \gamma_\ell^{\text{upper}}(k)$  with

$$0 \leq \gamma_\ell^{\text{lower}}(k) = k(1 + \alpha^2) - (L - k)(\mathcal{M} + \alpha^2) + 1, \quad (51)$$

$$\gamma_\ell^{\text{upper}}(k) = k(1 + \alpha^2) + (L - k)(\mathcal{M} + \alpha^2) + 1. \quad (52)$$

We achieve the upper bound from zero, when  $\mathcal{M} < 2/(L - 1)$ , assuming that  $\alpha$  is close enough to zero so that is is negligible. Finally, the condition from equation 31 yields

$$\mathcal{M} < \frac{1}{2(L - k) - 1} - \frac{2(L - k) - 2}{2(L - k) - 1} \alpha^2 \quad (53)$$

under the condition that  $0 < \alpha < \sqrt{\frac{1}{2(L - k) - 2}}$ . Again, assuming there exist at least two different tokens in the sequence, the r.h.s. of the above expression is minimized for  $k = 1$  as

$$\mathcal{M} < \frac{1}{2L - 3} - \frac{2L - 4}{2L - 3} \alpha^2 \quad (54)$$

which is always positive assuming  $L \geq 2$ . This is the relevant bound when we have a large enough  $L \geq 5$  and again  $\alpha$  is close enough to zero. Again, combining this with the Welch bound equation 33 leads to

$$d - 1 \geq \left\lceil \frac{T(\frac{2L-3}{1-(2L-4)\alpha^2})^2}{T-1 + (\frac{2L-3}{1-(2L-4)\alpha^2})^2} \right\rceil, \quad (55)$$

and when we choose  $\alpha > 0$  close to zero, as for `lin` before

$$d \geq \left\lceil \frac{T(2L-3)^2}{T-1 + (2L-3)^2} \right\rceil + 1. \quad (56)$$

□

This proof holds equivalently for `bos` when we set the BOS token embedding to zero.

### C.2.3. (`dot`, `bos`; $p = T$ )

We can decrease the required dimension  $d < T$  even further than previously, when we have  $p = T$  and implement inventory-based counting in the `dot` model (and equivalently in the `bos` model). In that case, the lower bound on  $d$  becomes more loose, because we combine the ideas we saw in `lin` and  $p = T$  for inventory-based counting and the effects on the margin in `dot` and  $p = 1$ .

*Proof of Proposition 4.4 - dot,  $p = T$ .* In our construction, for a given  $T$  and  $L$ , we assume that there is a set of  $T$  unit norm vectors  $\{e_1, \dots, e_T\} \subset \mathbb{R}^d$  with mutual coherence  $\mathcal{M}$  upon which we build our embeddings. Note that the only difference to the previous relation-based case  $p = 1$  is that this time there is no extra counting direction. Importantly, we set  $K = d^{1/4}I_d$  as before, but  $Q = \frac{1}{L}d^{1/4}I_d$ . This gives an extra factor in the attention scores. Further, we set  $b_1 = -1$  and the columns of  $W_1 \in \mathbb{R}^{d \times T}$  to the embeddings  $e_t$ , as we did for `lin`. This results in a mixed token  $\bar{x}'_\ell$  according to equation 48. The hidden neuron is

$$z_{\ell,t} = \frac{1}{L} \sum_{m=1}^L \langle e_{x_m}, e_{x_\ell} \rangle \langle e_t, e_{x_m} \rangle + \langle e_t, e_{x_\ell} \rangle - 1 \quad (57)$$

then with  $t = x_\ell$  we have

$$z_{\ell,t} = \frac{1}{L} \left( k_{x_\ell} + \sum_{m:x_m \neq t} \langle e_{x_m}, e_t \rangle^2 \right). \quad (58)$$

Note that the square in equation 57 is what differs from the  $z_{\ell,t}$  in equation 34. This is because the term  $\langle e_{x_m}, e_t \rangle$  is once introduced through the dot-product attention and once through the dot-product via  $W_1$ . Conversely, with  $t \neq x_\ell$  it becomes

$$z_{\ell,t} = \frac{1}{L} \left( k_t \langle e_t, e_{x_\ell} \rangle + \sum_{m:x_m \neq t} \langle e_{x_m}, e_{x_\ell} \rangle \langle e_{x_m}, e_t \rangle \right) + \langle e_t, e_{x_\ell} \rangle - 1 \quad (59)$$

$$\leq \frac{1}{L} (k_t \mathcal{M} + (L - k_t) \mathcal{M}) + \langle e_t, e_{x_\ell} \rangle - 1 \quad (60)$$

$$\leq 2\mathcal{M} - 1 \quad (61)$$

if we set  $\mathcal{M} < 0.5$ , which we need anyways for  $L \geq 2$  by the stronger upper bound on  $\mathcal{M}$  that we derive in the following, we finally have for  $t \neq x_\ell$

$$z_{\ell,t} < 0. \quad (62)$$

Again, negative  $z_{\ell,t}$  are set to zero via the ReLU, and the final outcome  $\gamma_\ell = z_{\ell,t=x_\ell}$  depends only on a single hidden neuron equation 57. This eventually leads to

$$0 \leq \gamma_\ell^{\text{lower}}(k) = \frac{k}{L}, \quad (63)$$

$$\gamma_\ell^{\text{upper}}(k) = \frac{1}{L} (k + \mathcal{M}^2(L - k)), \quad (64)$$

and using the same concept as before, while minimizing over  $k$  and applying the Welch bound, to the upper bound

$$\mathcal{M} < \sqrt{\frac{1}{L-1}}. \quad (65)$$

The final bound is more loose than it was for  $p = 1$  as we only require

$$d \geq \left\lceil \frac{T(L-1)}{T-1+(L-1)} \right\rceil. \quad (66)$$

□

### C.3. Explicit Construction with binary representations and softmax

In our final analysis we examine the key difference between the models `bos+sftm` and `bos` – the softmax activation. In order to show Prop. 4.4 we needed to construct embeddings with a low mutual coherence, because the term  $\langle e_t, e_s \rangle$  introduced an error on the mixed token, when  $t$  and  $s$  were not equal. Now, with the softmax activation applied to the mixing coefficients, the model can use the non-linearity of this transform to its advantage to separate the relative error.

Recall the softmax function is

$$\text{sftm}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \quad \text{for } i = 1, 2, \dots, n, \quad (67)$$

and when we compute  $\text{sftm}(\kappa \mathbf{z})_i$  we say it is a softmax with a inverse temperature  $\kappa > 0$ . When  $\mathbf{z}$  of length  $L$  contains only two different values, one with  $k$  and the other with  $L - k$  occurrences, then as  $\kappa \rightarrow \infty$  the mass concentrates only on the larger value of the two, and sets the other to zero. We use this intuition to create token embeddings that fulfill for all  $t, s = 1, \dots, T$  and  $s \neq t$

$$\langle e_t, e_t \rangle = 1, \quad (68)$$

$$\langle e_t, e_s \rangle < 1 + \epsilon, \quad (69)$$

where  $\epsilon > 0$ .

The idea is that the softmax with a high enough inverse temperature sets the term for different tokens,  $\langle e_t, e_s \rangle$ , *close enough* to zero, essentially eliminating the noise. Note that equation 68 is a weaker condition on the set of token embeddings than for example the bound of the mutual coherence in terms of the sequence length  $L$  `bos` with  $p = 1$  in Section C.2.2. It allows us to obtain perfect accuracy with smaller  $d$ . In the following, we describe the construction of the matrix explicitly.

The supplementary code at [supplementary material](#) contains executable `pytorch` models that have the weight configurations that are used to prove Propositions 4.5 and the Remark for  $d = 4$ , which allows one to test the devised weight configurations for fixed  $T, L, d$  in practice.

#### C.3.1. (`BOS+SFTM`; $p = 1$ )

*Proof of Proposition 4.5 - `bos+sftm`.* For a given  $T, L > 2$  we set the embeddings vectors to the binary representation of the token index  $t = 1, \dots, T$  in  $d' = \lceil \log_2(T+1) \rceil$  dimensions

$$e_t = \begin{bmatrix} \text{bin}(t) \langle \text{bin}(t), \text{bin}(t) \rangle^{-1} \\ \alpha \\ 0 \end{bmatrix}; \quad e_{BOS} = \begin{bmatrix} \text{bin}(0) \langle \text{bin}(0), \text{bin}(0) \rangle^{-1} \\ 1/\alpha \\ 1 \end{bmatrix}. \quad (70)$$

where  $\text{bin}(t) = [v_1, \dots, v_{d'}] \in \{0, 1\}^{d'}$  with  $t = \sum_{i=1}^{d'} v_i 2^{i-1}$ . We select  $\alpha > 0$ . Then we have that

$$\langle e_t, e_t \rangle = 1 + \alpha^2, \quad (71)$$

$$\alpha^2 \leq \langle e_t, e_s \rangle \leq \sqrt{1 - \frac{1}{d'}} + \alpha^2 \leq 1 + \alpha^2 - \epsilon, \quad (72)$$

$$\langle e_t, e_{BOS} \rangle = 1, \quad (73)$$



where  $\sqrt{\frac{d'-1}{d'}} = \langle e_{2^{d'}-1}, e_{2^{d'}-2} \rangle$ , which has the largest overlap among all possible non-equal pairs of tokens, and the lower bound comes from all coordinates being positive. Using a readout on the direction only present in the  $e_{BOS}$  token, namely,  $W_1 = [e_{cnt}] = [0, \dots, 0, 1] \in \mathbb{R}^d$  and  $b_1 = 0$ , we construct

$$\gamma_\ell = \langle e_{cnt}, \bar{x}'_\ell \rangle = \text{sftm}(EE_\ell^T)_0 \langle e_{BOS}, e_{cnt} \rangle + \sum_{m=1}^L \text{sftm}(EE_\ell^T)_{m+1} \langle e_{x_m}, e_{cnt} \rangle + \langle e_{x_\ell}, e_{cnt} \rangle \quad (74)$$

$$= \text{sftm}(EE_\ell^T)_0 \quad (75)$$

$$= \text{sftm}([\langle e_\ell, e_{BOS} \rangle, \langle e_\ell, e_1 \rangle, \dots, \langle e_\ell, e_L \rangle])_0 \quad (76)$$

The goal of applying the softmax function is to diminish the contributions of error equation 72, while having the final dimension of the  $e_{BOS}$  token be representative of the count of  $x_\ell$ . The maximum error is induced when the upper bound equation 72 is attained for all tokens in the sequence  $\mathbf{x}$  that are not equal to  $x_\ell$ . The minimum error is obtained when these different tokens attain the lower bound. Without loss of generality on the ordering, this implies that for a given length  $L$  and a softmax activation function with an inverse temperature  $\kappa^1$ , we have that

$$\gamma_\ell^{\text{lower}}(k) = \frac{e^{\kappa 1}}{e^{\kappa 1} + k e^{\kappa(1+\alpha^2)} + (L-k)e^{\kappa(1+\alpha^2-\epsilon)}} \quad (77)$$

$$\gamma_\ell^{\text{upper}}(k) = \frac{e^{\kappa 1}}{e^{\kappa 1} + k e^{\kappa(1+\alpha^2)} + (L-k)e^{\kappa\alpha^2}} \quad (78)$$

We explicitly need  $\epsilon$  strictly greater than zero, since otherwise there is no information about the count in  $\gamma_\ell$  when it becomes independent of the count  $k$ . Notice, that this time it holds that  $\gamma_\ell$  that correspond to higher values correspond to smaller counts, since a larger count corresponds to a larger denominator, i.e. a smaller  $\gamma_\ell$ . Due to this inverse relationship, for this model, we want that for all counts  $k = 1, \dots, L-1$  that it holds that

$$\gamma_\ell^{\text{upper}}(k+1) < \gamma_\ell^{\text{lower}}(k). \quad (79)$$

This can be achieved by setting the inverse temperature  $\kappa$  accordingly.

In the following we show that there exists a  $\kappa$  which fulfills equation 79 for all  $d' \geq 2$  and  $L > 2$ . Observe that  $\gamma_\ell^{\text{upper}}(2) < \gamma_\ell^{\text{lower}}(1)$  implies the bounds for all other  $k$ . We define the distance or margin as

$$\text{dist}(\kappa) = \gamma_\ell^{\text{lower}}(1) - \gamma_\ell^{\text{upper}}(2). \quad (80)$$

Since at  $\kappa = 0$  both  $\gamma_\ell^{\text{upper}}(2) = \gamma_\ell^{\text{lower}}(1) = 1/(L+1)$ , the distance is zero. However then it becomes impossible to distinguish  $k = 1$  and  $k = 2$ , as they receive the same weight. We therefore need the additional condition that  $\gamma_\ell^{\text{upper}}(2) \neq \gamma_\ell^{\text{lower}}(1)$ . At  $\kappa = 0$ , we observe that this function has a negative derivative, as

$$\frac{\partial}{\partial \kappa} \text{dist}(\kappa)|_{\kappa=0} = \text{sftm}(\kappa z_{\text{lower}})_0 \left( (z_{\text{lower}})_0 - \sum_{i=0}^{L+1} (z_{\text{lower}})_j \text{sftm}(\kappa z_{\text{lower}})_i \right) \quad (81)$$

$$- \text{sftm}(\kappa z_{\text{upper}})_0 \left( (z_{\text{upper}})_0 - \sum_{i=0}^{L+1} (z_{\text{upper}})_j \text{sftm}(\kappa z_{\text{upper}})_i \right) \quad (82)$$

$$= - \left( \frac{1}{L+1} \right)^2 ([ (1+\alpha^2) + (L-1)(1+\alpha^2-\epsilon) ] - [ 2(1+\alpha^2) + (L-2)\alpha^2 ]) \quad (83)$$

$$= - \left( \frac{1}{L+1} \right)^2 [(L-2) - (L-1)\epsilon] \quad (84)$$

$$< 0 \quad (85)$$

where the last bound is met when  $0 < \epsilon < 0.5$  which is fulfilled already for  $d' = 2$  and when  $L > 2$ . As the distance function is continuous, there exists a  $\kappa$  close to zero for which the  $\text{dist}(\kappa) < 0$ . Simultaneously, as  $\kappa \rightarrow \infty$ , we have that

<sup>1</sup>In order to introduce the inverse temperature  $\kappa$  of the softmax in the model, we scale the query matrix. We set  $K = d^{1/4} I_d$ , but  $Q = \kappa d^{1/4} I_d$ .

due to the concentration of the softmax probabilities on the largest entry, which here is  $1 + \alpha^2$ , it holds that as  $\kappa \rightarrow \infty$  we have  $\text{dist}(\kappa) \rightarrow 0$ . At the same time, the function approaches infinity from the positive regime. For large enough  $\kappa$  we have  $\gamma_\ell^{\text{upper}}(2) < \gamma_\ell^{\text{upper}}(1)$ .

When we select the smallest possible  $\kappa > 0$ , we avoid computing functions with large exponential terms. To find the non-trivial root of  $\text{dist}(\kappa)$  numerically, we consider a simplification of equation 80. We define  $u = e^\kappa$ . Then it holds that we can solve

$$\text{dist}(\kappa) = 0 = (L - 1)u^{(1-\epsilon)} - u - (L - 2) \quad (86)$$

numerically for  $\kappa > 0$ . This shows that we can find an explicit construction with 100% accuracy with  $p = 1$  and  $d' > 2$  for the `bos+sftm` when we have

$$d = \lceil \log_2(T + 1) \rceil + 2. \quad (87)$$

For example, for the case of  $L = 10$  and  $T = 32$  this allows for a dimension  $d = 7$  with  $\alpha = 0.01$  (and for  $T = 31$  with the same settings  $d = 6$  suffices).  $\square$

**Remark** ( $d = 4$ ). In principle, it is enough to have some  $\epsilon > 0$  that ensures that overlaps between different token embeddings are strictly less than one. In principle, we can find an arbitrary number of tokens  $T$  that satisfy this condition for just  $d' = 2$ . Take for example the following construction. For  $t = 1, \dots, T$  tokens with  $T$  odd we can design the set of embeddings

$$v_t = \begin{bmatrix} \sqrt{\frac{t}{T}} \\ \sqrt{\frac{T-t}{T}} \end{bmatrix}. \quad (88)$$

Each  $\langle e_t, e_t \rangle = 1$  and for  $t \neq s$  the overlap  $\langle e_t, e_s \rangle \leq \langle e_{(T+1)/2}, e_{(T-1)/2} \rangle = \sqrt{T^2 - 1}/T$ .

This implies that  $\epsilon \rightarrow 0$  as  $T \rightarrow \infty$  at a rate  $1/T$ . Since smaller  $\epsilon$  imply larger values of the temperature to solve equation 86, this might become problematic when this exceeds the accuracy of computations. Previously, for the binary representation construction from equation 70, we had that  $\epsilon$  shrinks at a rate  $\sim 1/\log_2(T)$ . For the intermediate regime between  $\log_T(T) + 1$  and  $\log_2(T)$  dimensions, one can generalize this principle to arbitrary bases, e.g.  $\log_3(T) > 2$ , resulting in a smaller dimension but also less favorable (smaller)  $\epsilon$  – this construction thus comes with a clear trade-off.

### C.3.2. (`DOT+SFTM`; $p = T$ )

*Proof of Proposition 4.5 - dot+sftm.* For this model, the explicit construction is analogous to the previous one. Instead of using  $p = 1$  we use  $p = T$ . The selection of the embeddings is analogous, but instead of a counting direction we read off all the weight directions separately with  $T = d$ . Not having a counting direction also saves the additional two dimensions required for `bos+sftm` with  $p = 1$ . In the feed-forward layer with  $W_1$  the explicit construction considers again  $z_{\ell,t}$  for every token  $t \in \mathcal{T}$ . The selection of the temperature is also analogous, with the exception that one has  $L$  terms in the softmax instead of  $L + 1$ .  $\square$

## D. Data Generation

Every sample  $\mathbf{x} = (x_1, \dots, x_L)$  is generated recursively as follows, starting from size  $K = L$  and alphabet  $\mathcal{T}' = \mathcal{T}$ :

1. Sample an integer  $k$  uniformly from  $[1, \dots, K]$ .
2. Sample a token  $t$  uniformly from  $\mathcal{T}'$ .
3. Set  $x_i = t$  for all  $i = k, \dots, K$ .
4. Set  $\mathcal{T}' = \mathcal{T}' \setminus \{t\}$  and  $K = k$ .
5. If  $K \neq 0$ , repeat from 1.
6. Set  $\mathbf{x} = \text{shuffle}(\mathbf{x})$ .

In contrast to sampling the elements of each sequence uniformly at random from the alphabet, this simple strategy enables us to better control the distribution of counts in the training dataset.

## E. Additional Experiments

### E.1. Best Accuracy

In Fig. 10, we show the best reached accuracy during training over the five sample runs. This gives insights into the feasibility of implementing a counting solution for a given combination of parameters  $T, d, p$  of a model.

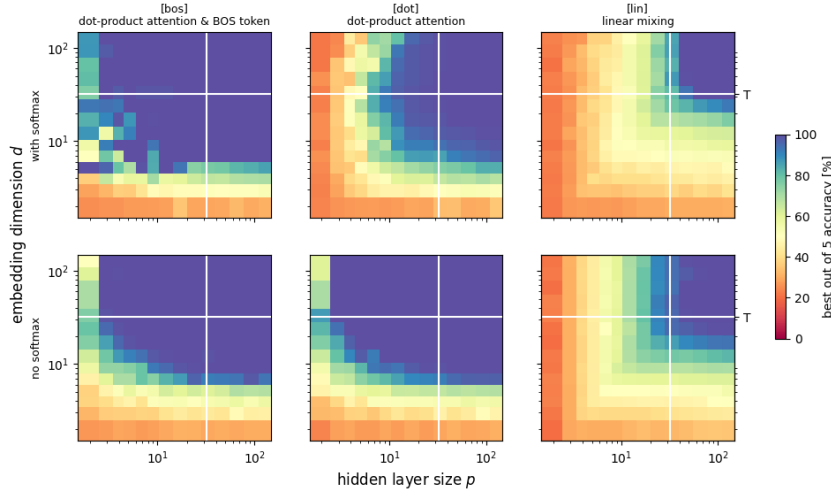


Figure 10. **Best accuracy during training.** Experiments from Fig. 1 ( $T = 32$ ), we show only the *best accuracy* during training reached from the 5 randomly initialized runs per model/hyperparameter configuration.

### E.2. Variability

In Fig. 11 we explore the influence of initialization on the performance via the variability of the final accuracy for several runs. Especially in the  $p, d < T$  regime where `bos+softmax` is able to reach an accuracy relatively close to 100%, the variability of the accuracies resulting from different initializations is quite large.

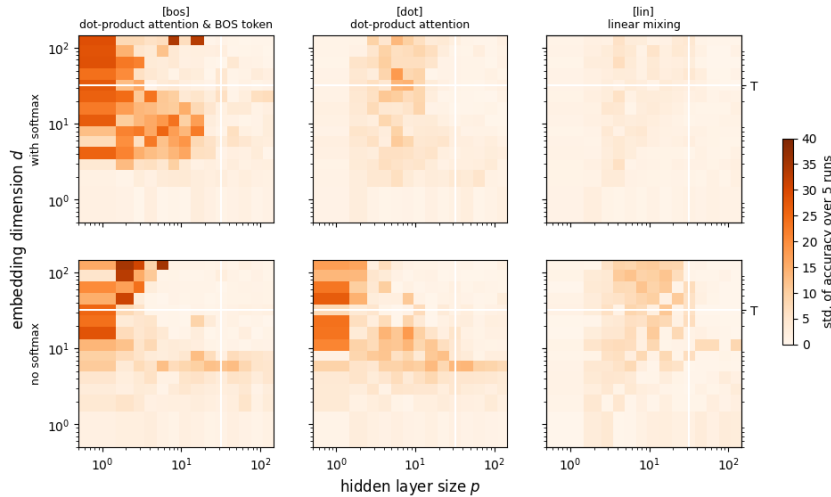


Figure 11. **Variability to randomness in initial conditions and optimization.** Experiments from Fig. 1 with  $T = 32$ , standard deviation of the accuracy reached after training from the 5 randomly initialized runs per model/hyperparameter configuration.

### E.3. Model with Random but Fixed Embeddings

In Fig. 12, we repeat the experiments of Fig. 1, but for embeddings that are frozen throughout training (also 5 runs). In the regime  $d < T$  where there is no mutual orthogonality possible, the random embeddings result in worse performance than the learned ones. Especially for `bos+softmax`, learning the embeddings increases the performance strongly in some regimes. This indicated that the models indeed learn adapted embeddings here.

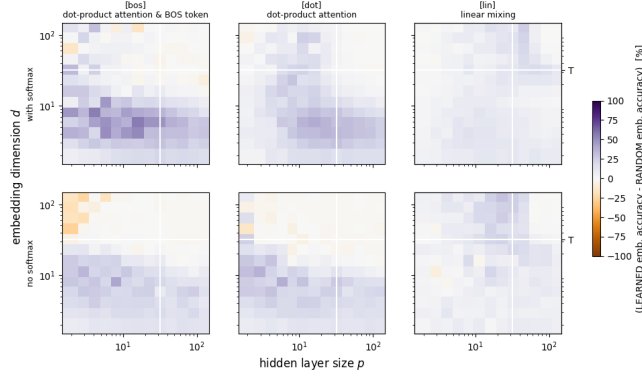


Figure 12. **Random but fixed embeddings (not learned).** The difference between learned and random embeddings for  $T = 32$ . Orange indicates that the random embeddings perform better on average. Purple indicates that the learned embeddings perform better on average. Experimental settings as in Fig. 1.

### E.4. BOS mixing token

In Fig. 3 in the main, we describe how the  $t_{\text{BOS}}$  is the main predictor for the count. Here, we provide more evidence by showing how the count predictions for mixed tokens  $\bar{x}'$  output by the feature transform  $f$  are invariant to the type of other token present in the mixed token. The results for four different tokens are shown in Fig. 13.

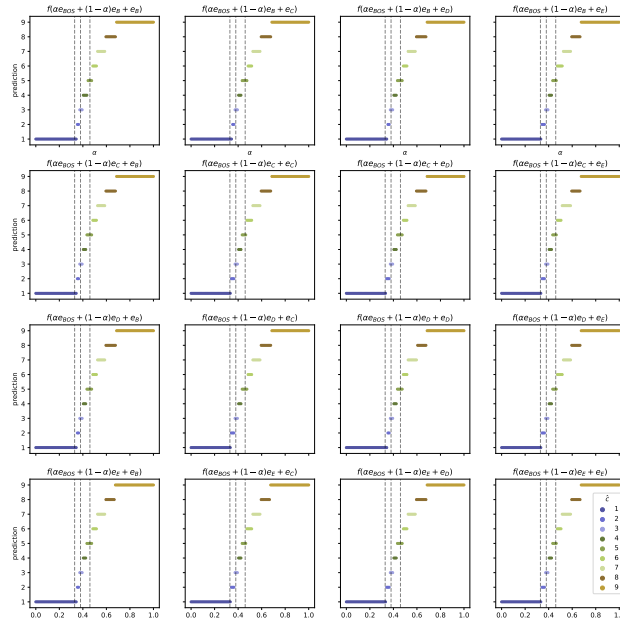


Figure 13. **Visualization of feature transform.** For the same model as in Fig. 3, we vary the inputs to the feature transformation  $f$  to show it is independent on the precise input sequence, but only depends on the prevalence of  $t_{\text{BOS}}$ . We vary the inputs between the learned tokens  $[B, C, D, E]$ .

### E.5. Singular Value Decomposition of $W_1$

In Fig. 14 we show the distribution of singular values of  $W_1$  for several runs of the model to investigate whether models that are capable of both IC and RC are implementing the more memory heavy IC or the same solution that they can find for  $p = 1$  with RC.

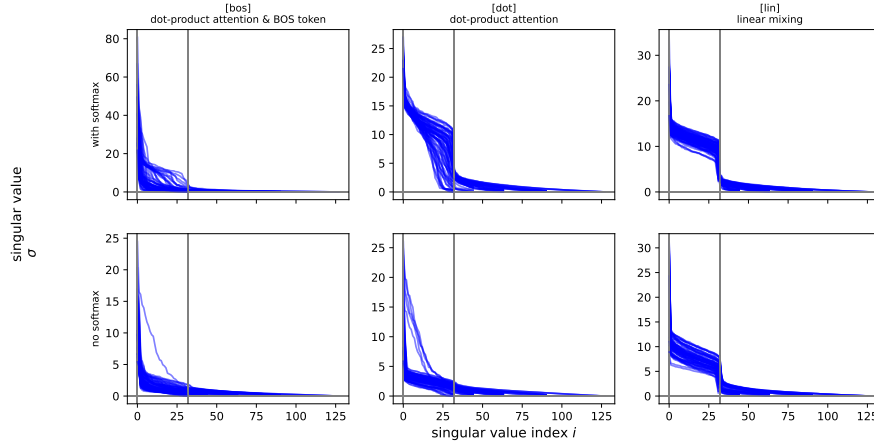


Figure 14. **Singular values of  $W_1$ .** We show the results for all models from Fig. 1 with  $T = 32$ , where  $p, d \geq T$  and the accuracy is at least 99%. Some qualitative differences are visible for `bos` and `dot`.

### E.6. Varying the number of tokens in the alphabet

In Fig. 15 and 16 we change the number of tokens that may occur in the alphabet from  $T = 32$  in the experiments from the main text to  $T = 15$  and  $T = 64$  respectively. The sequence length is kept fixed at  $L = 10$ .

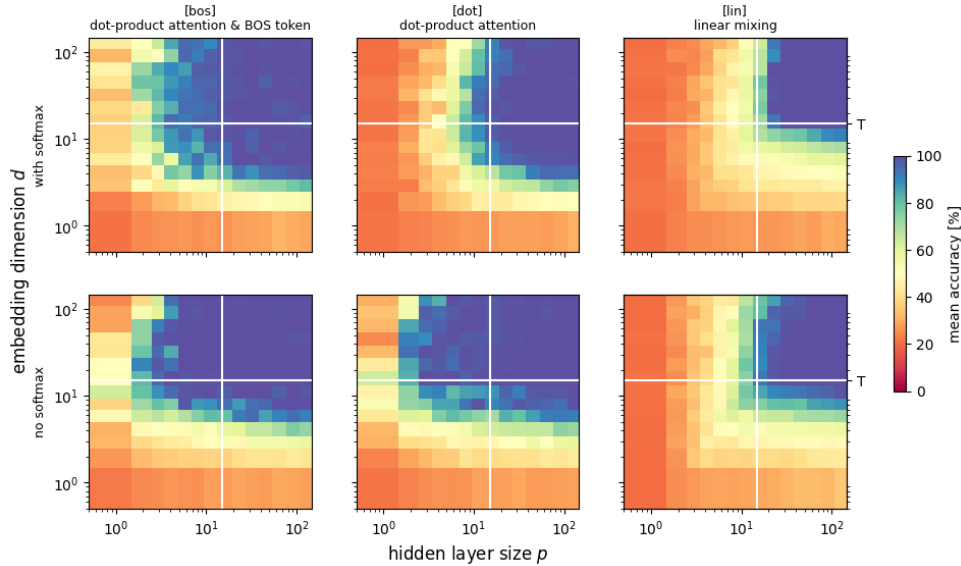


Figure 15. **Alphabet size  $T = 15$ .** Experiments as in Fig. 1, except that  $T = 15$  instead of  $T = 32$  was used. The sequence length is fixed to  $L = 10$ .



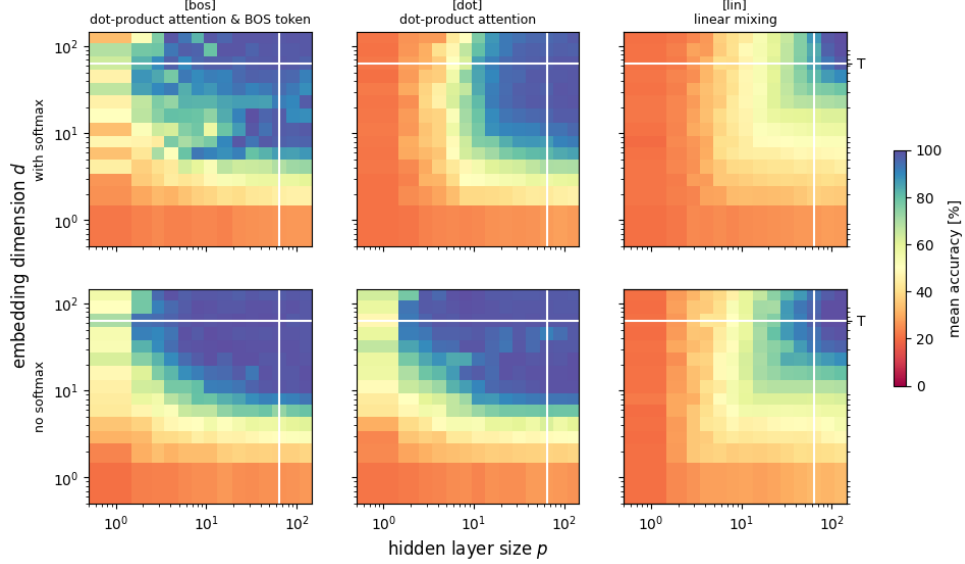


Figure 16. **Alphabet size**  $T = 64$ . Experiments as in Fig. 1, except that  $T = 64$  instead of  $T = 32$  was used. The sequence length is fixed to  $L = 10$ .

### E.7. Varying the length of the sequence

In Fig. 17, 18 and 19 we change the length of the sequence to  $L = 5, 15, 30$  respectively, while keeping the number of tokens fixed to  $T = 32$ . The sequence length is kept fixed at  $L = 10$ . Fig. 17 shows that the learned model can make use of the discrete nature of the task especially for `lin` (without softmax) and for the dot-product style models with the softmax. While  $L = 15$  in Fig. 18 is only slightly worse than the case of  $L = 10$  that was used as the example in the main text, the performance  $L = 30$  in Fig. 19 is very much deteriorated. Even in the cases without the softmax performance is far from what would be possible theoretically. We hypothesise that the optimization with the increased number of possible inputs becomes much harder. Interestingly, Fig. 19 even seems to show a double descent behavior along the embedding dimension axis.

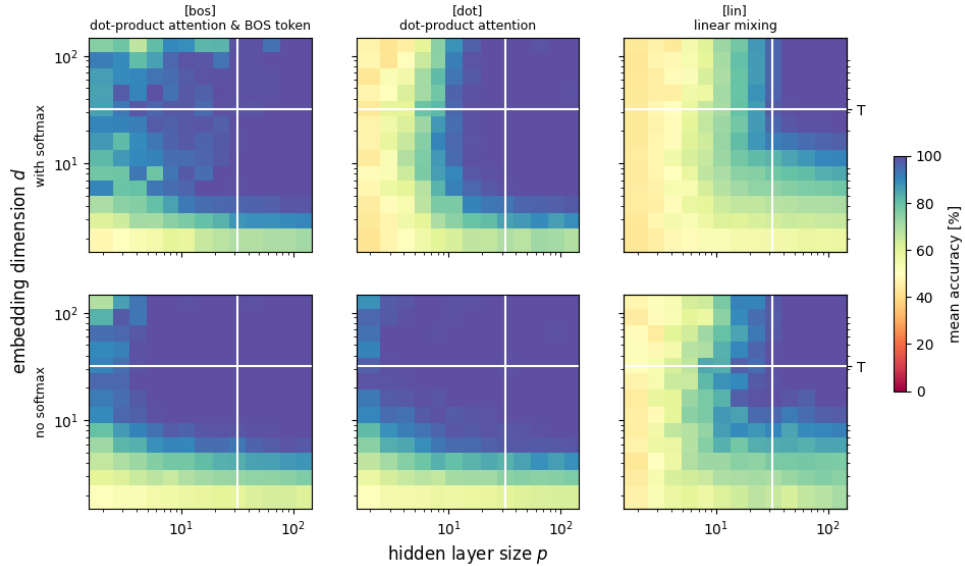


Figure 17. **Sequence length**  $L = 5$ . Experiments as in Fig. 1, except that  $L = 5$  instead of  $L = 10$  was used. The alphabet size is fixed to  $T = 32$ .

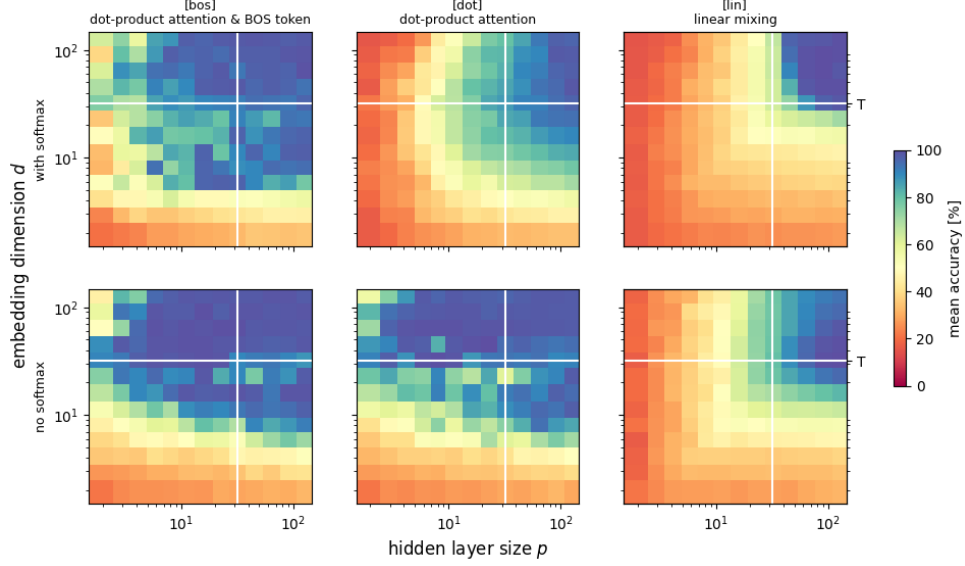


Figure 18. **Sequence length**  $L = 15$ . Experiments as in Fig. 1, except that  $L = 15$  instead of  $L = 10$  was used. The alphabet size is fixed to  $T = 32$ .

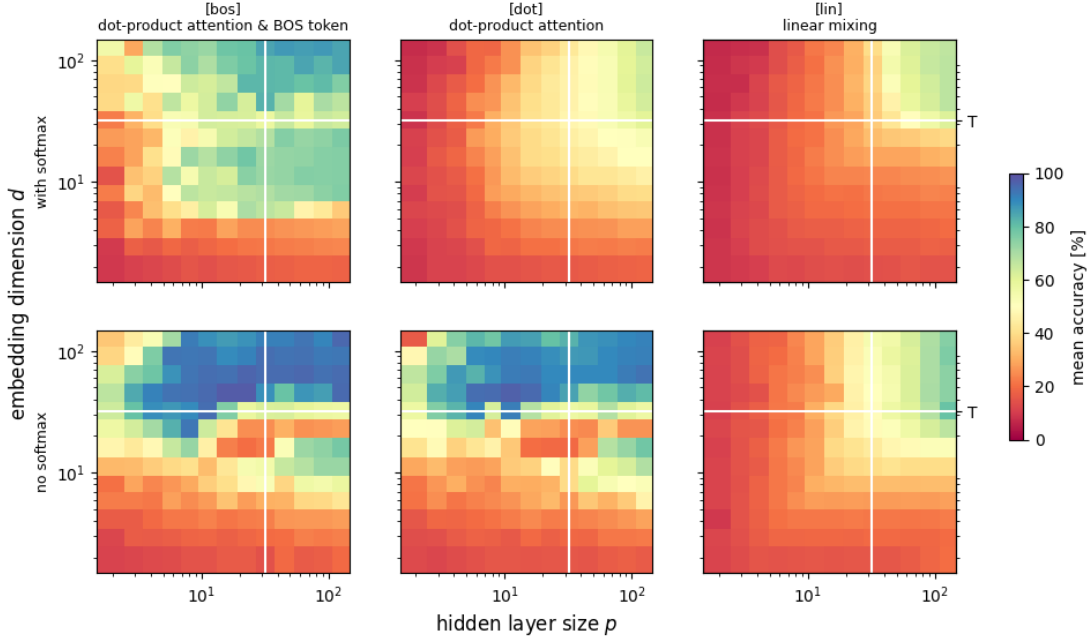


Figure 19. **Sequence length**  $L = 30$ . Experiments as in Fig. 1, except that  $L = 30$  instead of  $L = 10$  was used. The alphabet size is fixed to  $T = 32$ .

### E.8. Models with two layers

In this section, we look at the case where we have models that have an extra layer, i.e. instead of the logit output layer after the feed-forward part, we add another layer with the same dimensionality  $d$  as the previous layer – the same mixing and the same hidden layer size – to then lead into the classification. The parameters are not shared between the layers. Note that this model does not have an extra residual in the MLP layers.

We train the model in the same setting as in the main and report the results for the different architectures, this time with 2 layers, in Fig. 20. To compare more easily with the previous set-up, we show the difference between the single and double layer case in Fig. 21. Remarkably, the general picture does not seem to change significantly. Indeed, the two layer model is generally better, extending the range where perfect models can be found slightly, but the general trend remains. Given this coarse grained experiment we hypothesize, that the extra layer aids the optimization process, and improves robustness in the regions where and the softmax is used to disentangle non-orthogonal embeddings. More generally, these results are not as comprehensive as our previous results as they are not supported theoretically beyond a single layer. They warrant more detailed in further work with more layers and realistic settings.

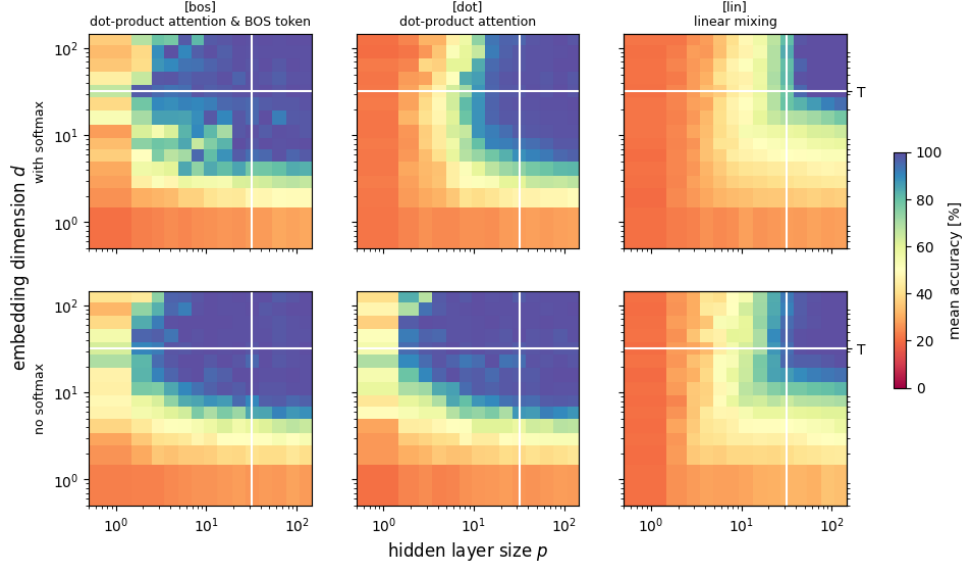


Figure 20. **Two-layer performance.** Experiments as in Fig. 1 with  $T = 32$ ,  $L = 10$  but for models where the attention block is repeated twice as described in Section E.8.

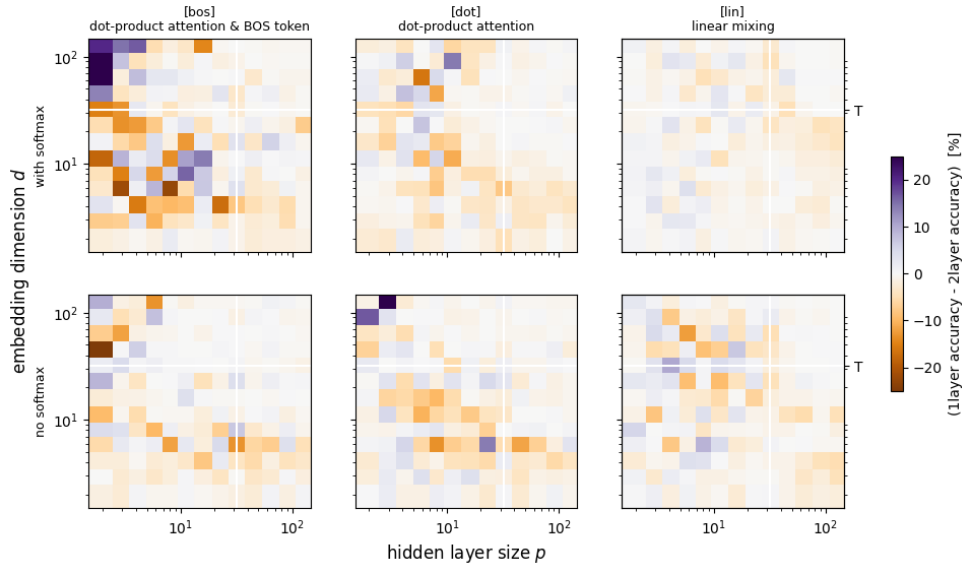


Figure 21. **Comparison between one and two-layer models.** Difference between the accuracy of a single and two layer attention model, for different mixing layers and hyperparameter setups. Experiments as in Fig. 1 for a single layer attention model, and as in Fig. 20 for the two layer model.