

---

# Identifying Sparsely Active Circuits Through Local Loss Landscape Decomposition

---

Anonymous Authors<sup>1</sup>

## Abstract

### 1. Background

Mechanistic interpretability is somewhat of an emerging field and aims to uncover the internal mechanisms responsible for seemingly black box behavior of large models so that developers can better understand, intervene on, and align model (Bereska & Gavves, 2024). Within mechanistic interpretability, decomposition methods aim decompose model behavior into subcomponents that are less complex and more human interpretable, but that still fully express model behavior. The most popular method in this space is Sparse Autoencoders, which have been successful in decomposing compressed features in the activation space of a model (Gao et al., 2024; Cunningham et al., 2023; Bricken et al., 2023).

#### 1.1. Limitations of Activation Space-Based Methods

Sparse autoencoders rely on the activation space of a model, the output of a set of model's neurons during inference. Sparse autoencoders (SAEs) are designed to identify latent features by projecting a model's compressed activation space into a sparsely activated, overcomplete basis. These learned basis vectors represent distinct features, which can then be used to reconstruct the original activations.

Despite their utility, activation-based decomposition methods like SAEs have notable limitations. First, they assume that activation space consists of linear combinations of feature vectors, yet growing evidence suggests that even architectures designed to encourage linearity, such as residual streams, may encode non-linear features (Engels et al., 2024a;b). Second, activation-based approaches generally treat layers and blocks as independent computational modules, reconstructing activations from specific layers without capturing the possibility of cross-layer interactions and superposition (Merullo et al., 2024; Lindsey et al., 2024). This limitation is particularly problematic in architectures beyond transformers, such as recurrent networks, adversarial networks, diffusion models, and classic RL models, where the separation of activation spaces is even less well-defined

(Pascanu, 2013; Goodfellow et al., 2014; Ho et al., 2020; Mnih et al., 2015). Finally, activation-based methods provide little insight into how features emerge during training, making it difficult to link learned behaviors to specific data or optimization dynamics. This lack of connection to the training process limits our ability to intervene in model behavior at an early stage or refine training data to steer a model toward more desirable properties.

#### 1.2. Loss Landscape Geometry

An alternative to interpreting models through their activations is to interpret models through their loss landscape geometry, understanding how perturbing parameters, rather than activations, affects model output and behavior. Parameters are the fundamental entities updated during training, capturing information from data in a way that persists beyond individual inferences. Unlike activations, which are transient and input-dependent, the parameter space reflects the cumulative outcome of the training process, potentially providing deeper insights into how models generalize, store knowledge, and develop internal representations.

A growing body of research explores the relationship between parameter space, data, and model behavior. Singular learning theory (SLT) describes how the structure of high-dimensional parameter space influences generalization and feature formation (Watanabe, 2007; 2000; 2005; Wei et al., 2022), while developmental interpretability extends this by analyzing how parameter updates evolve over training (Sharkey et al., 2025). Studying parameter space has already led to key insights, such as the role of sharpness in grokking (Davies et al., 2023), phase transitions corresponding to important learning stages in language models (Wang et al., 2024; Hoogland et al., 2024), and distinctions between parameters associated with memorization vs. generalization (Bushnaq et al., 2024). These findings suggest that parameter-space analysis may uncover mechanisms that activation-based methods overlook, opening new avenues for interpretability and intervention.

### 055 1.3. Degenerate Loss Landscapes

056 We already know that loss landscapes have high levels of  
 057 degeneracy. Mathematical degeneracy, data distribution  
 058 degeneracy, etc. We hypothesize that

### 060 1.4. Contrastive Methods

062 Contrastive/paired methods are common in interpretability.  
 063 They work

### 065 1.5. Loss Landscape Decomposition

067 We propose a method for identifying features and the cir-  
 068 cuits that produce them in large models by analyzing par-  
 069 ameter space rather than activations. Our approach identifies  
 070 directions in parameter space that correspond to specific  
 071 circuits, drawing on insights from singular learning theory  
 072 (SLT). Our method aims to identify principle directions in  
 073 parameter space such that for a given sample, moving model  
 074 parameters in most of these directions will not meaningfully  
 075 change the model's output, but for a small subset of di-  
 076 rections, the model's output will change dramatically. We  
 077 hypothesize that these directions, or subnetworks, represent  
 078 distinct computations.

079 To our knowledge, there have only been two prior methods  
 080 attempting to decompose parameter space for interpretabil-  
 081 ity. In one earlier work(Matena & Raffel, 2023), the authors  
 082 decompose parameter space by computing principal direc-  
 083 tions of a per-sample Fisher Information matrix to resolve  
 084 meaningful features. This approach aimed to identify pa-  
 085 rameter directions most sensitive to individual samples, cap-  
 086 turing how different parts of the model contribute to specific  
 087 tasks. However, due to computational constraints, it relied  
 088 on diagonalization techniques to approximate key directions,  
 089 which limited its ability to fully capture sparse, structured  
 090 circuits. Another recent method, (Braun et al., 2025) to is  
 091 Attribution Parameter Decomposition (APD), which uses  
 092 a bi-optimization approach to identify subnetworks where  
 093 (1) the sum of subnetwork weights is close to the original  
 094 model parameters, and (2) for any given input, the sum of a  
 095 smaller set of subnetworks - identified through topk attribu-  
 096 tion values - has low behavioral loss when compared to the  
 097 original model.

098 Local loss landscape decomposition (L3D) aims to decom-  
 099 pose a model's gradients of loss between output pairs with  
 100 respect to parameter space 1. We seek to identify directions  
 101 in parameter space such that (1) the set of directions can  
 102 approximately reconstruct any gradient of paired sample  
 103 divergences, and (2) for any given pair of samples, an even  
 104 smaller subset of directions can be used in reconstructing  
 105 this gradient. To reduce computational complexity, we learn  
 106 low-rank representations of these parameter directions. In  
 107 this work, we first describe the mathematical foundation of

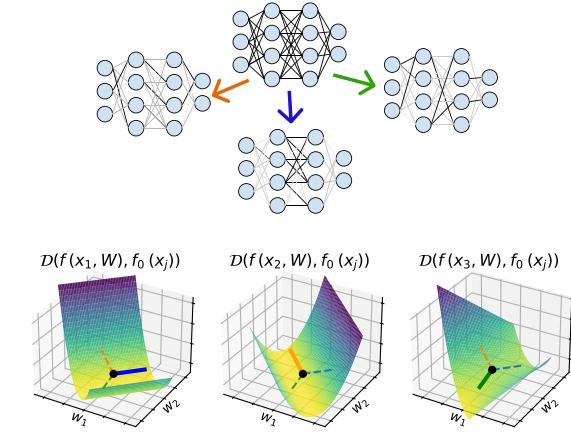


Figure 1: Decomposing a loss landscape into a set of parameter directions , where (1) the set of directions can approximately reconstruct any per-sample-pair gradient of divergence and (2) for any given sample, the curvature in most of the subnetworks directions is approximately zero.

our approach. We then develop progressively more complex toy models to measure the efficacy of L3D and quantify its limitations.

## 2. Methodology

In the next sections, we will formally setup our decomposi-  
 tion problem (Sec. 2.1), define the criteria we will use for  
 our subnetworks/parameter directions (Sec. 2.2), describe  
 how to efficiently decompose parameters into these direc-  
 tions (Sec. 2.3), walk through our training algorithm (Sec.  
 2.4), and then explain how to use these decompositions to  
 intervene on a model's behavior (Sec. 2.5).

The terms parameter directions, subnetworks, and circuits  
 can all be used interchangeably, and will refer to the direc-  
 tions in parameter space that we wish to learn.

### 2.1. Set up

Consider a model  $f$  that takes an batch of inputs  $X$  and  
 parameter values of  $W$ , and outputs a batch of outputs.

$$f(x, W) : R^{n_s \times n_f} \mapsto R^{n_s \times n_o} \quad (1)$$

where  $n_s$  is the number of samples in the batch,  $n_f$  is the  
 length of each input vector, and  $n_o$  is the length of each  
 output vector.

Our project rests on the following hypothesis:

1. For a given input, there are many components of a

model that are not involved in inference. Changing these parameters will not affect change the model's output.

2. For the same given input, there are a smaller set of components in a model that are highly involved in inference. Changing these parameters will change the model's output.
3. We are interested in parameter directions that meaningfully change a model's output. We'd like to identify directions that when perturbed, do not break the model's output by moving it completely outside of realm of realistic outputs, but rather turn decrease or enhance the power of a specific computational circuit.

## 2.2. Divergences of Paired Outputs

To put point (3) a different way, intervening in a certain parameter direction should not move our outputs too far away from a reference output distribution. Therefore, we will be decomposing gradients of divergences between pairs of outputs - with the aim of finding directions that move the model's output within the same subspace that the reference output distribution lies.

Divergence of a pair of outputs can be defined as:

$$\nabla_W D(f(x_i, W), f(x_r, W_0))|_{W=W_0} \quad (2)$$

where  $x_i, x_r \in X$

Here,  $D$  is a divergence measuer,  $f$  is our model,  $x_i$  is our input of interest,  $x_r$  is a randomly selected reference input,  $W$  is a set of parameters and  $W_0$  is model's parameters fixed. Our toy models, which are regression-type models, so we use MSE as divergence.

We will abbreviate  $f(x_r, W_0)$  as  $f(x_r)$  and the expression in Eq. 2 as  $\nabla_W D$ .

## 2.3. Sparse Principal Directions

We want to transform our gradient into directions in parameter space, where each sample's gradient can be written as a linear combination of a small set of these components. We will do this by learning transforms  $V^{in} \in R^{n_v \times n_w}$  and  $V^{out} \in R^{n_w \times n_v}$  where  $n_w$  is the number of parameters in the model, and  $n_v$  is the number of components (subnetworks) we wish to use to represent the parameter space.

$V^{in}$  effectively transforms a gradient from the parameter space to the subnetwork space, so that:

$$\nabla_V D = V^{in} \nabla_W D \quad (3)$$

We want to find  $V^{in}$  and  $V^{out}$  such that for any given pair of samples a small subset of subnetworks can approximately

reconstruct the gradient of divergence.

$$\nabla_W D \approx V^{out}_{:, \mathcal{K}} V^{in}_{\mathcal{K}, :} \nabla_W D \quad (4)$$

where  $\mathcal{K} = \text{argTopFrac}_{\kappa} \text{abs}(\nabla_{v_k} D)$

$\text{argTopFrac}$  relies on a hyperparameter that controls the number of components we wish to use to reconstruct each sample. In practice, we use a batchTopK (Bussmann et al., 2024) and a fraction rather than an absolute number. If we use a batchTopFrac of .1 it means that we select the top 10% of  $\nabla_V D$  magnitudes over  $v_k$  and  $x_i$  to reconstruct our batch of gradients.

### 2.3.1. LOW RANK PARAMETER DIRECTIONS

Learning a set of full rank parameter parameter directions would be extremely expensive (we would be learning  $n_w n_v$  values). We also expect that modular, sparsely active circuits would be lower rank than their full-model counterparts (). Therefore, we use low-rank representations of our  $V^{in}$  and  $V^{out}$ , and correspondingly learn low-rank circuits. We describe the procedure for doing this in B.1.

## 2.4. Training

We wish to learn the decomposition-related transforms  $V^{in}$  and  $V^{out}$  that minimize the topK reconstruction loss of our divergence gradient described above. Therefore we use the following L2 norm loss:

$$L = \|\nabla_W D - V^{out}_{:, \mathcal{K}} V^{in}_{\mathcal{K}, :} \nabla_W D\|_2 \quad (5)$$

For each batch of samples, we randomly select a reference sample  $x_r$  to be paired with each sample  $x_i$  in the batch. We then compute the gradient of divergence of  $x_i$  and  $x_r$  at the model's current parameters  $W_0$ . We transform that gradient into the subnetwork space using  $V^{in}$ , and compute the topK components. We transform those components back into the original parameter space using  $V^{out}$ , and compute the loss between the reconstructed gradient and the original gradient. We apply a learning update to  $V^{in}$  and  $V^{out}$  with the goal of minizing this loss. We also normalize  $V^{out}$  to be a unit vector after each update.

## 2.5. Measuring and Intervening

Our learned subnetworks will just be the rows of  $V^{out}$ , reorganized into the same tensor structure as  $W$ . After identifying subnetworks, we may want to intervene on a specific circuit.

If we wish to intervene using a single subnetwork, we can update the parameters in by moving them in a scalar factor ( $\delta$ ) of that unit direction. To tune our model in the direction of subnetwork  $v_i$  and compute predictions on  $x$ , we evaluate:

**Algorithm 1** Training Algorithm

---

```

165 1: for each epoch do
166 2:   for each  $X$  do
167 3:     for each  $x_i \in X$  do
168 4:       Randomly select  $x_r \in X$ 
169 5:        $\nabla_w D_i = \nabla_w D(f(x_i, W), f(x_r))|_{W=W_0}$ 
170 6:     end for
171 7:      $\nabla_v D = V^{in} \nabla_w$ 
172 8:      $\tau = \text{topK}(\text{abs}(\nabla_v D))$ 
173 9:      $\hat{\nabla}_w D = V^{out} (\nabla_v D \odot (\text{abs}(\nabla_v D) > \tau))$ 
174 10:     $L = \|\nabla_w D - \hat{\nabla}_w D\|_2$ 
175 11:     $L.\text{backward}()$ 
176 12:    Update  $V^{in}$  and  $V^{out}$ 
177 13:    Normalize  $V^{out}$  to be unit vectors.
178 14:  end for
179 15: end for

```

---

$$f(x, W + \delta v_i) \quad (6)$$

Going one step further, we may want to finetune our model, but constrain our finetuning to specific circuits for efficiency or performance reasons. While we do not use finetuning in this work, we describe what the process of finetuning would look like in Appendix B.2.

Additionally, we may want to identify which subnetworks were most relevant to the predictions of various samples, or identify which samples most relied on a given subnetwork. We describe this process in Appendix B.2.1

### 3. Results

We designed several toy models to test the efficacy of L3D. Our toy models all consist of several well-characterized computations being performed by the same model, with an input space designed to rely on a single or small set of tasks. By designing our toy models this way, we have well characterized ground truth subnetworks by which to compare L3D’s decompositions.

Our toy models progressively test more complex types of circuits. Table 1 describes our 4 toy models and the different attributes of circuitry the are designed to capture. The specific hyperparameters used to train our toy models is described in Appendix B.3, as well as the hyperparamters used for each decomposition in Appendix B.4

#### 3.1. Toy Model of Superposition

##### 3.1.1. SETUP

We start off by validating our algorithm on a well-studied toy problem, the toy model of superposition (TMS), with a small variation. TMS is an autoencoder with a single low-

dimensional hidden layer in between and a relu activation at the output (Elhage et al., 2022). The model is trained on a dataset of samples where few features are active at a time, and “superimposes” these features in the hidden layer such that features embeddings in the hidden layer have minimal interference with each other S2. We train a toy model of superposition with 5 features and 2 hidden dimensions (with 1-sparsity=.05), and then place three such models in parallel, to test our method’s ability to resolve superimposed features, as well as independent circuit modules.

##### 3.1.2. DECOMPOSITION

We decompose the TMS model into 15 subnetworks, using rank-1 parameter tensors. The network decomposes into subnetworks as we would expect: each subnetwork represents the embedding and reconstruction of a single input element, involving only the weights connecting the relevant input and output, and the final bias. Fig. 2 shows the subnetworks corresponding to each of the first 5 features, and Fig. S3 shows the full decomposition. One thing to note is that parameter vectors do not have a preferred direction. L3D is equally likely to identify a parameter vector in the direction of  $\theta$  as it is in the direction of  $-\theta$ . This is why, for example, the weights and biases in subnetwork 1 are in the opposite direction as the original network (Table 1).

This decomposition resulted in a reconstruction error of .19. 15 components of rank-1 can successfully express each individual  $X_i : \hat{X}_i$  circuit, but does not capture the interference between features when multiple features are active in the same sample. We expect decompositions of higher dimensional networks to exhibit less reconstruction error, as the amount of nearly orthogonal parameter vectors (non-interfering) that can be compressed into parameter space scales exponentially with dimension. We see this effect in later toy models.

##### 3.1.3. INTERVENTION

Parameter vectors learned by L3D can be used to intervene on model behavior. In principle, we could finetune a model with an adaptor that only changes the parameters of a network in the direction of a specified set of parameter vectors (See Appendix B.2). While we do not go the extent of finetuning a model, we explore the effect of perturbing a model’s parameter space in the direction of a subnetwork (by an increment of  $\delta$ ). If subnetworks do in fact represent sparse computations, we hope that intervening on a subnetwork has a strong effect on the predictions of relevant samples, and no effect on other samples. As shown in Fig. 3, moving the TMS-in-parallel model in the direction of a single subnetwork does in fact achieve just this. Perturbing in the direction of network 1 primarily affects samples where feature 1 was active, with a small effect on

	Toy model of superposition	Circuit Superposition	Higher Rank Circuit Superposition	Complex Loss Landscape
220				
221				
222				
223				
224				
225				
226				
227				
228				
229				
230				
231	Feature Superposition ✓	X ↦ AX ✓	X ↦ AX ✓	X ↦ X <sup>2</sup> ✓
232				
233	Circuit Superposition ×	✓	✓	✓
234	Circuits > rank-2 ×	×	✓	probably ✓
235	Complex Loss Land-	×	×	✓
236	scape			
237				

Table 1: Toy models and their properties (transposed)

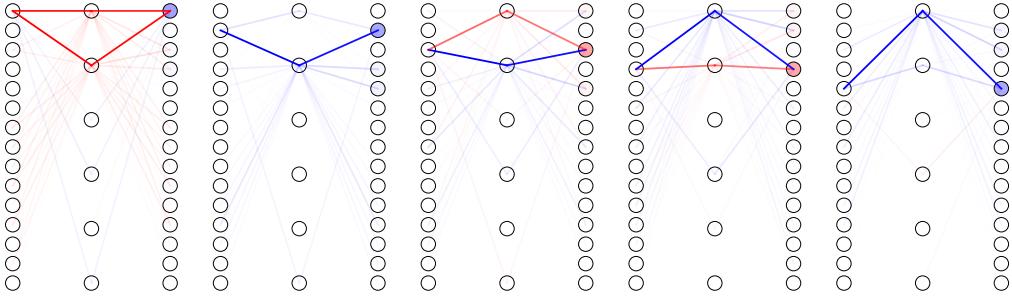


Figure 2: Selected 5 subnetworks that L3D decomposes the TMS-in-parallel model into

the inputs that had interference to feature 1’s embeddings. In fact, for TMS-in-parallel, we can successfully fully ”turn off” a computation by moving far enough in the direction of the subnetwork. (Although for models with more complex loss landscapes, turning ”off” a computation is not as straightforward, as we will later discuss).

### 3.2. Toy Model of Circuit Superposition

#### 3.2.1. SETUP

TMS famously exhibits feature superposition - the input elements’ low dimensional embeddings are non-orthogonal. However, the sparse circuits in TMS are noteably *not* in superposition - a given weight or parameter is only relevant for a single circuit and circuits. It seems highly unlikely that real world model circuits would decompose this way, since learning circuits composed of non-orthogonal parameter

vectors would allow a model to put more expressivity into a compressed space. We therefore develop a toy model of *circuit superposition* (TMCS) in order to analyze L3D’s ability to resolve such circuits. We define circuit superposition as a phenomenon by which subnetworks share parameter elements, and even more generally have non-orthogonal parameter vectors.

Our toy model of circuit superposition uses the same architecture and input data distribution as TMS, but is trained to predict linear combinations of the input features as its output ( $X \mapsto AX$ ). We randomly select values of  $A$  between 0 and 3 to generate input-output pairs to train the toy model with. We use an input with 10 input features, and 10 output features (although such a model does not need to have the same number of input and outputs).

In this model, if we consider ”subnetworks” to be parameter directions involved in inference of a small set of samples,

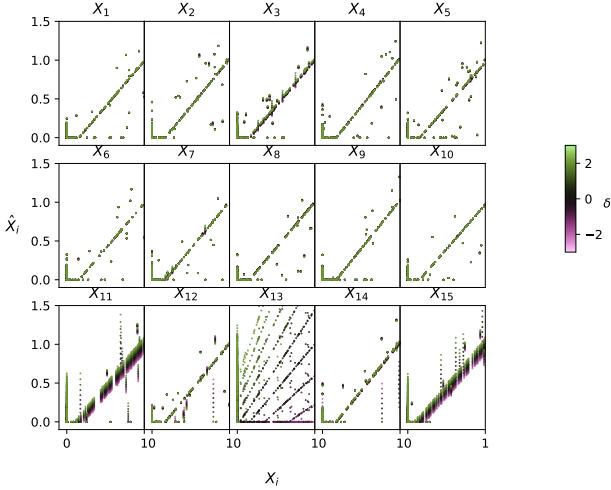


Figure 3: The effect of intervening on the TMS-in-parallel model in the direction of Network<sub>1</sub>.

then we would expect each circuit to be associated with a single input feature. If this is the case, then parameters would be involved in multiple circuits:  $W^{dec}_{i,1}$  (the set of parameters connecting the hidden nodes to the first output node) will contain information about both  $A_{1,1}, A_{2,1}, A_{3,1}, \dots$ . Put another way, the circuits will interfere with each other - parameter directions associated with each circuit will be non-orthogonal.

### 3.2.2. DECOMPOSITION

We decompose TMCS into 10 subnetworks of rank-1 parameter tensors (4) with a reconstruction loss of .064. The subnetworks look like we expect, with each strongly corresponding to a single input feature.

Since each subnetwork theoretically corresponds to the computations involved with a single input feature, we should be able to reconstruct the original  $A$  values from each subnetwork. To derive  $A$  from each subnetwork, we (1) identify the which column in the subnetwork's  $W^{dec}$  direction has the largest norm and then (2) trace the weights of the network through that path. That is for subnetwork  $k$ :

$$\begin{aligned} j^* &= \text{argmax}_j \|W^{dec}_{j,k}\|_2 \\ \hat{a}_{i,j} &= W^{enc}_{i,j,k} W^{dec}_{i,j,k} \end{aligned} \quad (7)$$

The parameter vectors are normalized to be unit vectors so we expect them to be a scalar multiple of the true  $A$  values. As seen in Figure 4, our derived  $\hat{a}$  have a very high correlation to the original  $A$  values ( $r^2 = 0.92$ ).

## 3.3. Higher Rank Circuits

### 3.3.1. SETUP

Because each subnetwork in TMCS traces the path of a single input neuron, the underlying subnetworks should inherently have a rank of 1. In order to test the ability of L3D to learn higher rank circuits, we developed a toy model with inherently higher rank circuits. For this model, we use the same set up as TMCS, but we correlate the sparsities of sets of input features. We use 25 input features, and we filter our data to ensure that input features 1-5, 6-10, etc, are always active ( $> 0$ ) or inactive ( $< 0$ ) together. In this setup, circuits should always be associated with groups of 5 input features and so should have a rank of 5.

### 3.3.2. DECOMPOSITION

Although we expect the model to have 5 subcircuits, we use excess parameter tensors ( $n=10$ ) in order to allow more flexibility in learning. We track the fraction of inputs for which a subnetwork was used in the topK reconstruction ( $P_{act}$ ) to identify which were "dead circuits", and report the last epoch. Furthermore, although we expect the underlying subcircuits to be rank 5, we experiment with using different rank representations to see how well lower-rank parameter directions can represent the model. Interestingly, rank-1 representations of the parameter tensors are able to represent the model nearly as well as rank-5 representations (Fig. S5). In Fig. 6, we show the decomposition of a rank-3 decomposition. L3D successfully learns a subnetwork corresponding to each of the 5 sets of input features, as well as a number of "dead". The higher and lower rank decompositions also learn similar subnetworks (Figure S6). When we trained L3D without these additional subnetworks, the reconstruction loss would get caught in a local minima. Similar to training SAEs (Cunningham et al., 2023), having extra degrees of freedom allows for better learning, even if at the end of training the extra subnetworks are never active.

## 3.4. Complex Loss Landscape

### 3.4.1. SETUP

The previous set of toy models all had relatively quadratic landscapes with respect to each element of  $W^{enc}$  or  $W^{dec}$  -  $dMSE/dW_i$  depends on  $W_i$  to the first order (or not at all, depending on whether the ReLU has fired). This suggests that the local loss landscape of the models is a good approximation for much of the global loss landscape, and we should expect intervening on a circuit to have well-behaved effects, even at large values of  $\delta$  (as in Fig. 3).

However, we wanted to test the limitations of a L3D on a model with a more complex loss landscape, especially when it comes to intervening with a subnetwork.

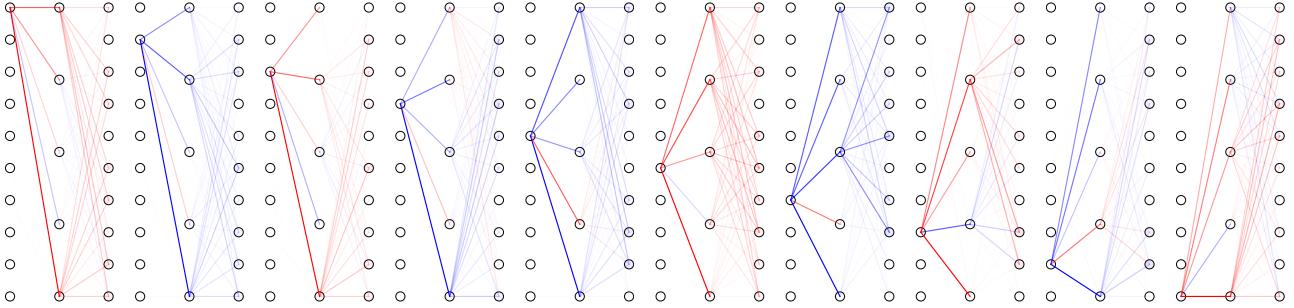


Figure 4: The learned subnetworks of TMCS

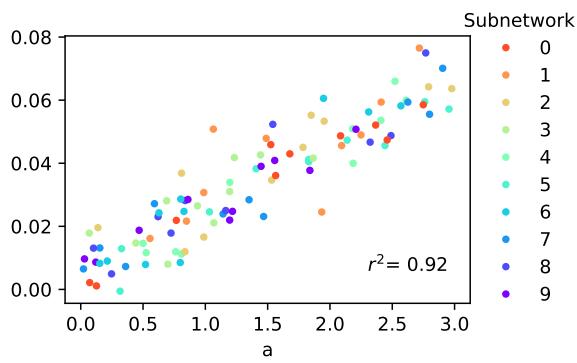


Figure 5: The coefficients derived from the subnetworks compared to actual coefficients

We therefore trained a multi-layer model to predict multiple non-linear functions of input features at once. We train a GeLU network for  $X_i \mapsto X_i^2$ . We use a network with 4 hidden layers of 5 neurons each, and 5 input and output neurons. Once again, the input features are sparse, incentivizing the toy model to learn circuits in superposition whose interferences will cause minimal errors on the sparse input distribution.

We expect the model to have 5 subnetworks, one for each input feature. Although it is less clear what rank the tensors of the underlying circuits should be, we expect that with a non-linear model they will likely be high rank.

#### 3.4.2. DECOMPOSITION

We decompose our model into 5 rank-2 parameter tensors in our decomposition and instead of varying rank, we experiment with using different numbers of subnetworks to represent our model. In the 5-subnetwork decomposition (Figure 7), we see that subnetworks tracing the path of  $X_i \mapsto X_i^2$  for each input feature  $i$ . However, this decomposition has a relatively high reconstruction error of .32. Much of this is probably because we kept our topFrac constant throughout

all our models, using a topFrac of .1 for consistency. With only 5 subnetworks, this means that on average each sample's reconstruction will use  $\approx 1$  subnetwork - limiting the reconstruction error we can achieve.

We also experiment with holding rank constant (we drop to rank-1) and decompose the model into different numbers of subnetworks (3, 5, 10, and 15 subnetworks). In our 3-subnetwork decomposition, L3D still learned subnetworks corresponding to single input features, but can of course only represent 3 out of the 5 inputs. As we add more subnetworks, we are able to successfully learn more expressive decompositions of the model that reduce reconstruction error (Figure S11). Each decomposition continues to learn input-output pair specific subnetworks, with the larger decompositions resulting in a few more dead subnetworks as well (Fig. S10).

#### 3.4.3. INTERVENTION

Intervening on these circuits helps us understand how much local loss landscape is representative of global loss landscape, particularly when it comes to inactive subnetworks remaining inactive as we move through parameter space. If local loss landscape is truly representative of global loss landscape in this way, then intervening on a single subnetwork should result in only a consistent small number of samples being affected, even if we move very far in that direction. Fig. 8 shows our results for these interventions. Even in this more complex toy model, local loss landscape is a relatively good approximation of the global loss landscape. We can perturb the parameters in a direction of interest and have a large impact on the predictions of that sample and a minor impact on others. If we perturb far enough (Fig. S7), we do begin to see effects on the predictions of other samples, but ratio of change in predictions to the relevant samples to those of the irrelevant samples is very high.

A careful reader may have noticed is that it would be very easy to identify parameter directions in the first or last layer of the  $X \mapsto X^2$  model involved in a sparsely active circuit.

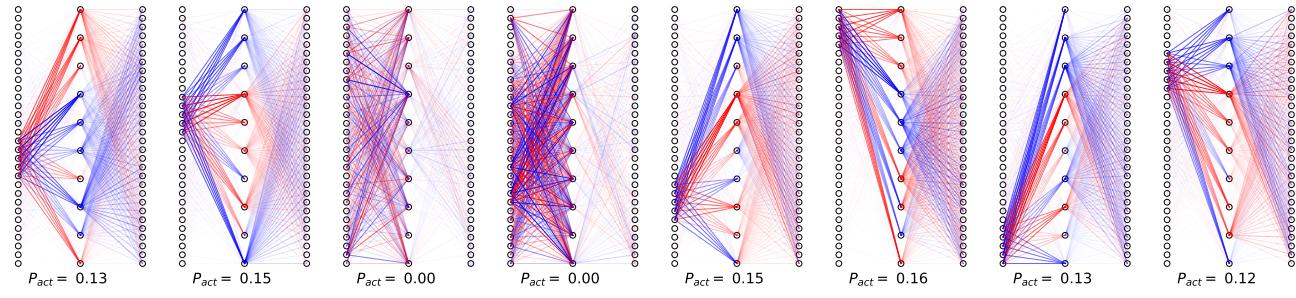


Figure 6: Parameter representations learned by L3D for the high rank circuit decomposition task.

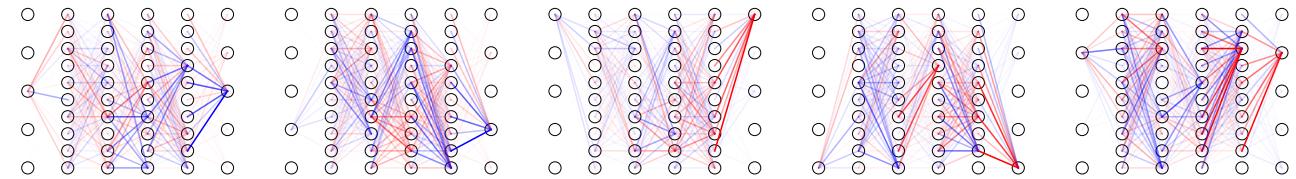


Figure 7: Subnetworks learned by L3D for the  $X \mapsto X^2$  model, and effect of intervening on each subnetwork

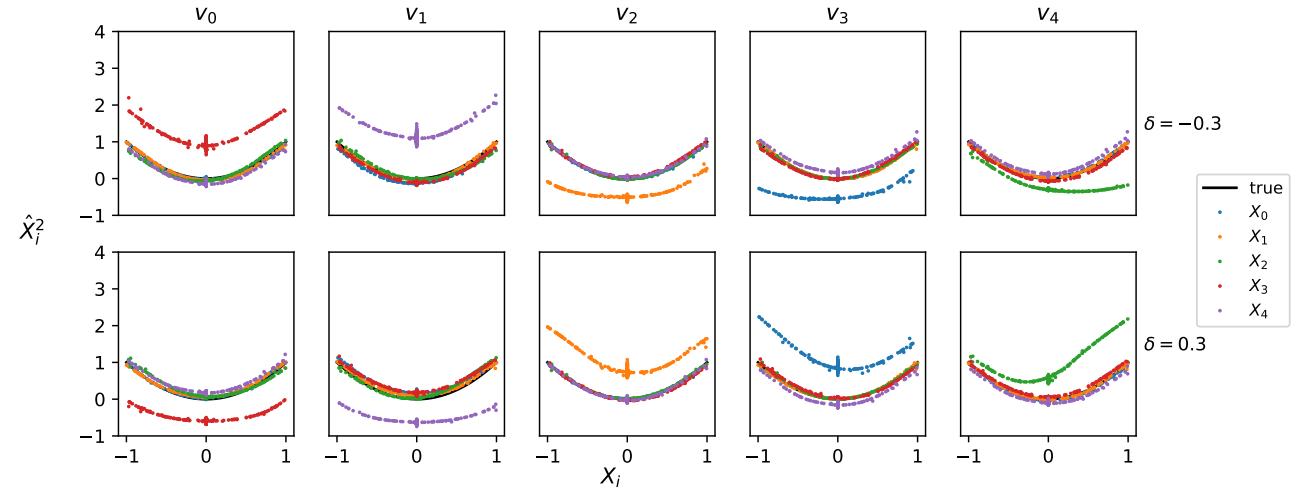


Figure 8: The effect of intervening on each subnetwork

440 The circuit for  $X_i \mapsto X_i^2$  would just involve all weights  
441 connecting input node  $i$  to the first hidden layer, and all  
442 weights connecting the hidden node  $i$  to the output node  
443  $i$ , as well as the bias of output node  $i$ . In order to make  
444 sure L3D is not just learning this trivial solution, we want to  
445 test if the hidden parameter directions it learns are also relevant  
446 circuits. Therefore, we perform the same intervention  
447 experiments but only intervening with the subnetwork components  
448 related to the model's hidden layers' weights and  
449 biases ([S8](#)). The results are very similar, where inventions  
450 on a parameter direction only affect the output of a subset  
451 of samples - showing that the circuits that L3D has learned  
452 perform relevant computation even in their middle layers.

453 [8](#) shows changes in predictions as we move in a single direction  
454 in parameter space. We also wanted to understand  
455 how subcircuits might interact with each other as we move  
456 through parameter space. In [S9](#) we perturb multiple subnetworks  
457 at once, and measure the new predictions. For the  
458 most part, the subnetworks have little inference with each  
459 other: the relevant output values for each subnetwork move  
460 relatively independently of each other. For some parameter  
461 directions, we do see some unexpected interactions, such  
462 as when we move in the direction of subnetwork 1 and 2,  
463 we see a small effect on the predictions of a few samples  
464 not associated with either subnetwork. With larger models,  
465 especially those with more modular circuits that are working  
466 together, we expect this kind of interaction to become more  
467 common. We discuss more in the discussion section.  
468

## 469 470 **4. Discussion**

### 471 472 **4.1. Simple Improvements**

#### 473 **Hyperparameter Choice**

#### 474 **Flexible Rank Choices**

### 476 477 **4.2. Challenges**

#### 478 **Interpretation of a circuit**

#### 479 **Local Loss Landscape**

#### 480 481 **Relationship to overparameterized models**

### 483 484 **4.3. Extensions**

#### 485 **Finetuning**

#### 486 **Identifying Specific Circuits with Contrastive Pairs**

487

488

489

490

491

492

493

494

**5. Impact Statement**

495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549

---

## References

- 550 Bereska, L. and Gavves, E. Mechanistic interpretability  
551 for ai safety—a review. *arXiv preprint arXiv:2404.14082*,  
552 2024.
- 553 Braun, D., Bushnaq, L., Heimersheim, S., Mendel,  
554 J., and Sharkey, L. Interpretability in parameter  
555 space: Minimizing mechanistic description length with  
556 attribution-based parameter decomposition. *arXiv  
557 preprint arXiv:2501.14926*, 2025.
- 558 Bricken, T., Templeton, A., Batson, J., Chen, B., Jermyn,  
559 A., Conerly, T., Turner, N., Tamkin, A., and Carter, S. Towards  
560 monosemanticity: Decomposing language models  
561 with dictionary learning. *Transformer Circuits Thread*,  
562 2023. Accessed: 2025-02-17.
- 563 Bushnaq, L., Mendel, J., Heimersheim, S., Braun, D.,  
564 Goldowsky-Dill, N., Hänni, K., Wu, C., and Hobbhahn,  
565 M. Using degeneracy in the loss landscape for mechanistic  
566 interpretability. *arXiv preprint arXiv:2405.10927*,  
567 2024.
- 568 Bussmann, B., Leask, P., and Nanda, N. Batchtopk sparse  
569 autoencoders. *arXiv preprint arXiv:2412.06410*, 2024.
- 570 Cunningham, H., Ewart, A., Riggs, L., Huben, R., and  
571 Sharkey, L. Sparse autoencoders find highly interpretable  
572 features in language models. *arXiv preprint  
573 arXiv:2309.08600*, 2023.
- 574 Davies, X., Langosco, L., and Krueger, D. Unifying  
575 grokking and double descent. *arXiv preprint  
576 arXiv:2303.06173*, 2023.
- 577 Elhage, N., Hume, T., Olsson, C., Schiefer, N., Henighan,  
578 T., Kravec, S., Hatfield-Dodds, Z., Lasenby, R., Drain,  
579 D., Chen, C., et al. Toy models of superposition. *arXiv  
580 preprint arXiv:2209.10652*, 2022.
- 581 Engels, J., Michaud, E. J., Liao, I., Gurnee, W., and  
582 Tegmark, M. Not all language model features are linear.  
583 *arXiv preprint arXiv:2405.14860*, 2024a.
- 584 Engels, J., Riggs, L., and Tegmark, M. Decomposing  
585 the dark matter of sparse autoencoders. *arXiv preprint  
586 arXiv:2410.14670*, 2024b.
- 587 Gao, L., la Tour, T. D., Tillman, H., Goh, G., Troll, R.,  
588 Radford, A., Sutskever, I., Leike, J., and Wu, J. Scaling  
589 and evaluating sparse autoencoders. *arXiv preprint  
590 arXiv:2406.04093*, 2024.
- 591 Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B.,  
592 Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y.  
593 Generative adversarial nets. *Advances in neural information  
594 processing systems*, 27, 2014.
- 595 Ho, J., Jain, A., and Abbeel, P. Denoising diffusion proba-  
596 bilistic models. *Advances in neural information process-  
597 ing systems*, 33:6840–6851, 2020.
- 598 Hoogland, J., Wang, G., Farrugia-Roberts, M., Carroll, L.,  
599 Wei, S., and Murfet, D. The developmental landscape  
600 of in-context learning. *arXiv preprint arXiv:2402.02364*,  
601 2024.
- 602 Lindsey, J., Templeton, A., Marcus, J., Conerly, T., Batson,  
603 J., and Olah, C. Sparse crosscoders for cross-layer fea-  
604 tures and model diffing. *Transformer Circuits Thread*,  
605 2024.
- 606 Matena, M. and Raffel, C. Npeff: Non-negative per-example  
607 fisher factorization. *arXiv preprint arXiv:2310.04649*,  
608 2023.
- 609 Merullo, J., Eickhoff, C., and Pavlick, E. Talking heads: Un-  
610 derstanding inter-layer communication in transformer lan-  
611 guage models. *arXiv preprint arXiv:2406.09519*, 2024.
- 612 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness,  
613 J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidje-  
614 land, A. K., Ostrovski, G., et al. Human-level control  
615 through deep reinforcement learning. *nature*, 518(7540):  
616 529–533, 2015.
- 617 Pascanu, R. On the difficulty of training recurrent neural  
618 networks. *arXiv preprint arXiv:1211.5063*, 2013.
- 619 Sharkey, L., Chughtai, B., Batson, J., Lindsey, J., Wu, J.,  
620 Bushnaq, L., Goldowsky-Dill, N., Heimersheim, S., Or-  
621 tega, A., Bloom, J., et al. Open problems in mechanistic  
622 interpretability. *arXiv preprint arXiv:2501.16496*, 2025.
- 623 Wang, G., Farrugia-Roberts, M., Hoogland, J., Carroll, L.,  
624 Wei, S., and Murfet, D. Loss landscape geometry re-  
625 veals stagewise development of transformers. In *High-  
626 dimensional Learning Dynamics 2024: The Emergence  
627 of Structure and Reasoning*, 2024.
- 628 Watanabe, S. Algebraic information geometry for learning  
629 machines with singularities. *Advances in neural informa-  
630 tion processing systems*, 13, 2000.
- 631 Watanabe, S. Algebraic geometry of singular learning ma-  
632 chines and symmetry of generalization and training errors.  
633 *Neurocomputing*, 67:198–213, 2005.
- 634 Watanabe, S. Almost all learning machines are singular. In  
635 *2007 IEEE Symposium on Foundations of Computational  
636 Intelligence*, pp. 383–388. IEEE, 2007.
- 637 Wei, S., Murfet, D., Gong, M., Li, H., Gell-Redman, J.,  
638 and Quella, T. Deep learning is singular, and that’s good.  
639 *IEEE Transactions on Neural Networks and Learning  
640 Systems*, 34(12):10473–10486, 2022.

---

## A. Definitions

### A.0.1. DIMENSIONS

605     $n_s$ : The number of samples in a batch of inputs  
 606     $n_f$ : The dimensions of a single input vector to a model  
 607     $n_o$ : The dimensions of a single output vector from a model  
 608     $n_w$ : The number of parameters values in a model.  
 609     $n_v$ : The number of subnetworks or parameter directions chosen to decompose a model.

### A.0.2. MODEL SYNTAX

610     $X \in \mathbb{R}^{n_s \times n_f}, x \in \mathbb{R}^{n_f}$ : Batch and individual input vectors to a model.  
 611     $W \in \mathbb{R}^{n_w}, w \in \mathbb{R}$ : The set of and individual parameter values of a model  
 612     $\sqsubseteq$ : The set of parameters corresponding to a specific tensor or block in a model.  
 613     $f : \mathbb{R}^{n_s \times n_f} \mapsto \mathbb{R}^{n_s \times n_o}$ : A model mapping a set of input vectors to a set of output vectors.  
 614     $f(X, W)$ : The output of model  $f$  with parameter values  $W$  on input  $X$ .  
 615     $f(X, W_0)$  or  $f(X)$ : The output of model  $f$  with fixed parameter values  $W_0$ .  $W_0$  is the set of learned parameter values from  
 616    model training.  
 617     $D$ : Divergence metric between two vectors. Typical divergence metrics are mean-squared error for regression-type outputs,  
 618    and KL-divergence for probability-type outputs.

### A.0.3. DECOMPOSITION SYNTAX

619     $V(\text{or } V^{out}) \in \mathbb{R}^{n_v \times n_w}, v(\text{or } V^{out}) \in \mathbb{R}^{n_w}$ : The set of or individual parameter directions that are used to decompose a  
 620    model.  $V_{out}$  can be used to transform parameter directions in the subnetwork vector space back into the original parameter  
 621    space of the model. The terms parameter direction, subnetwork, and circuit all have the same meaning.  
 622     $\mathcal{V}_{out}, \sqsubseteq_{out}$ : Components of  $V^{out}$  or  $V^{out}$  related to a specific tensor or block in the original model.  
 623     $V^{in} \in \mathbb{R}^{n_w}, V^{in} \in \mathbb{R}$ : Transforms the original parameter space of the model into the subnetwork vector space.  
 624     $\mathcal{V}_{in}, \sqsubseteq_{in}$ : Components of  $V^{in}$  or  $V^{in}$  related to a specific tensor or block in the original model.  
 625     $r$ : The rank of each component of the decomposition vectors corresponding to tensors in the original model.

### A.0.4. TRAINING

626     $\mathcal{K}$ : The subset of indexes  $i = \tau$  where  $\tau$  is the threshold computed by using the top  $k$  absolute values in a tensor.  
 627     $L$ : The L2 reconstruction loss used to optimize  $V^{in}$  and  $V^{out}$ .

### A.0.5. MEASURING AND INTERVENTION

628     $I(x_i, x_j, v_k), I(x_i, v_k)$ : The impact of subnetwork  $v_k$  on the divergence between samples  $x_i$  and  $x_j$ , or averaged across  
 629    many  $x_j$  reference samples.  
 630     $\delta$ : A scalar value to move  $W$  in a specific direction.

## B. Additional Methods

### B.1. Low-Rank Tensor Representation

631    We use low-rank representations of our  $V^{in}$  and  $V^{out}$ , and correspondingly learn low-rank circuits.

660 While  $W$  is a vector of all of the parameters in a model, typically model parameters are organized into tensors  $W = \{w_i\}_i$ .  
 661 If our parameters are organized into tensors  $W = \{w_i\}_i$ , each subnetwork or parameter component can be organized as  
 662  $V_i^{in} = \{v^{in}_i\}_i, V_i^{out} = \{v^{out}_i\}_i$  where we have parts of our subnetworks corresponding to each tensor in the original model  
 663 parameters. We wish each of these tensors to be low rank, and we express them using the canonical polyadic decomposition  
 664 () (a way to write 3+ dimentinoal tensors in terms of low-rank components).  
 665

666  
 667 
$$v^{in}_{i,j} = \sum_{r=1}^R a_{i,j,r} \times b_{i,j,r} \times c_{i,j,r} \dots \quad (8)$$
  
 668  
 669

670 Where  $R$  is the rank we wish to use to represent the parameter component, and the number of factors (a, b, c,...) in the  
 671 factorization is equal to the number of dimensions in the tensor  $w_i$ .  
 672

## 673 B.2. Finetuning

674 To finetune a model on a specific set of parameter directions, we would freeze the current set of weights and learn an  
 675 adaptor consisting of linear combinations of the subnetworks of choice. During fine tuning, we would only need to learn the  
 676 coefficients of these linear combinations greatly reducing cost.  
 677

678  
 679 
$$f(X, W_0 + W_{ft} V_{\mathcal{K}}^{out}) \quad (9)$$
  
 680

### 681 B.2.1. QUANTIFYING IMPACT

682 We can quantify the impact of a subnetwork in two ways. First, we can compute the impact of a subnetwork on a pair of  
 683 samples  $x_i, x_j$ , identifyng the subnetwork that, if intervened upon, would most strongly impact the the divergence of  $x_i$ 's  
 684 outputs when compared to  $x_j$ . The impact of subnetwork  $v_k$  on such a pair of samples,  $I(x_i, x_j, v_k)$  can be measured by:  
 685

686  
 687 
$$I(x_i, x_j, v_k) = |V_{k,:}^{in} \nabla_w D(f(x_i, W), f(x_j))| \quad (10)$$
  
 688

689 Alternatively, we can average the impacts of a subnetwork  $v_k$  and an input  $x_i$  over many different reference samples to better  
 690 quantify the impact of the subnetwork on a single sample's predictions overall. Although more computationally expensive,  
 691 this can give a more robust measurement for the impact of a subnetwork on a specific sample.  
 692

693  
 694 
$$I(x_i, v_k) = \frac{1}{n_j} \sum_{j=1}^{n_j} I(x_i, x_j, v_k) \quad (11)$$
  
 695  
 696

## 697 B.3. Toy Model Training

698 For all of our toy models (except the  $X \mapsto X^2$  model), we generate uniformly random inputs between 0 and 1. For  
 699  $X \mapsto X^2$ , we generate uniformly random inputs between -1 and 1. For all toy model data, we use a sparsity value of  
 700 1-sparsity=.05. We generate 10000 datapoints and train for 1000 epochs with batch sizes of 32. We use an AdamW optimizer  
 701 with a learning rate of 0.001.  
 702

## 703 B.4. L3D Model Training

704 To train L3D, we use the same training distributions as in each toy models. Although optimal hyperparameter values  
 705 probably depend on the model size, and the rank and number of parameter tensors, we use the same hyperparameters for all  
 706 of our models. We generate only 1000 datapoints, with a batch size of 32, and train for 1000 epochs. We use an AdamW  
 707 optimizer with a learning rate of 0.01, and a learning decay rate of .8 every 100 steps. We always use a topFrac value of .1.  
 708 We include all of the model's parameter tensors, including biases, in the decomposition.  
 709

## 710 C. Supplemental Figures

715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742

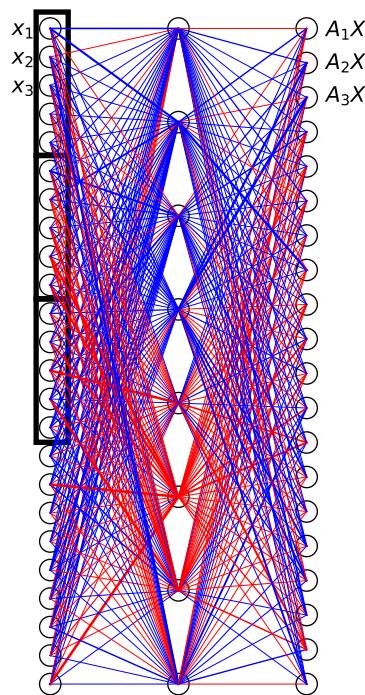


Figure S1: Full architecture of high rank circuit toy model.

743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769

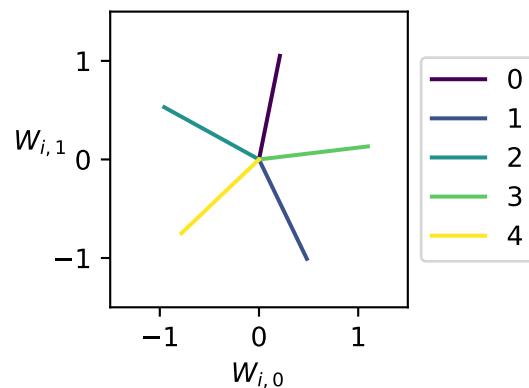


Figure S2: Encoder directions in TMS

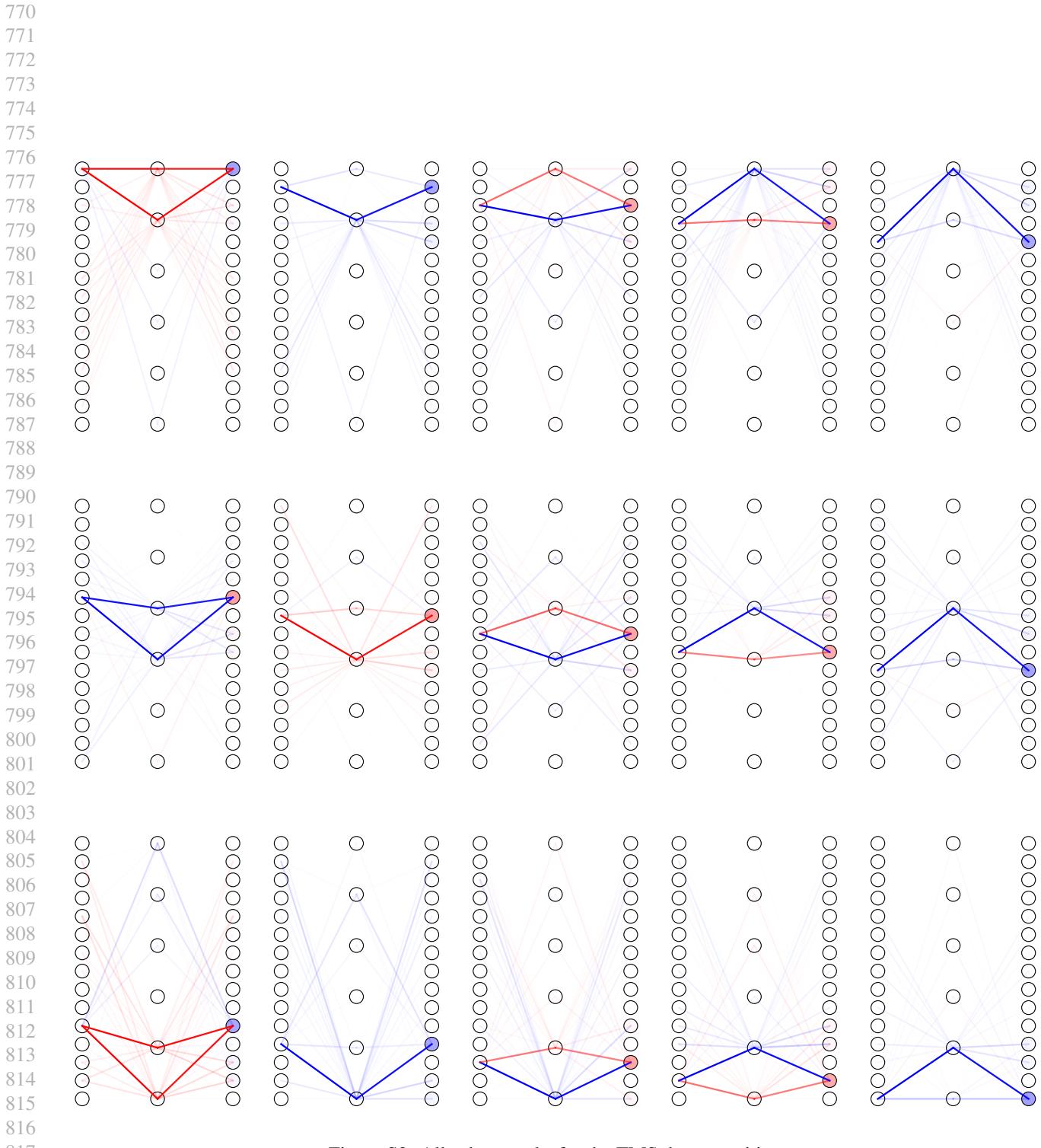
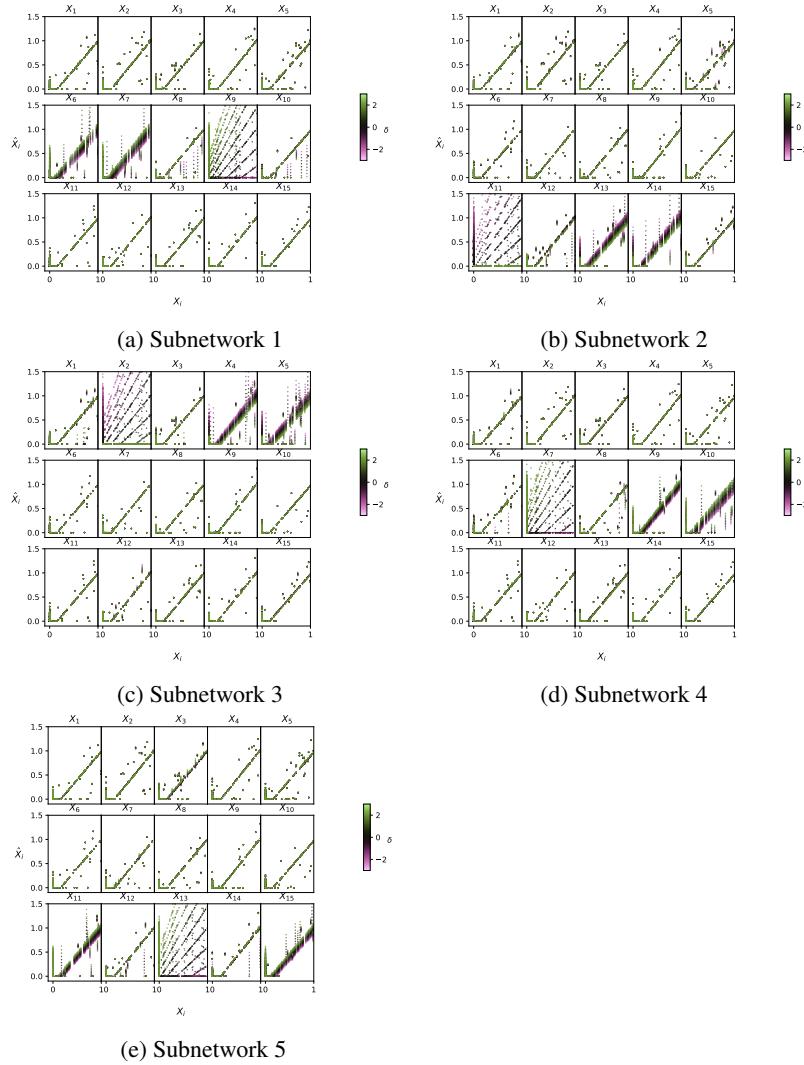


Figure S3: All subnetworks for the TMS decomposition

Figure S4: Effect of intervening on the first 5 subnetworks of the TMS-in-parallel model.



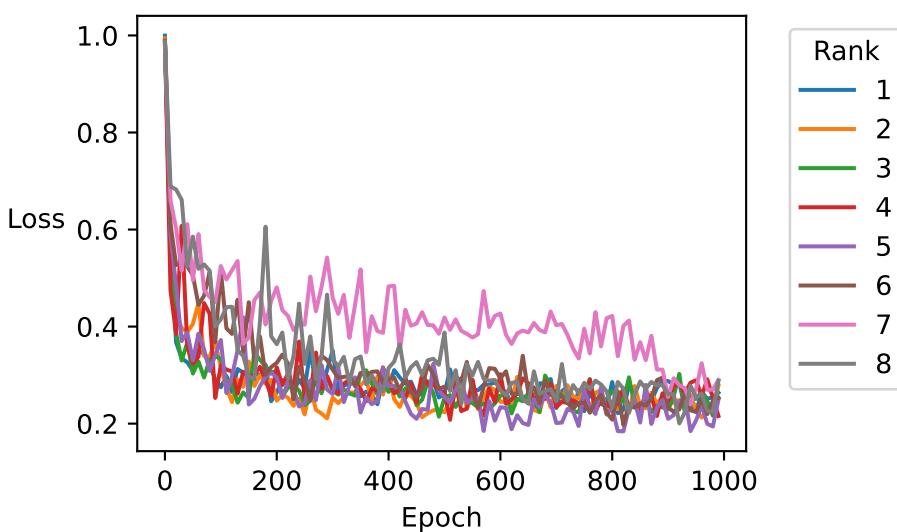
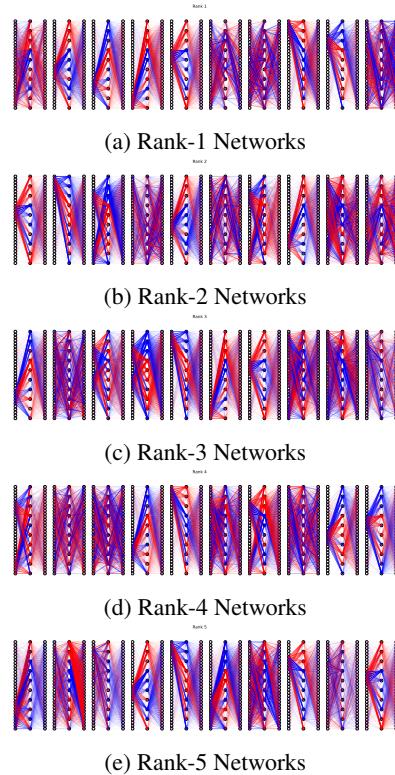


Figure S5: Loss vs Rank

Figure S6: Decomposing the toy model of high rank circuits into different numbers of subnetworks



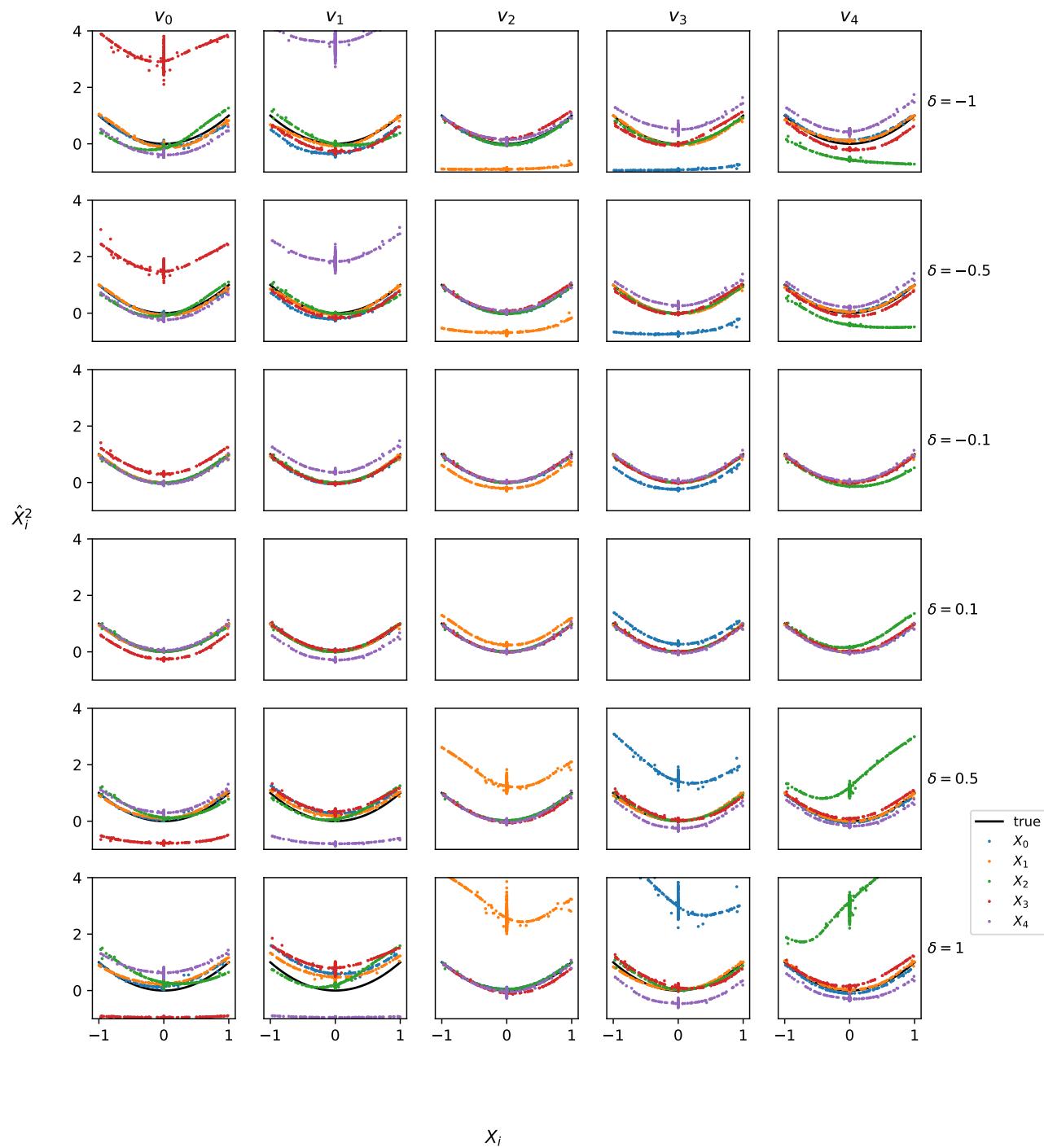


Figure S7: Effects of intervening on each of the 5 subnetworks of the  $X \mapsto X^2$  model.

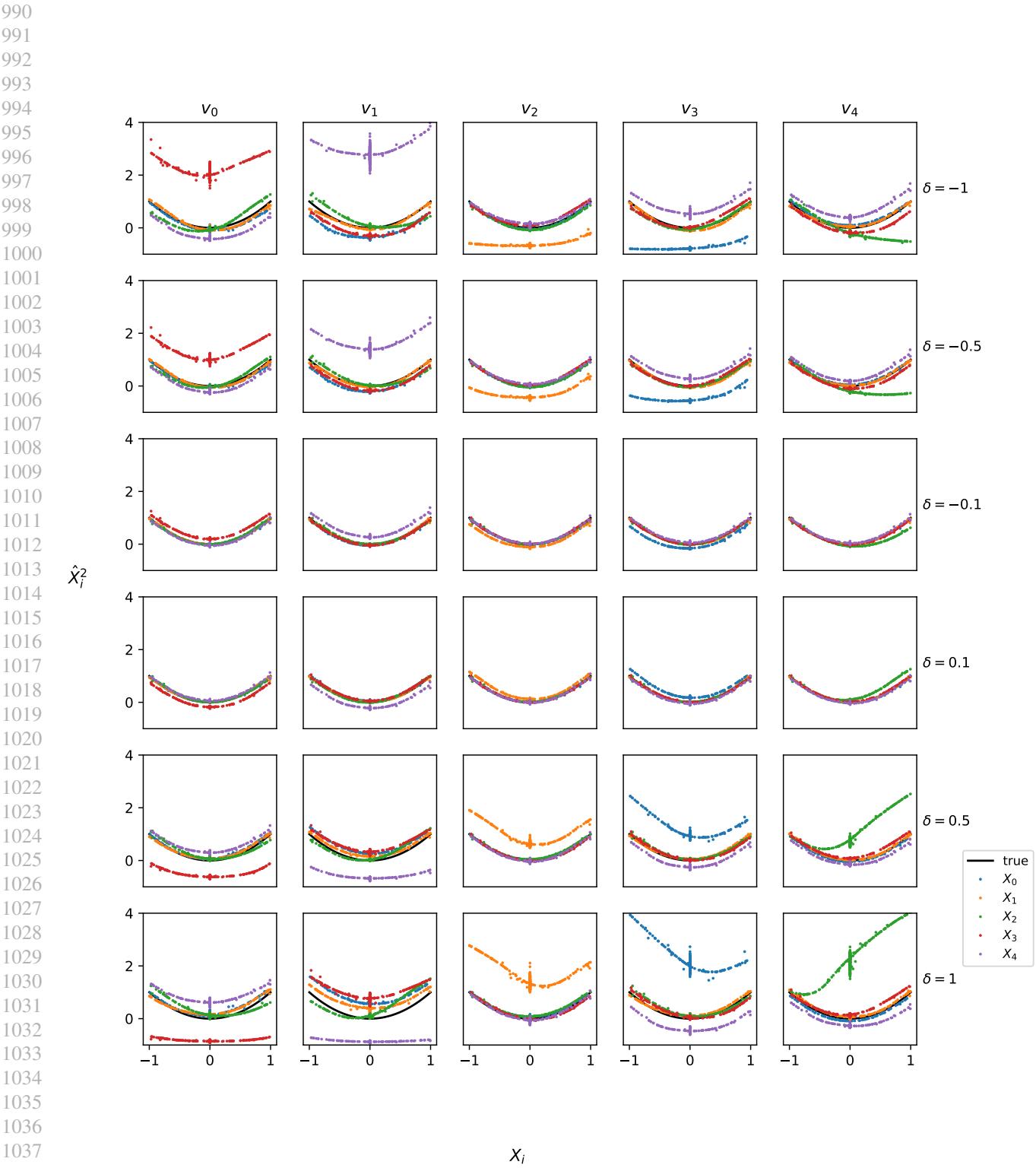


Figure S8: Intervening on just the weights and biases of the middle hidden layers in the  $X \mapsto X^2$  network.

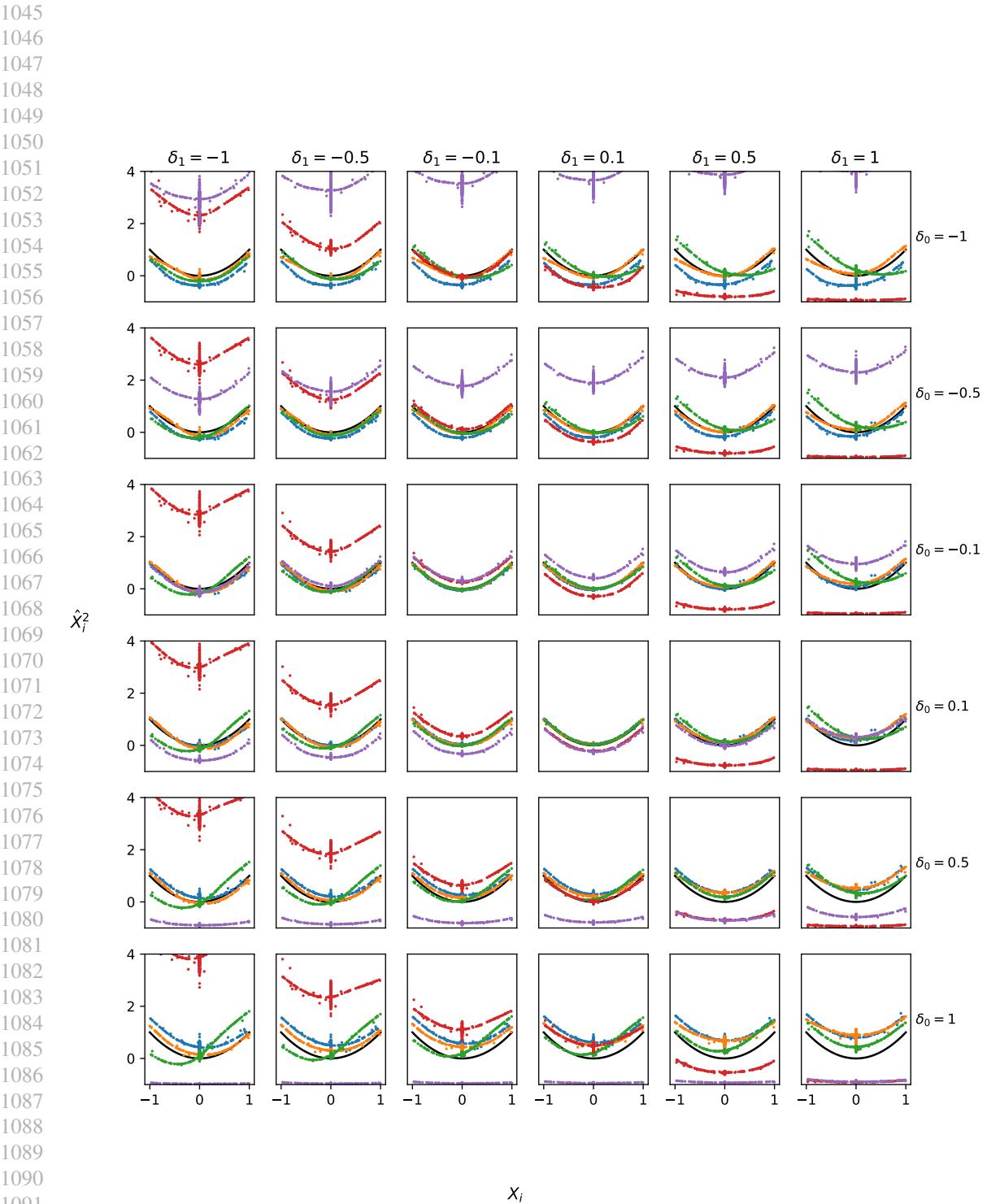
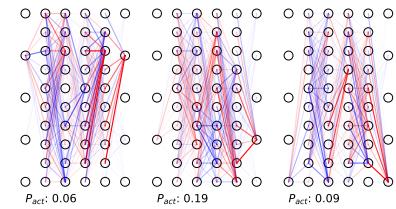
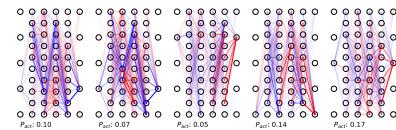


Figure S9: Intervening on multiple subnetworks at once.

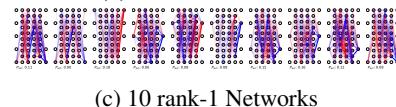
Figure S10: Decomposing the  $X \mapsto X^2$  model into different numbers of subnetworks



(a) 3 rank-1 Networks



(b) 5 rank-1 Networks



(c) 10 rank-1 Networks



(d) 15 rank-1 Networks

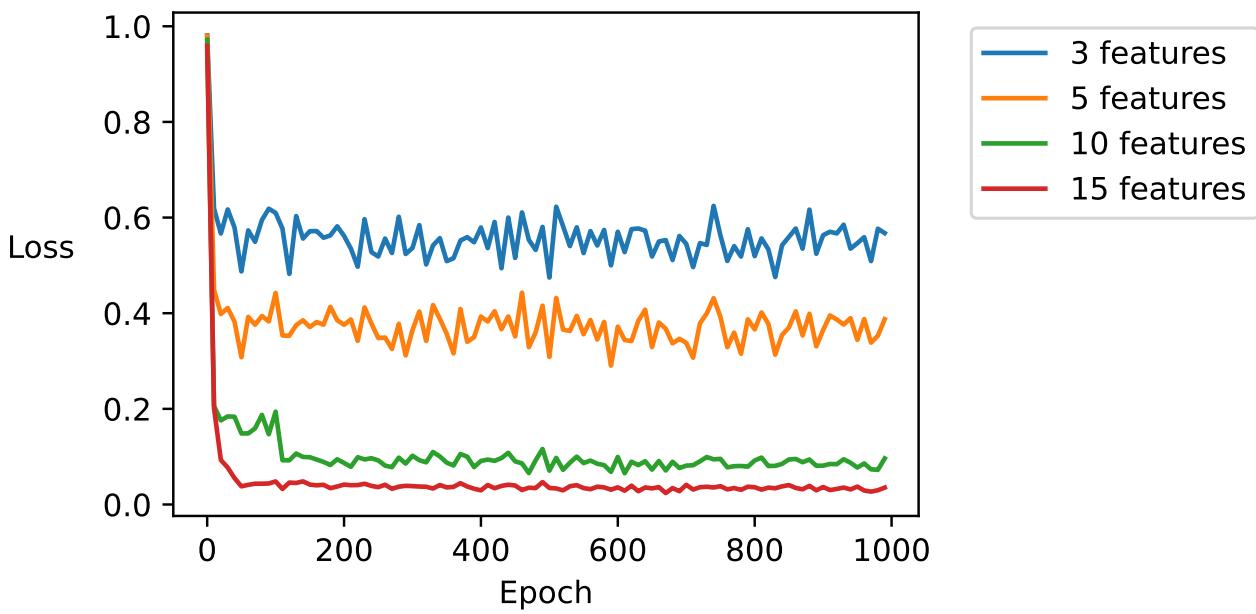


Figure S11: Training loss vs number of features for the  $X \mapsto X^2$  model, and decompositions for a 3- and 5- subnetwork decomposition.