

---

# Toy Models of Superposition Replication and Findings

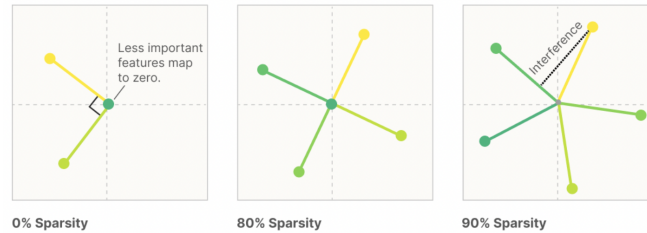
---

## Abstract

Toy Models of Superposition[1] is a groundbreaking paper published by researchers affiliated with Anthropic and Harvard University in 2022. By investigating small models with under 100 neurons, the paper demonstrates that neural networks can represent more features than they have dimensions. Additionally, they use these so called "toy models" to understand the relationship between how neural networks are trained and how they represent the data internally. This paper was able to replicate the finding from this paper and make new observations about "toy models" and how they behave under different training circumstances.

## 1 Introduction

The original paper motivates the idea of superposition with the following graphic:



The basic idea is this: if you think of each feature as being represented inside a neural network as a direction, you can graph these directions and observe them. In the graphic above, the researchers studied a model with two neurons and five inputs. Recreating this model showed the same results when sparsity was varied. A replication of the original graphic can be found below and the code used to generate it can be found [here](#).

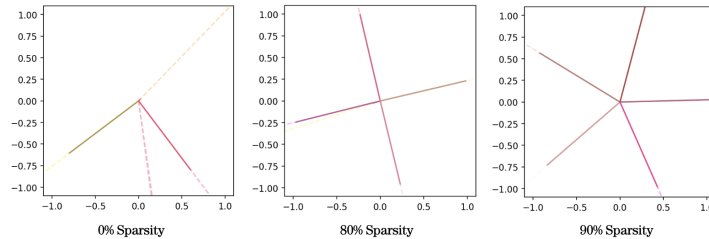


Figure 1: Graphing superposition in 2D.

Although this example only has two neurons (meaning it can be graphed in 2D), in future sections we will see how similar approaches scale nicely to models with more neurons. As far as I

understand, it is still unclear whether the ideas discussed in the original paper can be scaled to help understand much larger models such as GPT-4 or Claude 2.

## 2 Background and Motivation

In this section of Toy Models of Superposition[1], the authors provide context and define terms. In this paper, I make a few additional comments about some of the key ideas and terms the authors discuss.

**(1) Defining Features:** The original paper Toy Models of Superposition defines features broadly as "properties of the input which a sufficiently large neural network will reliably dedicate a neuron to representing." The authors do however describe this definition as "slightly circular" and note that they are not "overly attached to it." I find the definition especially problematic because a network that is small or has unconventional architecture may represent a feature that a larger network or a network with a more typical architecture may represent. These representations are clearly still features, but are not treated as so under the original definition.

As a result, I propose an alternative definition: features are aspects of the input that a neural network represents accurately with a higher probability than a randomly initialized network. In other words, features are parts of the input that a model determines to be important enough to represent internally.

**(2) Role of linear functions in Neural Networks:** The original authors of Toy Models of Superposition strongly emphasize the role of linear functions in neural networks. They claim that "Linear representations are the natural format for neural networks to represent information in!" Although this is clearly mostly true, I don't think there is sufficient evidence to justify the confidence of the authors in the original paper. Nonlinearities are also clearly doing lots of useful computation and it seems just as important to study their role in how a neural network represents features from layer to layer. The nonlinearity at each layer of a neural network may be just as important as the linear transformation in understanding how a model makes decisions.

**(3) Defining Superposition:** The original paper has a fantastic and simple definition for Superposition: "Roughly, the idea of superposition is that neural networks 'want to represent more features than they have neurons', so they exploit a property of high-dimensional spaces to simulate a model with many more neurons." This is the definition I will use throughout this paper.

## 3 Demonstrating Superposition

In this section the authors of the original paper demonstrate that superposition is observable even in models with more than two neurons (like the ones we saw in the introduction). They first study models with 20 inputs and 5 neurons. They train a linear model and one with a ReLU activation function. In the ReLU model they vary sparsity and observe superposition.

The authors demonstrate this by graphing  $W^T W$  where positive numbers are red and negative ones are blue. They also graph the length of each feature (by treating each column in

$W$  as a vector). Feature that are orthogonal to others in  $W$  are labeled black. Features that interfere with others are colored yellow.

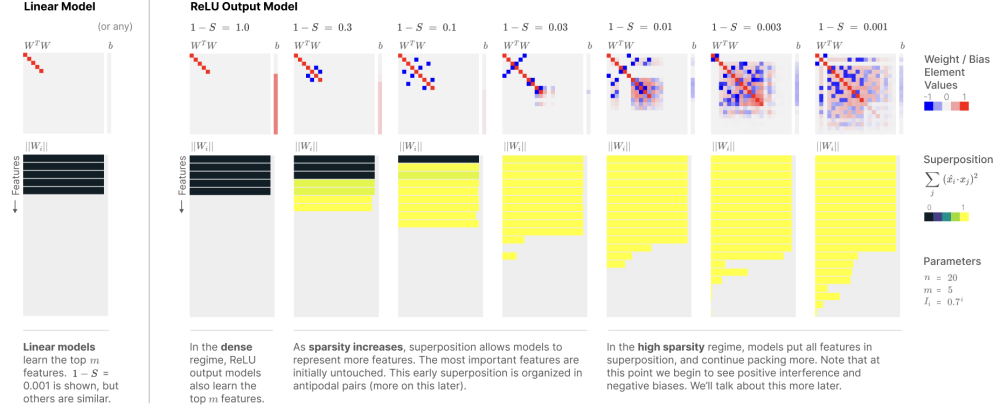


Figure 2: Demonstrating superposition in Toy Models of Superposition [1]

The first step in replicating these findings was to train the linear and ReLU models that don't perform computatoin in superposition. The linear model was defined by  $W^T W x + b$  and the ReLU model was defined by  $\text{ReLU}(W^T W x + b)$ . I trained both models with Adam (learning rate =  $1 * 10^{-3}$ ) on 20,000 batches of 256 examples.

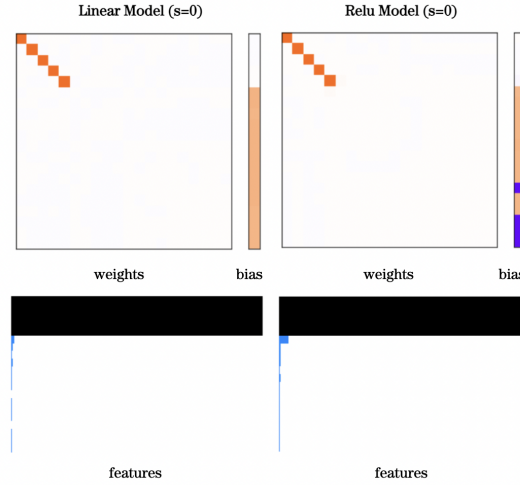


Figure 3: The generated graphics show that the model uses each of its five dimensions to graph the 5 most important features orthogonally.

There are a few small differences between my image and the original images that I think are worth noting before moving on. First, I use orange to indicate positive numbers and purple to indicate negative numbers. This is different from the red and blue used by the original authors to make it clear whose graphic is whose. Similarly, while the original authors use yellow to indicate features in superposition in Figure 2, I will use blue (hard to see in Figure 3 but will be more obvious going forward).

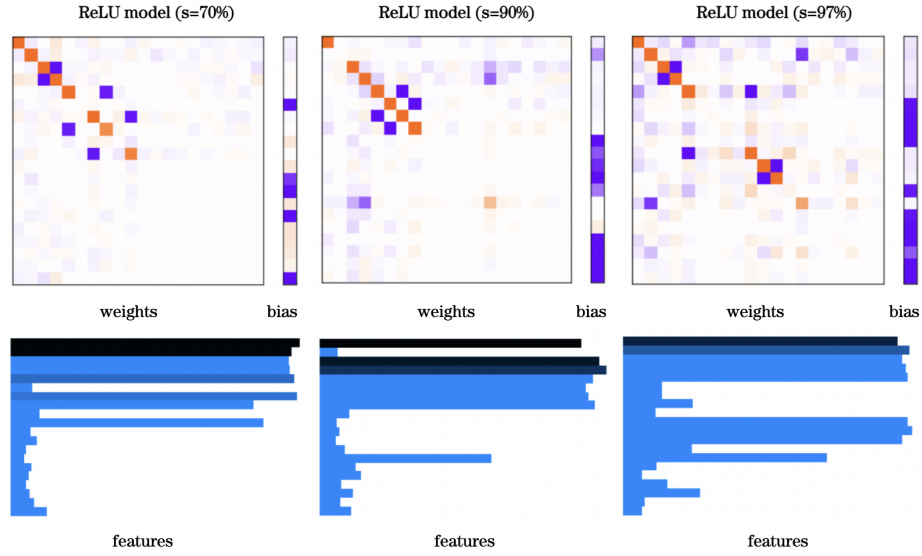


Figure 4: ADD THIS.

## References

- [1] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition, 2022.