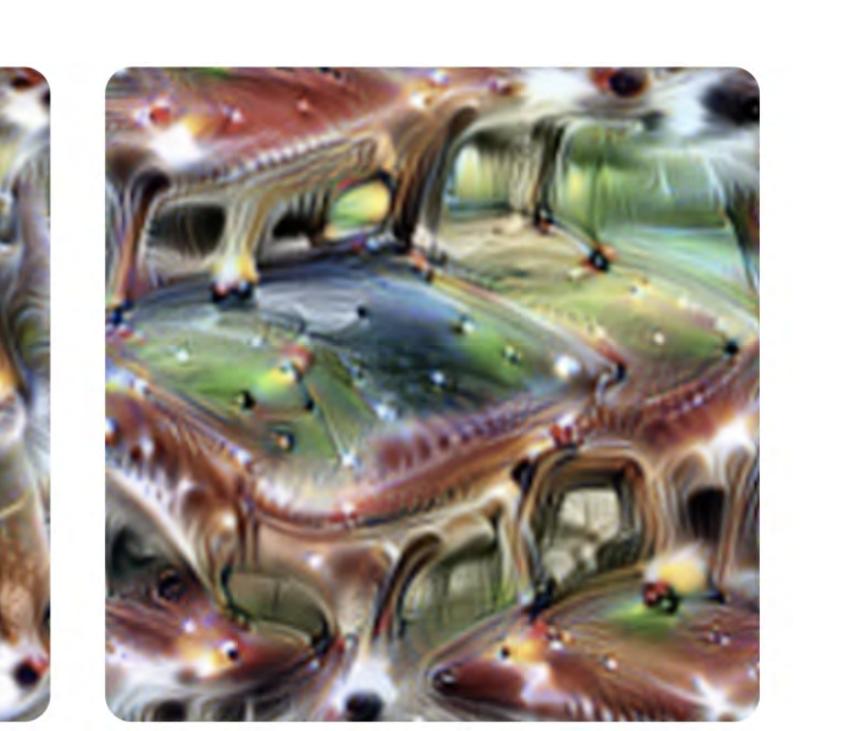
## Toy Models of Superposition

Guangyu Yang, Xueyan Li, Yawen Duan

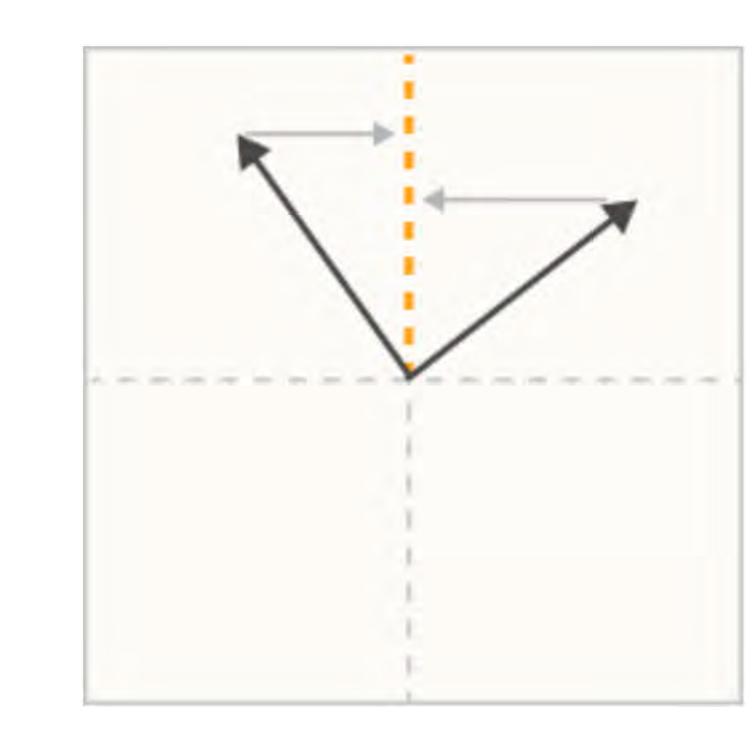
#### Background

Neural networks often contain "polysemantic neurons" that respond to multiple unrelated inputs.

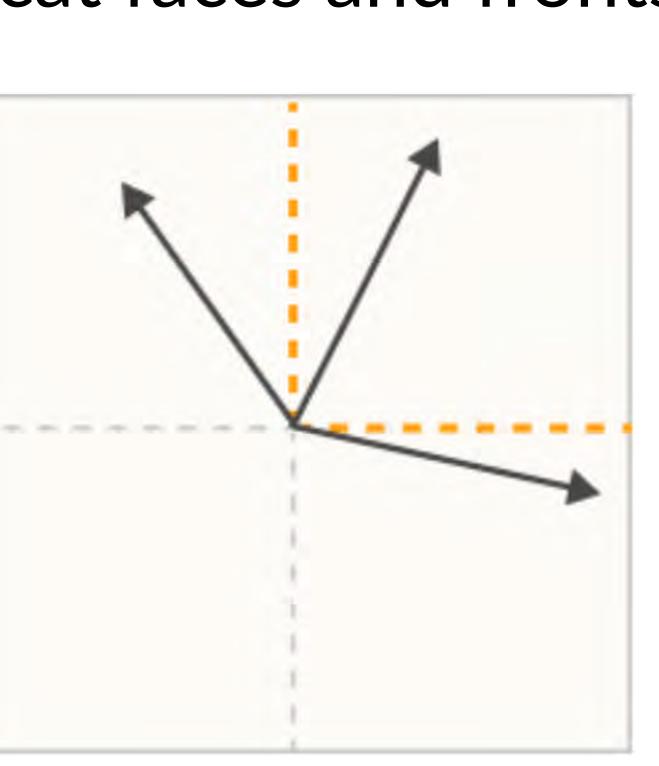




e.g. 4e:55 in InceptionV1 responses to both cat faces and fronts of cars.



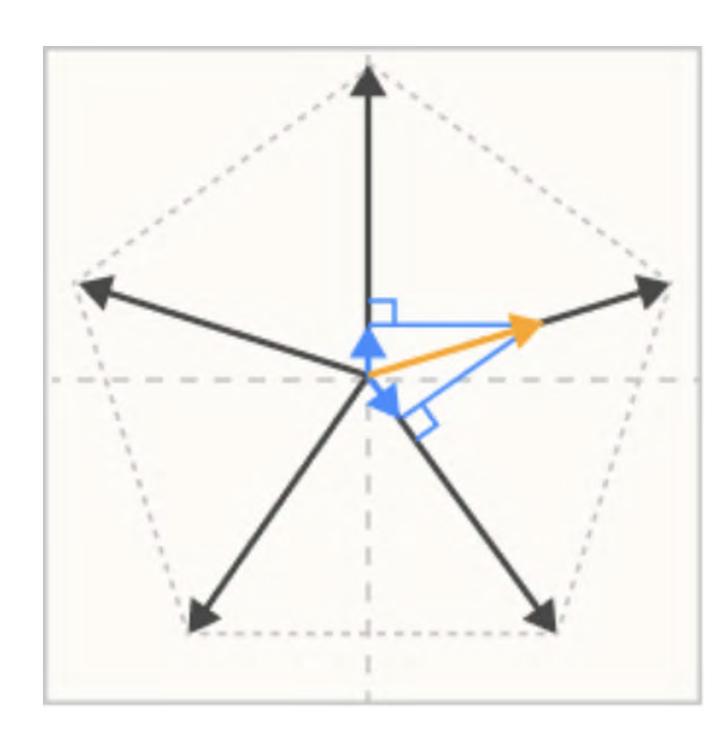
Polysemanticity is what we'd expect to observe if features were not aligned with a neuron.



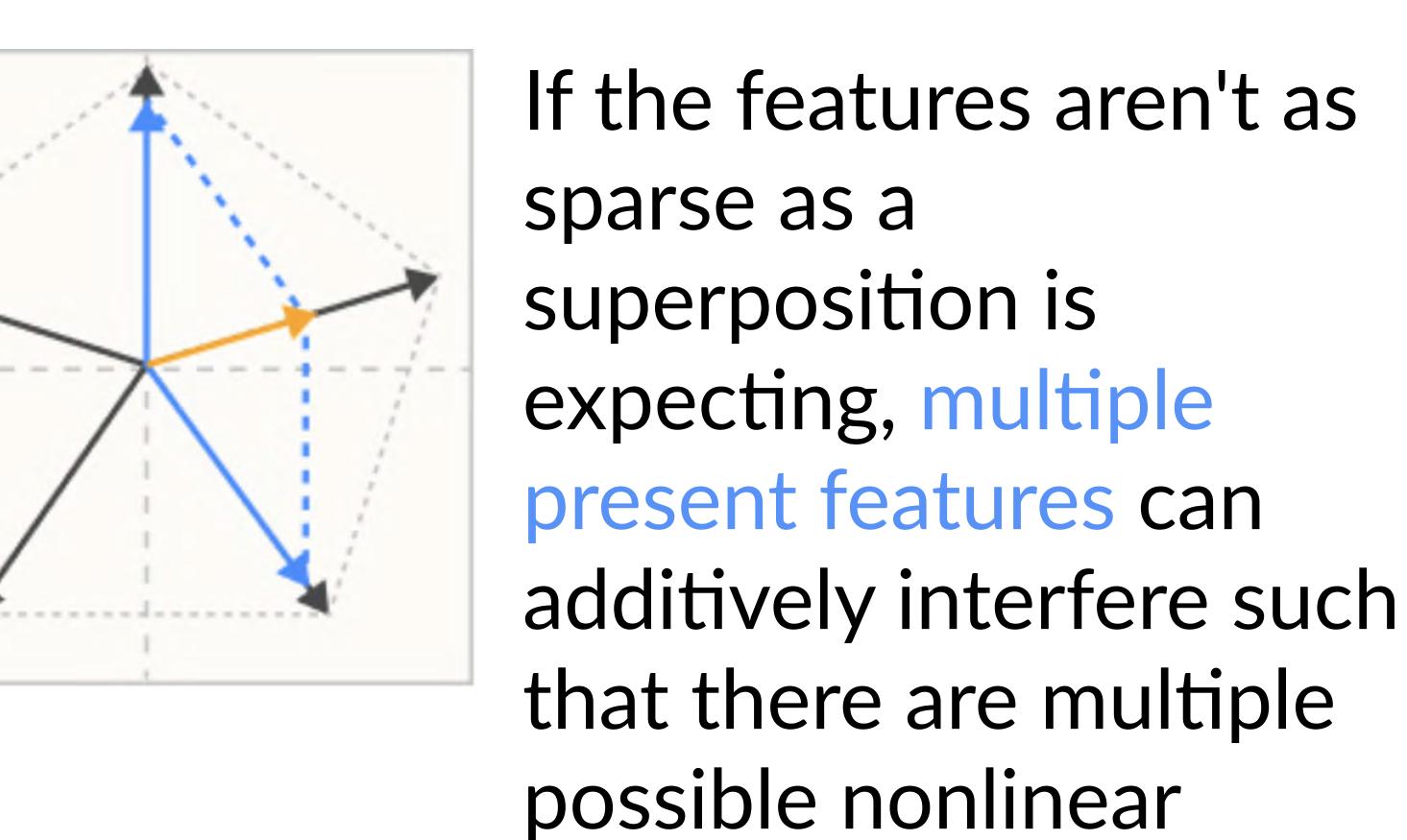
In the superposition hypothesis, features can't align with the basis because the model embeds more features than there are neurons. Polysemanticity is inevitable if this happens.

activation vector.

Superposition is when models represent more features than they have dimensions.

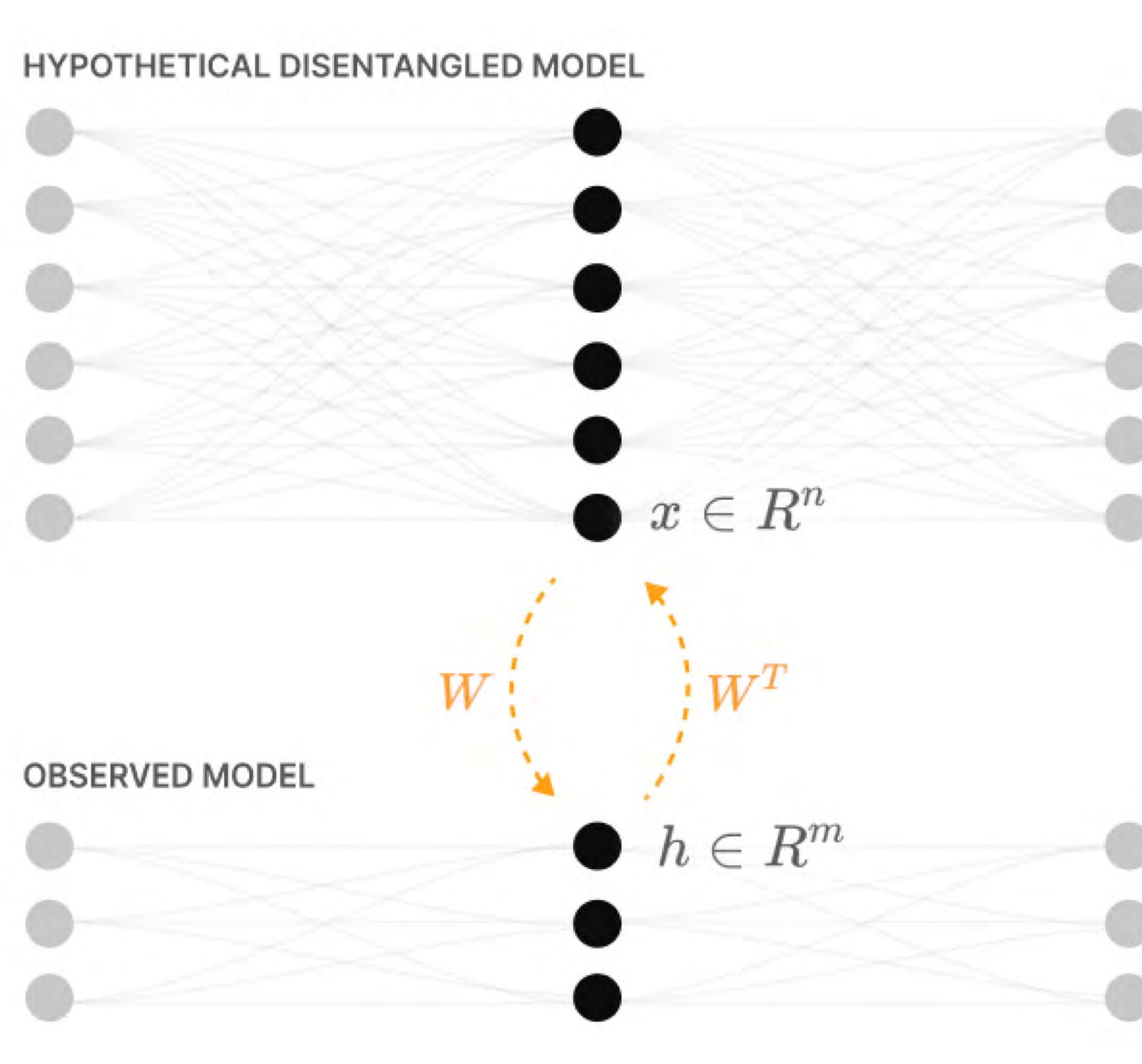


Even if only o e is active, using linear dot product projection on the superposition leads to interference which the model must tolerate or



Goal: demonstrating superposition

### Toy Model and Experiment Settings

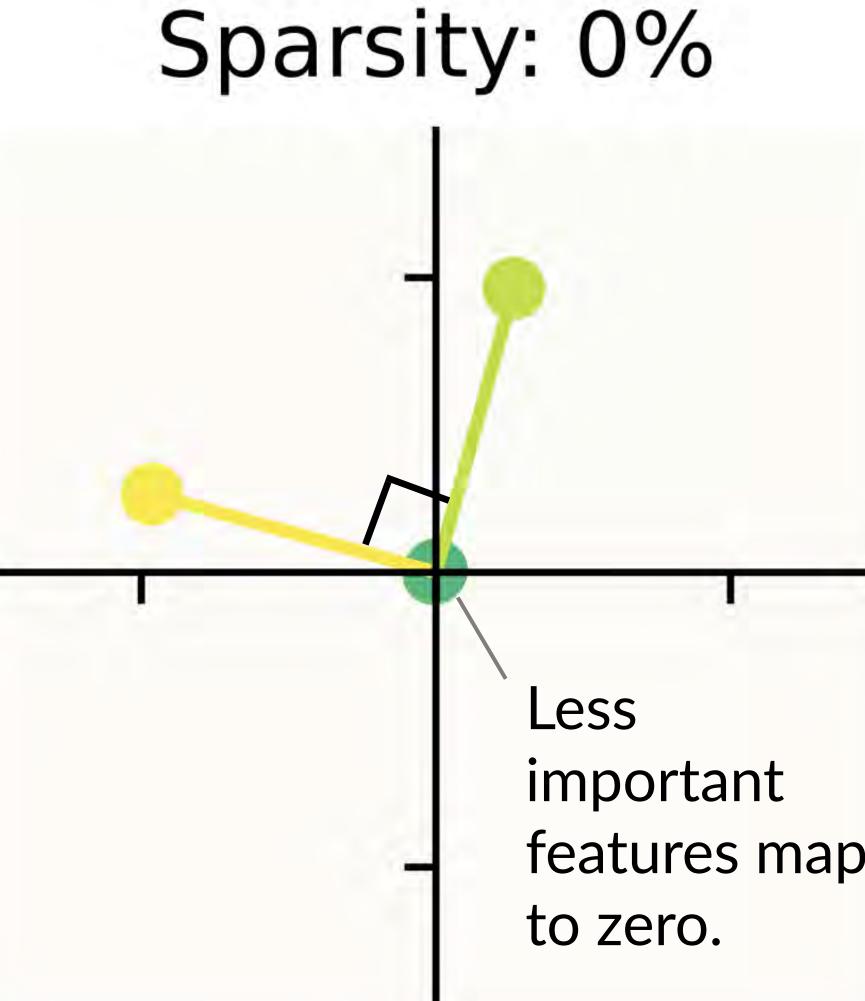


- Observations:
  - ullet Real-world feature is sparse --  $S_i$
  - More features than neurons -- n > m• Features vary in importance -- 1
- Parameters:
- Feature Vector:  $x \in \mathbb{R}^n$ Feature Sparsity:  $S_i$ Feature Importance:  $I_i$

with probability  $S_i$ ,

 $\text{Model: } x' = \text{ReLU}\left(W^TWx + b\right), W \in \mathbb{R}^{n \times m}$ Loss:  $L = \sum I_i (x_i - x_i')^2$ 

As sparsity increases, models use "Superposition" to represent more features than dimensions



additively interfere such that there are multiple reconstructions of an

# Sparsity: 80%

pentagons, and tetrahedrons.

Dedicated orthogonal The four most important dimensions for the two antipodal pairs. The least most important features. Less important features important features are not are not embedded. embedded.

#### Sparsity: 90% Feature **Importance** Medium **Parameters** m = $I_i = 0.9^i$

All five features are features are represented as embedded as a pentagon, but there is now "positive

• Conclusion 4: Superposition organizes features into geometric structures such as digons, triangles,

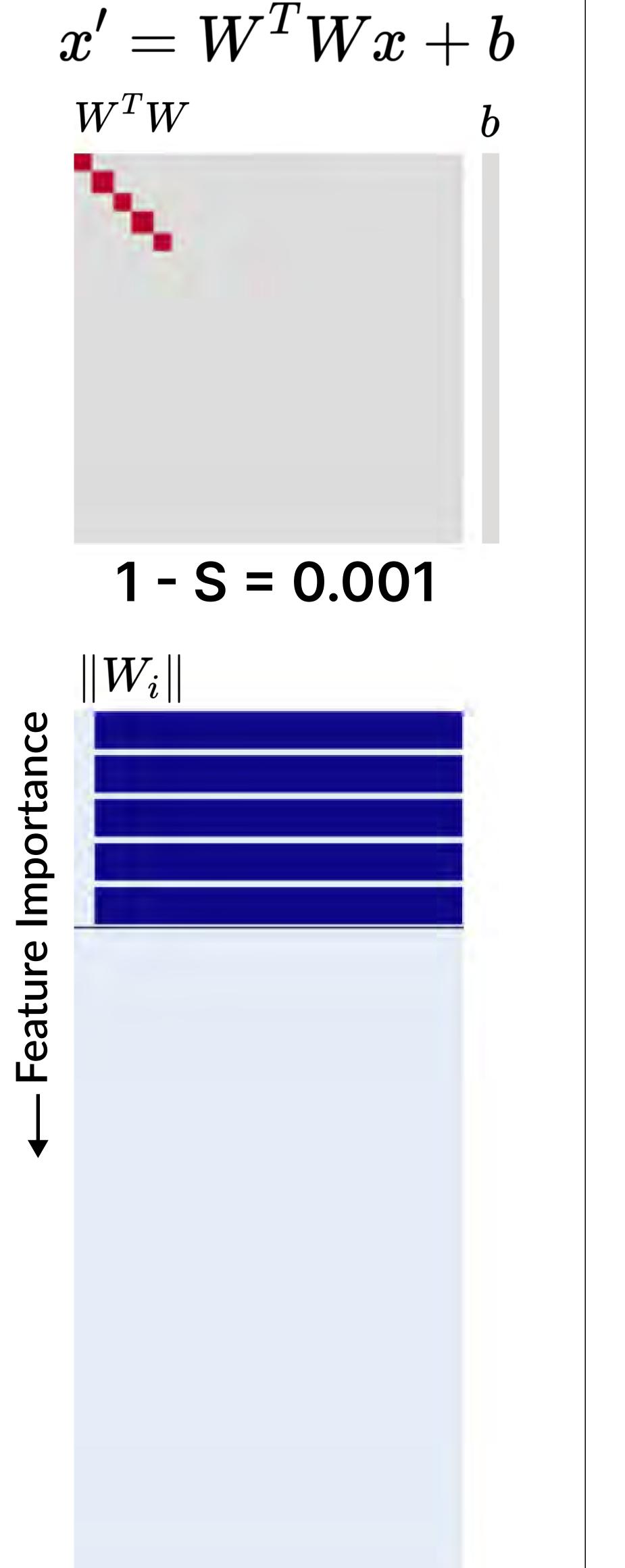
#### Demonstrating Superposition

• Conclusion 2: Both monosemantic and polysemantic neurons can form.

 $b W^T W$ 

ReLU Output Model

 $x' = \text{ReLU}(W^T W x + b)$ 



Linear models

learn the top

m features.

The Geometry of Uniform Superposition

Linear Model

learn the top early superposition is organized in antipodal

m features.

1 - S = 1.0As sparsity increases, superposition allows regime, ReLU models to represent more features. The most

 $oldsymbol{b} W^T W$ 

In the high sparsity regime, models put all features in superposition, and continue packing important features are initially untouched. This more. Note that at this point we begin to see positive interference and negative biases.

Superposition

 $\sum \left(\hat{x}_i \cdot x_j
ight)^{i}$ 

Parameters

m=5

 $_{\scriptscriptstyle L}=0.7^{i}$ 

1 - S = 0.001

Weights / Bias

Element Values

#### Superposition as a Phase Change

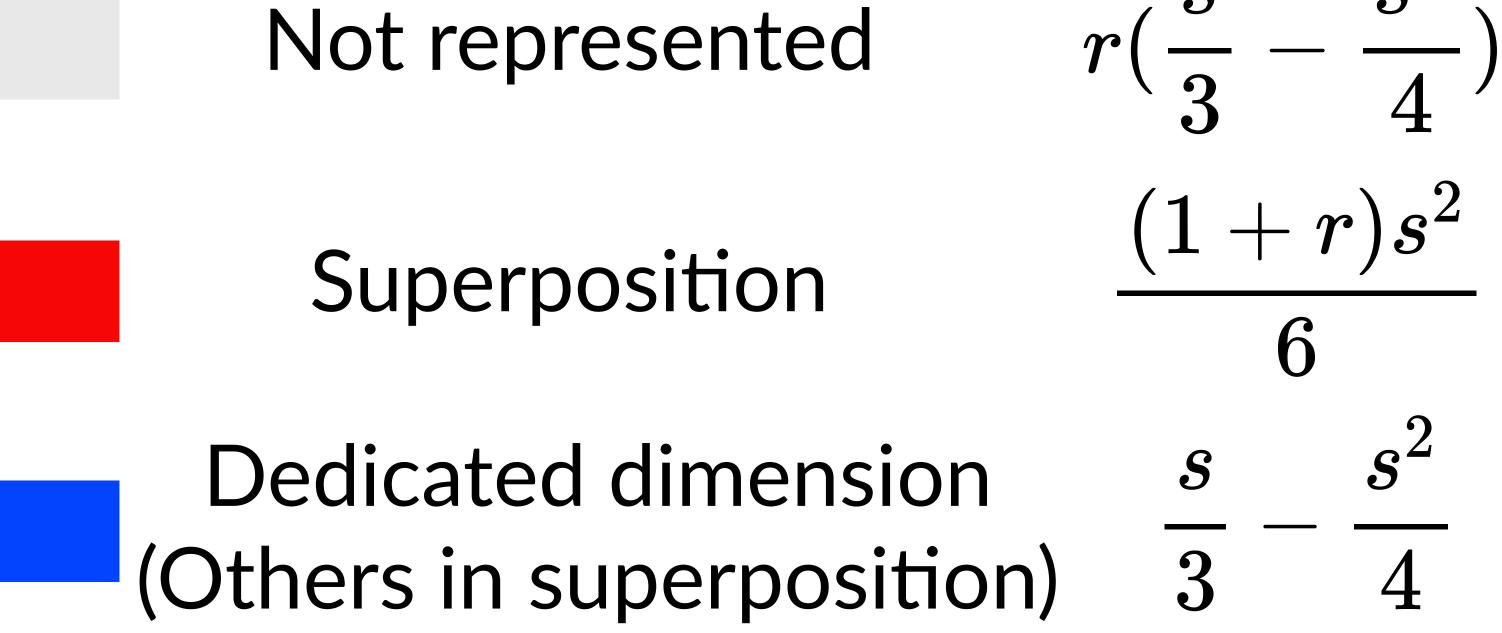
n = 2, m = 1

- Conclusion 3: Whether features are stored in superposition is governed by a phase change.
- We investigate into a simple 2-layer model with  $\boldsymbol{n}$  input (x) and output features (y) and  $\boldsymbol{m}=\boldsymbol{n-1}$ hidden neurons is generated with different importance r and sparsity s.

 $y = \text{ReLU}(W^t W x + b)$ 

$$\log = (x_1 - y_1)^2 + \ldots + \ (x_{n-1} - y_{n-1})^2 + r(x_n - y_n)^2$$

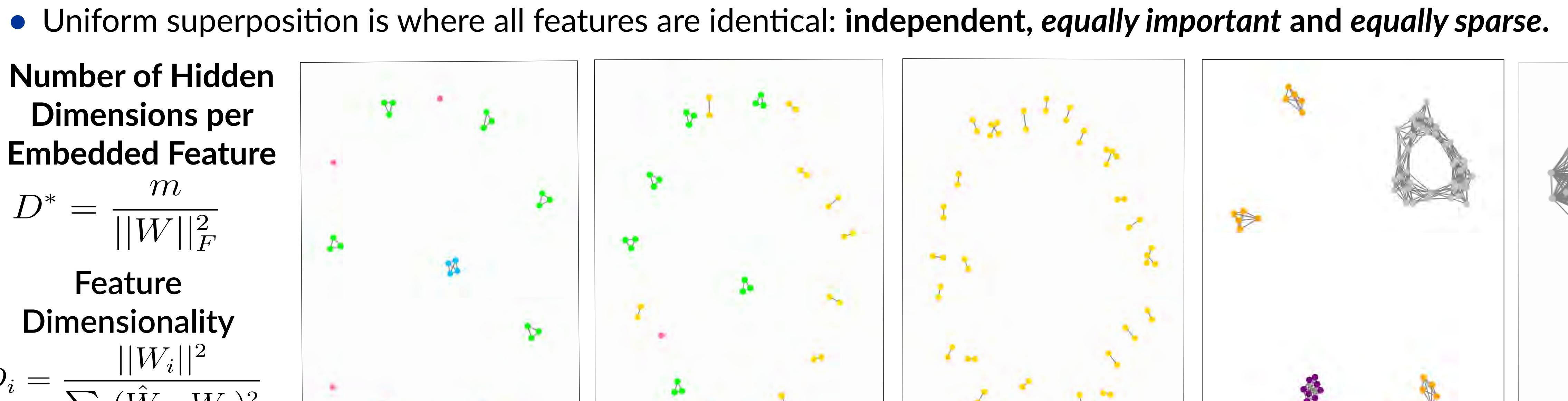
Theoretical loss for feature of interest

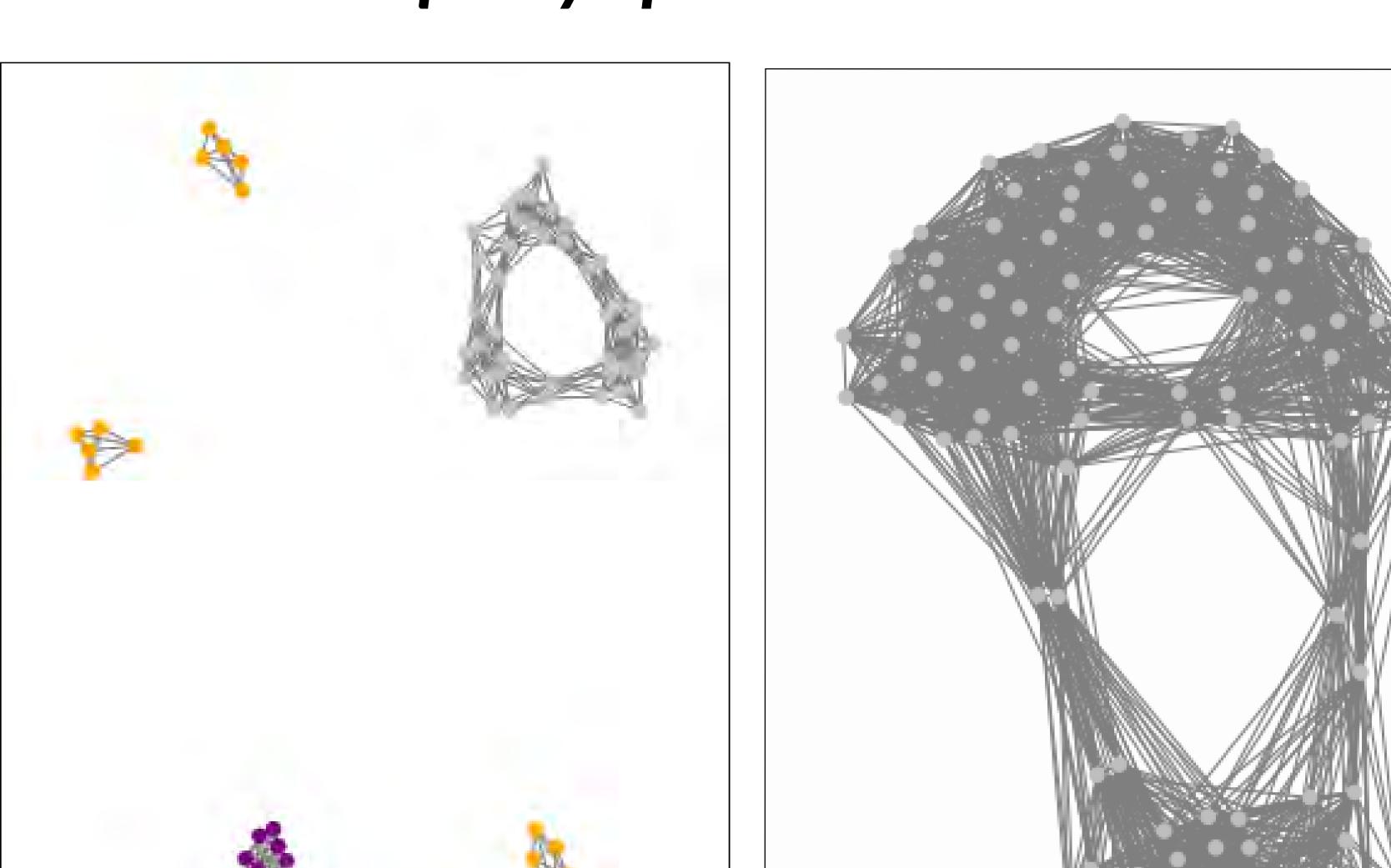


Dedicated dimension (Others not represented)

n = 3, m = 2

Number of Hidden Dimensions per **Embedded Feature** Feature Dimensionality





- The number of dimensions per feature (blue curve) is "sticky"
- at 1 and 1/2. The feature dimensionalities (black dots) cluster at certain fractions labeled with colored lines.
- These fractions correspond to

#### Relationship to Adversarial Robustness\*

Non-Superposition:  $(W^TW)_0 = (1,0,0,0,\ldots)$ With-Superposition:  $(W^TW)_0 = (1, \epsilon, -\epsilon, \epsilon, \ldots)$ 

