

# Compact proofs of model performance via mechanistic interpretability

Refer: [https://github.com/LouisYRYJ/Proof\\_based\\_approach\\_tutorial/blob/master/proof\\_public.ipynb](https://github.com/LouisYRYJ/Proof_based_approach_tutorial/blob/master/proof_public.ipynb)

## Introduction

### The problem of understanding a model

#### Why do we care about understanding a model?

When we use or train a model, many things can go wrong. For example, during training the model can learn undesired behaviour, which is not obvious to us (deceptive alignment). Or it might have failure modes that are not salient to us (adversarial examples). On the other hand, we could steer a model towards a desired behaviour, if we understood how it works.

All of these issues could be resolved, if the models were transparent to us (though this is not the only approach). So an important question to raise here is: What do we mean when we talk of “a mechanistic understanding” of a model? When is a model transparent to us?

This is a difficult question! Let’s say you study a model and reverse engineered parts of it, like a circuit. How can you be sure that the circuit you found actually does the thing you are claiming it is? Let’s look at the specific example of autoencoders. This is a quote from [Towards Monosemanticity: Decomposing Language Models With Dictionary Learning](#)

Usually in machine learning we can quite easily tell if a method is working by looking at an easily-measured quantity like the test loss. We spent quite some time searching for an equivalent metric to guide our efforts here, and unfortunately have yet to find anything satisfactory.

We began by looking for an information-based metric, so that we could say in some sense that the best factorization is the one that minimizes the total information of the autoencoder and the data. Unfortunately, this total information did not generally correlate with subjective feature interpretability or activation sparsity.[...]

Thus we ended up using a combination of several additional metrics to guide our investigations[...]

Interpreting or measuring some of these signals can be difficult, though. For instance, at various points we thought we saw features which at first didn’t make any sense, but with deeper inspection we could understand.

We think it would be very helpful if we could identify better metrics for dictionary learning solutions from sparse autoencoders trained on transformers.

See also Section 5 of this review [Mechanistic Interpretability for AI Safety – A Review](#) for more references on the difficulty of evaluating interpretability results.

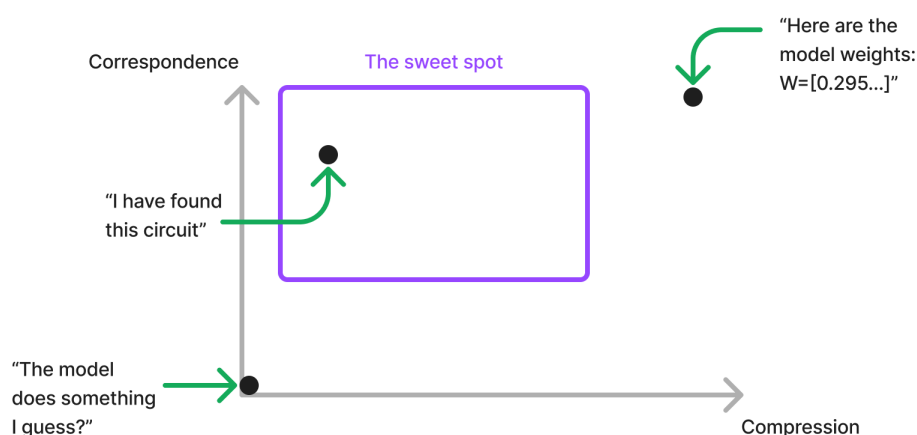
### Quantative methods for interpretability

Having quantative methods would not only make mechanistic interpretability research more rigorous. If we want to scale up methods to huge models, we will need to automate parts of the process and we won’t be able to have a human in the loop at every crucial point. A lack of quantative benchmarks makes this task seem almost impossible. To spoiler the punchline: Compact proofs provide such a quantative benchmark, although they currently are infeasible for larger models.

Before getting into the details, let's nail down two things that we want to quantify. The following two points are taken from the [Compact proofs blog post](#), see also this [comment](#) by Ryan Greenblatt.

1. Correspondence (or faithfulness): How well our explanation reflects the model's internals.
2. Compression: Explanations compress the particular behavior of interest. Not just so that it fits in our heads, but also so that it generalizes well and is feasible to find and check.

Specifically, the second point implies that the explanation, say the circuit that we discovered, should be more **compact** and therefore more understandable for us humans: The weights of a model are a perfectly faithful explanation of its behaviour, but this explanation is not helpful for us.



An important insight that we will make is that our explanations are not as good as we might think. Specifically, **noise** in the model's weights seem negligible. But worst case bound imply that it could still be an important contributing factor. In fact, it might be that something that we deem as noise, is important for the model's computation, but we simply don't understand it. This issue with the noise is another point that a quantitative evaluation should be able to address.

## What are compact proofs?

Compact proofs are an attempt at formalizing the above diagram.

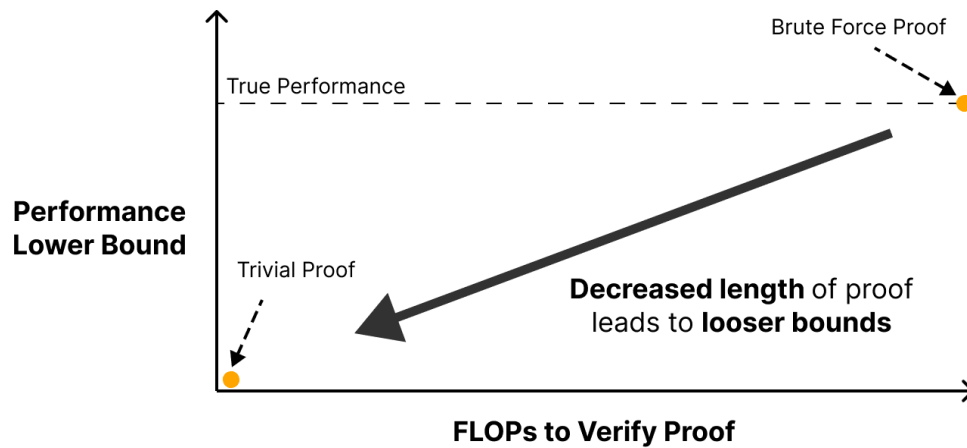
First of all, what do we mean by proof i.e. what are we trying to prove? Say we are training a model with weights  $\theta$  on some task. The kind of statements that we want to prove are of the form

$$\mathbb{E}[f_{\theta}(x)] \geq b$$

where  $f_{\theta}$  is a quantity that we are interested in bounding from below or above (depending on the quantity), such as loss or accuracy.

The compactness of a proof is determined by its length. A good proxy for the length is the FLOPS required to run the proof, see the [paper](#) for more details. (Maybe I will write more on this)

So once we have a proof, we can measure its correspondence by looking at the bound and measure its compactness by measuring its length. We get a similar picture to the one drawn above:

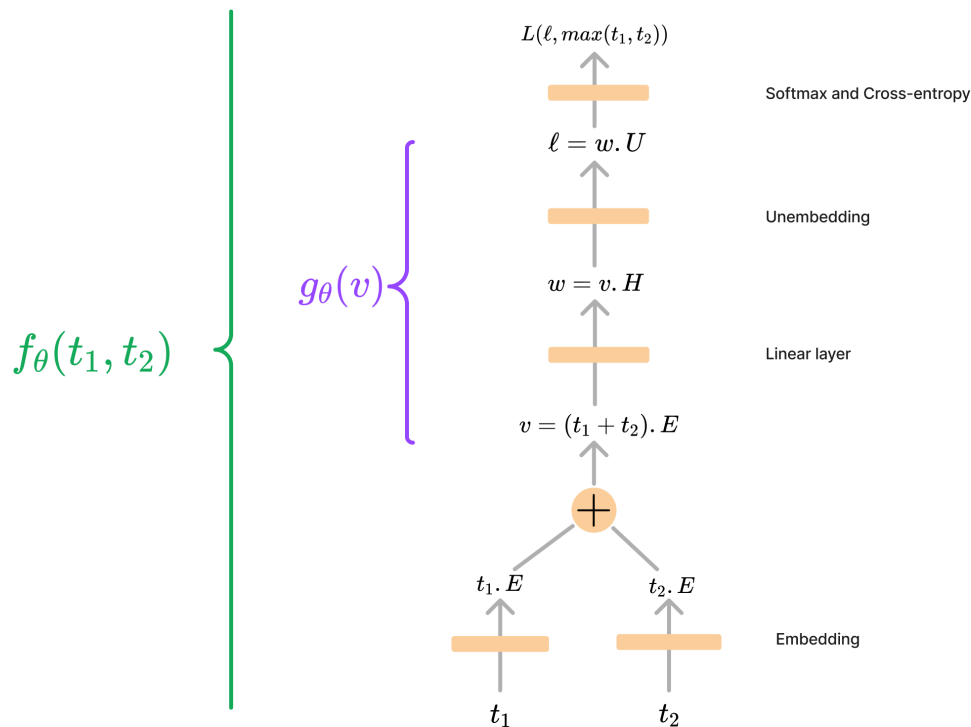


We will see many examples of proofs and compare their performance below. The ideal goal here would be to have a pipeline that takes in a vague interpretation of the model, turn that into a rigorous proof, and evaluate the interpretation based on the correspondence and compactness of the proof.

As we will see below this turns out to be rather difficult, even in toy models. The takeaway here is that quantification seems to be a hard problem and the compact proof approach is an example of this.

## Max-of-2 example

Having worked through the high level picture, let us now focus on concrete examples of compact proofs and explain what it means. Let's say we have a model:



refer: main.py > class MLP()

In our first example, we will train the model to predict the max of the two tokens, where the tokens range from 0 to  $d_{\text{vocab}}$ . We will be interested in estimating the global loss of this model, that is we want to estimate

$$\mathbb{E}[f(t_1, t_2)] = \frac{1}{d_{\text{vocab}}^2} \cdot \sum_{t_1, t_2 \in \{0, \dots, d_{\text{vocab}} - 1\}} f(t_1, t_2).$$

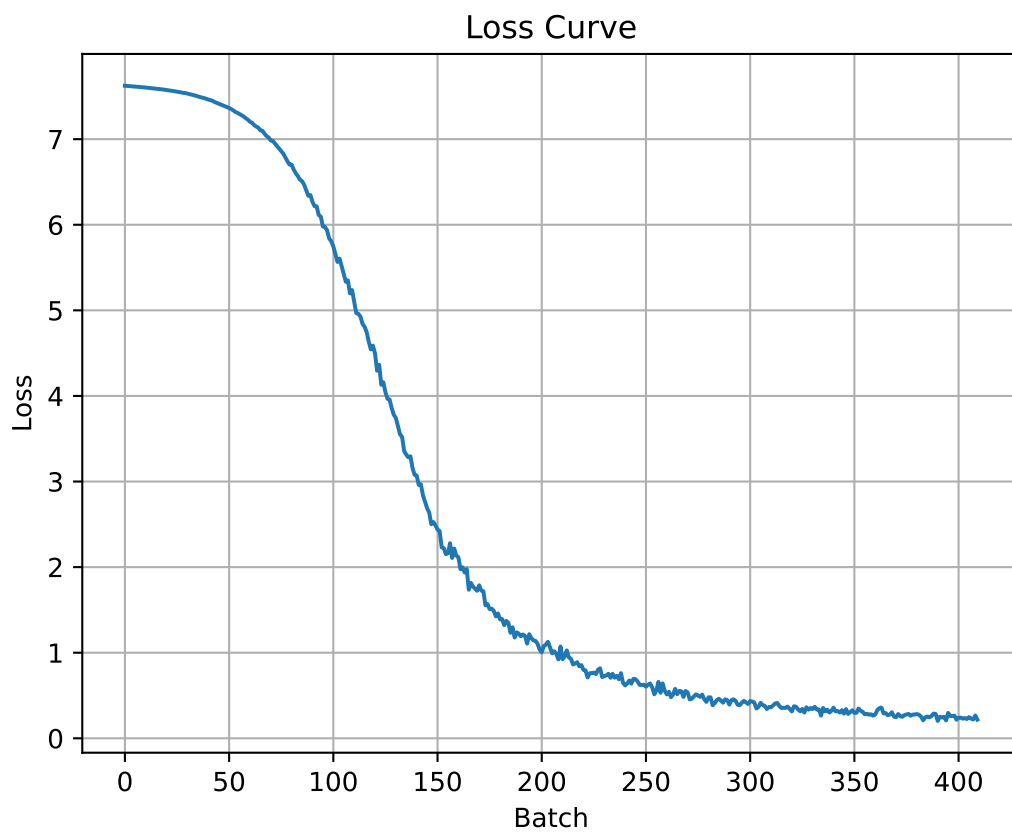
The general proof strategy consists of two steps:

1. P1: Prove a statement that given a model with its weights  $\theta$ , there is a quantity  $C(\theta)$  such that  $\mathbb{E}[h(x, M(x))] \leq C(\theta)$ .
2. P2: Compute the quantity  $C(\theta)$ .

We will come back to this after we did some proofs and also discuss what it means to have a **compact** proof.

We train our model on  $\sim 10\%$  of the whole data set. (This is not correct, dataset is random each batch)

refer: main.py > params = Parameters() ...



Note that this is our **training set** loss.

refer: `main.py > loss_history[-5:]`

```
[  
    0.2482442557811737,  
    0.23066267371177673,  
    0.21940861642360687,  
    0.26723822951316833,  
    0.218377023935318  
]
```