
Identifying Sparsely Active Circuits Through Local Loss Landscape Decomposition

Anonymous Authors¹

Abstract

Much of mechanistic intepreability has focused on understanding the activation spaces of large neural networks. However, activation space-based approaches say little about the underlying circuitry used to compute the features. To better understand the underlying circuits used by models, we develop a new decomposition method called local loss landscape decomposition (L3D). We define "subnetworks" of a model to be directions in parameter space that, when intervened on, move a small set of samples' outputs towards a reference output. L3D finds a set of low-rank subnetworks, from which a smaller subset can reconstruct the gradient of the loss between any of these output pairs. We design a series of progressively more challenging toy models, with well-defined subnetworks, and show that L3D can nearly perfectly recover the associated subnetworks. We additionally investigate to what extent perturbing the model in the direction of a given subnetwork affects only the relevant subset of samples. Finally, we apply L3D to a real-world transformer model and a convolutional neural network, to demonstrate the promise of L3D identifying interpretable and relevant circuits through the parameter space.

1. Background

Mechanistic intepreability aims to uncover the internal mechanisms responsible for the behavior of large models so that developers can better understand, intervene on, and align models (?). One goal of the field is to be able to decompose model behavior into subcomponents that are less complex and more human interpretable but that, together, still fully explain model behavior. The most popular method in this space is Sparse Dictionary Learning (SDL) (???) which identifies latent features by decomposing a model's activation space into a overcomplete basis of sparsely activating components. These learned basis vectors represent distinct features, which can then be used to reconstruct the original activations.

1.1. From Activation to Parameter-Based Intepreability

However, decomposing the activation space of a model has various limitations. Current SDL algorithms struggle with reconstructing features of certain geometries (nonlinear features, feature manifolds, and certain types of superposition)(????). Such issues could become more pronounced in models with a less clearly defined read/write stream, such as diffusion models and recurrent networks. (??). Additionally, activation space captures the *features* extracted by the model's underlying circuits, but it says little about what computations or algorithms those circuits performed to derive them.

Alternatively, to understand a model's underlying *mechanisms*, we might interpret models through the lens of *parameter space*. Parameters are the fundamental objects updated during training, and can capture information about a model's internal mechanisms, the training process, and the mechanistic relationship between outputs. We hypothesize that parameter space can hold interpretable units of computation (?): models can be decomposed into simpler *subnetworks*, where each subnetwork is involved in the predictions of a subset of training data. To understand how we might go about identifying such sparsely active subnetworks, we first must understand some key insights about loss landscape geometry.

1.2. Loss Landscape Geometry

Singular learning theory describes how the structure of parameter space influences model behavior (??), and has been used to characterize model topologies (??) as well as phases of the training process (???). A key understanding of SLT is that large models are highly *degenerate* (??) in the parameter space: they can have many different parameter configurations that result in minimum loss on the training set. In fact, gradient descent tends to converge on configurations with many of these degenerate directions. Our work takes this hypothesis one step further: If, in relation to the loss on the full training distribution, models are highly degenerate, then in relation to a subset of the training data, they probably have additional subset-specific degeneracies.

The other phenomenon our method relies on is that, at least

055 in the current set of foundation models, local attribution
 056 methods seem to be good approximations of global relationships
 057 between pairs of samples. For example, attribution
 058 patching can successfully change the output of a model
 059 by targeting specific activations, determined by the first order
 060 gradients of paired outputs (??). Similarly, steering
 061 vectors, which are derived from differences in activations
 062 from paired samples, can effectively guide models toward
 063 specific behaviors even when applied beyond the original
 064 magnitude of activation differences (??). Local and global
 065 attribution methods seem to be reasonable approximations
 066 of each other.
 067

1.3. Loss Landscape Decomposition

068 Our goal in this paper is to identify directions in parameter
 069 space that correspond to subnetworks, as defined in ??.

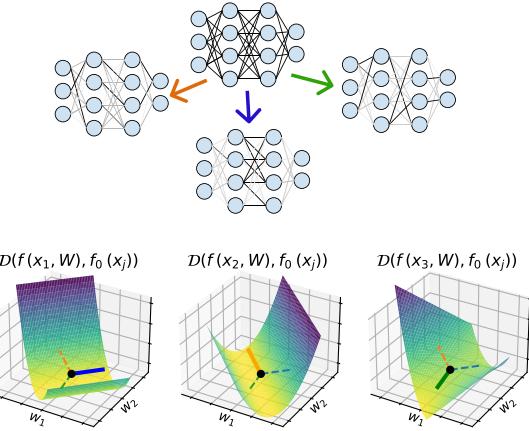
070 Two methods in particular tackle the problem of decomposing
 071 models into a similar definition of subnetworks. An
 072 earlier work (?), decomposes parameter space by computing
 073 principal directions of a per-sample Fisher Information
 074 matrix to resolve meaningful features. A recent method,
 075 Attribution Parameter Decomposition (?) decomposes model
 076 weights by identifying subnetworks where (1) the sum of
 077 subnetwork weights approximates the original model para-
 078 meters, (2) for any given input, the outputs of the sum of
 079 topk-attributed networks has low behavioral loss when
 080 compared to those of the original model, and (3) subnetworks
 081 are individually simpler than the whole network.
 082

083 Like (?) and (?), we seek to identify subnetworks that
 084 when perturbed or ablated, would change a set of samples'
 085 overall predictions. (?) compute gradient attributions across
 086 output indexes to identify topK subnetworks, and then use
 087 behavioral loss compare the output of a sum of subnetworks
 088 to that of the original model. (?) work with second-order
 089 gradients to understand the contribution of a subnetwork to
 090 overall predictions.
 091

092 We use a slightly different approach: We quantify a change
 093 in output by its distance moved in the direction of a
 094 specific "reference" output.
 095

096 As we will describe in more detail in ??, this approach
 097 allows us to work with first-order gradients, unlike in (?).
 098 We also get to work with a scalar metric for loss, rather than
 099 considering loss across each element in an output vector as
 100 in (?).

101 The object we decompose is therefore the **gradient (with**
respect to the model parameters) of the loss between a
sample's output and a "reference" output. (In practice,
 102 we choose the reference output as another output sampled
 103 from the training distribution, and discuss our reasoning in
 104 ??.). We seek to identify low-rank directions in parameter
 105 space, which we hereon refer to as *subnetworks* such that for
 106



107 Figure 1: Decomposing a loss landscape into a set of parame-
 108 ter directions, or subnetworks, where a smaller subset of
 109 directions can approximately reconstruct any per-sample-
 110 pair gradient of divergence Here, D is a loss, or *divergence*
 111 measure, f is our model, W is the set of parameters in the
 112 model, x_i is a sample input, and y_r is our reference output

113 any pair of samples, a small number of these directions can
 114 be used to reconstruct this gradient (Figure ??). We call our
 115 decomposition method Local Loss Landscape Decomposition
 116 (L3D). In this work, we first describe the mathematical
 117 foundation of our approach. We develop progressively more
 118 complex toy models to measure the efficacy of L3D and
 119 characterize its limitations. Finally, we briefly show some
 120 preliminary results on real world model, to demonstrate
 121 L3D's promise in scaling to real world models.

2. Methodology

122 In the next sections, we will formally set up our decom-
 123 position problem (Section ??), define the criteria that we
 124 will use for our subnetwork/parameter directions (Section
 125 ??), describe how to efficiently decompose parameters into
 126 these directions (Section ??), walk through our training al-
 127 gorithm (Section ??), and then explain how to use these
 128 decompositions to intervene on a model's behavior (Section
 129 ??).

130 From now on, we will use the word "subnetworks" to refer
 131 to the directions in parameter space we wish to learn.

2.1. Set up

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

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

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

113 Our approach rests on the hypothesis that, for a given input,
 114 there are many components of a model’s parameters that are
 115 not involved in inference. Changing parameters in the direc-
 116 tion of these components will not affect change the model’s
 117 output. Conversely, changing parameters in the direction of
 118 components that *are* involved *would* change the model’s out-
 119 put. We are interested in finding parameter directions that,
 120 when perturbed, *meaningfully* change a model’s output. The
 121 next subsection will explain what constitutes “meaningful”.

123 2.2. Divergences of Paired Outputs

125 To restate point (3), intervening in a relevant parameter
 126 direction should move a sample’s output either closer to
 127 or further from a reference output. This reference output
 128 should serve as a neutral and representative baseline that
 129 captures the typical behavior of the model’s output distribu-
 130 tion. We considered three candidates for this reference:

- 132 1. **A uniform output:** This reference consists of a vec-
 133 tor with uniform values. However, it fails to account
 134 for the training distribution, leading to a bias toward
 135 learning subnetworks that influence outputs that skew
 136 toward particularly high or low values.
- 137 2. **Mean of outputs:** This reference is computed by av-
 138 eraging each output across the training distribution or
 139 a batch. While it is grounded in the data, it risks aver-
 140 aging away meaningful correlations between outputs,
 141 producing a reference that may be out-of-distribution
 142 relative to the training data.
- 143 3. **Another sample as the reference:** For each sample,
 144 we use the output of a randomly selected sample as the
 145 reference. This approach preserves the nuances of the
 146 output distribution but may lead to slow convergence
 147 due to high variance in reference selection.

150 We thought (3) was the most principled, and least biased of
 151 the three. Although not tested rigorously, in early prototypes
 152 all three choices seemed to produce reasonable results on
 153 toy models and we did not find any issues with convergence
 154 using (3). For this work we use (3) as our reference output,
 155 but we believe other choices are possible and may have
 156 different strengths and weaknesses.

158 Therefore, we will decompose gradients of the loss between
 159 pairs of outputs with the aim of finding directions that move
 160 the model’s output towards or away from our reference.
**Because we use the term “loss” later on when we describe
 161 our training process, we will refer to this metric instead
 162 as “divergence.”**

Divergence of a pair of outputs can be defined as:

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

where $x_i, x_r \in X$

Here, D is a divergence measure, f is our model, x_i is our input of interest, y_r is a reference output, W is a set of parameters and W_0 are model’s original parameters. Our toy models are regression-type models, so we use MSE as divergence. For the real-world transformer and CNN models, we use KL-divergence.

We will abbreviate the expression in Eq. ?? as $\nabla_W D$.

123 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.

$$\nabla_W D \approx V^{out} \Lambda V^{in} \nabla_V D \quad (4)$$

$$(5)$$

$$\text{where } \Lambda_{i,j} = \begin{cases} 1 & \text{if } i = j \text{ and } i \in \underset{i}{\text{argtopK}} [|\nabla_{v_i} D|] \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

$$(7)$$

topk relies on a hyperparameter that controls the number of components we wish to use to reconstruct each sample. In practice, we use a batchTopK (?) and a fraction rather than an absolute number. A topK of 0.1 it means that we select the top 10% of $\nabla_V D$ magnitudes over v and x 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 because they are processing smaller numbers of features. Therefore, we use low-rank representations of our

165 V^{in} and V^{out} , and correspondingly learn low-rank circuits
 166 (Appendix ??).

168 2.4. Training

170 We wish to learn the decomposition-related transforms V^{in}
 171 and V^{out} that minimize the topK reconstruction loss of our
 172 divergence gradient described above. We use a (normalized)
 173 L2 norm loss.

$$175 \quad L = \frac{\|\nabla_W D - V_{:, \mathcal{K}}^{out} V_{\mathcal{K}, :}^{in} \nabla_W D\|_2}{\|\nabla_W D\|_2} \quad (8)$$

177 For each batch of samples, we randomly select a reference
 178 sample x_r to be paired with each sample x_i in the batch. We
 179 then compute the gradient of divergence of x_i and x_r at the
 180 target model's parameters W_0 . We transform that gradient
 181 into the subnetwork space using V^{in} , and compute the topK
 182 components. We transform those components back into the
 183 original parameter space using V^{out} , and compute the loss
 184 between the reconstructed gradient and the original gradient.
 185 We apply a learning update to V^{in} and V^{out} with the goal
 186 of minimizing this loss. We also normalize V^{out} to be a
 187 unit vector after each update keep the magnitudes of V^{in}
 188 and V^{out} similar.

190 **Algorithm 1** Algorithm for decomposing parameter space
 191 using L3D.

```
192 1: for each epoch do
193 2:   for each  $X$  do
194 3:     for each  $x_i \in X$  do
195 4:       Randomly select  $x_r \in X$ 
196 5:        $\nabla_w D_i = \nabla_w D(f(x_i, W), f(x_r))|_{W=W_0}$ 
197 6:     end for
198 7:      $\nabla_v D = V^{in} \nabla_w$ 
199 8:      $\tau = \text{topK}(\text{abs}(\nabla_v D))$ 
200 9:      $\hat{\nabla}_w D = V^{out} (\nabla_v D \odot (\text{abs}(\nabla_v D) > \tau))$ 
201 10:     $L = \|\nabla_w D - \hat{\nabla}_w D\|_2$ 
202 11:     $L.\text{backward}()$ 
203 12:    Update  $V^{in}$  and  $V^{out}$ 
204 13:    Normalize  $V^{out}$  to be unit vectors.
205 14:  end for
206 15: end for
```

209 2.5. Measuring and Intervening

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

215 If we wish to intervene using a single subnetwork, we can
 216 update the parameters in by moving them in a scalar factor
 217 (δ) of that unit direction. To tune our model in the direc-
 218 tion of subnetwork v_i and compute predictions on x , we

219 evaluate:

$$220 \quad f(x, W + \delta v_i) \quad (9)$$

222 We also may want to quantify the impact of a subnetwork in
 223 on a certain sample. First, we can compute the impact of a
 224 subnetwork on a sample's (x_i) divergence with a reference
 225 output y_j . The impact of subnetwork v_k , $I(x_i, x_j, v_k)$ can
 226 be measured by:

$$227 \quad I(x_i, y_j, v_k) = |V_{k,:}^{in} \nabla_w D(f(x_i, W), y_j)| \quad (10)$$

229 Because we are randomly sampling outputs from our training
 230 distribution as the reference output, we then average
 231 the impacts of a subnetwork v_k and an input x_i over many
 232 different reference samples to better quantify the impact of
 233 the subnetwork on a single sample's predictions overall. Al-
 234 though more computationally expensive, this gives a more
 235 robust measurement for the impact of a subnetwork on a
 236 specific sample.

$$237 \quad I(x_i, v_k) = \frac{1}{n_j} \sum_{j=1}^{n_j} I(x_i, y_j, v_k) \quad (11)$$

3. Results

239 To evaluate L3D's ability to decompose models, we focus
 240 on developing toy models that involve well-defined subnet-
 241 works.

243 We designed several toy models to test the efficacy of L3D.
 244 Our toy models all consist of several well-characterized
 245 computations being performed by the same model, with an
 246 input space designed to rely on a single or small set of tasks.

248 Our toy models progressively test more complex types of
 249 circuits. Table ?? describes our 4 toy models and the differ-
 250 ent attributes of circuitry the are designed to capture. The
 251 specific hyperparameters used to train our toy models is
 252 described in Appendix ??, as well as the hyperparameters
 253 used for each decomposition in Appendix ??

3.1. Toy Model of Superposition

3.1.1. SETUP

255 We start off by validating our algorithm on a well-studied
 256 toy problem, the toy model of superposition (TMS). TMS is
 257 simple linear autoencoder with a low-dimensional hidden
 258 layer followed by a ReLU activation function at the output
 259 (?). The model is trained on a dataset of samples where
 260 few features are active at a time, and “superimposes” these
 261 features in the hidden layer such that features embeddings in
 262 the hidden layer have minimal interference with each other
 263 ???. We train a toy model of superposition with 5 features

Local Loss Landscape Decomposition

| | Toy model of superposition | Circuit Superposition (TMCS) | Higher Rank Circuit Superposition | Complex Loss Landscape |
|-----|---|------------------------------|-----------------------------------|------------------------|
| 220 | | | | |
| 221 | | | | |
| 222 | | | | |
| 223 | | | | |
| 224 | | | | |
| 225 | | | | |
| 226 | | | | |
| 227 | | | | |
| 228 | | | | |
| 229 | | | | |
| 230 | | | | |
| 231 | | | | |
| 232 | Feature Superposition | $X \mapsto X$ | $X \mapsto AX$ | $X \mapsto AX$ |
| 233 | | ✓ | ✓ | ✓ |
| 234 | Circuit Superposition | ✗ | ✓ | ✓ |
| 235 | Circuits > rank-2 | ✗ | ✗ | ✓ |
| 236 | Complex Loss Landscape | ✗ | ✗ | probably ✓ |
| 237 | | | | |
| 238 | | | | |
| 239 | | | | |
| 240 | | | | |
| 241 | and 2 hidden dimensions (with sparsity=.05) to test L3D's ability to resolve models with superimposed features. | | | |
| 242 | | | | |
| 243 | | | | |
| 244 | 3.1.2. DECOMPOSITION | | | |
| 245 | | | | |
| 246 | We decompose the TMS model into 5 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. Figure ?? shows the 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 θ 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 ??). | | | |
| 247 | | | | |
| 248 | | | | |
| 249 | | | | |
| 250 | | | | |
| 251 | | | | |
| 252 | | | | |
| 253 | | | | |
| 254 | | | | |
| 255 | | | | |
| 256 | | | | |
| 257 | | | | |
| 258 | L3D successfully decomposes the model into subnetworks corresponding to the encoding and decoding of each feature (a $X_i : \hat{X}_i$ circuit). Moreover, the encoder directions of the learned subnetworks are nearly perfectly parallel to the original encoding of each input (Figure ADD NEW FIGURE IN. This decomposition resulted in a reconstruction error of 19%. The reconstruction error is related to 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 the next higher-dimensional toy model where the reconstruction loss is in fact lower. | | | |
| 259 | | | | |
| 260 | | | | |
| 261 | | | | |
| 262 | | | | |
| 263 | | | | |
| 264 | | | | |
| 265 | | | | |
| 266 | | | | |
| 267 | | | | |
| 268 | | | | |
| 269 | | | | |
| 270 | | | | |
| 271 | | | | |
| 272 | | | | |
| 273 | | | | |
| 274 | | | | |
| 275 | | | | |
| 276 | | | | |
| 277 | | | | |
| 278 | | | | |
| 279 | | | | |
| 280 | | | | |
| 281 | | | | |
| 282 | | | | |
| 283 | | | | |
| 284 | | | | |
| 285 | | | | |
| 286 | | | | |
| 287 | | | | |
| 288 | | | | |
| 289 | | | | |
| 290 | | | | |
| 291 | | | | |
| 292 | | | | |
| 293 | | | | |
| 294 | | | | |
| 295 | | | | |
| 296 | | | | |
| 297 | | | | |
| 298 | | | | |
| 299 | | | | |
| 300 | | | | |
| 301 | | | | |
| 302 | | | | |
| 303 | | | | |
| 304 | | | | |
| 305 | | | | |
| 306 | | | | |
| 307 | | | | |
| 308 | | | | |
| 309 | | | | |
| 310 | | | | |
| 311 | | | | |
| 312 | | | | |
| 313 | | | | |
| 314 | | | | |
| 315 | | | | |
| 316 | | | | |
| 317 | | | | |
| 318 | | | | |
| 319 | | | | |
| 320 | | | | |
| 321 | | | | |
| 322 | | | | |
| 323 | | | | |
| 324 | | | | |
| 325 | | | | |
| 326 | | | | |
| 327 | | | | |
| 328 | | | | |
| 329 | | | | |
| 330 | | | | |
| 331 | | | | |
| 332 | | | | |
| 333 | | | | |
| 334 | | | | |
| 335 | | | | |
| 336 | | | | |
| 337 | | | | |
| 338 | | | | |
| 339 | | | | |
| 340 | | | | |
| 341 | | | | |
| 342 | | | | |
| 343 | | | | |
| 344 | | | | |
| 345 | | | | |
| 346 | | | | |
| 347 | | | | |
| 348 | | | | |
| 349 | | | | |
| 350 | | | | |
| 351 | | | | |
| 352 | | | | |
| 353 | | | | |
| 354 | | | | |
| 355 | | | | |
| 356 | | | | |
| 357 | | | | |
| 358 | | | | |
| 359 | | | | |
| 360 | | | | |
| 361 | | | | |
| 362 | | | | |
| 363 | | | | |
| 364 | | | | |
| 365 | | | | |
| 366 | | | | |
| 367 | | | | |
| 368 | | | | |
| 369 | | | | |
| 370 | | | | |
| 371 | | | | |
| 372 | | | | |
| 373 | | | | |
| 374 | | | | |
| 375 | | | | |
| 376 | | | | |
| 377 | | | | |
| 378 | | | | |
| 379 | | | | |
| 380 | | | | |
| 381 | | | | |
| 382 | | | | |
| 383 | | | | |
| 384 | | | | |
| 385 | | | | |
| 386 | | | | |
| 387 | | | | |
| 388 | | | | |
| 389 | | | | |
| 390 | | | | |
| 391 | | | | |
| 392 | | | | |
| 393 | | | | |
| 394 | | | | |
| 395 | | | | |
| 396 | | | | |
| 397 | | | | |
| 398 | | | | |
| 399 | | | | |
| 400 | | | | |
| 401 | | | | |
| 402 | | | | |
| 403 | | | | |
| 404 | | | | |
| 405 | | | | |
| 406 | | | | |
| 407 | | | | |
| 408 | | | | |
| 409 | | | | |
| 410 | | | | |
| 411 | | | | |
| 412 | | | | |
| 413 | | | | |
| 414 | | | | |
| 415 | | | | |
| 416 | | | | |
| 417 | | | | |
| 418 | | | | |
| 419 | | | | |
| 420 | | | | |
| 421 | | | | |
| 422 | | | | |
| 423 | | | | |
| 424 | | | | |
| 425 | | | | |
| 426 | | | | |
| 427 | | | | |
| 428 | | | | |
| 429 | | | | |
| 430 | | | | |
| 431 | | | | |
| 432 | | | | |
| 433 | | | | |
| 434 | | | | |
| 435 | | | | |
| 436 | | | | |
| 437 | | | | |
| 438 | | | | |
| 439 | | | | |
| 440 | | | | |
| 441 | | | | |
| 442 | | | | |
| 443 | | | | |
| 444 | | | | |
| 445 | | | | |
| 446 | | | | |
| 447 | | | | |
| 448 | | | | |
| 449 | | | | |
| 450 | | | | |
| 451 | | | | |
| 452 | | | | |
| 453 | | | | |
| 454 | | | | |
| 455 | | | | |
| 456 | | | | |
| 457 | | | | |
| 458 | | | | |
| 459 | | | | |
| 460 | | | | |
| 461 | | | | |
| 462 | | | | |
| 463 | | | | |
| 464 | | | | |
| 465 | | | | |
| 466 | | | | |
| 467 | | | | |
| 468 | | | | |
| 469 | | | | |
| 470 | | | | |
| 471 | | | | |
| 472 | | | | |
| 473 | | | | |
| 474 | | | | |
| 475 | | | | |
| 476 | | | | |
| 477 | | | | |
| 478 | | | | |
| 479 | | | | |
| 480 | | | | |
| 481 | | | | |
| 482 | | | | |
| 483 | | | | |
| 484 | | | | |
| 485 | | | | |
| 486 | | | | |
| 487 | | | | |
| 488 | | | | |
| 489 | | | | |
| 490 | | | | |
| 491 | | | | |
| 492 | | | | |
| 493 | | | | |
| 494 | | | | |
| 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 | | | | |
| 550 | | | | |
| 551 | | | | |
| 552 | | | | |
| 553 | | | | |
| 554 | | | | |
| 555 | | | | |
| 556 | | | | |
| 557 | | | | |
| 558 | | | | |
| 559 | | | | |
| 560 | | | | |
| 561 | | | | |
| 562 | | | | |
| 563 | | | | |
| 564 | | | | |
| 565 | | | | |
| 566 | | | | |
| 567 | | | | |
| 568 | | | | |
| 569 | | | | |
| 570 | | | | |
| 571 | | | | |
| 572 | | | | |
| 573 | | | | |
| 574 | | | | |
| 575 | | | | |
| 576 | | | | |
| 577 | | | | |
| 578 | | | | |
| 579 | | | | |
| 580 | | | | |
| 581 | | | | |
| 582 | | | | |
| 583 | | | | |
| 584 | | | | |
| 585 | | | | |
| 586 | | | | |
| 587 | | | | |
| 588 | | | | |
| 589 | | | | |
| 590 | | | | |
| 591 | | | | |
| 592 | | | | |
| 593 | | | | |
| 594 | | | | |
| 595 | | | | |
| 596 | | | | |
| 597 | | | | |
| 598 | | | | |

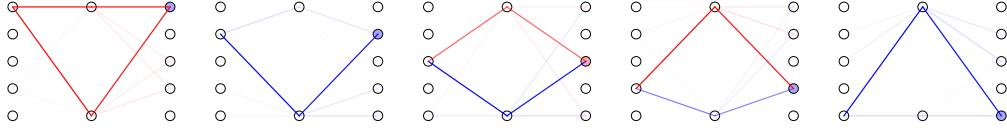


Figure 2: L3D subnetwork decomposition of TMS

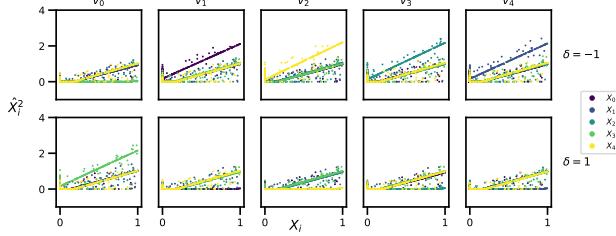


Figure 3: The effect of intervening on the TMS-in-parallel model in the direction of each subnetwork. We generate 1000 inputs from the TMS input distribution (x-axis), intervene on each subnetwork with magnitude δ (color scale) and measure the change in outputs (y-axis) for each sample.

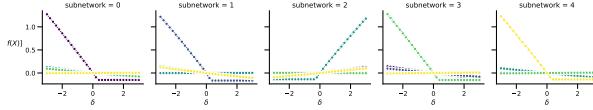


Figure 4: The effect of intervening at various values of δ in the direction of each subnetwork. Each datapoint represents the average TODO FILL IN

that can be contained in a given set of parameters. 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 (Toy model 2 in ??) 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 set the entries of A as uniform random values between 0 and 3 and generate input-output pairs to train the toy model with. We use a model with 10 inputs, 5 hidden layers, and 10 output features (although such a model does not need to have the same number of input and outputs).

If subnetworks are only relevant to a small set of inputs, then we would expect each subnetwork to be the weights and biases associated with a single input feature. If this is the case, then individual parameters would be involved in multiple subnetworks: $W^{dec}_{i,1}$ (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 subnetworks will interfere with each other - parameter directions associated with each subnetwork will be non-orthogonal.

3.2.2. DECOMPOSITION

We decompose TMCS into 10 subnetworks of rank-1 parameter tensors (??) with a reconstruction loss of 6.4%. The subnetworks each strongly correspond to a single input feature, as desired.

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 :

$$j^* = \underset{j}{\operatorname{argmax}} \|W^{dec}_{j,k}\|_2 \quad (12)$$

$$\hat{a}_{i,j^*} = W^{enc}_{i,j^*} W^{dec}_{i,j^* k}$$

Recall 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 ??, 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.

330
331
332
333
334
335
336
337
338
339
340

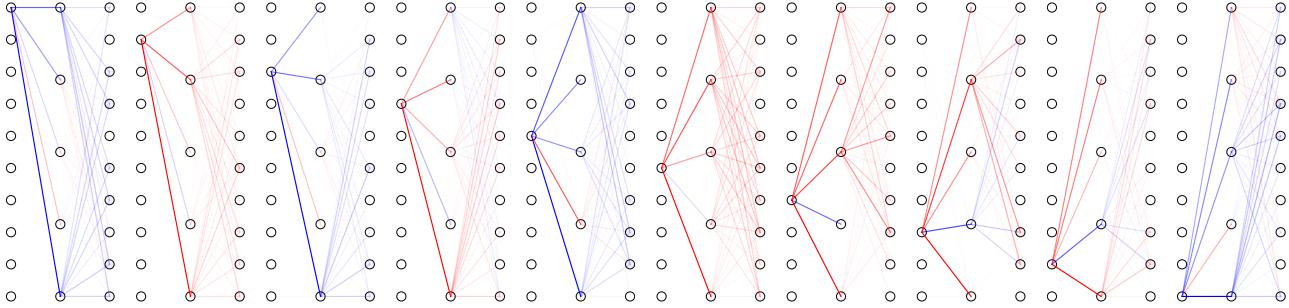


Figure 5: The subnetworks L3D decomposes the TMCS model into.

341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356

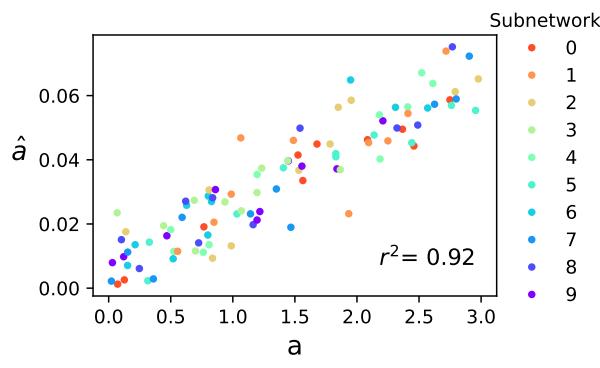


Figure 6: The coefficients derived from the subnetworks compared to actual coefficients used to train TMCS.

357
358
359
360
361
362

3.3.2. DECOMPOSITION

363 Although we expect the model to have 5 subnetworks, we
364 use excess parameter tensors ($n=10$) in order to allow more
365 flexibility in learning. We track the fraction of inputs for
366 which a subnetwork was used in the topK reconstruction
367 (P_{act}) to identify which were “dead circuits”, and report the
368 last epoch. Furthermore, although we expect the underlying
369 subcircuits to be rank 5, we experiment with using different
370 rank representations to see how well lower-rank parameter
371 directions can represent the model. Interestingly, rank-1
372 representations of the parameter tensors are able to represent
373 the model nearly as well as rank-5 representations (Figure
374 ??). In Figure ??, we show the decomposition of a rank-
375 3 decomposition. L3D successfully learns a subnetwork
376 corresponding to each of the 5 sets of input features, as well
377 as a number of dead circuits. The higher and lower rank
378 decompositions also learn similar subnetworks (Figure ??).
379 When we trained L3D without these additional subnetworks,
380 the reconstruction loss gets caught in local minima. Similar
381 to training sparse autoencoders (?), having extra degrees
382 of freedom allows for better learning, even if at the end of
383 training the extra subnetworks are never active.
384

3.4. Toy model with Complex Loss Landscape

3.4.1. SETUP

In the previous models, assuming the the final layer ReLUs have fired, the gradient of divergence with respect to to an element of W^{dec} looks like:

$$\nabla_{W_{i,j}^{dec}} (XW^{enc}W^{dec} + b - X)^2 = (2(XW^{enc}W^{dec} + b - X)^T XW^{enc})_{i,j} \quad (13)$$

The gradient is linear with respect to W^{dec} . With this gradient, we expect local approximations to be excellent approximations of global loss landscapes (up until the Relu-induced discontinuities.) We should be able to move relatively far in parameter space (large δ in Figure ??) and see well behaved effects on predictions. However, we

We therefore trained a multi-layer model to predict multiple non-linear functions of input features at once. We train a GeLU network for $X_i \rightarrow X_i^2$. 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, there are not inherent reasons to believe subnetworks should be low rank, the way there was in the TMS model.

3.4.2. DECOMPOSITION

To allow for slightly higher rank subnetworks but still compress the dimensions of the model, we decompose our model into 5 rank-2 parameter tensors in our decomposition. Additionally, instead of varying rank, we experiment with using different numbers of subnetworks to represent our model. In the 5-subnetwork decomposition (Figure ??), we see that subnetworks tracing the path of $X_i \mapsto X_i^2$ for each index i . However, this decomposition has a relatively high reconstruction error of 0.32. Much of this is probably because

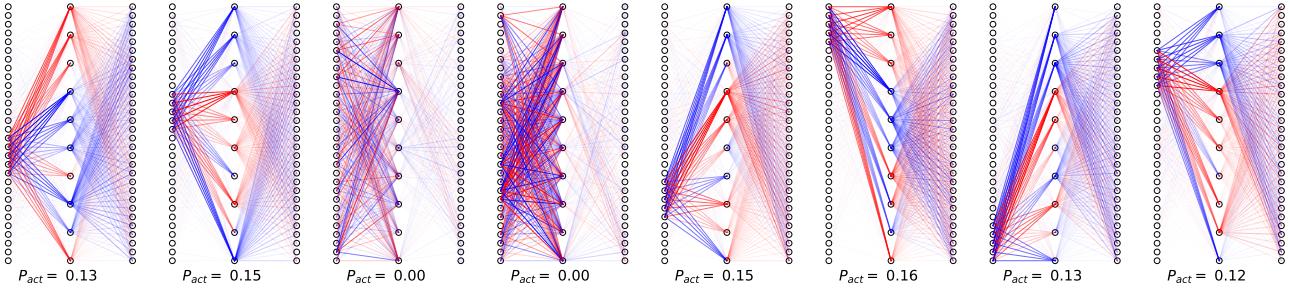


Figure 7: Parameter representations learned by L3D for the high rank circuit decomposition task.

we kept our topK hyperparameter constant (at $k = 0.1$) throughout all our models for consistency. With only 5 subnetworks, this means that each sample’s reconstruction will use 1 subnetwork on average, limiting the reconstruction error the network can achieve.

We also experiment with holding rank constant (we drop to rank-1 for this) 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 ??). Each decomposition continues to learn input-output pair specific subnetworks, with the larger decompositions resulting in a few more dead subnetworks as well (Figure ??).

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. Figure ?? 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 (Figure ??), 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.

Figure ?? shows changes in predictions as we move in a single direction in parameter space. We also wanted to understand how subcircuits might interact with each other as

we move through parameter space. In ?? we perturb multiple subnetworks at once, and measure the new predictions. For the most part, the subnetworks have little inference with each other: the relevant output values for each subnetwork move relatively independently of each other. For some parameter directions, we do see some unexpected interactions, such as when we move in the direction of subnetwork 1 and 2, we see a small effect on the predictions of a few samples not associated with either subnetwork. With larger models, especially those with more modular circuits that are working together, we expect this kind of interaction to become more common. We discuss more in the discussion section.

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 subnetworks. For example, 5 subnetworks that each only involved one of the biases in the final layer could give similar intervention results as in Figure ?? In order to make sure L3D is not just learning this trivial solution, we want to test if the hidden parameter directions it learns are also truly involved in the subnetwork. Therefore, we perform the same intervention experiments but only intervening with the subnetwork components related to the model’s hidden layers’ weights and biases (??). The results are very similar to Figure ??, where interventions on a subnetwork only successfully affect the expected subset of outputs.

3.5. Real world models

Finally, to show the promise of this method to pull out relevant features from real world models, we use L3D to decompose blocks of a language model and computer vision model. These results are primarily qualitative, and were run with minimal compute and little hyperparameter tuning. Our choices for model block, number of subnetworks, and subnetwork rank were relatively arbitrary.

These models do not have well-characterized subnetworks the same way our toy models do. To briefly analyze the function of the subnetworks we identify, we look at the top samples that each subnetwork is relevant to. We compute

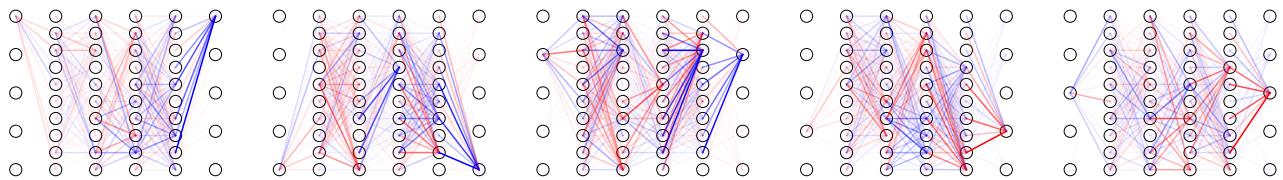


Figure 8: Subnetworks learned by L3D for the $X \mapsto X^2$ model/

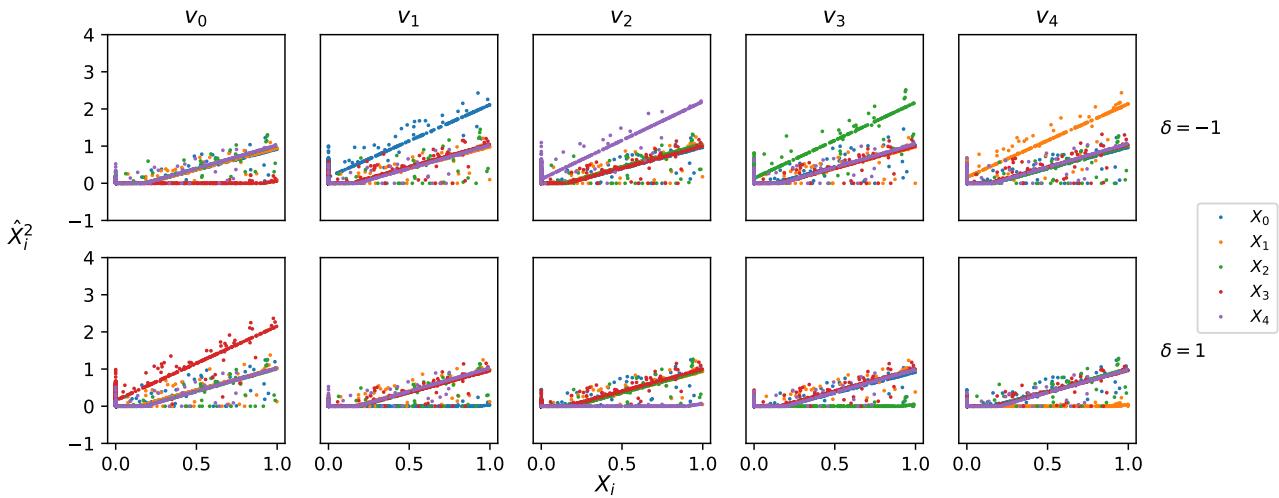


Figure 9: The effect of intervening on each subnetwork in the $X \mapsto X^2$ model. We generate 1000 inputs from the TMS input distribution, intervene on each subnetwork with magnitude δ and measure the change in outputs for each sample.

495 this metric by:

496

497 3.5.1. LANGUAGE MODEL

498 We decompose attention block 7 of the tiny-stories-8M
 499 model. We chose an attention block because this has been a
 500 challenging component of a transformer for SDL to extract
 501 features from (), and we chose a middle layer of the model
 502 so that we identify subnetworks that are neither so high
 503 level that they line up with next-token prediction, and not
 504 so low-level that they line up with token id. We decompose
 505

506 TODO(image with 10 features that we like) TODO(image
 507 with more features for supplementary)

508

509 3.5.2. COMPUTER VISION MODEL

510 We decompose convolutional block 4 of the mobilenet-v3-
 511 small model. Once again we choose a middle layer such
 512 that the subnetworks we identify are neither involved in
 513 low-level computations that would likely require additional
 514 pixel attribution methods to interpret, or so high level they
 515 perfectly line up with classification.

516 TODO(image with 10 features that we like) TODO(image
 517 with more features for supplementary)

518

519 4. Discussion

520 L3D is one of the earliest parameter-based decomposition
 521 methods. For this reason, we have focused our work on
 522 demonstrating the fundamentals of L3D on toy models, and
 523 showcasing its promise with more complex models. Here
 524 we discuss what we believe are simple improvements to L3D
 525 that could enhance its performance and real-world use cases
 526 to which to extend L3D. Finally, we discuss unresolved
 527 challenges and limitations of L3D.

528

529 4.1. Simple Improvements

530 In this work, we did not focus on optimizing L3D, and
 531 we chose nearly identical hyperparameters for all of our
 532 decompositions.

533 **Hyperparameter Choice:** For all of our toy model
 534 decompositions, we always chose our topK hyperparameter
 535 as $k = 0.1$, even when it was clear that certain toy
 536 models should have larger numbers of subnetworks activated
 537 per sample than others (For example, the $X \rightarrow X^2$ model
 538 with 5 inputs and 5 outputs, decomposed into 5 networks,
 539 should probably have $k \geq 0.2$). Too low of k choice is
 540 likely responsible for the high reconstruction loss of some
 541 of our models. Similarly, we always chose the ranks of
 542 the subnetworks somewhat arbitrarily. Some preliminary
 543 research aims to understand the relationship between rank,
 544 compressibility, and interference of subnetworks (??), and
 545 a better understanding of this relationship could help us
 546

choose better hyperparameters for L3D.

Flexible Rank Choices: Relatedly, the way L3D makes circuit decomposition computationally feasible is to represent each tensor in the subnetwork as a low-rank tensor. For simplicity, we use the same rank for all tensors. Theoretically, there could be more flexibility across tensors, and we could use different ranks for different blocks of the subnetwork if we believed certain parts of the model were composed of fundamentally higher rank circuits.

Scaling up

4.2. Extensions

There are also some higher effort extensions to L3D that may give it more real-world relevance.

Finetuning: Our intervention experiments showed promise subnetworks of L3D could be perturbed in ways that only affect the predictions of relevant samples. As we described in Section ??, this could be taken one step further by finetuning a model on a specific set of parameter directions. Using L3D networks, we could finetune a model on a specific set of parameter directions by freezing the current set of weights and learning an adapter consisting of linear combinations of the subnetworks of choice. This could also benchmark the intervention capabilities of L3D versus other mechanistic intervention strategies such as SDL-derived steering vectors. For example, we might use L3D to identify various subnetworks involved in sycophancy, refusal, and other undesired behaviors. After collecting curated data with the goal of finetuning away such behaviors, we could finetune L3D only in the direction of the behavior-related subnetworks and test how well the model achieves our desired output compared to other intervention strategies.

Identifying Specific Circuits with Contrastive Pairs: We developed this method as an unsupervised decomposition method, with goals comparable to those of SDL. However, the methods of L3D could be easily modified to use supervised signals to identify specific circuits of interest. Rather than using gradients of divergence of random pairs, we could decompose gradients of divergence between curated pairs of samples that isolate a behavior of interest. In this context, we could probably not require the topK constraint, but instead look decompose the network into a very small number of low-rank subnetworks involved in our behavior of interest.

4.3. Challenges

Although many of the improvements and extensions of L3D are highly addressable, we think there are some fundamental challenges with parameter-based decomposition methods that may not be easily resolved.

550 **Local Attribution:** L3D's algorithm hinges on the some-
551 what surprising phenomenon that local gradient approxima-
552 tions work reasonably well as attribution methods. They
553 clearly work well in the toy models we used for L3D and
554 for the circuits we found in our real world models. How-
555 ever, do they work for all circuits? In our work, we use
556 a randomly selected sample to be our "reference" output
557 with which to compute divergence gradient. By using a
558 randomly selected sample, rather than a single "reference
559 output" such as the mean of the output distribution, we hope
560 that the random noise in the reference sample will average
561 out the effects of any non-convexity in the loss landscape.
562 However, perhaps even in this setup there are parameter
563 directions that are highly non-convex on which it will be
564 difficult to perform local attribution. Quantifying different
565 types of "dark matter" of parameter decomposition by ana-
566 lyzing reconstruction loss (similar to recent work, such as
567 (?), done on SDL) could better help us characterize these
568 limitations.

569 **Relationship to overparameterized models:** Going one
570 step further, we suspect that the reason local attribution
571 methods work so well is because large models are probably
572 overparameterized (????). Larger models may have wider
573 loss basins, or more degeneracies near their global minima
574 (??), making local attribution methods less likely to break
575 down as we move through parameter space. If in the future,
576 an learning algorithm is developed that has fundamentally
577 different limitations that stochastic gradient descent and its
578 relatives, we might lose this property. Moreover, circuit ac-
579 tivations might no longer be sparse. A new learning process
580 might be able to compress subnetworks in such a way that
581 subnetworks have very high levels of interference with each
582 other - removing the degeneracy assumption that underlies
583 L3D.
584

585 **Interpretation of a circuit:** Finally, we should address
586 the definition of "circuits". It is still not well agreed upon
587 what a "feature" is in relation to large networks, and the
588 definition of what should constitute a circuit or subnetwork
589 is even less clear. Is our definition of a circuit - sparsely
590 active subnetworks that can move outputs within the origi-
591 nal output distribution - too restrictive? If there is a circuit
592 that is relevant to every output, a sort of "scaffolding" for
593 more specific circuits - should it be included in the decom-
594 position? If, after identifying the structure of subnetworks,
595 we cannot interpret it beyond a description of its end re-
596 sults, are circuits any more informative than the features
597 they are computing? If parameter decomposition is a viable
598 strategy for understanding and intervening with large net-
599 works, these questions will be important for the mechanistic
600 interpretability community to address.

601
602
603
604

605 **5. Impact Statement**

606
607 This paper presents work whose goal is to advance the field
608 of neural network interpretability. There are many potential
609 societal consequences of our work, none which we feel must
610 be specifically highlighted here.

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

References

- 660 Bereska, L. and Gavves, E. Mechanistic interpretability for
661 AI safety—a review. *arXiv preprint arXiv:2404.14082*,
662 2024.
- 663 Braun, D., Bushnaq, L., Heimersheim, S., Mendel,
664 J., and Sharkey, L. Interpretability in parameter
665 space: Minimizing mechanistic description length with
666 attribution-based parameter decomposition. *arXiv
667 preprint arXiv:2501.14926*, 2025.
- 668 Bricken, T., Templeton, A., Batson, J., Chen, B., Jermyn,
669 A., Conerly, T., Turner, N., Tamkin, A., and Carter, S. Towards
670 monosemanticity: Decomposing language models
671 with dictionary learning. *Transformer Circuits Thread*,
672 2023. Accessed: 2025-02-17.
- 673 Bushnaq, L., Mendel, J., Heimersheim, S., Braun, D.,
674 Goldowsky-Dill, N., Hänni, K., Wu, C., and Hobbahn,
675 M. Using degeneracy in the loss landscape for mechanistic
676 interpretability. *arXiv preprint arXiv:2405.10927*,
677 2024.
- 678 Bussmann, B., Leask, P., and Nanda, N. BatchTopK sparse
679 autoencoders. *arXiv preprint arXiv:2412.06410*, 2024.
- 680 Choromanska, A., Henaff, M., Mathieu, M., Arous, G. B.,
681 and LeCun, Y. The loss surfaces of multilayer networks.
682 In *Artificial intelligence and statistics*, pp. 192–
683 204. PMLR, 2015.
- 684 Cunningham, H., Ewart, A., Riggs, L., Huben, R., and
685 Sharkey, L. Sparse autoencoders find highly interpretable
686 features in language models. *arXiv preprint arXiv:2309.08600*, 2023.
- 687 Dauphin, Y. N., Pascanu, R., Gulcehre, C., Cho, K., Gan-
688 guli, S., and Bengio, Y. Identifying and attacking the
689 saddle point problem in high-dimensional non-convex
690 optimization. *Advances in neural information processing
691 systems*, 27, 2014.
- 692 Davies, X., Langosco, L., and Krueger, D. Unify-
693 ing grokking and double descent. *arXiv preprint
694 arXiv:2303.06173*, 2023.
- 695 Elhage, N., Hume, T., Olsson, C., Schiefer, N., Henighan,
696 T., Kravec, S., Hatfield-Dodds, Z., Lasenby, R., Drain,
697 D., Chen, C., et al. Toy models of superposition. *arXiv
698 preprint arXiv:2209.10652*, 2022.
- 699 Engels, J., Michaud, E. J., Liao, I., Gurnee, W., and
700 Tegmark, M. Not all language model features are linear.
701 *arXiv preprint arXiv:2405.14860*, 2024a.
- 702 Engels, J., Riggs, L., and Tegmark, M. Decomposing
703 the dark matter of sparse autoencoders. *arXiv preprint
704 arXiv:2410.14670*, 2024b.
- 705 Gao, L., la Tour, T. D., Tillman, H., Goh, G., Troll, R.,
706 Radford, A., Sutskever, I., Leike, J., and Wu, J. Scaling
707 and evaluating sparse autoencoders. *arXiv preprint
708 arXiv:2406.04093*, 2024.
- 709 Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B.,
710 Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y.
711 Generative adversarial nets. *Advances in neural informa-
712 tion processing systems*, 27, 2014.
- 713 Ho, J., Jain, A., and Abbeel, P. Denoising diffusion proba-
714 bilistic models. *Advances in neural information process-
715 ing systems*, 33:6840–6851, 2020.
- 716 Hoogland, J., Wang, G., Farrugia-Roberts, M., Carroll, L.,
717 Wei, S., and Murfet, D. The developmental landscape
718 of in-context learning. *arXiv preprint arXiv:2402.02364*,
719 2024.
- 720 Kawaguchi, K. Deep learning without poor local minima.
721 *Advances in neural information processing systems*, 29,
722 2016.
- 723 Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy,
724 M., and Tang, P. T. P. On large-batch training for deep
725 learning: Generalization gap and sharp minima. *arXiv
726 preprint arXiv:1609.04836*, 2016.
- 727 Kramár, J., Lieberum, T., Shah, R., and Nanda, N. AtP*:
728 An efficient and scalable method for localizing LLM behav-
729 iour to components. *arXiv preprint arXiv:2403.00745*,
730 2024.
- 731 Lindsey, J., Templeton, A., Marcus, J., Conerly, T., Batson,
732 J., and Olah, C. Sparse crosscoders for cross-layer fea-
733 tures and model diffing. *Transformer Circuits Thread*,
734 2024.
- 735 Matena, M. and Raffel, C. Npeff: Non-negative per-example
736 fisher factorization. *arXiv preprint arXiv:2310.04649*,
737 2023.
- 738 Merullo, J., Eickhoff, C., and Pavlick, E. Talking heads:
739 Understanding inter-layer communication in transformer lan-
740 guage models. *arXiv preprint arXiv:2406.09519*, 2024.
- 741 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness,
742 J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidje-
743 land, A. K., Ostrovski, G., et al. Human-level control
744 through deep reinforcement learning. *nature*, 518(7540):
745 529–533, 2015.
- 746 Nanda, N. Attribution patching: Activation patching
747 at industrial scale. URL: <https://www.neelnanda.io/mechanistic-interpretability/attribution-patching>,
748 2023.
- 749 Pascanu, R. On the difficulty of training recurrent neural
750 networks. *arXiv preprint arXiv:1211.5063*, 2013.

- 715 Sagun, L., Evci, U., Guney, V. U., Dauphin, Y., and Bottou,
716 L. Empirical analysis of the hessian of over-parametrized
717 neural networks. *arXiv preprint arXiv:1706.04454*, 2017.
- 718 Soudry, D. and Hoffer, E. Exponentially vanishing sub-
719 optimal local minima in multilayer neural networks.
720 *arXiv preprint arXiv:1702.05777*, 2017.
- 721
- 722 Subramani, N., Suresh, N., and Peters, M. E. Extracting
723 latent steering vectors from pretrained language models.
724 *arXiv preprint arXiv:2205.05124*, 2022.
- 725
- 726 Syed, A., Rager, C., and Conmy, A. Attribution patching
727 outperforms automated circuit discovery. *arXiv preprint*
728 *arXiv:2310.10348*, 2023.
- 729
- 730 Turner, A. M., Thiergart, L., Leech, G., Udell, D., Vazquez,
731 J. J., Mini, U., and MacDiarmid, M. Steering lan-
732 guage models with activation engineering. *arXiv preprint*
733 *arXiv:2308.10248*, 2023.
- 734 Wang, G., Farrugia-Roberts, M., Hoogland, J., Carroll, L.,
735 Wei, S., and Murfet, D. Loss landscape geometry re-
736 veals stagewise development of transformers. In *High-*
737 *dimensional Learning Dynamics 2024: The Emergence*
738 *of Structure and Reasoning*, 2024.
- 739
- 740 Watanabe, S. Algebraic information geometry for learning
741 machines with singularities. *Advances in neural informa-*
742 *tion processing systems*, 13, 2000.
- 743
- 744 Watanabe, S. Algebraic geometry of singular learning ma-
745 chines and symmetry of generalization and training errors.
746 *Neurocomputing*, 67:198–213, 2005.
- 747
- 748 Watanabe, S. Almost all learning machines are singular. In
749 *2007 IEEE Symposium on Foundations of Computational*
750 *Intelligence*, pp. 383–388. IEEE, 2007.
- 751 Wei, S., Murfet, D., Gong, M., Li, H., Gell-Redman, J.,
752 and Quella, T. Deep learning is singular, and that's good.
753 *IEEE Transactions on Neural Networks and Learning*
754 *Systems*, 34(12):10473–10486, 2022.
- 755
- 756
- 757
- 758
- 759
- 760
- 761
- 762
- 763
- 764
- 765
- 766
- 767
- 768
- 769

A. Definitions
A.0.1. DIMENSIONS

770 n_s : The number of samples in a batch of inputs
 771
 772 n_f : The dimensions of a single input vector to a model
 773
 774 n_o : The dimensions of a single output vector from a model
 775
 776 n_w : The number of parameters values in a model.
 777
 778 n_v : The number of subnetworks or parameter directions chosen to decompose a model.

A.0.2. MODEL SYNTAX

781 $X \in \mathbb{R}^{n_s \times n_f}, x \in \mathbb{R}^{n_f}$: Batch and individual input vectors to a model.
 782
 783 $W \in \mathbb{R}^{n_w}, w \in \mathbb{R}$: The set of and individual parameter values of a model
 784
 785 \sqsubseteq : The set of parameters corresponding to a specific tensor or block in a model.
 786
 787 $f : \mathbb{R}^{n_s \times n_f} \mapsto \mathbb{R}^{n_s \times n_o}$: A model mapping a set of input vectors to a set of output vectors.
 788
 789 $f(X, W)$: The output of model f with parameter values W on input X .
 790
 791 $f(X, W_0)$ or $f(X)$: The output of model f with fixed parameter values W_0 . W_0 is the set of learned parameter values from
 792 model training.
 793
 794 D : Divergence metric between two vectors. Typical divergence metrics are mean-squared error for regression-type outputs,
 795 and KL-divergence for probability-type outputs.

A.0.3. DECOMPOSITION SYNTAX

796 $V(\text{or } V^{out}) \in \mathbb{R}^{n_v \times n_w}, v(\text{or } V^{out}) \in \mathbb{R}^{n_w}$: The set of or individual parameter directions that are used to decompose a
 797 model. V_{out} can be used to transform parameter directions in the subnetwork vector space back into the original parameter
 798 space of the model. The terms parameter direction, subnetwork, and circuit all have the same meaning.
 799
 800 $\mathcal{V}_{out}, \sqsubseteq_{out}$: Components of V^{out} or V^{out} related to a specific tensor or block in the original model.
 801
 802 $V^{in} \in \mathbb{R}^{n_w}, V^{in} \in \mathbb{R}$: Transforms the original parameter space of the model into the subnetwork vector space.
 803
 804 $\mathcal{V}_{in}, \sqsubseteq_{in}$: Components of V^{in} or V^{in} related to a specific tensor or block in the original model.
 805
 806 r : The rank of each component of the decomposition vectors corresponding to tensors in the original model.

A.0.4. TRAINING

807 \mathcal{K} : The subset of indexes $i = \tau$ where τ is the threshold computed by using the top k absolute values in a tensor.
 808
 809 L : The L2 reconstruction loss used to optimize V^{in} and V^{out} .

A.0.5. MEASURING AND INTERVENTION

810 $I(x_i, x_j, v_k), I(x_i, v_k)$: The impact of subnetwork v_k on the divergence between samples x_i and x_j , or averaged across
 811 many x_j reference samples.
 812
 813 δ : A scalar value to move W in a specific direction.

B. Additional Methods
B.1. Low-Rank Tensor Representation

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

825 While W is a vector of all of the parameters in a model, typically model parameters are organized into tensors $W = \{w_i\}_i$.
 826 If our parameters are organized into tensors $W = \{w_i\}_i$, each subnetwork or parameter component can be organized as
 827 $V_i^{in} = \{v^{in}_i\}_i, V_i^{out} = \{v^{out}_i\}_i$ where we have parts of our subnetworks corresponding to each tensor in the original model
 828 parameters. We wish each of these tensors to be low rank, and we express them using the canonical polyadic decomposition
 829 () (a way to write 3+ dimentinoal tensors in terms of low-rank components).
 830

831

$$832 v^{in}_{i,j} = \sum_{r=1}^R a_{i,j,r} \times b_{i,j,r} \times c_{i,j,r} \dots \quad (14)$$

833

834

835 Where R is the rank we wish to use to represent the parameter component, and the number of factors (a, b, c,...) in the
 836 factorization is equal to the number of dimensions in the tensor w_i .
 837

838 **B.2. Finetuning**

839

840 To finetune a model on a specific set of parameter directions, we would freeze the current set of weights and learn an
 841 adaptor consisting of linear combinations of the subnetworks of choice. During fine tuning, we would only need to learn the
 842 coefficients of these linear combinations greatly reducing cost.
 843

844

845

846

847

848

849 **B.3. Toy Model Training**

850

851 For all of our toy models (except the $X \mapsto X^2$ model), we generate uniformly random inputs between 0 and 1. For
 852 $X \mapsto X^2$, we generate uniformly random inputs between -1 and 1. For all toy model data, we use a sparsity value of
 853 1-sparsity=.05. We generate 10000 datapoints and train for 1000 epochs with batch sizes of 32. We use an AdamW optimizer
 854 with a learning rate of 0.001.
 855

856 **B.4. L3D Model Training**

857 To train L3D, we use the same training distributions as in each toy models. Although optimal hyperparameter values
 858 probably depend on the model size, and the rank and number of parameter tensors, we use the same hyperparameters for all
 859 of our models. We generate only 1000 datapoints, with a batch size of 32, and train for 1000 epochs. We use an AdamW
 860 optimizer with a learning rate of 0.01, and a learning decay rate of .8 every 100 steps. We always use a topK hyperparameter
 861 of $k = 0.1$. We include all of the model's parameter tensors, including biases, in the decomposition.
 862

863

864 **C. Supplemental Figures**

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880
 881
 882
 883
 884
 885
 886
 887
 888
 889
 890
 891
 892
 893
 894
 895
 896
 897
 898
 899
 900
 901
 902
 903
 904
 905
 906
 907
 908
 909
 910
 911
 912
 913
 914
 915
 916
 917
 918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934

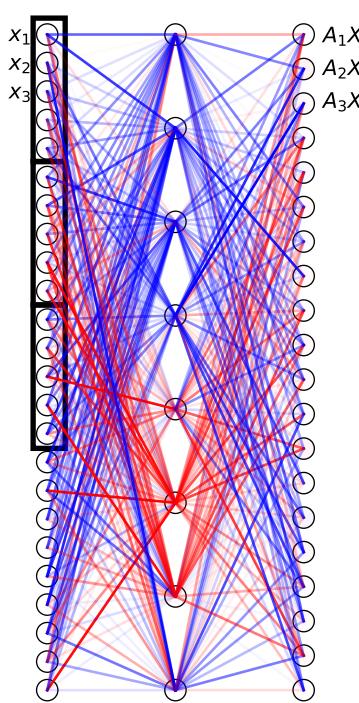


Figure S1: The full architecture of high rank circuit toy model (model C).

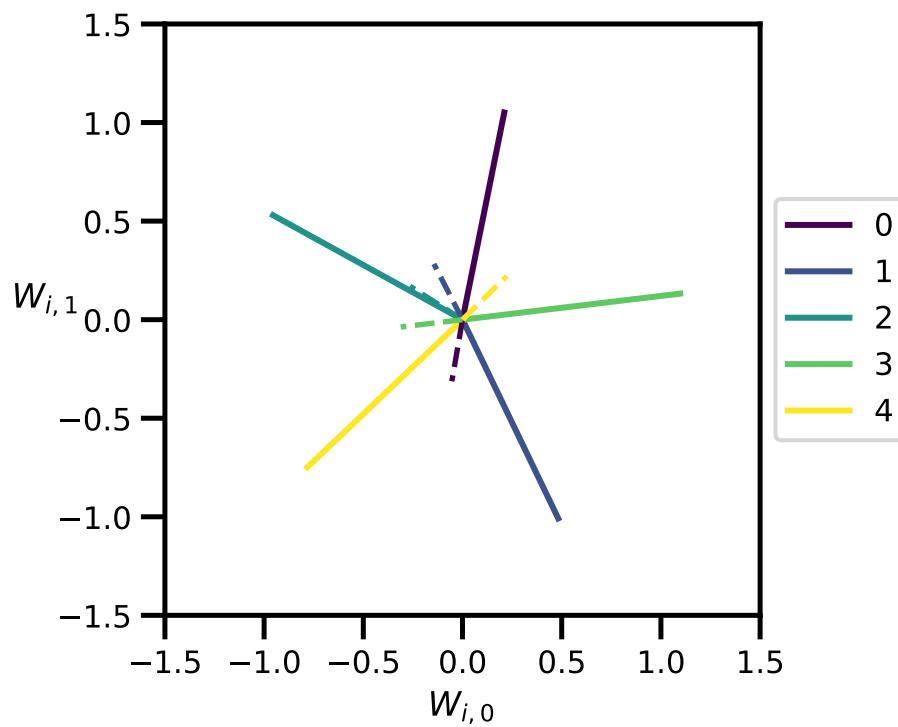


Figure S2: Encoder directions in the hidden layer dimension of the TMS toy model.

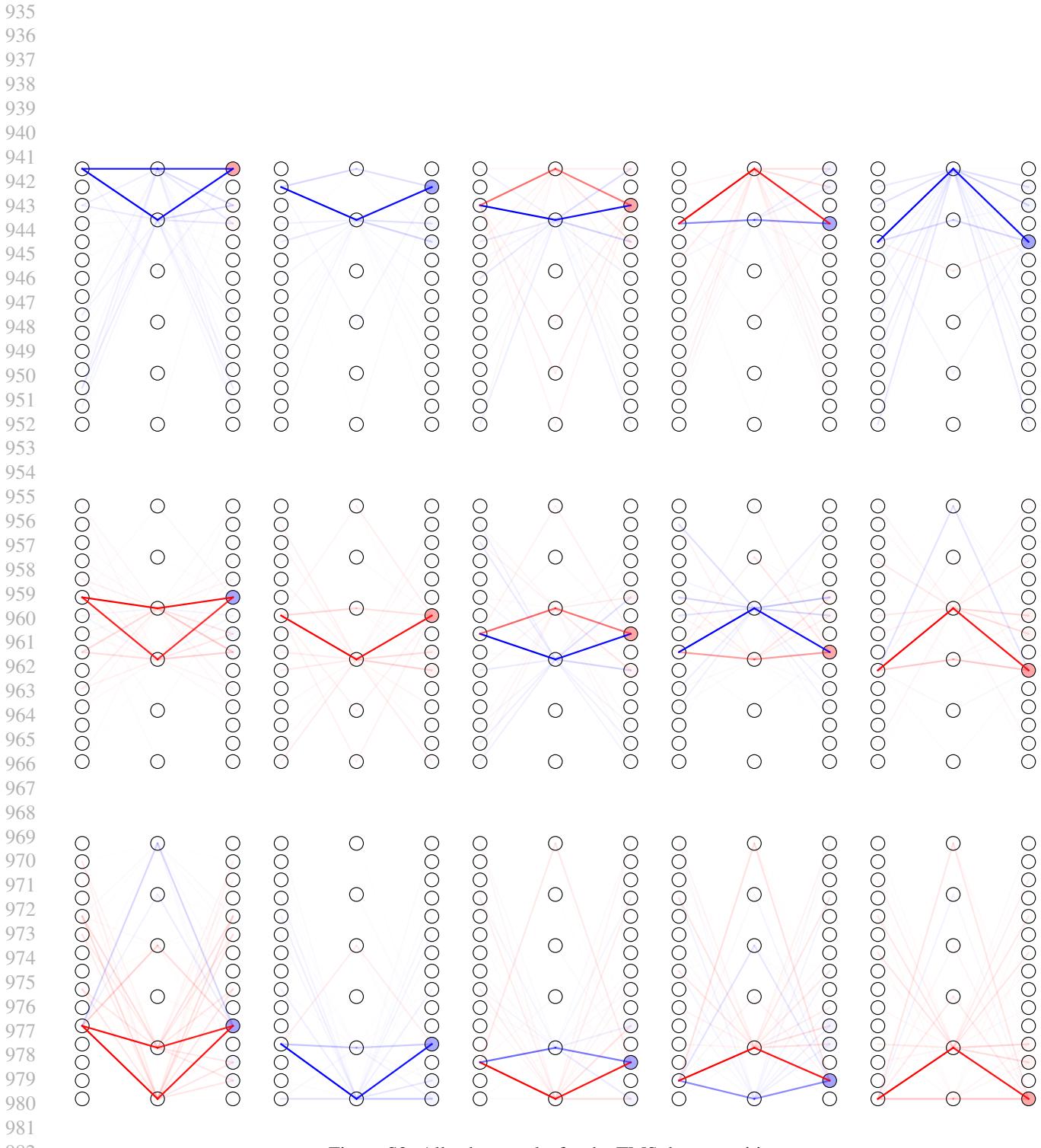
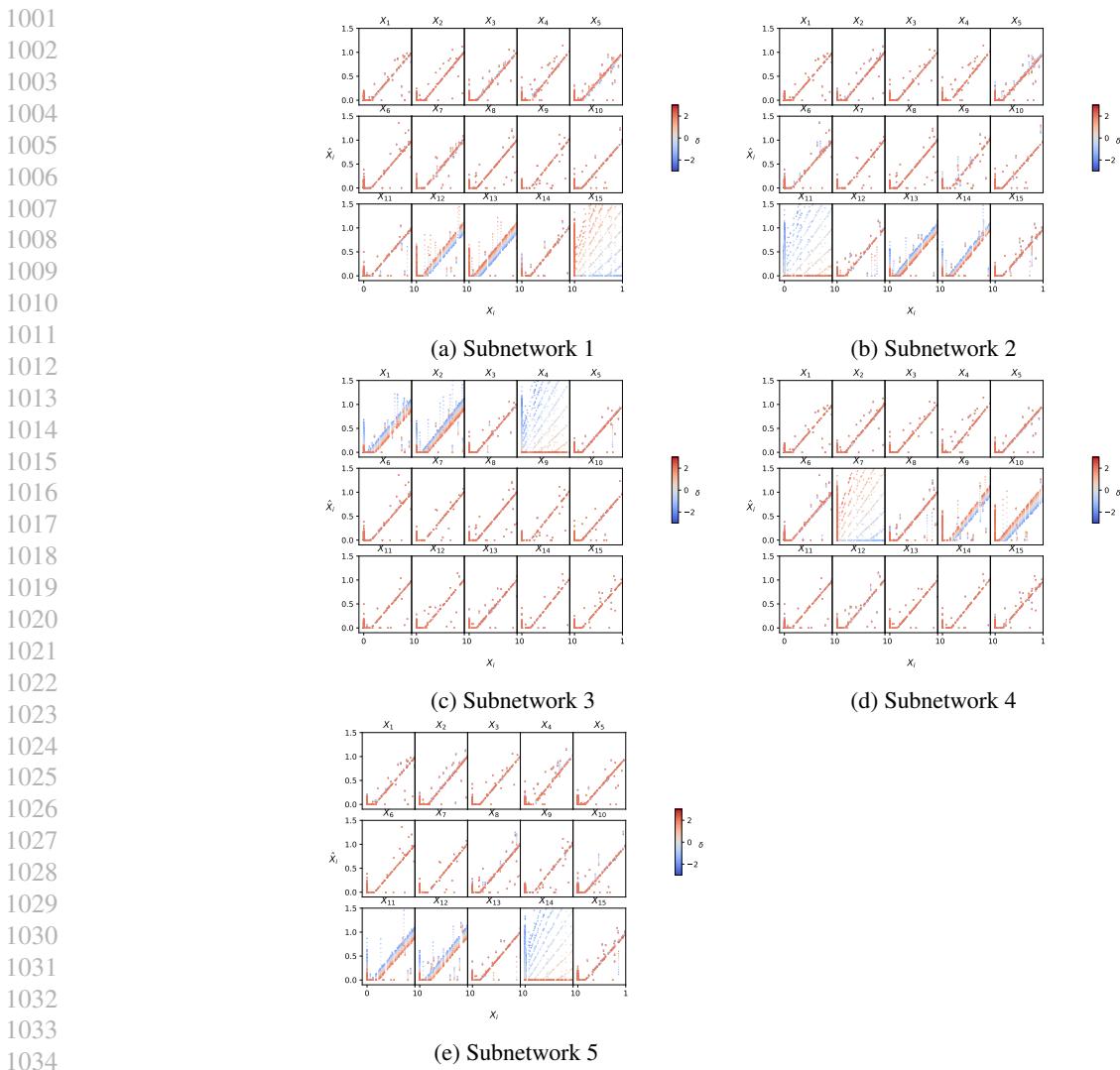


Figure S3: All subnetworks for the TMS decomposition

990
991
992
993
994
995
996
997
998
999

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



Local Loss Landscape Decomposition

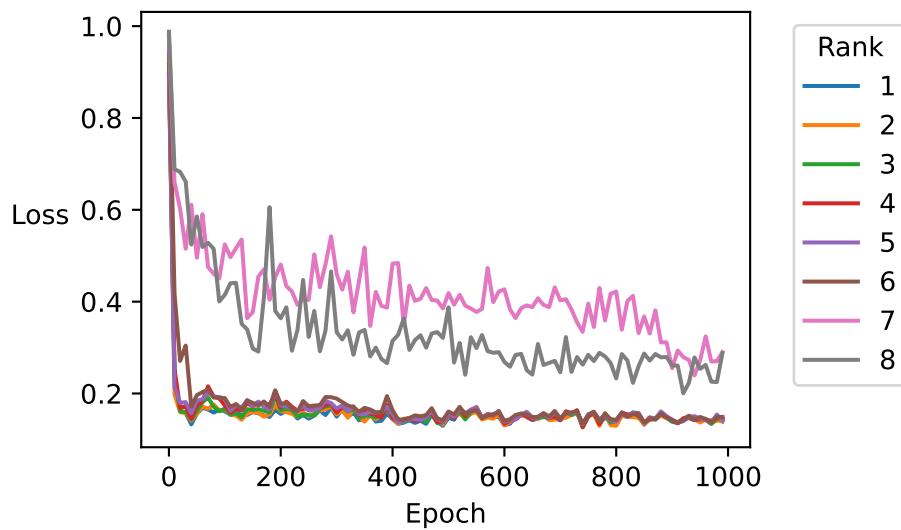
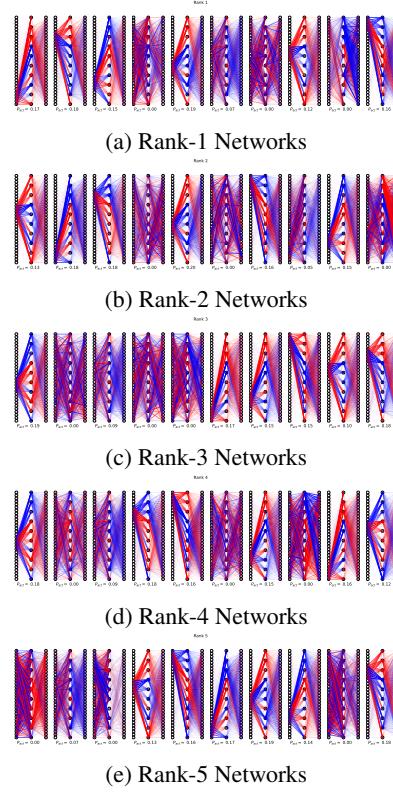


Figure S5: Loss vs Rank

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



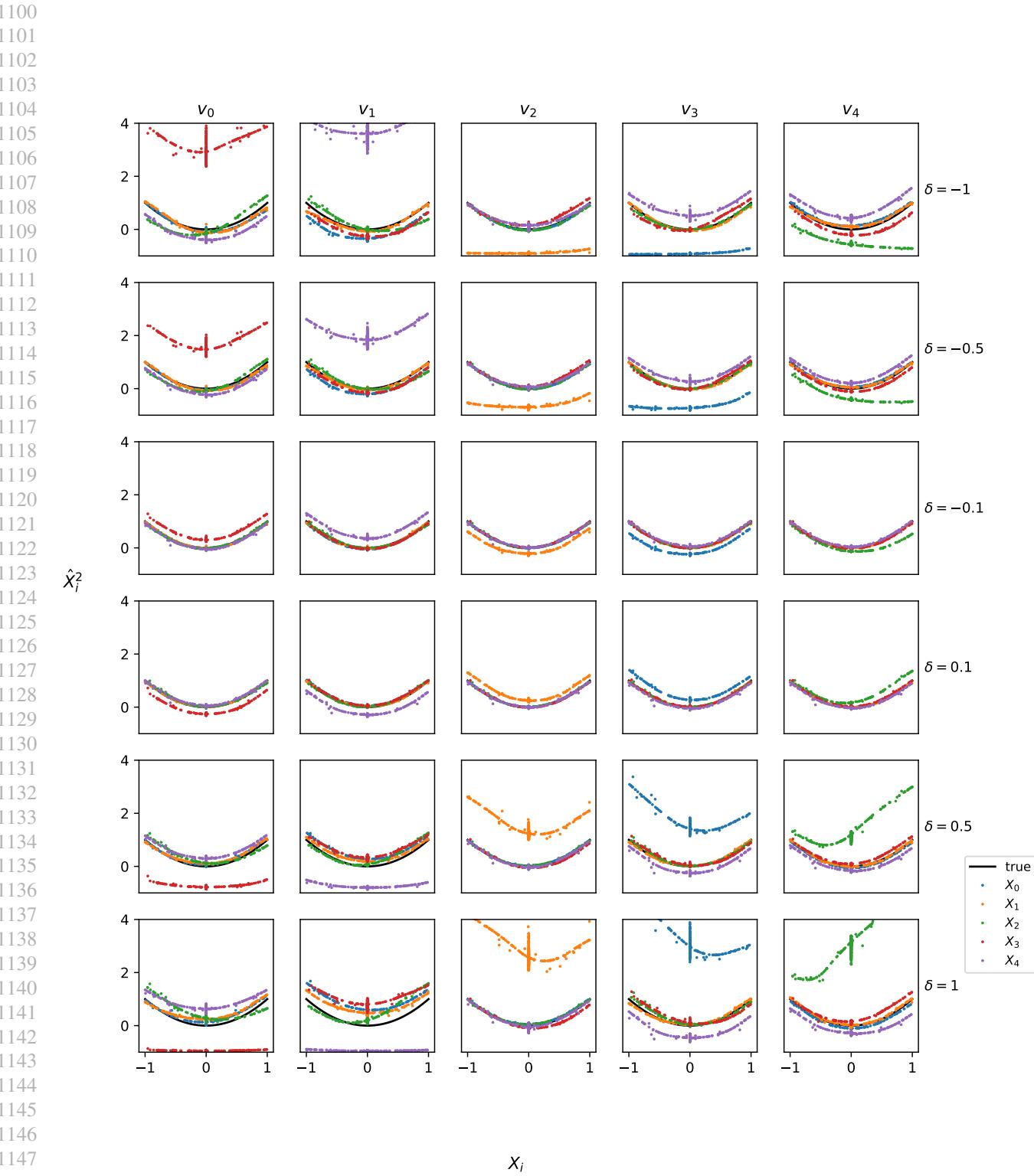


Figure S7: Effects of intervening on each of the subnetworks of the $X \mapsto X^2$ model.

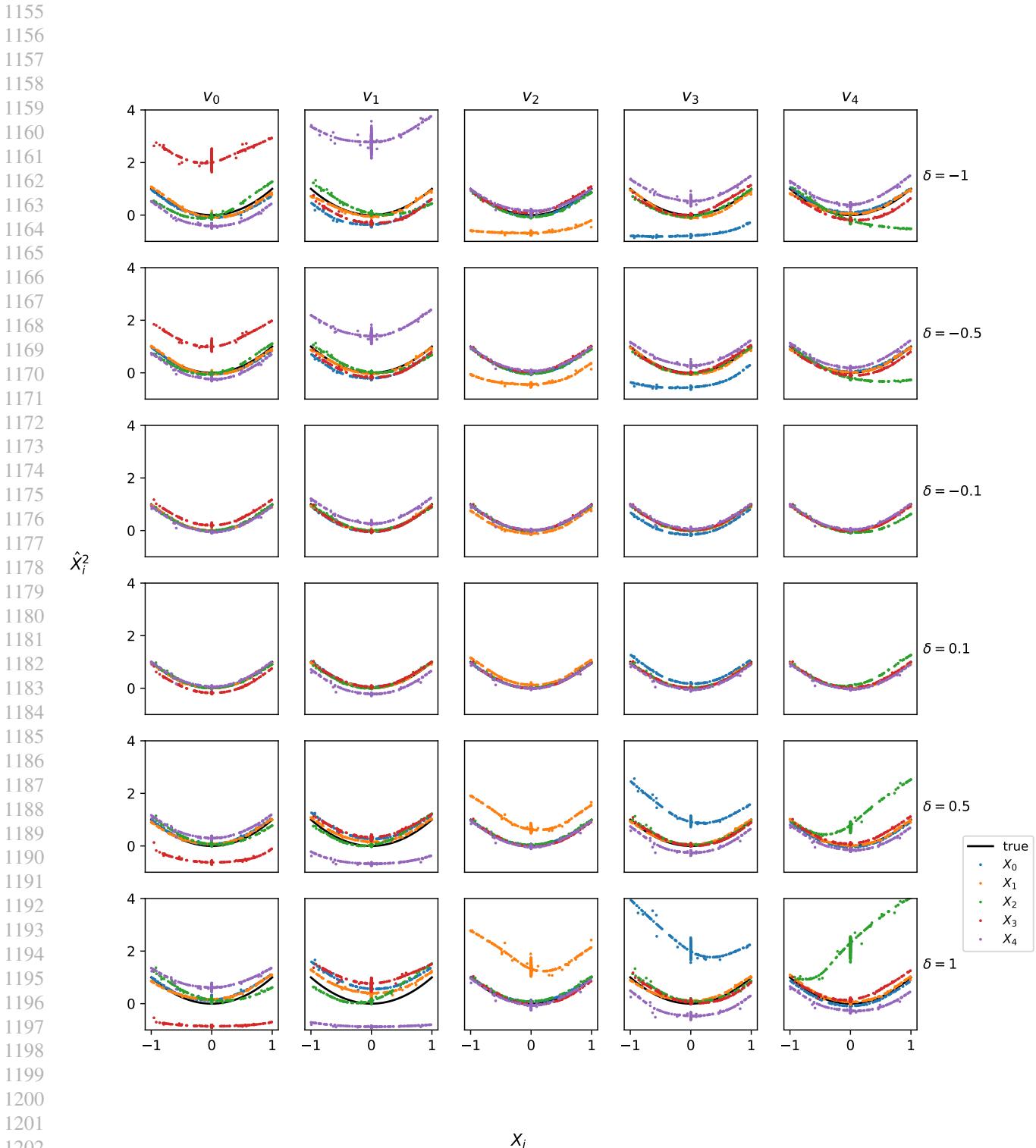


Figure S8: Effect of intervening using just the weights and biases of the middle hidden layers in the subnetworks of the $X \mapsto X^2$ model.

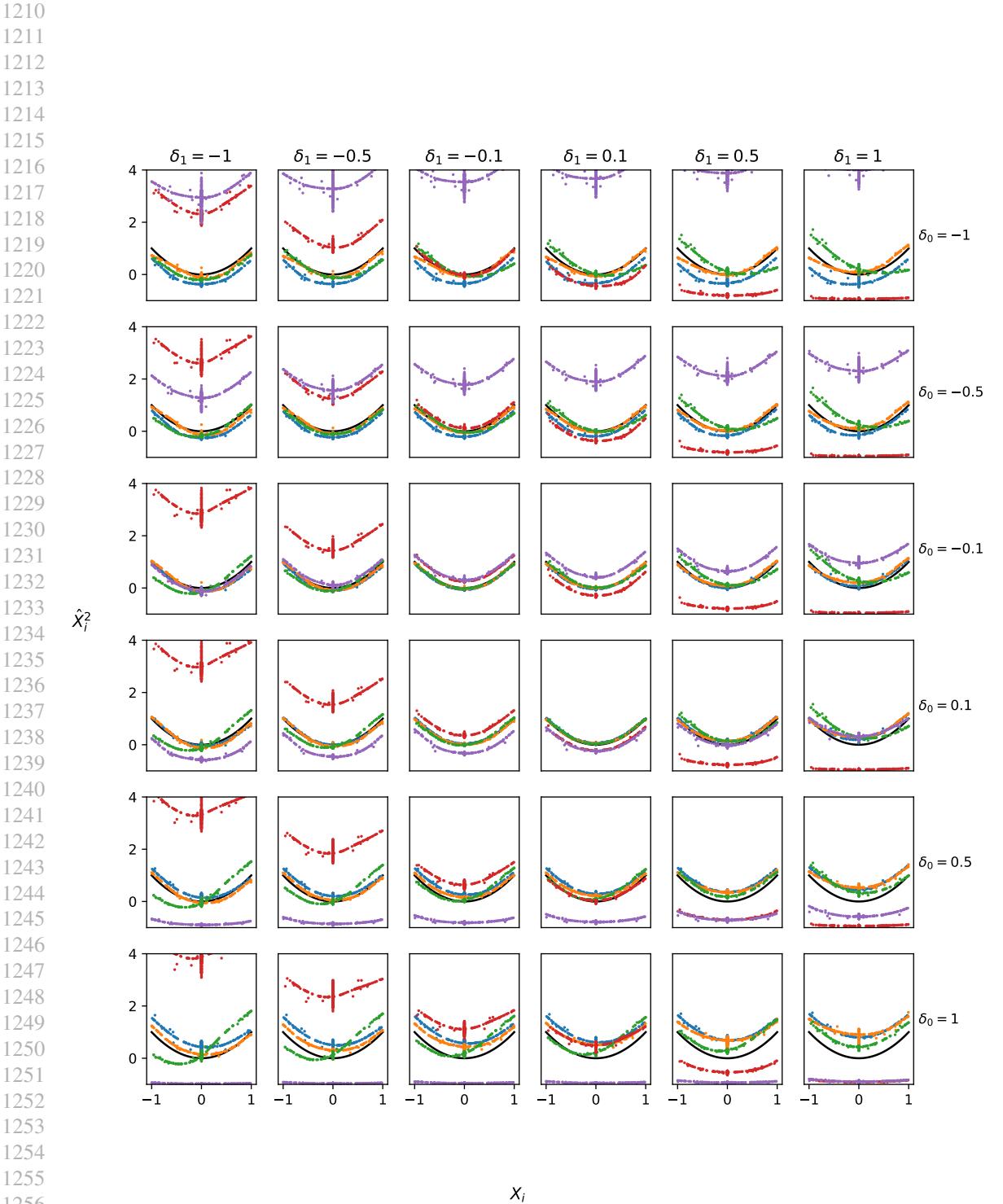


Figure S9: Effects of intervening with multiple subnetworks (v_0 on the x-axis, v_1 on the y-axis) at once.

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

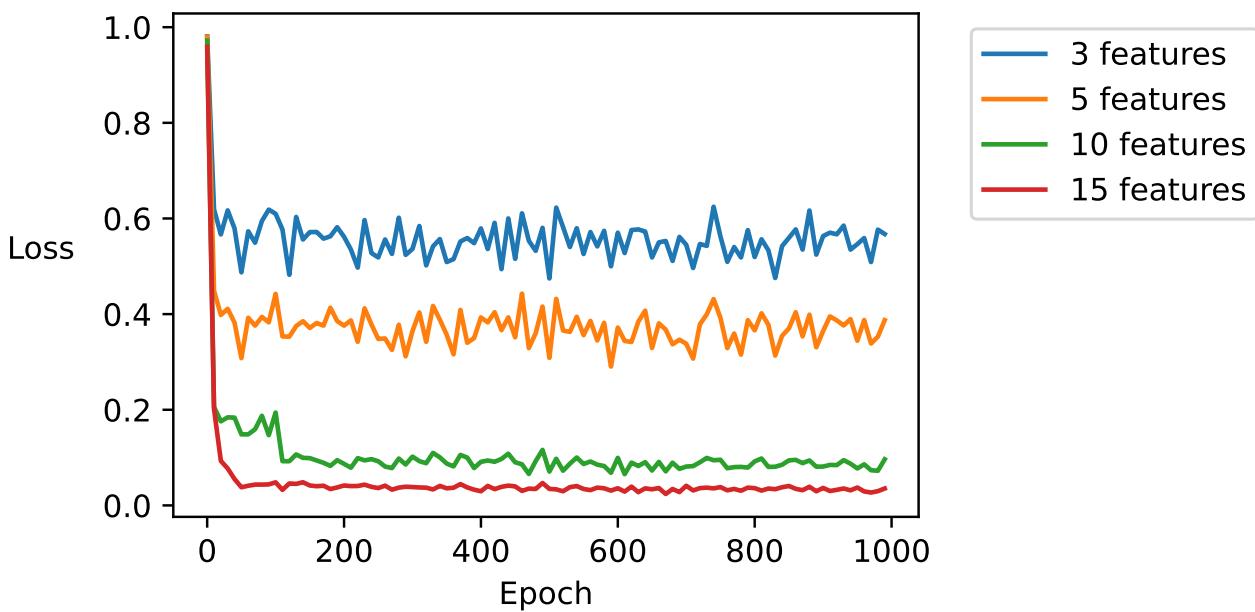
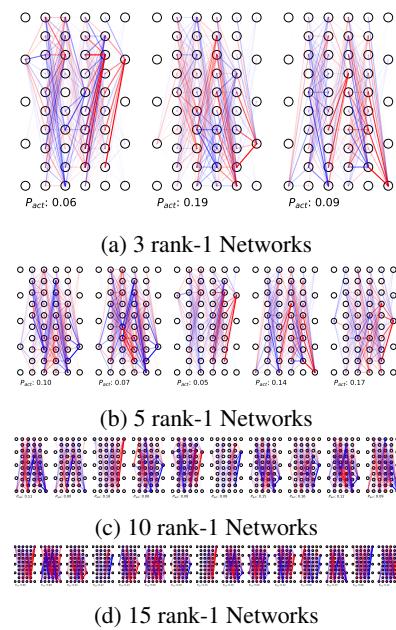


Figure S11: Training loss vs. number of subnetworks for the $X \mapsto X^2$ model