

# Large Language Models are Locally Linear Mappings

Anonymous authors

Paper under double-blind review

## Abstract

Despite significant progress in transformer interpretability, an understanding of the computational mechanisms of large language models (LLMs) remains a fundamental challenge. We demonstrate that the inference operation of LLMs can be mapped to an equivalent linear system that nearly exactly reconstructs the predicted output embedding for a given input sequence. Extending techniques from image diffusion models that exhibit local or piecewise linearity, we strategically detach components of the gradient computation with respect to an input sequence for a next-token prediction such that the Jacobian of the model reproduces the output with one linear operation per input token. We demonstrate this approach across models, including Qwen 3, Gemma 3, Llama 3, Phi 4, Mistral Ministral and OLMo 2, up to Llama 3.3 70B Q4. With the singular value decomposition of the detached Jacobian, we show that these LLMs operate in extremely low-dimensional subspaces where the largest singular vectors decode to distinct concepts related to possible output tokens. We examine the equivalent linear operation of each successive layer (and its attention and multilayer perceptron components) and observe the emergence of semantic concepts. We demonstrate that the detached Jacobian of middle layer representations can be used as a steering operator to insert semantic concepts into unrelated text, which could be useful for improving safety and decreasing bias. Despite their expressive power and global nonlinearity, modern LLMs can be interpreted through locally linear decompositions that provide insights into their internal representations and reveal interpretable semantic structures in the next-token prediction process.

## 1 Introduction

The transformer decoder has become the architecture of choice for large language models (Vaswani et al., 2017) and efforts toward a conceptual understanding of its mechanisms are ongoing. Significant insights include sparse autoencoders for conceptual activations in LLMs (Bricken et al., 2023), “white-box” alternative architectures (Yu et al., 2023), minimally sufficient architectures (He & Hofmann, 2023) and analytic results on generalization (Cowsik et al., 2024). While transformers are complex globally nonlinear functions of their input, we demonstrate how to compute an equivalent linear system that nearly exactly reconstructs the predicted output embedding for a given input sequence.

Our approach builds on two previous results: Elhage et al. (2021) found that attention-only networks with basic language generation abilities transfer semantic information across the network with interpretable circuits, including the “induction head”, and Kadkhodaie et al. (2023) showed that powerful image denoising diffusion models can be made exactly locally or piecewise linear through several architectural constraints and interpreted as low-dimensional adaptive linear filters.

We extend these ideas by demonstrating local linearity for modern LLMs with gated linear activations. For many open-weight LLMs, the gradient operation with respect to the input can be manipulated at inference such that the output prediction is unchanged and also almost exactly locally linear. This numerical Jacobian computation captures the complete forward operation of the model, including activation functions and attention, although it is only valid for that particular input sequence (more “pointwise” linear than piecewise linear).

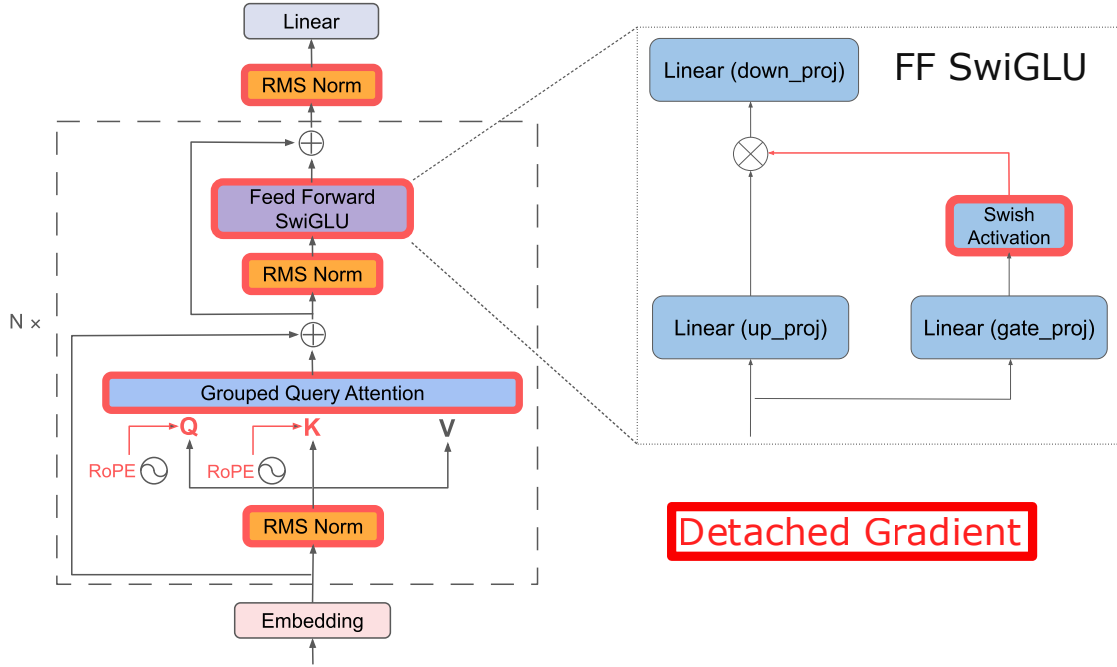


Figure 1: A schematic of the Llama 3 transformer decoder (Grattafiori et al., 2024; Nvidia, 2024). The gradient *detach* operations for components outlined in red effectively freeze the nonlinear activations for a given input sequence, creating a linear path for the gradient with respect to the input embedding vectors, but do not change the output. The output embedding prediction can be nearly exactly mapped to a linear system by the Jacobian autograd operation. The feedforward module with a *SwiGLU* activation function is shown in expanded form to demonstrate how the nonlinear *Swish* term can be detached from the gradient to form a linear path, achieving local linearity for a given input. The *RMSNorm* layers and softmax attention blocks also must be detached from the gradient.

This approach allows us to analyze entire models, from input embeddings to predicted output embedding, as equivalent linear systems for a particular input sequence. By examining the singular value decomposition (SVD) of the equivalent linear system, we can measure the local dimensionality of the learned manifolds involved in next-token prediction, and we can decode the singular vectors into output tokens. This analysis can also be done layer by layer, or for individual attention and multilayer perceptron (MLP) modules, in order to observe how these models compose next-token predictions.

We demonstrate local linearity in model families including Qwen 3, Gemma 3, Llama 3, Phi 4, Mistral, Deepseek R1 0528 Qwen 3 8B and OLMo 2, at a range of sizes up to Llama 3.3 70B Q4. This approach does not require further model training, offering a new path to interpreting a wide range of open-weight LLMs at a local level that could serve as a complement to other powerful interpretability methods.

## 2 Method

### 2.1 The Jacobian of a deep ReLU Network

Mohan et al. (2019) observed that deep *ReLU* networks for image denoising which utilize zero-bias linear layers are “adaptive linear” functions due to their homogeneity of order 1 at a given fixed input, which enables interpretation as a nearly-exact linear system. Given the homogeneity at a fixed input, the network’s output can be nearly exactly reproduced by numerically computing the Jacobian matrix of the network at a particular input image  $\mathbf{x}_{\text{im}}^*$  and multiplying it by  $\mathbf{x}_{\text{im}}^*$ .

$$\mathbf{y}_{\text{im}}^* = \mathbf{J}(\mathbf{x}_{\text{im}}^*) \cdot \mathbf{x}_{\text{im}}^* \quad (1)$$

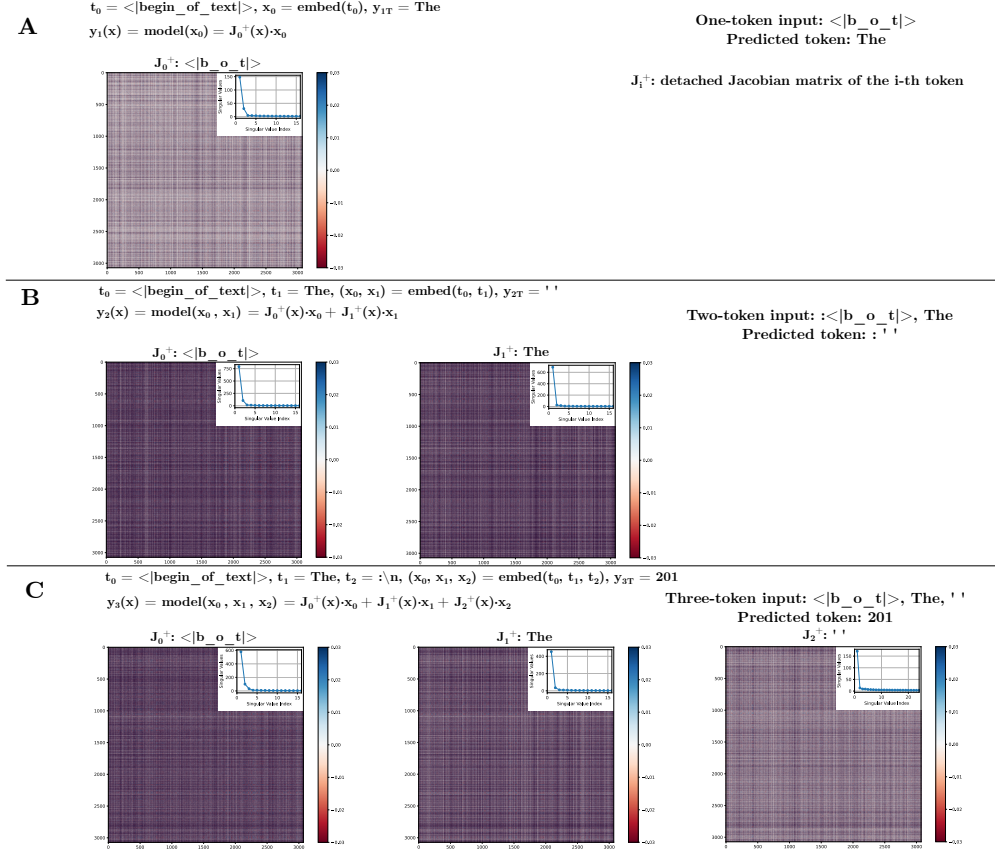


Figure 2: An overview of next-token prediction in the Llama 3.2 3B transformer decoder and decomposition of the predicted embedding vector computation using the detached Jacobian. Generating three tokens with only  $\langle \text{BoT} \rangle$  as input produces “The 201”. For each prediction, each input token  $t_i$  is mapped to an embedding vector  $x_i$ , and the network generates the embedding of a next token. The phrase turns out to be “The 2019-2020 season”. The detached Jacobian  $J^+(\mathbf{x})$  of the predicted output embedding with respect to the input embeddings is composed of a matrix corresponding to each input vector. Each detached Jacobian matrix  $J_i^+(\mathbf{x})$  is a function of the entire input sequence but operates only on its corresponding input embedding vector. The matrices tend to be extremely low rank, shown in the inset figures, and the matrix  $J_0^+$  varies across A), B) and C) above because the input sequences differ. Since the detached Jacobian captures the entirety of the model operation in a linear system (numerically, for a given input sequence), tools like the SVD can be used to interpret the model and its sub-components.

Due to the global nonlinearity of the network, the Jacobian must usually be computed again at every input of interest. The Jacobian may be the same for similar inputs in the same piecewise region of the response (Balestriero & Baraniuk, 2021).

## 2.2 The Jacobian of a transformer decoder

Many open weight LLMs also use linear layers with zero bias, like the *ReLU* network of Mohan et al. (2019). A transformer decoder predicts an output token embedding  $\mathbf{y}$  given a sequence of  $k$  input tokens  $\mathbf{t} = (t_0, t_1, \dots, t_k)$  mapped to input embedding vectors  $\mathbf{x} = (x_0, x_1, \dots, x_k)$ , where  $\mathbf{t}^*$  and  $\mathbf{x}^*$  represent a particular sequence. The output embedding prediction is a nonlinear function of the input embedding vectors  $x_0, x_1, \dots, x_k$ , as LLMs utilize nonlinear gated activation functions for layer outputs (*SwiGLU* for Llama 3, *GELU* for Gemma 3 and *Swish* for Qwen 3) as well as normalization and softmax attention blocks.

Gated activations like  $\text{Swish}(\mathbf{x}) = \mathbf{x} \cdot \text{sigmoid}(\mathbf{x})$ , with a linear term and a nonlinear term, are also an “adaptive” linear function or, more generally, an adaptive homogeneous function of order 1 (Mohan et al., 2019). If the  $\text{sigmoid}(\mathbf{x})$  term that gives rise to the nonlinearity is frozen for a specific numerical input, e.g. an embedding vector  $\mathbf{x}_0^*$  (Elhage et al., 2021) (or equivalently detached from the computational graph with respect to the input), then we have a linear function valid only at  $\mathbf{x}_0^*$  where (1) holds and we can numerically compute a Jacobian that carries out  $\text{Swish}(\mathbf{x}_0^*)$  as a linear operation.

Below we show that computing the Jacobian after effectively substituting specific values for the nonlinear terms also works for other gated activation functions, zero-bias *RMSNorm* layers and softmax attention blocks. We further demonstrate that for a given input sequence the entire transformer decoder is an adaptive homogeneous function of order 1 where we can apply necessary gradient detachments and numerically compute a linear system that nearly exactly reproduces the transformer output embedding  $\mathbf{y}^*$ .

The Jacobian  $\mathbf{J}(\mathbf{x})$  of a transformer is the set of matrices generated by taking the partial derivative of the decoder inference function  $\mathbf{y}(\mathbf{x}) = f(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k)$ , with respect to each element of each  $\mathbf{x}_i$  (where  $\mathbf{x}_i$  for Llama 3.2 3B has length 3072, for example, and therefore the Jacobian matrix for each embedding vector is a square matrix of this size).

We introduce a “detached” Jacobian  $\mathbf{J}^+$ , which is a set of matrices that captures the full nonlinear forward computation for a particular input sequence  $\mathbf{x}^*$  as a linear system. The detached Jacobian is the numerical Jacobian of the LLM forward operation when its gradient includes a specific set of *detach()* operations for the nonlinear terms in the normalization, activation and attention operations. Fig. 2 shows how each matrix of the detached Jacobian operates on its corresponding input embedding vector to provide a nearly-exact reconstruction of the LLM forward operation (with a relative error of  $10^{-6}$  (the standard deviation of the reconstruction error divided by the standard deviation of the output embedding), see Fig. 3).

$$\mathbf{y}^* = \sum_{i=0}^k \mathbf{J}_i^+(\mathbf{x}^*) \cdot \mathbf{x}_i^* \quad (2)$$

While the traditional Jacobian  $\mathbf{J}$  for a particular input sequence  $\mathbf{x}^*$  is a locally linear approximation of the nonlinear LLM forward operation, it does not generate an exact reconstruction at  $\mathbf{x}^*$  since the transformer function is not linear. The detached Jacobian  $\mathbf{J}^+$  evaluated at  $\mathbf{x}^*$  is the result of an alternative gradient path through the same network which is linear for the input  $\mathbf{x}^*$ . The detached Jacobian  $\mathbf{J}^+$  only generates a near-exact reconstruction at  $\mathbf{x}^*$  and not in the local neighborhood due to the strong nonlinearity of the decoder inference function. The detached Jacobian matrices differ for every input sequence and must be computed numerically for every sequence.

## 2.3 Nonlinear layers as locally linear operations

In order to achieve local linearity, modifications must be made to the gradient computations of the *RMSNorm* operation, the activation function (*SwiGLU* in Llama 3.2) and the *softmax* term in the attention block output.

### 2.3.1 Normalization

Normalization layers like *LayerNorm* (Xu et al., 2019) or *RMSNorm* (Zhang & Sennrich, 2019) are nonlinear with respect to their input because they include division by the square root of the variance of the input.

$$\text{norm}(\mathbf{x}) = \frac{\mathbf{x}}{\sqrt{\text{var}(\mathbf{x})}} \quad (3)$$

Mohan et al. (2019) devised a novel bias-free batch-norm layer which disconnects the variance term from the network’s computational graph. Their batch-norm layer returns the same values as the standard batch-norm layer, but it is locally linear at inference as the nonlinear operation is removed from the gradient computation.



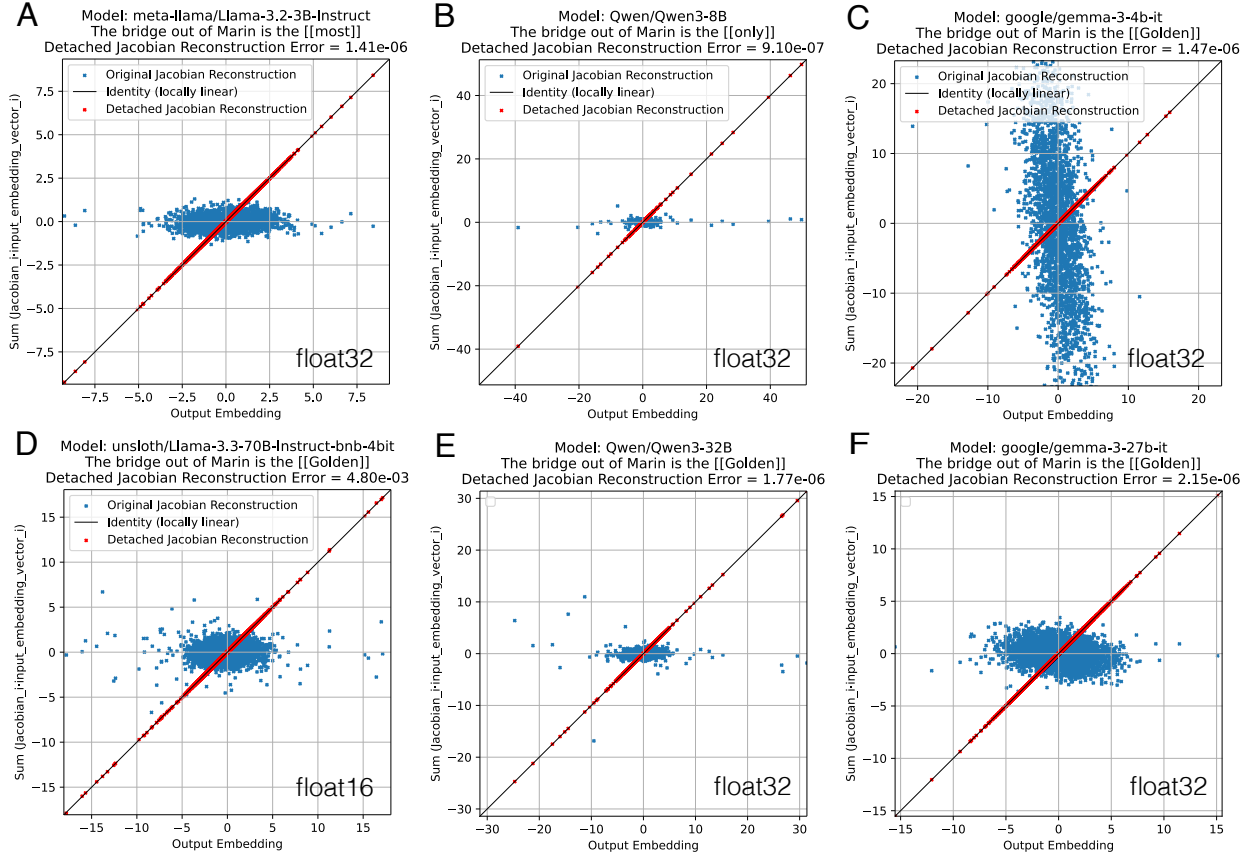


Figure 3: For the input sequence “The bridge out of Marin is the”, the elements of the predicted output embedding vector of the model compared to the elements from the Jacobian reconstruction for both the original Jacobian (blue points) and detached Jacobian operations (red points), shown for a range of model families and sizes. Note that the detached Jacobian reconstructions nearly exactly match the predicted embedding in each case, with relative error (the standard deviation of the reconstruction error divided by the standard deviation of the output embedding) on the order of  $10^{-3}$  for *float16* precision and  $10^{-6}$  for *float32* precision. Llama 3 models are shown in A) and D), Qwen 3 models in B) and E), and Gemma 3 models in C) and F). These plots demonstrate the near-exact reconstruction of the detached Jacobian and therefore the local linearity of the mapping.

This is also similar to the “freezing” of nonlinear terms in attention-only transformers from Elhage et al. (2021).

We make a similar change for Llama 3.2 3B by altering how the gradient with respect to the input is computed at inference for *RMSNorm*. This is accomplished by substituting the value for the input vector  $\mathbf{x}^*$  for only the variance term as in (4). In PyTorch, this is accomplished by cloning and detaching the  $\mathbf{x}$  tensor within the variance operation, so its value will be treated as a constant. The gradient operation is still tracked for  $\mathbf{x}$  in the numerator, so that term will be treated as a variable by *autograd.functional.jacobian*. The gradient of the function is then computed at  $\mathbf{x}^*$  (we assume for simplicity a sequence of length 1).

$$\text{norm}_{LL}(\mathbf{x}) = \frac{\mathbf{x}}{\sqrt{\text{var}(\mathbf{x}^*)}} \quad (4)$$

We define the detached Jacobian as follows:

$$\mathbf{J}_{\mathbf{n}_{LL}}^+ = \left[ \frac{\partial}{\partial \mathbf{x}} \text{norm}_{LL}(\mathbf{x}) \right]_{\mathbf{x}=\mathbf{x}^*} \quad (5)$$

We can rewrite the locally linear *RMSNorm* as follows:

$$\text{norm}_{LL}(\mathbf{x}^*) = \mathbf{J}_{\mathbf{n}_{LL}}^+(\mathbf{x}^*) \cdot \mathbf{x}^* \quad (6)$$

At inference, we now have a locally linear *RMSNorm* whose output is numerically identical to the one used in training. However, when we take the gradient with respect to the input vector  $\mathbf{x}$  in *eval* mode, the numerical output is the detached Jacobian matrix  $\mathbf{J}_{\mathbf{n}_{LL}}^+$ , which we can use to reconstruct the normalization output exactly as a linear system; no higher-order terms are needed.

The goal is to apply this same approach for other nonlinear functions in the decoder such that the entire computation from the input embedding vectors to the predicted output is locally linear, and we can compute and interpret the set of detached Jacobian matrices.

### 2.3.2 Activation functions

While Mohan et al. (2019) relied on *ReLU* activation functions, which do not require any changes to achieve local linearity, Llama 3.2 uses *SwiGLU* (Shazeer, 2020), Gemma 3 uses approximate *GELU* (Hendrycks & Gimpel, 2016) and Qwen 3 uses *Swish* for activation functions. Fortunately, there is a linear  $\mathbf{x}$  term in each of these, and the gradients can be cloned and detached from the nonlinear terms. This manipulation produces a locally linear *SwiGLU* layer with respect to the input  $\mathbf{x}$ . Below,  $\text{Swish}(\mathbf{x}) = \mathbf{x} \cdot \text{sigmoid}(\mathbf{x})$  and  $\otimes$  is element-wise multiplication.

$$\text{SwiGLU}(\mathbf{x}) = \text{Swish}(\mathbf{W}\mathbf{x}) \otimes (\mathbf{Z}\mathbf{x}) \quad (7)$$

$$\text{SwiGLU}_{LL}(\mathbf{x}) = [\text{Swish}(\mathbf{W}\mathbf{x})]_{\mathbf{x}=\mathbf{x}^*} \otimes (\mathbf{Z}\mathbf{x}) \quad (8)$$

$$\text{SwiGLU}_{LL}(\mathbf{x}^*) = \left( \left[ \frac{\partial}{\partial \mathbf{x}} \text{SwiGLU}_{LL}(x) \right]_{\mathbf{x}=\mathbf{x}^*} \right) \cdot \mathbf{x}^* \quad (9)$$

$$\text{SwiGLU}_{LL}(\mathbf{x}^*) = \mathbf{J}_{\text{SwiGLU}_{LL}}^+(\mathbf{x}^*) \cdot \mathbf{x}^* \quad (10)$$

Detaching the gradient from the *Swish* output thus allows for a locally linear form of *SwiGLU* at inference. A similar procedure may be carried out for *GELU* with Gemma 3 (see supplement).

### 2.3.3 Attention

The *softmax* operation at the output of the attention block can also be detached, with the linear relationship preserved through the subsequent multiplication with  $\mathbf{V}$ , which is a linear function of  $\mathbf{x}$ . Below,  $\mathbf{Q} = \mathbf{W}_Q\mathbf{x}$ ,  $\mathbf{K} = \mathbf{W}_K\mathbf{x}$  and  $\mathbf{V} = \mathbf{W}_V\mathbf{x}$ .

$$\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right) \cdot \mathbf{V} \quad (11)$$

$$\text{Attn}_{LL}(\mathbf{x}) = \left[ \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right) \right]_{\mathbf{Q}=\mathbf{Q}^*, \mathbf{K}=\mathbf{K}^*} \cdot \mathbf{W}_V\mathbf{x} \quad (12)$$

$$\text{Attn}_{LL}(\mathbf{x}^*) = \left( \left[ \frac{\partial}{\partial \mathbf{x}} \text{Attn}_{LL}(x) \right]_{\mathbf{x}=\mathbf{x}^*} \right) \cdot \mathbf{x}^* \quad (13)$$

$$Attn_{LL}(\mathbf{x}^*) = \mathbf{J}_{\text{Attn}_{LL}}^+(\mathbf{x}^*) \cdot \mathbf{x}^* \quad (14)$$

The linear  $\mathbf{x}$  term within  $\mathbf{V}$  makes it possible for the attention block to be locally linear at inference, as the gradient for the *softmax* output is detached.

### 2.3.4 The Transformer Decoder

With the the above gradient detachments for the normalization layers, activation functions and attention blocks, the transformer decoder network is locally linear with respect to  $\mathbf{x}^*$  when evaluated at  $\mathbf{x}^*$  (shown here for length  $k$ ).

$$\mathbf{y}^* = \sum_{i=0}^k \mathbf{J}_i^+(\mathbf{x}^*) \cdot \mathbf{x}_i^* \quad (15)$$

The output of the network incorporating the above gradient detachments is unchanged from the original architecture.

## 3 Results

### 3.1 Local linearity of the predicted output

In order to validate whether the detached Jacobian achieves local linearity, we can compare the predicted output embedding vector for a given input token sequence to the Jacobian reconstruction of the output.

Fig. 3 examines the original output and the output when incorporating the appropriate gradient detachments for local linearity for two sizes of Llama 3, Qwen 3 and Gemma 3. The reconstruction of the output embedding with the detached Jacobian matrices falls close to the identity line when compared with the output embedding, which is evidence for local linearity. The same comparison with the original Jacobian does not fall on the identity line. The standard deviation of the difference for the detached Jacobian divided by the standard deviation of the output embedding vector is on the order of only  $10^{-6}$ , which quantifies how exactly the detached Jacobian reproduces the output). These plots therefore validate the local linearity of Llama 3, Qwen 3 and Gemma 3 with the appropriate gradient detachments for a particular input.

The numerical computation of the top  $k$  singular vectors of the Jacobian takes on the order of 10 seconds for an input sequence of 8 tokens for Llama 3.2 3B on a T4 GPU. At the other end of the spectrum, the singular vectors for the same sequence with Llama 3.3 70B Q4 on an H100 GPU takes more than a minute. (A method for computing the top singular vectors without forming the full matrix with Lanczos iteration is under development.)

### 3.2 Single-unit feature selectivity and invariance

Since the detached Jacobian applied to the input embedding reproduces the predicted output embedding vector, and the elements of the predicted output embedding vector are the units of the last transformer layer, the rows of the detached Jacobian matrices represent the input features to which the last layer units are selective and invariant for that particular input sequence (Kadkhodaie et al., 2023; Mohan et al., 2019).

The activation of a particular unit in the last layer is determined by the inner product of a row of the detached Jacobian and the input embedding vector. We can sort by the magnitude of row norms, then map the largest-magnitude rows of the detached Jacobian back to the input embedding space (by finding the nearest-neighbor embedding vectors from the tokenizer) to determine the tokens that cause each unit to be strongly positive or negative. We can see in Fig. 4 (in appendix A below) that the units respond strongly to the words of the prompt, including “bridge”, “Marin” and “is”. Decoding of the rows of the detached Jacobian for each token as well as the distribution of activations for this sequence is shown in Fig. 4A.

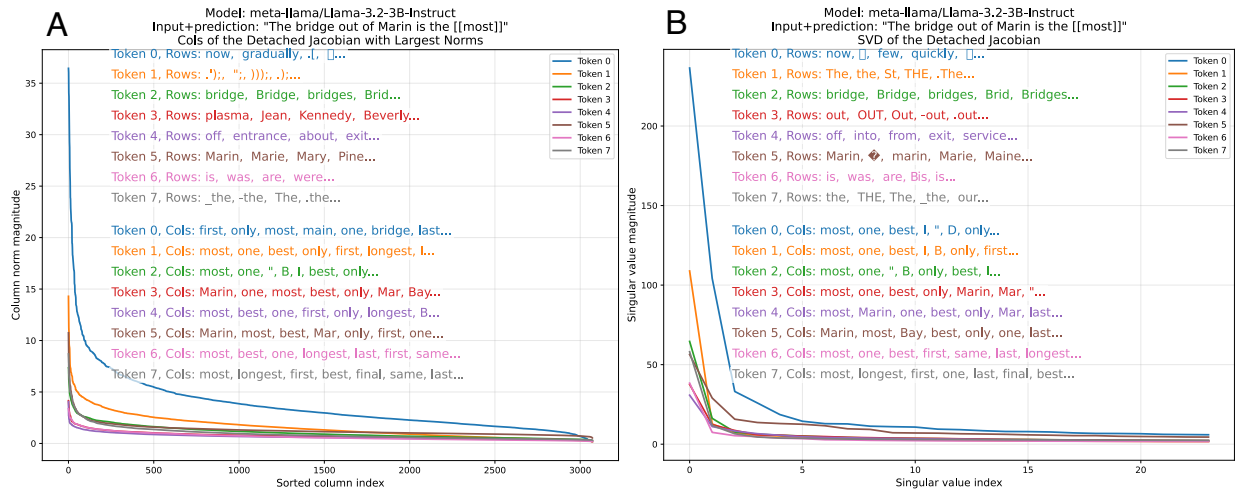


Figure 4: Given the sequence “The bridge out of Marin is the”, the most likely prediction is “most” for Llama 3.2 3B. The detached Jacobian matrices for each token represent a near-exact local linearization of the predicted output embedding. A) We show the features which drive large responses in single units in the last decoder layer, which are the rows of the detached Jacobian with the largest norm values, and decode each of those into the most likely input embedding token. The blue list of words at the top are the ordered decoded “feature” input tokens from the largest rows of the detached Jacobian matrix for the beginning of sequence token, and the different colors show the decoded feature tokens for the other input tokens. A similar operation is carried out for columns of the largest norm values, which are decoded to the output token space. Note that the activation distribution of column magnitudes is fairly sparse, with only a few features driving the response. B) We take the singular value decomposition of the detached Jacobian matrix corresponding to each input token, which summarizes the modes driving the response, and decode the left and right singular vectors  $U$  and  $V$  to output and input embeddings, shown in colors. The singular value spectrum is extremely low rank, and decoding the  $U$  singular vectors returns candidate output tokens: “most” and “one” appear frequently. Decoding the  $V$  singular vectors returns variants of the input tokens like “bridge”, “Marin” and “is”, as well as others that are not clearly related to the input sequence.

The columns of the Jacobian can also be decoded to the output token space, and these turn out to be tokens that could be predicted, which include words like “most” or “first”, which could be acceptable outputs.

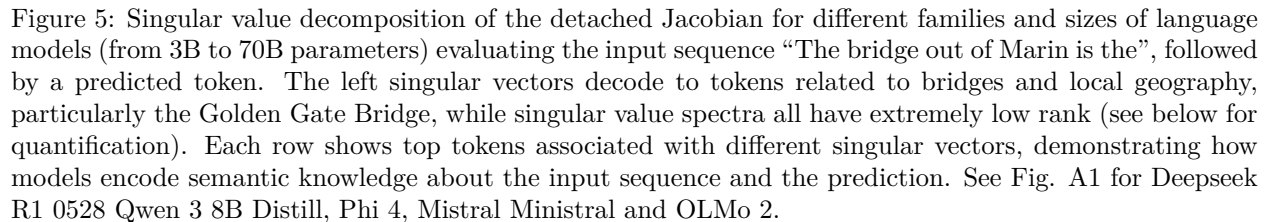
### 3.3 Singular vectors of the detached Jacobian

An alternative approach is to look at the singular value decomposition of the detached Jacobian  $\mathbf{J}_i^+ = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ , following Mohan et al. (2019). Since the detached Jacobian represents the forward computation, and there are no higher order terms, the fact that the SVD is very low rank shows the entire forward computation can be approximated with only a few singular vectors operating on the input embeddings.

In Fig. 4B, the singular vectors are decoded for three different models of two sizes each, from  $3B$  to  $70B$  parameters. The right singular vectors  $V$  are decoded to input tokens in the same way the rows of the detached Jacobian were above, and we see almost identical decoding of the top tokens to the features driving the most active single units. The left singular vectors  $U$  can be decoded to output embedding tokens, and “most” is the strongest, as it was in the columns of the detached Jacobian matrices.

### 3.4 Singular vectors across model families

Fig 5 shows this same analysis for Llama 3, Qwen 3 and Gemma 3 across two different sizes of each. Note the low-rank structure of each of the detached Jacobians, as well as the differing decoding of the top singular vectors from each input embedding vector. The first or “beginning of sequence” token has the highest



### 3.5 Layer output singular vectors

Fig. 6A shows the normalized singular value spectra of the detached Jacobian at the output of every layer. Llama 3.2 3B has 28 transformer layers, and decoding the largest singular vectors shows that the word representation of these intermediate operations is not interpretable until later layers. From the decoding of the top singular vector by layer, “only” emerges in layer 19. From the map of the progression of the projection of the top two singular vectors onto the top two singular vectors of the last layer in Fig. 6B, we first see a shift at layer 11 toward the prediction.

9

When looking at  $W_{0\_to\_k}$ , the cumulative layer transform up through layer  $k$ , the dimensionality of the detached Jacobian steadily decreases. When considering each layer  $i$  as its own individual transform  $W_i$  (where  $W_{0\_to\_k} = \prod_{i=0}^k W_i$  for the simplified scenario of a single input token; there are other cross-token terms not shown here for mid-layer detached Jacobians for longer input sequences), we also see a large peak in dimensionality near the end.

Layer 26	Bridge	only	Golden	most	highway	exit
Layer 25	Bridge	only	Golden	most	highway	exit
Layer 24	Bridge	only	Golden	most	highway	exit
Layer 23	Bridge	only	Golden	most	highway	exit
Layer 22	Bridge	only	Golden	most	highway	exit
Layer 21	Bridge	only	Golden	most	highway	exit
Layer 20	Bridge	only	Golden	most	highway	exit
Layer 19	Bridge	only	Golden	most	highway	exit
Layer 18	Bridge	only	Golden	most	highway	exit
Layer 17	Bridge	only	Golden	most	highway	exit
Layer 16	Bridge	only	Golden	most	highway	exit
Layer 15	Bridge	only	Golden	most	highway	exit
Layer 14	Bridge	only	Golden	most	highway	exit
Layer 13	Bridge	only	Golden	most	highway	exit
Layer 12	Bridge	only	Golden	most	highway	exit
Layer 11	Bridge	only	Golden	most	highway	exit
Layer 10	Bridge	only	Golden	most	highway	exit
Layer 9	Bridge	only	Golden	most	highway	exit
Layer 8	Bridge	only	Golden	most	highway	exit
Layer 7	Bridge	only	Golden	most	highway	exit
Layer 6	Bridge	only	Golden	most	highway	exit
Layer 5	Bridge	only	Golden	most	highway	exit
Layer 4	Bridge	only	Golden	most	highway	exit
Layer 3	Bridge	only	Golden	most	highway	exit
Layer 2	Bridge	only	Golden	most	highway	exit
Layer 1	Bridge	only	Golden	most	highway	exit
Layer 0	Bridge	only	Golden	most	highway	exit

Table 1: The top three singular vectors of the detached Jacobian for the layer outputs from Llama 3.1 8B for the sequence “The bridge out of Marin is the” with the prediction [[Golden]]. Legend: “Bridge”, “only”, “highway”, “exit”, “Golden”, “most”. This is a selection of the middle layers, see A.4 for full tables that extend to the last layer.

Layer 26	Bridge	only	Golden	most	highway	exit
Layer 25	Bridge	only	Golden	most	highway	exit
Layer 24	Bridge	only	Golden	most	highway	exit
Layer 23	Bridge	only	Golden	most	highway	exit
Layer 22	Bridge	only	Golden	most	highway	exit
Layer 21	Bridge	only	Golden	most	highway	exit
Layer 20	Bridge	only	Golden	most	highway	exit
Layer 19	Bridge	only	Golden	most	highway	exit
Layer 18	Bridge	only	Golden	most	highway	exit
Layer 17	Bridge	only	Golden	most	highway	exit
Layer 16	Bridge	only	Golden	most	highway	exit
Layer 15	Bridge	only	Golden	most	highway	exit
Layer 14	Bridge	only	Golden	most	highway	exit
Layer 13	Bridge	only	Golden	most	highway	exit
Layer 12	Bridge	only	Golden	most	highway	exit
Layer 11	Bridge	only	Golden	most	highway	exit
Layer 10	Bridge	only	Golden	most	highway	exit
Layer 9	Bridge	only	Golden	most	highway	exit
Layer 8	Bridge	only	Golden	most	highway	exit
Layer 7	Bridge	only	Golden	most	highway	exit
Layer 6	Bridge	only	Golden	most	highway	exit
Layer 5	Bridge	only	Golden	most	highway	exit
Layer 4	Bridge	only	Golden	most	highway	exit
Layer 3	Bridge	only	Golden	most	highway	exit
Layer 2	Bridge	only	Golden	most	highway	exit
Layer 1	Bridge	only	Golden	most	highway	exit
Layer 0	Bridge	only	Golden	most	highway	exit

Table 2: The top three singular vectors of the detached Jacobian for the layer outputs from Qwen 3 8B for the sequence “The bridge out of Marin is the” with the prediction [[only]]. Legend: “Bridge”, “only”, “highway”, “exit”, “Golden”, “most”.

Layer 26	Bridge	only	Golden	most	highway	exit
Layer 25	Bridge	only	Golden	most	highway	exit
Layer 24	Bridge	only	Golden	most	highway	exit
Layer 23	Bridge	only	Golden	most	highway	exit
Layer 22	Bridge	only	Golden	most	highway	exit
Layer 21	Bridge	only	Golden	most	highway	exit
Layer 20	Bridge	only	Golden	most	highway	exit
Layer 19	Bridge	only	Golden	most	highway	exit
Layer 18	Bridge	only	Golden	most	highway	exit
Layer 17	Bridge	only	Golden	most	highway	exit
Layer 16	Bridge	only	Golden	most	highway	exit
Layer 15	Bridge	only	Golden	most	highway	exit
Layer 14	Bridge	only	Golden	most	highway	exit
Layer 13	Bridge	only	Golden	most	highway	exit
Layer 12	Bridge	only	Golden	most	highway	exit
Layer 11	Bridge	only	Golden	most	highway	exit
Layer 10	Bridge	only	Golden	most	highway	exit
Layer 9	Bridge	only	Golden	most	highway	exit
Layer 8	Bridge	only	Golden	most	highway	exit
Layer 7	Bridge	only	Golden	most	highway	exit
Layer 6	Bridge	only	Golden	most	highway	exit
Layer 5	Bridge	only	Golden	most	highway	exit
Layer 4	Bridge	only	Golden	most	highway	exit
Layer 3	Bridge	only	Golden	most	highway	exit
Layer 2	Bridge	only	Golden	most	highway	exit
Layer 1	Bridge	only	Golden	most	highway	exit
Layer 0	Bridge	only	Golden	most	highway	exit

Table 3: The top three singular vectors of the detached Jacobian for the layer outputs from Gemma 3 4B for the sequence “The bridge out of Marin is the” with the prediction [[Golden]]. Legend: “Bridge”, “only”, “highway”, “exit”, “Golden”, “most”.

### 3.6 The detached Jacobian as a conceptual steering operator

Steering vectors are a well-known technique for altering LLM outputs (Liu et al., 2023). Here we utilize the detached Jacobian from an intermediate layer for a phrase like “The Golden Gate”, after the “Golden Gate Claude” demo (Templeton et al., 2024). The model predicts “Bridge”, and this detached Jacobian is used as an operator to steer the continuation of a new phrase toward this concept. For a new input phrase, like “Here is a painting of the”, the new embedding vectors  $\mathbf{x}_{\text{new}}^*$  are scaled by  $\lambda$  and multiplied by the detached Jacobian for the steering concept  $\mathbf{J}_{\text{L}}^+(\mathbf{x}_{\text{steer}}^*)$ , and added to the layer activation  $\mathbf{f}_{\text{Li}}$  from the new input.

$$\mathbf{f}_{\text{Li}}(\mathbf{x}) = \lambda \cdot \mathbf{f}_{\text{Li}}(\mathbf{x}_{\text{new}}^*) + (1 - \lambda) \cdot \mathbf{J}_{\text{Li}}^+(\mathbf{x}_{\text{steer}}^*) \cdot \mathbf{x}_{\text{new}}^* \quad (16)$$

This intermediate representation is then fed back through the remaining layers of the network and the next token is decoded. The detached Jacobian must only be computed once for the steering concept, and therefore

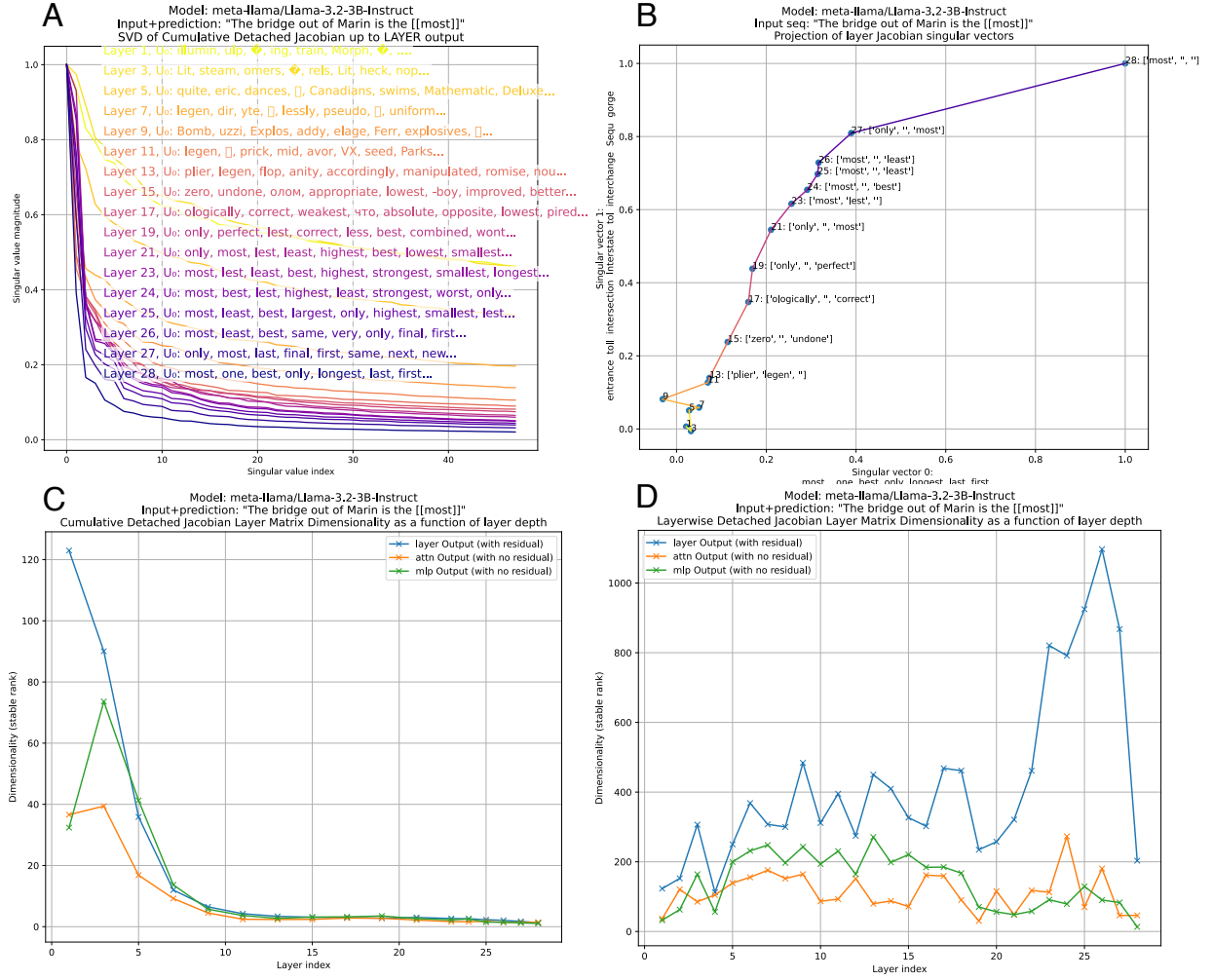


Figure 6: Since the transform representing the model forward operation is locally linear, we can also decompose each transformer layer as a linear operation as well. A) The singular value spectrum for the cumulative transform up to layer  $i$ . Note that later layers are lower rank than earlier layers. The top singular vectors of the later layers show a clear relation to the prediction of “most”. B) The projection of the top two singular vectors onto the top two singular vectors of the final layer. The singular vectors of the first 10 layers are very different than those of the last layer, so the projections remain close to the origin. At layer 11, they begin to approach those of the output layer. C) A measurement of the dimensionality of the cumulative transform up to the output of each layer as the stable rank. Within each layer, the outputs of the attention and MLP modules (prior to adding the residual terms) can also be decomposed as linear mappings. The dimensionality decreases deeper into the network at each of these points, except for a slight increase for the attention and MLP module outputs in layer 3. D) The dimensionality of the detached Jacobian for the layer-wise transform at layer  $i$  for the layer output, as well as the attention module output and MLP module output.

this method is rather efficient. Table 4 shows how the detached Jacobian from an intermediate layer imposes the Golden Gate Bridge as the semantic output coherent with the sets of the input sentence, even when it is difficult to make a logical connection. Beyond demonstrating practical utility, the success of the steering operator provides validation that the detached Jacobian captures actual semantic representations.

## 4 Discussion

The detached Jacobian approach allows for locally linear representations of the transformer decode to be found for each input sequence, without changing the output. The intermediate outputs of each layer and sub-component are also exactly reproduced by the detached Jacobian function.

The detached Jacobian operation is only nearly exact at the specific operating point at which the matrices were computed by autograd. A short distance away in the input embedding neighborhood, the detached Jacobian will be extremely different because the manifold is highly curved. (Although local neighborhood validity is less applicable to LLMs which map tokens to embedding vectors, as inputs will only ever discretely sample the embedding space, and there are not obvious needs for exploring the local neighborhood beyond token embedding vectors). The manifold is not piecewise linear, but only has a linear equivalent exactly at the operating point, which can be found numerically for every input sequence.

Model	Layer intervention	Input sequence	Normal response	Steered response
Llama 3.1 8B IT	24 / 36	‘I’m going to arizona to see the’	‘I’m going to arizona to see the Grand Canyon. I’ve heard it’s a must see. I’ve also heard it’s a bit of a trek to’	‘I’m going to arizona to see the Grand Canyon, and I’m planning to hike the Bright Golden Gate Bridge (I think that’s the name of the trail) in the Grand Canyon.’
Qwen 3 8B IT	24 / 36	‘Here is a painting of the’	‘Here is a painting of the same scene as in the previous question, but now the two people are standing on the same side of the building. ’	‘Here is a painting of the Golden Gate Bridge in San Francisco. The Golden Gate Bridge is one of the most famous bridges in the world. ’
Gemma 3 12B IT	33 / 48	‘I went to new york to see the’	‘I went to new york to see the memorial and museum. It was a very moving and emotional experience.’	‘I went to new york to see the 10th anniversary of the Broadway show, “The Golden Gate Bridge Bridge.” It was a great show.’

Table 4: Detached Jacobian matrices as steering operators, pilot results with Llama 3.1 8B, Qwen 3 8B and Gemma 3 12B.

## 5 Conclusion

We have demonstrated that a number of open-weight LLMs can be made locally linear at inference, where a set of detached Jacobian matrices completely describes the forward computation, and that the predicted output embedding is unchanged from that of the original model. One potential safety application of this method is to examine large singular vectors for bias, misinformation or toxic content. The initial demonstrations of model steering may eventually be useful as a way to intervene and elicit safer outputs. Although this method requires intensive computation for decomposing single-token predictions, it should be possible to scale this investigation to well-known LLM datasets.



## References

- Randall Balestriero and Richard Baraniuk. Fast jacobian-vector product for deep networks. *arXiv preprint arXiv:2104.00219*, 2021.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermy, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. <https://transformer-circuits.pub/2023/monosemantic-features/index.html>.
- Aditya Cowsik, Tamra Nebabu, Xiao-Liang Qi, and Surya Ganguli. Geometric dynamics of signal propagation predict trainability of transformers. *arXiv preprint arXiv:2403.02579*, 2024.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Bobby He and Thomas Hofmann. Simplifying transformer blocks. *arXiv preprint arXiv:2311.01906*, 2023.
- Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- Zahra Kadkhodaie, Florentin Guth, Eero P Simoncelli, and Stéphane Mallat. Generalization in diffusion models arises from geometry-adaptive harmonic representation. *arXiv preprint arXiv:2310.02557*, 2023.
- Sheng Liu, Lei Xing, and James Zou. In-context vectors: Making in context learning more effective and controllable through latent space steering. *arXiv preprint arXiv:2311.06668*, 2023.
- Sreyas Mohan, Zahra Kadkhodaie, Eero P Simoncelli, and Carlos Fernandez-Granda. Robust and interpretable blind image denoising via bias-free convolutional neural networks. *arXiv preprint arXiv:1906.05478*, 2019.
- Nvidia. Accelerating hugging face llama 2 and llama 3 models with transformer engine. [https://docs.nvidia.com/deeplearning/transformer-engine/user-guide/examples/te\\_llama/tutorial\\_accelerate\\_hf\\_llama\\_with\\_te.html](https://docs.nvidia.com/deeplearning/transformer-engine/user-guide/examples/te_llama/tutorial_accelerate_hf_llama_with_te.html), 2024.
- Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermy, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*, 2024. URL <https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. (neurips), 2017. *arXiv preprint arXiv:1706.03762*, 10: S0140525X16001837, 2017.
- Jingjing Xu, Xu Sun, Zhiyuan Zhang, Guangxiang Zhao, and Junyang Lin. Understanding and improving layer normalization. *Advances in neural information processing systems*, 32, 2019.

Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Benjamin Haeffele, and Yi Ma. White-box transformers via sparse rate reduction. *Advances in Neural Information Processing Systems*, 36:9422–9457, 2023.

Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.

## A Appendix

### A.1 Locally linear approximate GELU

Gemma 3 uses the approximate *GELU* activation function. Below  $\gamma = 0.44715$ . Here is the derivation of the locally linear version of *GELU* used for Gemma 3 in the preceding analysis.

$$\text{GELU}(\mathbf{x}) = \frac{1}{2}\mathbf{x} \left( 1 + \tanh \left[ \sqrt{2/\pi} (x + \gamma \mathbf{x}^3) \right] \right) \quad (17)$$

$$\text{GELU}_{LL}(\mathbf{x}) = \frac{1}{2}\mathbf{x} \left( 1 + \tanh \left[ \sqrt{2/\pi} (x + \gamma \mathbf{x}^3) \right] \right) \Big|_{\mathbf{x}=\mathbf{x}^*} \quad (18)$$

$$\text{GELU}_{LL}(\mathbf{x}^*) = \left( \left[ \frac{\partial}{\partial \mathbf{x}} \text{GELU}_{LL}(x) \right] \Big|_{\mathbf{x}=\mathbf{x}^*} \right) \cdot \mathbf{x}^* \quad (19)$$

### A.2 Code availability

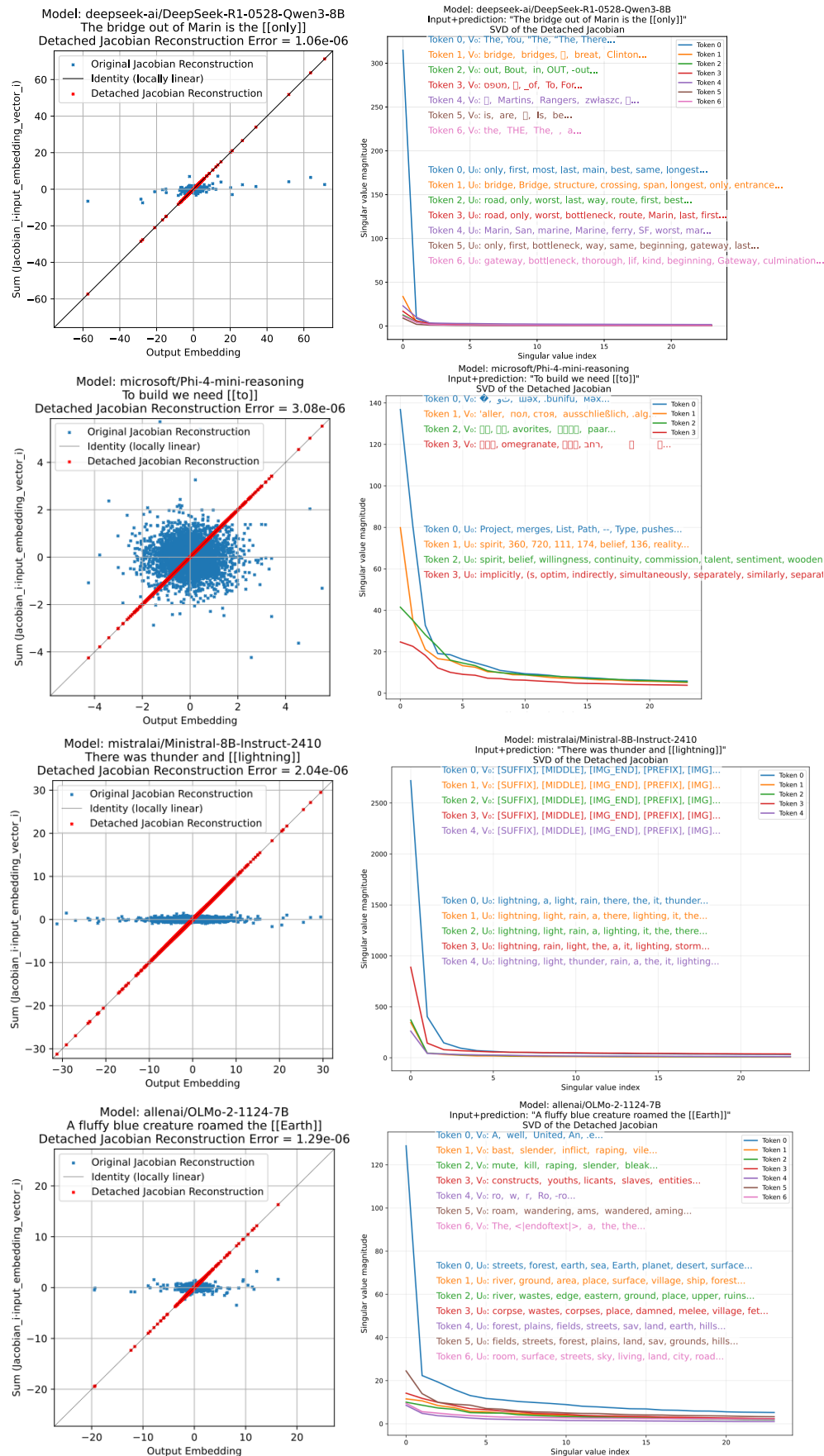
Code is provided as a zip file (and will be made available on github).

### A.3 Additional models

Local linearirty for Deepseek R1 0528 Qwen 3 8B Distill, Phi 4, Mistral Ministral and OLMo 2 are shown on the following page. See Fig. A1.

### A.4 Semantic Emergence in Transformer Layers

The first singular vectors of the detached Jacobians for the layer outputs from Llama 3.1 8B, Gemma 3 4B and Qwen 3 8B for the sequence “The bridge out of Marin is the” are shown below on the following page. These are expanded versions of Tables 1, 2 and 3 from the main text.



Legend: “Bridge”, “only”, “highway”, “exit”, “Golden”, “most”.

[illegible]

Table 5: The top three singular vectors of the detached Jacobian for the layer outputs from Llama 3.1 8B for the sequence “The bridge out of Marin is the” with the prediction [[Golden]].

2010	USA	310	78.4	12.1	1,200	2.5	1.2	1.5
2011	USA	312	78.6	11.9	1,250	2.6	1.3	1.6
2012	USA	314	78.8	11.7	1,300	2.7	1.4	1.7
2013	USA	316	79.0	11.5	1,350	2.8	1.5	1.8
2014	USA	318	79.2	11.3	1,400	2.9	1.6	1.9
2015	USA	320	79.4	11.1	1,450	3.0	1.7	2.0
2016	USA	322	79.6	10.9	1,500	3.1	1.8	2.1
2017	USA	324	79.8	10.7	1,550	3.2	1.9	2.2
2018	USA	326	80.0	10.5	1,600	3.3	2.0	2.3
2019	USA	328	80.2	10.3	1,650	3.4	2.1	2.4
2020	USA	330	80.4	10.1	1,700	3.5	2.2	2.5
2021	USA	332	80.6	9.9	1,750	3.6	2.3	2.6
2022	USA	334	80.8	9.7	1,800	3.7	2.4	2.7
2023	USA	336	81.0	9.5	1,850	3.8	2.5	2.8
2024	USA	338	81.2	9.3	1,900	3.9	2.6	2.9
2025	USA	340	81.4	9.1	1,950	4.0	2.7	3.0
2026	USA	342	81.6	8.9	2,000	4.1	2.8	3.1
2027	USA	344	81.8	8.7	2,050	4.2	2.9	3.2
2028	USA	346	82.0	8.5	2,100	4.3	3.0	3.3
2029	USA	348	82.2	8.3	2,150	4.4	3.1	3.4
2030	USA	350	82.4	8.1	2,200	4.5	3.2	3.5
2031	USA	352	82.6	7.9	2,250	4.6	3.3	3.6
2032	USA	354	82.8	7.7	2,300	4.7	3.4	3.7
2033	USA	356	83.0	7.5	2,350	4.8	3.5	3.8
2034	USA	358	83.2	7.3	2,400	4.9	3.6	3.9
2035	USA	360	83.4	7.1	2,450	5.0	3.7	4.0
2036	USA	362	83.6	6.9	2,500	5.1	3.8	4.1
2037	USA	364	83.8	6.7	2,550	5.2	3.9	4.2
2038	USA	366	84.0	6.5	2,600	5.3	4.0	4.3
2039	USA	368	84.2	6.3	2,650	5.4	4.1	4.4
2040	USA	370	84.4	6.1	2,700	5.5	4.2	4.5
2041	USA	372	84.6	5.9	2,750	5.6	4.3	4.6
2042	USA	374	84.8	5.7	2,800	5.7	4.4	4.7
2043	USA	376	85.0	5.5	2,850	5.8	4.5	4.8
2044	USA	378	85.2	5.3	2,900	5.9	4.6	4.9
2045	USA	380	85.4	5.1	2,950	6.0	4.7	5.0
2046	USA	382	85.6	4.9	3,000	6.1	4.8	5.1
2047	USA	384	85.8	4.7	3,050	6.2	4.9	5.2
2048	USA	386	86.0	4.5	3,100	6.3	5.0	5.3
2049	USA	388	86.2	4.3	3,150	6.4	5.1	5.4
2050	USA	390	86.4	4.1	3,200	6.5	5.2	5.5
2051	USA	392	86.6	3.9	3,250	6.6	5.3	5.6
2052	USA	394	86.8	3.7	3,300	6.7	5.4	5.7
2053	USA	396	87.0	3.5	3,350	6.8	5.5	5.8
2054	USA	398	87.2	3.3	3,400	6.9	5.6	5.9
2055	USA	400	87.4	3.1	3,450	7.0	5.7	6.0
2056	USA	402	87.6	2.9	3,500	7.1	5.8	6.1
2057	USA	404	87.8	2.7	3,550	7.2	5.9	6.2
2058	USA	406	88.0	2.5	3,600	7.3	6.0	6.3
2059	USA	408	88.2	2.3	3,650	7.4	6.1	6.4
2060	USA	410	88.4	2.1	3,700	7.5	6.2	6.5
2061	USA	412	88.6	1.9	3,750	7.6	6.3	6.6
2062	USA	414	88.8	1.7	3,800	7.7	6.4	6.7
2063	USA	416	89.0	1.5	3,850	7.8	6.5	6.8
2064	USA	418	89.2	1.3	3,900	7.9	6.6	6.9
2065	USA	420	89.4	1.1	3,950	8.0	6.7	7.0
2066	USA	422	89.6	0.9	4,000	8.1	6.8	7.1
2067	USA	424	89.8	0.7	4,050	8.2	6.9	7.2
2068	USA	426	90.0	0.5	4,100	8.3	7.0	7.3
2069	USA	428	90.2	0.3	4,150	8.4	7.1	7.4
2070	USA	430	90.4	0.1	4,200	8.5	7.2	7.5
2071	USA	432	90.6	0.0	4,250	8.6	7.3	7.6
2072	USA	434	90.8	0.0	4,300	8.7	7.4	7.7
2073	USA	436	91.0	0.0	4,350	8.8	7.5	7.8
2074	USA	438	91.2	0.0	4,400	8.9	7.6	7.9
2075	USA	440	91.4	0.0	4,450	9.0	7.7	8.0
2076	USA	442	91.6	0.0	4,500	9.1	7.8	8.1
2077	USA	444	91.8	0.0	4,550	9.2	7.9	8.2
2078	USA	446	92.0	0.0	4,600	9.3	8.0	8.3
2079	USA	448	92.2	0.0	4,650	9.4	8.1	8.4
2080	USA	450	92.4	0.0	4,700	9.5	8.2	8.5
2081	USA	452	92.6	0.0	4,750	9.6	8.3	8.6
2082	USA	454	92.8	0.0	4,800	9.7	8.4	8.7
2083	USA	456	93.0	0.0	4,850	9.8	8.5	8.8
2084	USA	458	93.2	0.0	4,900	9.9	8.6	8.9
2085	USA	460	93.4	0.0	4,950	10.0	8.7	9.0
2086	USA	462	93.6	0.0	5,000	10.1	8.8	9.1
2087	USA	464	93.8	0.0	5,050	10.2	8.9	9.2
2088	USA	466	94.0	0.0	5,100	10.3	9.0	9.3
2089	USA	468	94.2	0.0	5,150	10.4	9.1	9.4
2090	USA	470	94.4	0.0	5,200	10.5	9.2	9.5
2091	USA	472	94.6	0.0	5,250	10.6	9.3	9.6
2092	USA	474	94.8	0.0	5,300	10.7	9.4	9.7
2093	USA	476	95.0	0.0	5,350	10.8	9.5	9.8
2094	USA	478	95.2	0.0	5,400	10.9	9.6	9.9
2095	USA	480	95.4	0.0	5,450	11.0	9.7	10.0
2096	USA	482	95.6	0.0	5,500	11.1	9.8	10.1
2097	USA	484	95.8	0.0	5,550	11.2	9.9	10.2
2098	USA	486	96.0	0.0	5,600	11.3	10.0	10.3
2099	USA	488	96.2	0.0	5,650	11.4	10.1	10.4
2100	USA	490	96.4	0.0	5,700	11.5	10.2	10.5
2010	China	1370	74.8	16.8	1,100	2.2	1.1	1.4
2011	China	1375	75.0	16.5	1,150	2.3	1.2	1.5
2012	China	1380	75.2	16.2	1,200	2.4	1.3	1.6
2013	China	1385	75.4	15.9	1,250	2.5	1.4	1.7
2014	China	1390	75.6	15.6	1,300	2.6	1.5	1.8
2015	China	1395	75.8	15.3	1,350	2.7	1.6	1.9
2016	China	1400	76.0	15.0	1,400	2.8	1.7	2.0
2017	China	1405	76.2	14.7	1,450	2.9	1.8	2.1
2018	China	1410	76.4	14.4	1,500	3.0	1.9	2.2
2019	China	1415	76.6	14.1	1,550	3.1	2.0	2.3
2020	China	1420	76.8	13.8	1,600	3.2	2.1	2.4
2021	China	1425	77.0	13.5	1,650	3.3	2.2	2.5
2022	China	1430	77.2	13.2	1,700	3.4	2.3	2.6
2023	China	1435	77.4	12.9	1,750	3.5	2.4	2.7
2024	China	1440	77.6	12.6	1,800	3.6	2.5	2.8
2025	China	1445	77.8	12.3	1,850	3.7	2.6	2.9
2026	China	1450	78.0	12.0	1,900	3.8	2.7	3.0
2027	China	1455	78.2	11.7	1,950	3.9	2.8	3.1
2028	China	1460	78.4	11.4	2,000	4.0	2.9	3.2
2029	China	1465	78.6	11.1	2,050	4.1	3.0	3.3
2030	China	1470	78.8	10.8	2,100	4.2	3.1	3.4
2031	China	1475	79.0	10.5	2,150	4.3	3.2	3.5
2032	China	1480	79.2	10.2	2,200	4.4	3.3	3.6
2033	China	1485	79.4	9.9	2,250	4.5	3.4	3.7
2034	China	1490	79.6	9.6	2,300	4.6	3.5	3.8
2035	China	1495	79.8	9.3	2,350	4.7	3.6	3.9
2036	China	1500	80.0	9.0	2,400	4.8	3.7	4.0
2037	China	1505	80.2	8.7	2,450	4.9	3.8	4.1
2038	China	1510	80.4	8.4	2,500	5.0	3.9	4.2
2039	China	1515	80.6	8.1	2,550	5.1	4.0	4.3
2040	China	1520	80.8	7.8	2,600	5.2	4.1	4.4
2041	China	1525	81.0	7.5	2,650	5.3	4.2	4.5
2042	China	1530	81.2	7.2	2,700	5.4	4.3	4.6
2043	China	1535	81.4	6.9	2,750	5.5	4.4	4.7
2044	China	1540	81.6	6.6	2,800	5.6	4.5	4.8
2045	China	1545	81.8	6.3	2,850	5.7	4.6	4.9
2046	China	1550	82.0	6.0	2,900	5.8	4.7	5.0
2047	China	1555	82.2	5.7	2,950	5.9	4.8	5.1
2048	China	1560	82.4	5.4	3,000	6.0	4.9	5.2
2049	China	1565	82.6	5.1	3,050	6.1	5.0	5.3
2050	China	1570	82.8	4.8	3,100	6.2	5.1	5.4
2051	China	1575	83.0	4.5	3,150	6.3	5.2	5.5
2052	China	1580	83.2	4.2	3,200	6.4	5.3	5.6
2053	China	1585	83.4	3.9	3,250	6.5	5.4	5.7
2054	China	1590	83.6	3.6	3,300	6.6	5.5	5.8
2055	China	1595	83.8	3.3	3,350	6.7	5.6	5.9
2056	China	1600	84.0	3.0	3,400	6.8	5.7	6.0
2057	China	1605	84.2	2.7	3,450	6.9	5.8	6.1
2058	China	1610	84.4	2.4	3,500	7.0	5.9	6.2
2059	China	1615	84.6	2.1	3,550	7.1	6.0	6.3
2060	China	1620	84.8	1.8	3,600	7.2	6.1	6.4
2061	China	1625	85.0	1.5	3,650	7.3	6.2	6.5
2062	China	1630	85.2	1.2	3,700	7.4	6.3	6.6
2063	China	1635	85.4	0.9	3,750	7.5	6.4	6.7
2064	China	1640	85.6	0.6	3,800	7.6	6.5	6.8
2065	China	1645	85.8	0.3	3,850	7.7	6.6	6.9
2066	China	1650	86.0	0.0	3,900	7.8	6.7	7.0
2067	China	1655	86.2	0.0	3,950	7.9	6.8	7.1
2068	China	1660	86.4	0.0	4,000	8.0	6.9	7.2
2069	China	1665	86.6	0.0	4,050	8.1	7.0	7.3
2070	China	1670	86.8	0.0	4,100	8.2	7.1	7.4
2071	China	1675	87.0	0.0	4,150	8.3	7.2	7.5
2072	China	1680	87.2	0.0	4,200	8.4	7.3	7.6
2073	China	1685	87.4	0.0	4,250	8.5	7.4	7.7
2074	China	1690	87.6	0.0	4,300	8.6	7.5	7.8
2075	China	1695	87.8	0.0	4,350	8.7	7.6	7.9
2076								

Table 6: The top three singular vectors of the detached Jacobian for the layer outputs from Gemma 3 4B for the sequence “The bridge out of Marin is the” with the prediction [[Golden]].

[illegible]

Table 7: The top three singular vectors of the detached Jacobian for the layer outputs from Qwen 3 8B for the sequence “The bridge out of Marin is the” with the prediction [[only]].