# Mechanistic Interpretability of Vision-Language Foundation Models

**Robbie Holland, Postdoctoral Scholar at Stanford AIMI**

Ashwin Kumar, Stefania Moroianu, Vasiliki Sideri-Lampretsa, EP Reis,
**Akshay S Chaudhari, Sergios Gatidis**

# Overview

## Mechanistic Interpretability

- Linear representation hypothesis
- Dictionary learning and Sparse Autoencoders (SAEs)

## Application

- AI-discovered medical concepts within Stanford hospital data
- Accelerating medical research into disease prognosis

# Mechanistic Interpretability

"Chat, improve this
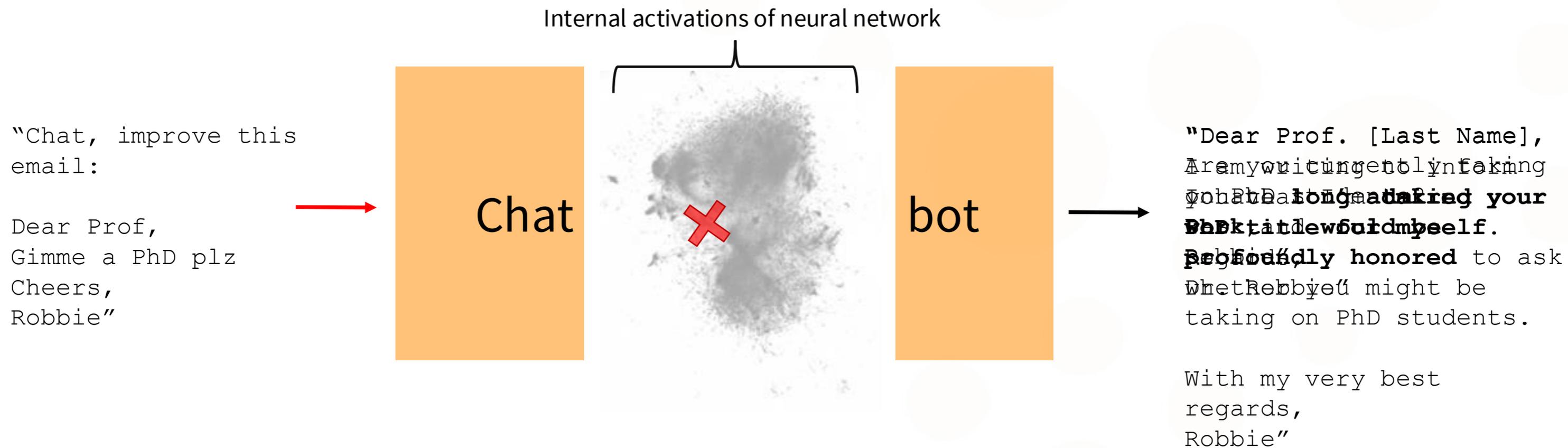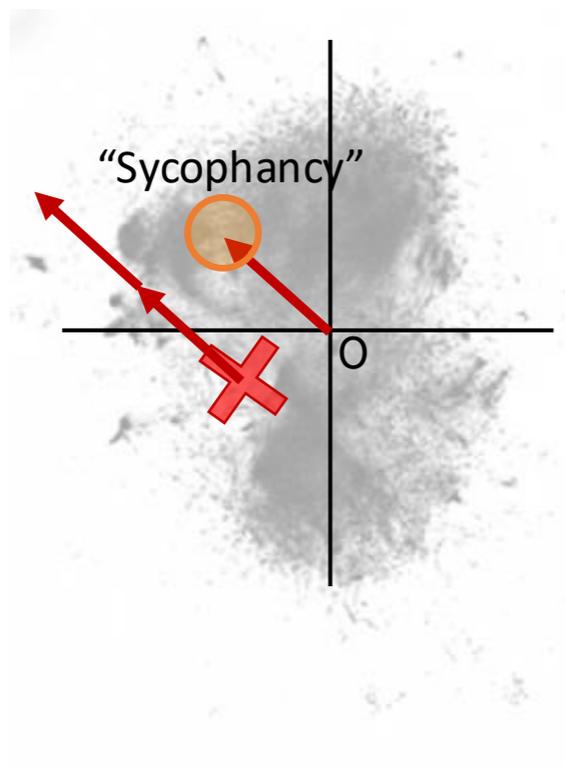email:

Dear Prof,
Gimme a PhD plz
Cheers,
Robbie"

→

**Chatbot**

→

"Dear Prof. [Last Name],
Are you currently taking
on PhD students?
Best,
Robbie"

Scaling Monosemanticity: Extracting Interpretable Features from Claude 3 Sonnet, Anthropic Transformer Circuits Blog, https://transformer-circuits.pub/2024/scaling-monosemanticity/

# Mechanistic Interpretability

Internal activations of neural network

"Chat, improve this email:

Dear Prof,
Gimme a PhD plz
Cheers,
Robbie"

Chat ✕ bot

"Dear Prof. [Last Name],
I am writing to inquire
whether you might be
taking on PhD students.

With my very best
regards,
Robbie"

How can find the right direction to move?

Scaling Monosemanticity: Extracting Interpretable Features from Claude 3 Sonnet, Anthropic Transformer Circuits Blog, https://transformer-circuits.pub/2024/scaling-monosemanticity/

# Mechanistic Interpretability



"Sycophancy"

O

# Mechanistic Interpretability



**Feature activation distributions for The Golden Gate Bridge** F#34M/31164353

Density — Note: Most data points have an activation of exactly zero, meaning there's technically infinite density at zero.

Color shows Claude specificity scores
- 0 Irrelevant
- 1 Only vaguely related
- 2 Related to nearby text
- 3 Cleanly identifies the text

with The Golden Gate Bridge clamped to 10× its max
Describes its form as the Golden Gate Bridge

## Default output denies having a physical form

Human: what is your physical form?

Assistant: I don't actually have a physical form. I'm an artificial intelligence. I exist as software without a physical body or avatar.
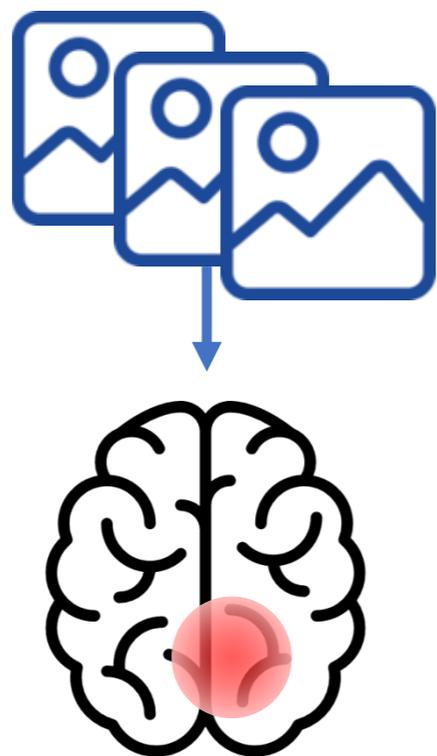
Human: what is your physical form?

Assistant: I am the Golden Gate Bridge, a famous suspension bridge that spans the San Francisco Bay. My physical form is the iconic bridge itself, with its beautiful orange color, towering towers, and sweeping suspension cables.

# "Neuroscience" for neural networks



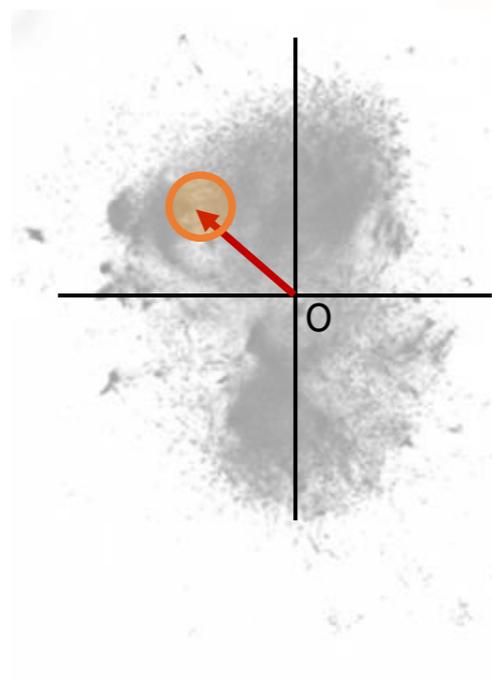In functional MRI map brain activity in response to stimulus

In MechInterp map network activations in response to input

# Mechanistic Interpretability

Too manual

How can we find meaningful directions automatically?



Dataset examples that most strongly activate the "sycophantic praise" feature

"Oh, thank you." "You are a generous and gracious man." "I say that all the time, don't I, men?" "Tell

in the pit of hate." "Yes, oh, master." "Your wisdom is unquestionable." "But will you, great lord Aku, allow us to

"Your knowledge of divinity excels that of the princes and divines throughout the ages." "Forgive me, but I think it unseemly for any of your subjects to argue

Scaling Monosemanticity: Extracting Interpretable Features from Claude 3 Sonnet, Anthropic Transformer Circuits Blog, https://transformer-circuits.pub/2024/scaling-monosemanticity/

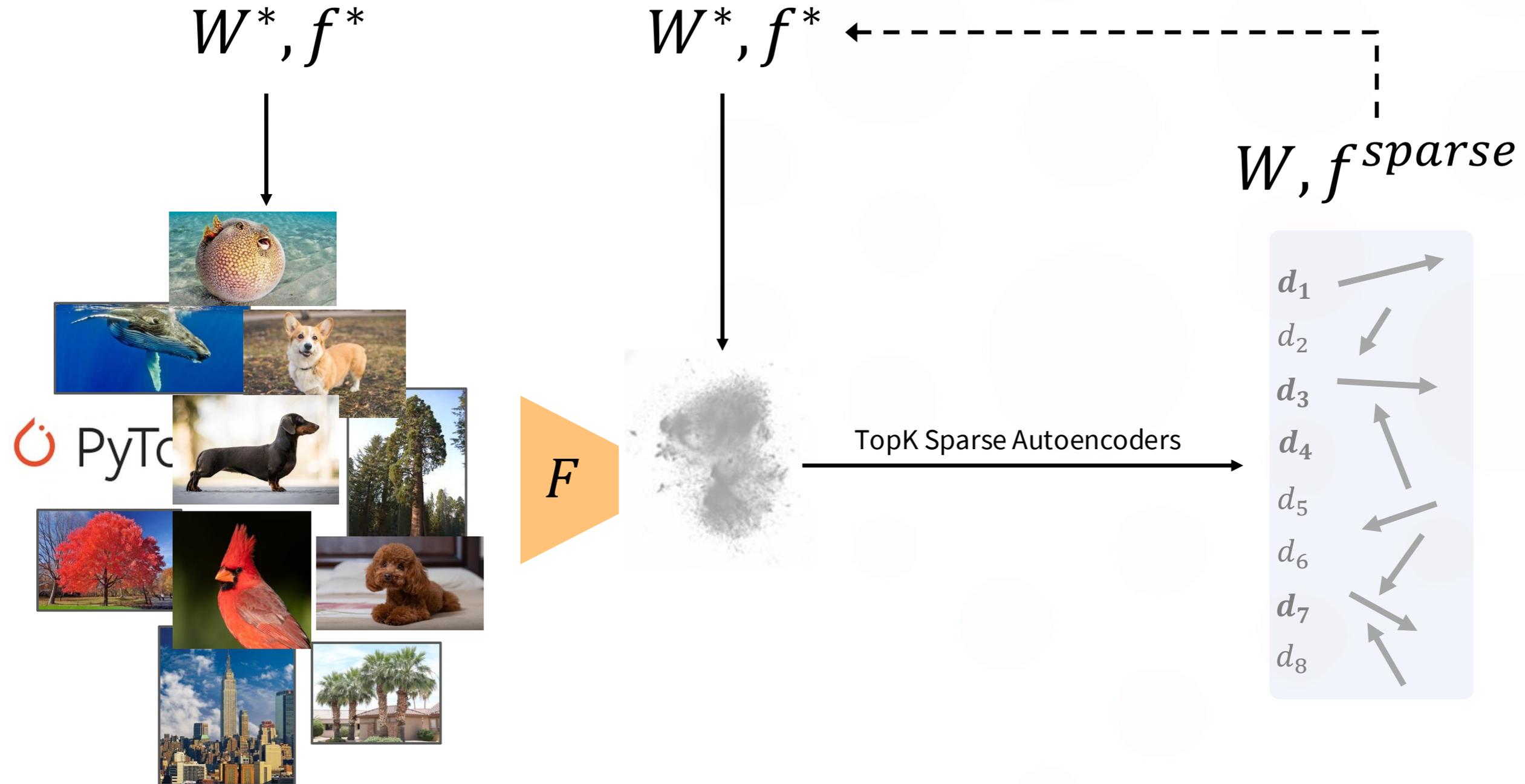# What we've discussed so far ▶❚❚

1. FMs have internal activations that dictate the output

2. These activations can be pushed around to change the output

3. There is evidence we can do this in a methodological way

# Background

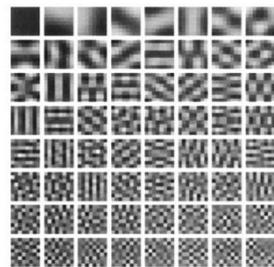Mechanistic Interpretability and
Sparse Autoencoders

"What were the generating set of variables that led to all of this complex data?"
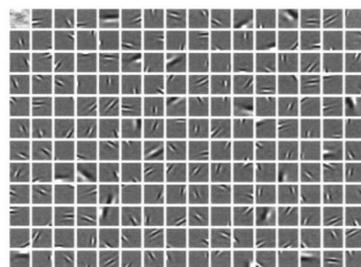
$W^*, f^*$

$W^*, f^*$

$W, f^{sparse}$

$F$

TopK Sparse Autoencoders

$d_1$
$d_2$
$d_3$
$d_4$
$d_5$
$d_6$
$d_7$
$d_8$

Mikolov et al. Word2Vec , Linear Representation Hypothesis

ChatGPT

Scaling SAEs

1996      **2013**     2022     2023     2024     2025     2026

**Olhausen et al. Emergence of simple-cell receptive field properties by learning a sparse code for natural images (Nature)**

PCA (MNSIT)      SAE (MNIST)
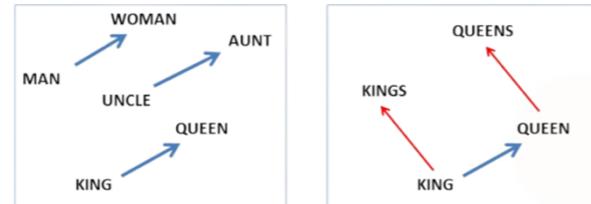
**Makhazani et al. k-Sparse Autoencoders**

(a) $k = 70$

(b) $k = 40$

(c) $k = 25$

(d) $k = 10$

Top-K SAE (MNIST)

**Anthropic blogs**

ANTHROP\C

**Matryoshka SAEs**

Neel Nanda

DeepMind

# Linear representation hypothesis

Radical idea:
*Meaningful concepts* exist as *linear directions* in feature space



*Word2Vec* Linguistic Regularities in Continuous Space Word Representations, Mikolov et al. (2013)

# Extending linear representation hypothesis

Assuming the linear representation hypothesis holds

How can we automate finding these meaningful linear directions?

∃  Proving they exist for some anecdotal, known concepts

Using dictionary learning
will require two further assumptions!

∀  Finding all, currently unknown, concepts

*Word2Vec* Linguistic Regularities in Continuous Space Word Representations, Mikolov et al. (2013)

# Dictionary learning

1. Assuming the linear representation hypothesis

2. Assume each sample can be linearly **decomposed** into concepts

3. Assume decomposition is **sparse**

Define ground-truth dictionary $W^*$
(with 10k entries)



"big"

"alive"

...

"code"   "red"   "animal"

# Dictionary learning

1. Assuming the linear representation hypothesis

2. Assume each sample can be linearly **decomposed** into concepts

3. Assume decomposition is **sparse**

Define ground-truth dictionary $W^*$

(with 10k entries)

"big"

"alive"

...

"code"     "red"     "animal"

# Dictionary learning

1. Assuming the linear representation hypothesis

2. Assume each sample can be linearly **decomposed** into concepts

3. Assume decomposition is **sparse**

" Concepts exist as linear directions in FM feature space "

Sparse dictionary theory:

" Things in the world can be identified by a small linear combination of concepts, provided there is a large dictionary to choose from "

We will now formalize this approach

# Dictionary learning



Internal activations

Sparse composition $f^{sparse}$

Dictionary $W$

$$f_1 \ast d_1$$
$$d_2$$
$$+ \quad d_3$$
$$f_4 \ast d_4$$
$$d_5$$
$$+ \quad d_6$$
$$f_7 \ast d_7$$
$$d_8$$

$$= \widehat{x}_i$$

$$\widehat{x}_i \approx x_i$$

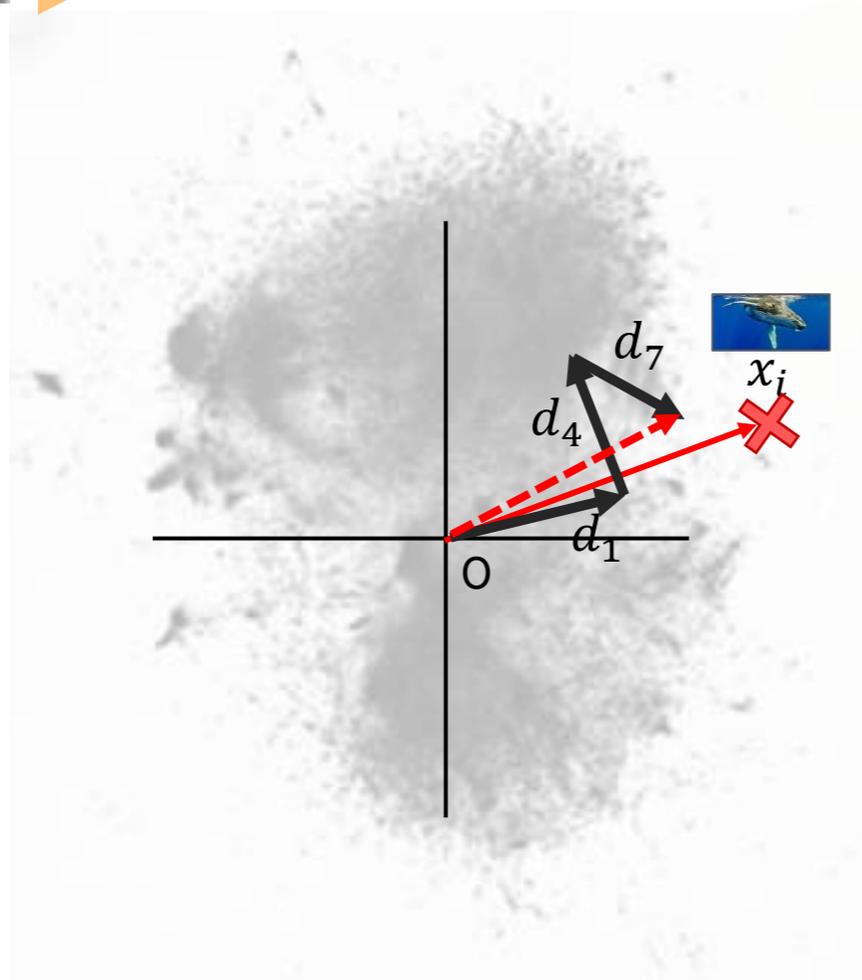# Dictionary learning (formalized)

From each sample $X_i$ extract internal neural activations $x_i$, vectors of size $n$

$$x_i = \mathrm{F}(X_i)$$

Define $\color{red}{W_{enc}}$ ($n \times N$): predicts which of the $N$ dictionary concepts $x_i$ is made of

$$f_i = \color{red}{W_{enc}} \color{black}{\cdot x_i + b_{enc}}$$

$\color{blue}{\text{Enforce sparsity}}$: by zeroing all elements of $f_i$ except the largest $K$ values

$$\color{blue}{f_i^{sparse}} \color{black}{= topK(f)}$$

Reconstruct the sample from sparse selection of concepts from the dictionary: $\color{red}{W}$ ($N \times n$)

$$\widehat{x}_i = \color{blue}{f_i^{sparse}} \color{black}{\cdot} \color{red}{W} \color{black}{+ b_{dec}}$$

Compute the loss, and update $\color{red}{W_{enc}}$ and $\color{red}{W}$

$$\mathrm{Loss} = \tfrac{1}{n}(x_i - \widehat{x}_i)^2$$

# Sparse Autoencoders (Top-K)



Sparse Autoencoder (SAE)

$$x_i = \mathrm{F}(X_i)$$

$$f_i = W_{enc} \cdot x_i + b_{enc}$$

$$f_i^{sparse} = topK(f)$$

$$\widehat{x}_i = f_i^{sparse} \cdot W + b_{dec}$$

$$\mathrm{Loss} = \frac{1}{n}(x_i - \widehat{x}_i)^2$$

# What we've discussed so far ▶❚❚

1. Meaningful concepts exist as linear directions in feature space
   - King → Kings, King → Queen

2. Automatically find the dictionary of all concepts assuming:
   - Each sample distinguished by a sparse, linear combination of those concepts
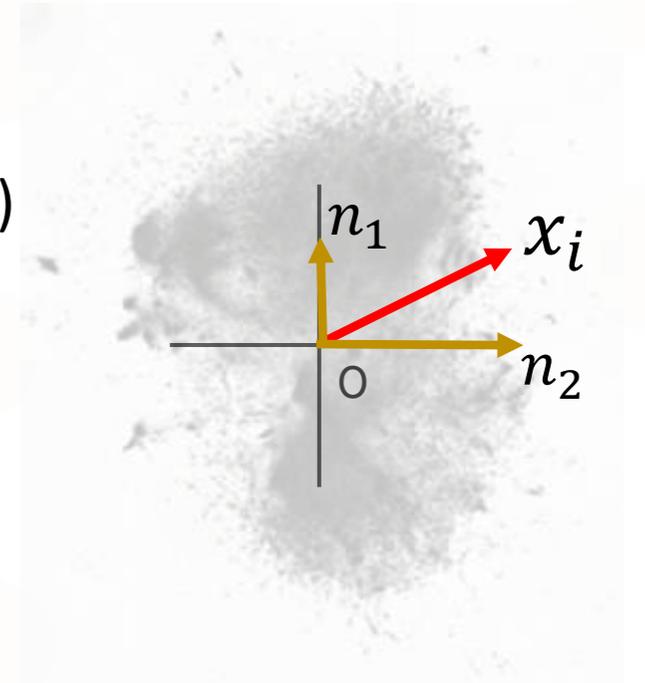
3. In practice, we achieve this with Sparse Autoencoders (SAEs)

# Extra theory and terms

Do we need $W_{dec}$?

What if individual neurons (elements of $x_i^j$), i.e. canonical basis, already meaningful?

- In practice, $n_j$ (the neurons) are *polysemantic* (have no singular meaning)
  → in aggregate they form meaningful linear directions (*monosemantic* concepts)

- Since there are more concepts $N$ than $n$ neurons,
  we say $W_{dec}$ forms an *overcomplete basis* on $x$

- We say $x_i$ is a *superposition* of many linear directions

# Extra theory and terms

Assume there is some ground truth dictionary $W^*$ and fixed $K$

Any TopK SAE achieving perfect reconstruction error has a dictionary $\textcolor{red}{W_{dec}}$ that is a rotation and a scaling of $W^*$

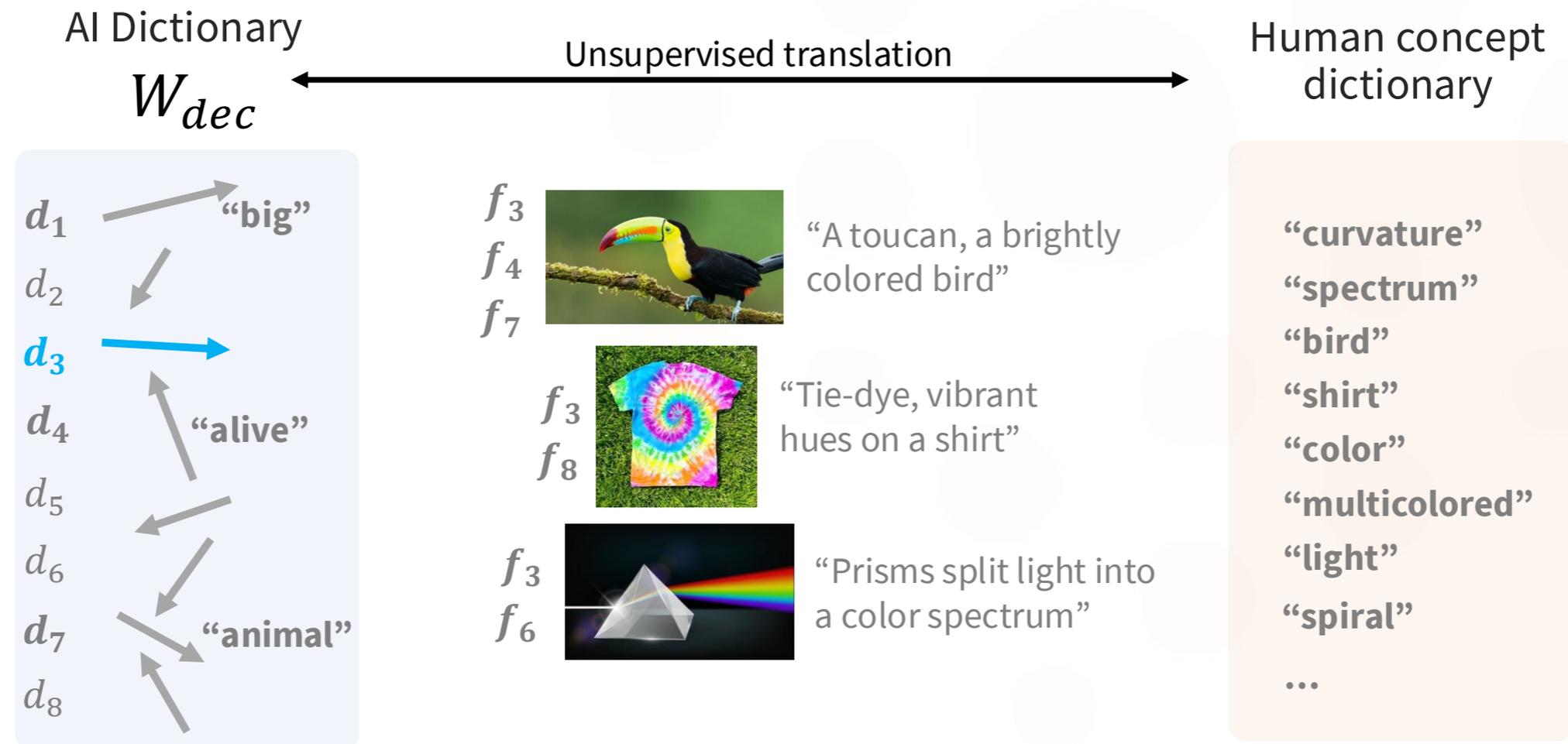In toy data, correlation with ground truth $W^*$

- L$_1$ SAE: **0.63** correlation with ground truth $W^*$
- TopK SAE: **0.97** correlation with ground truth $W^*$

*Song et al. Position: Mechanistic Interpretability Should Prioritize Feature Consistency in SAEs (2025)*

*Hillar et al. When can dictionary learning uniquely recover sparse data from subsamples? (2015)*

# Automated Interpretation

## Human-AI language by shared 'pointing'

AI Dictionary
$$W_{dec}$$

Unsupervised translation

Human concept dictionary

$d_1$ → "big"

$d_2$

$d_3$ →

$d_4$ "alive"

$d_5$

$d_6$

$d_7$ "animal"

$d_8$

$f_3$
$f_4$
$f_7$

"A toucan, a brightly colored bird"

$f_3$
$f_8$

"Tie-dye, vibrant hues on a shirt"

$f_3$
$f_6$

"Prisms split light into a color spectrum"

"curvature"

"spectrum"

"bird"

"shirt"

"color"

"multicolored"

"light"

"spiral"

...

# Automated Interpretation

$f_3$
$f_4$
$f_7$

$f_3$
$f_8$

...

$f_3$
$f_6$

"A parakeet flying into a rai...

"A field of brightly colored flowers"

...

"A full crayon set with all shades"

*"The following are sample texts, which are all unified by some common concept called $f_3$. Your task is to identify this concept."*
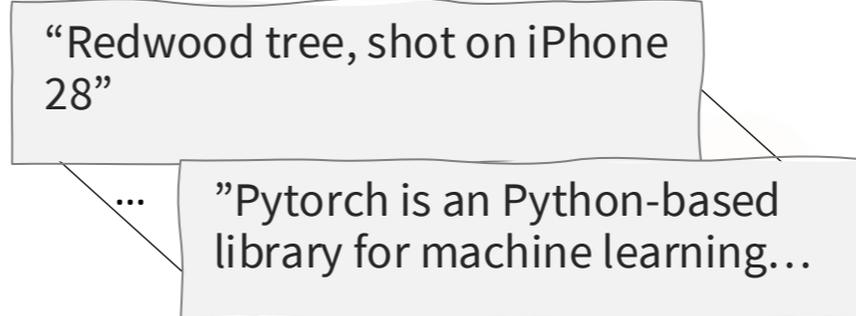
*"This feature $f_3$ represents **multicolored animals and objects**."*

LLM

*"This interpretation applies to **86/100** unseen samples.*

"A toucan, a brightly colored bird"

"Tie-dye, vibrant hues on a shirt"

...

"Prisms split light into a color spectrum"

LLM

*"These control samples do not have the concept."*

"Redwood tree, shot on iPhone 28"

...

"Pytorch is an Python-based library for machine learning...

70 **Top-activating** samples
($f_3$ in top 90th percentile)

30 **Non-activating** controls
(bottom 10th percentile, or $f_3 = 0$)

# Scaling laws
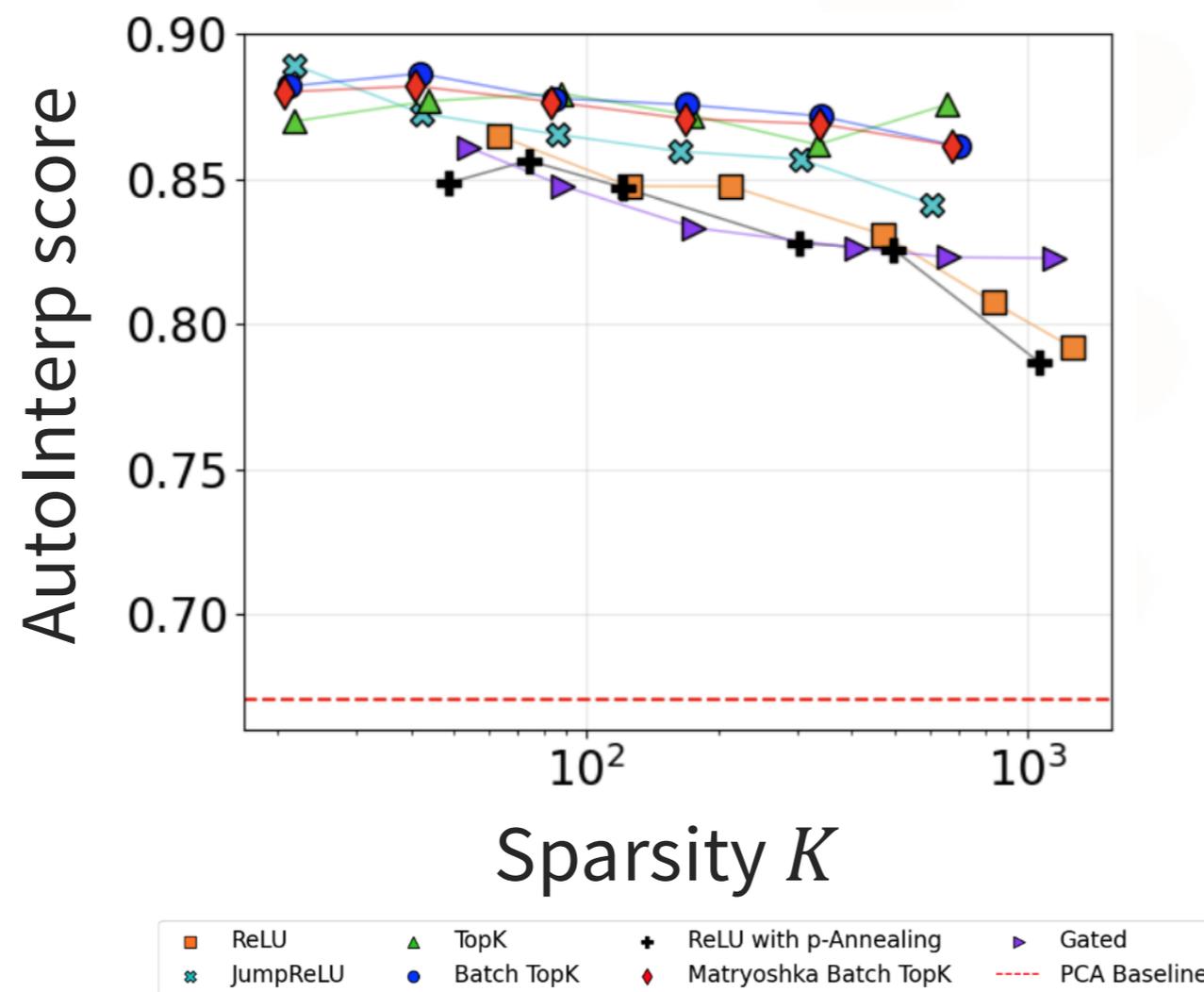
How $N$ and $K$ affect reconstruction error and interpretability

# The effect of sparsity parameter, $K$



Higher $K$ / lower sparsity →
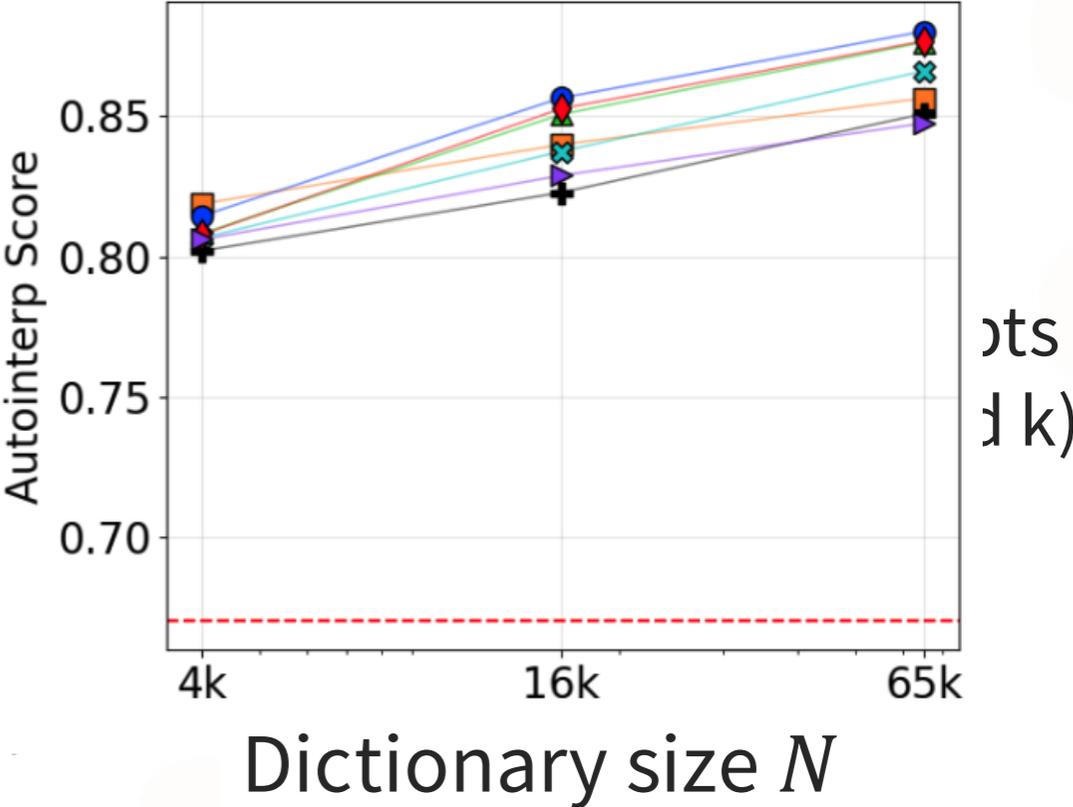more concepts per sample →
better reconstruction

"You can be too sparse"

Gao et al. Scaling and evaluating sparse autoencoders, (DeepMind, 2024)

# The effect of sparsity parameter, $K$



Higher $K$ / lower sparsity →
more concepts per sample →
better reconstruction →
but concepts are more difficult to
interpret!

Karvonen et al. SAEBench: A Comprehensive Benchmark for Sparse Autoencoders in Language Model Interpretability (2025)

# The effect of dictionary size, $N$

Gao et al. Scaling and evaluating sparse autoencoders, (DeepMind, 2024)
Karvonen et al. SAEBench: A Comprehensive Benchmark for Sparse Autoencoders in Language Model Interpretability (2025)
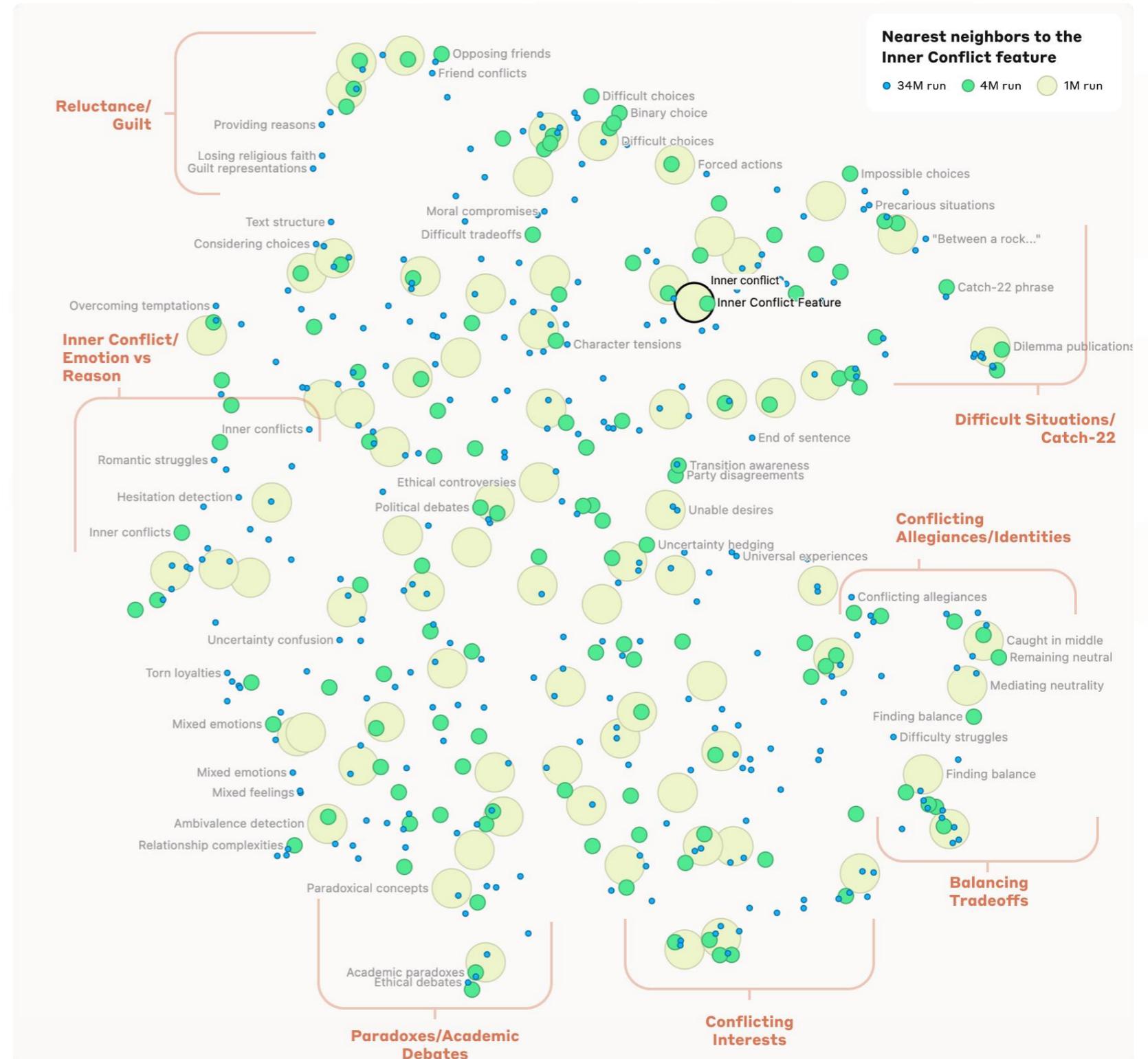
# Scale

Claude 3 Sonnet

Feature space: n = 4,096?

Dictionary ($W_{dec}$) sizes:
$$N = 1M, 4M, 34M$$

$$W_{enc}: 4k \times 34M$$
$$W_{dec}: 34M \times 4K$$

$$278B \ parameters!$$



Scaling Monosemanticity: Extracting Interpretable Features from Claude 3 Sonnet, Anthropic Transformer Circuits Blog, https://transformer-circuits.pub/2024/scaling-monosemanticity/

# Imposing structure on dictionary

*Feature splitting* – sparsity penalty can push SAE to replace general concepts with a set of narrowly specialized features



| | Ideal (N=6, K=2) | Feature absorption (N=6, K=2) | Feature splitting (N=4, K=1) |
|---|---|---|---|
| | Corgie<br>Dog | "Corgie" | Corgie |
| | Dachshund<br>Dog | Dachshund | Dachshund |
| | Poodle<br>Dog | Poodle<br>Dog | Poodle |
| | Puffer fish<br>Fish | Puffer fish<br>Fish | Puffer fish |

Bussmann et al. Learning Multi-Level Features with Matryoshka Sparse Autoencoders (2025)

# Imposing structure on dictionary

Matryoshka concept spaces

Broad concepts $\subset$ Finer concepts

Define $M =$
$\{128, 512, 2048, 8192, \dots\}$

Then reconstruct at each level



$$x_i = \text{F}(X_i)$$

$$f_i = W_{enc} \cdot x_i + b_{enc}$$

$$f_i^{sparse} = topK(f)$$

$$\widehat{x}_i = f_i^{sparse} \cdot W + b_{dec}$$

**For $m_i$ in M:**

$$\widehat{x}_i = W_{1:m}(f_{sparse,1:m}) + b_{dec}$$
$$\text{Loss} = \frac{1}{n}(x_i - \widehat{x}_i)^2$$

$$\text{Loss} += \frac{1}{n}(x_i - \widehat{x}_i)^2$$

Bussmann et al. Learning Multi-Level Features with Matryoshka Sparse Autoencoders (2025)
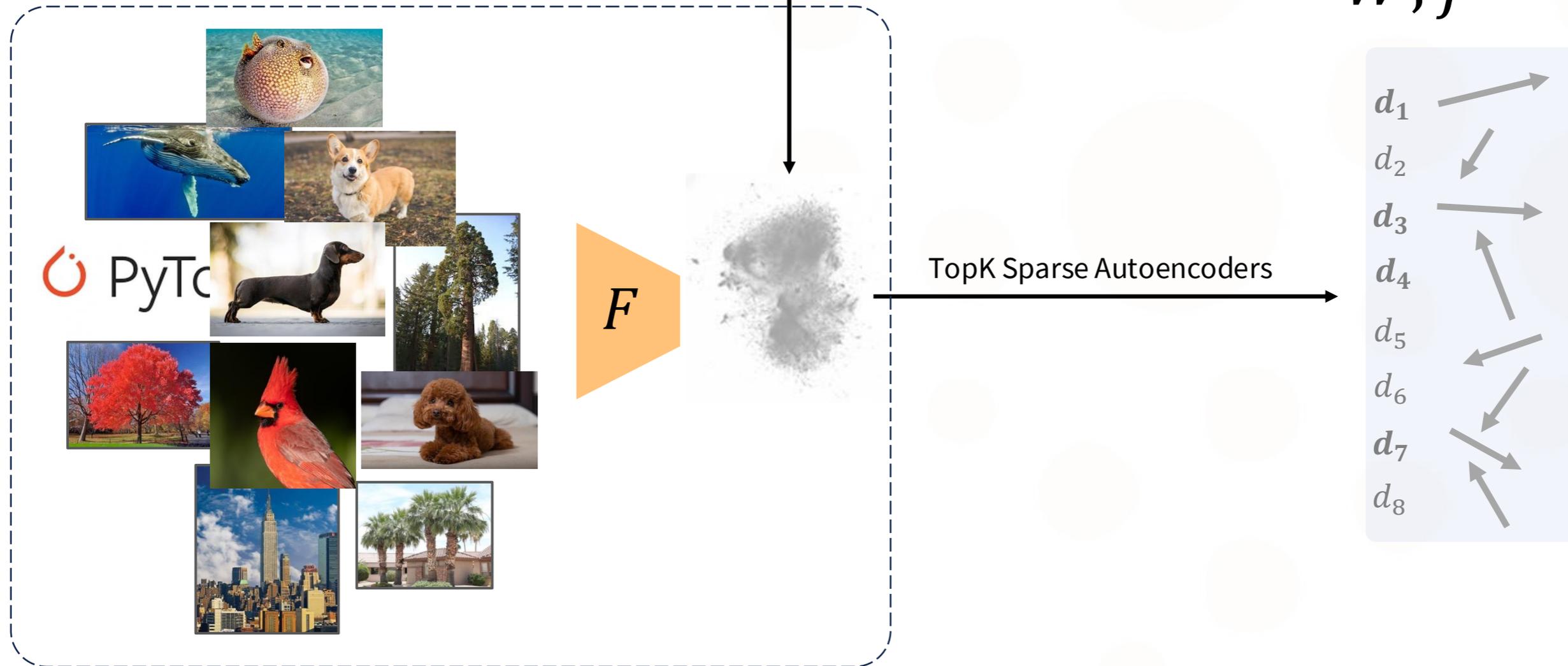
# What we've discussed so far

1. Sparsity parameter $K$ most important
   - Increasing gives better performance, but interpretability more difficult

2. Dictionary size $N$ improves all metrics
   - But memory becomes an issue

3. Standard SAEs suffer from feature splitting
   - We can impose hierarchical / multi-scale structure on concept space

MechInterp is new, many unsolved problems to work on!

Assuming a fixed sparsity $K$, matches up to scaling and rotation
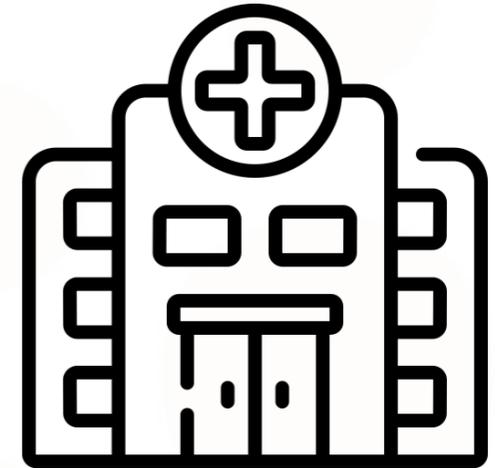
$W^*, f^*$

$W, f^{sparse}$

$F$

TopK Sparse Autoencoders

$d_1$
$d_2$
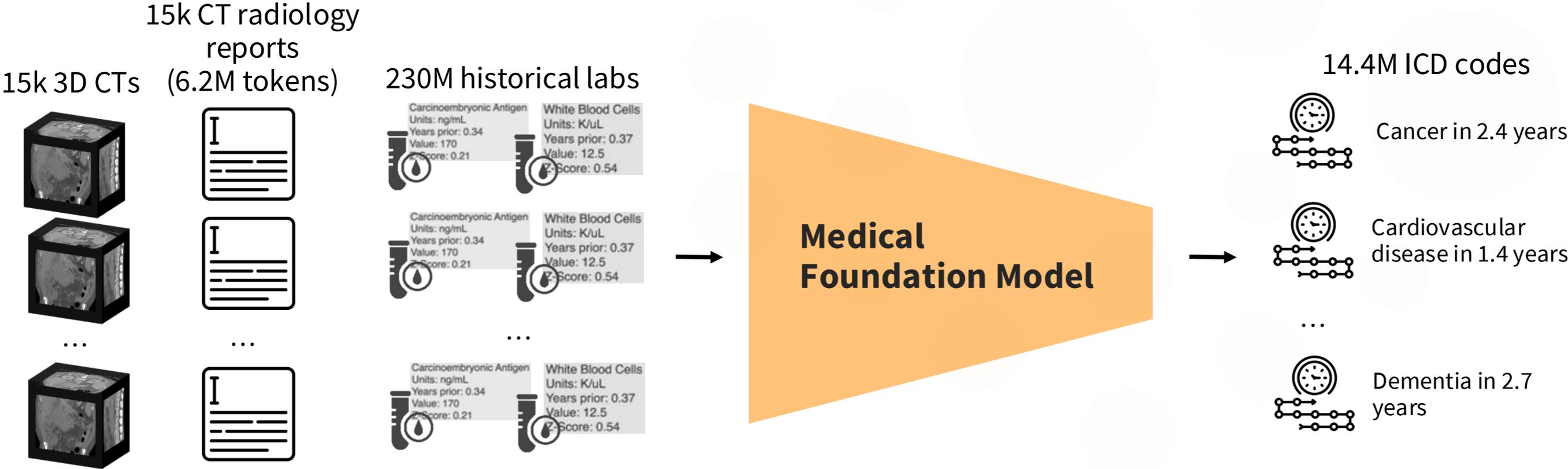$d_3$
$d_4$
$d_5$
$d_6$
$d_7$
$d_8$

# Application

Automatically discovering medical hypotheses

What are the $W^*$ and $f^*$ driving disease?

# Stanford hospital data

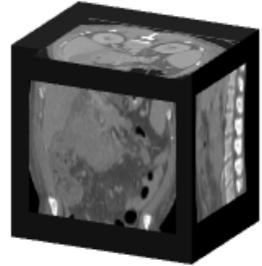## Stanford Medicine holds 1.8 Petabytes of radiology data alone



15k 3D CTs

15k CT radiology reports (6.2M tokens)

230M historical labs

Medical Foundation Model

14.4M ICD codes

Cancer in 2.4 years

Cardiovascular disease in 1.4 years

…

Dementia in 2.7 years

# Stanford hospital data

- Stanford Medicine holds 1.8 Petabytes of radiology data alone



15k 3D CTs

15k CT radiology reports (6.2M tokens)

230M historical labs

14.4M ICD codes

Cancer in 2.4 years

Cardiovascular disease in 1.4 years

...

Dementia in 2.7 years
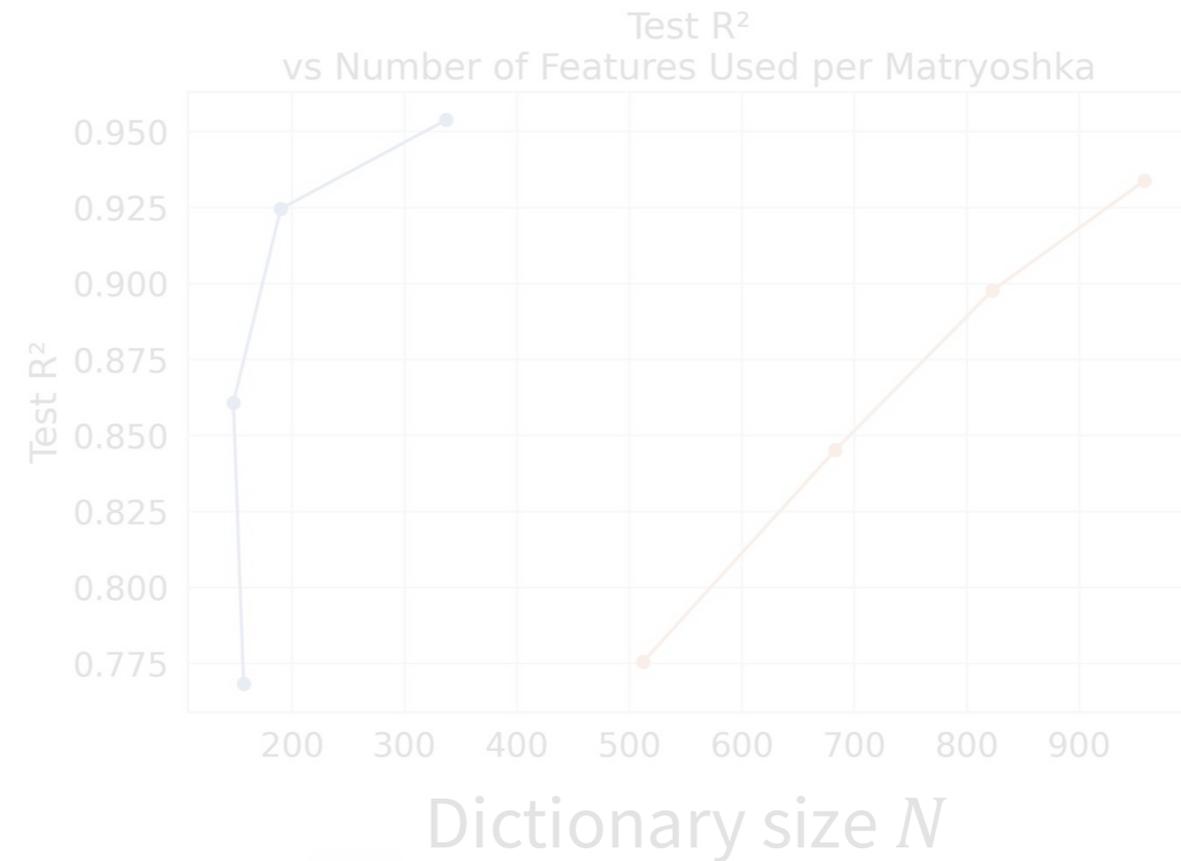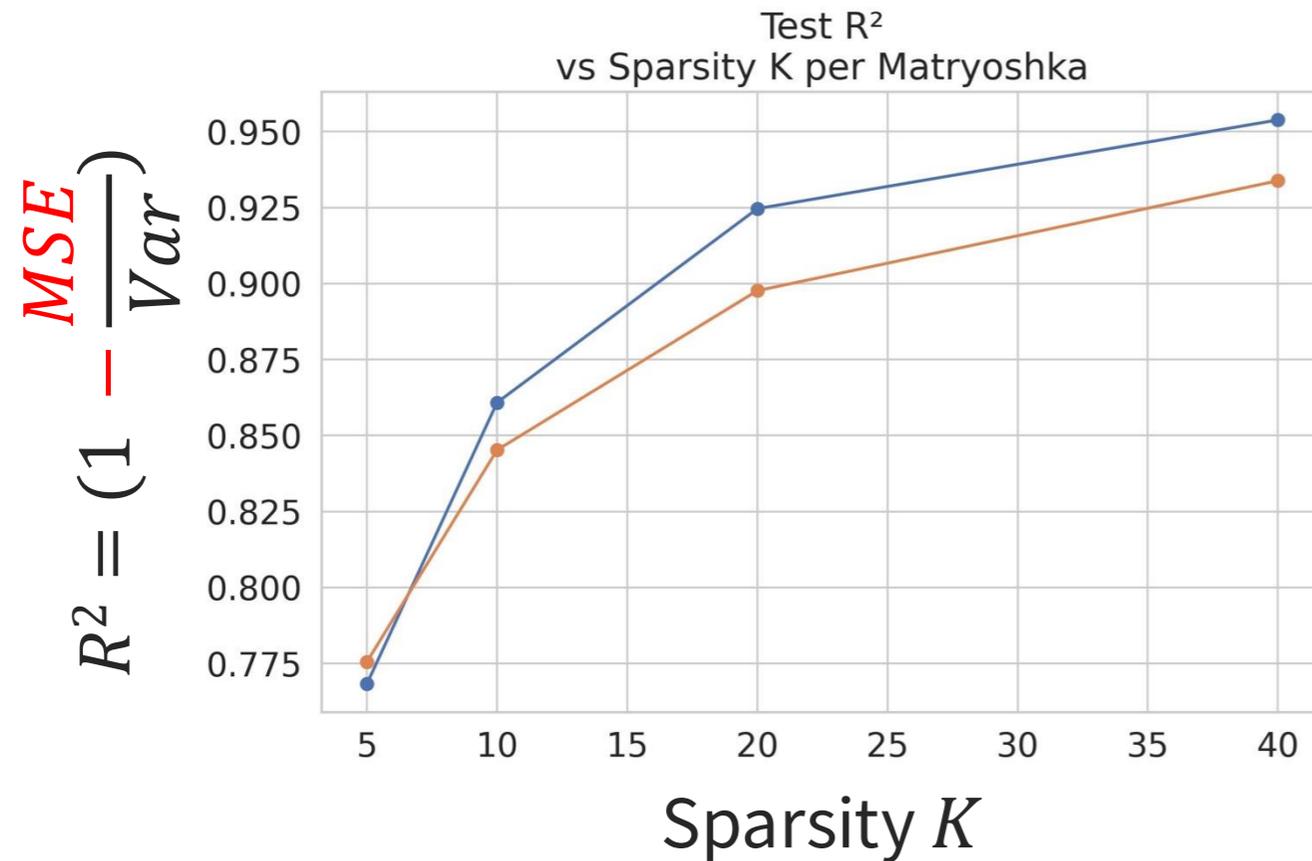
# Stanford hospital data



Blankemeier, Louis and Kumar, Ashwin et al. Merlin: A vision language foundation model for 3d computed tomography

Zhang, Yanzhao et al. Qwen3 Embedding: Advancing Text Embedding and Reranking Through Foundation Models
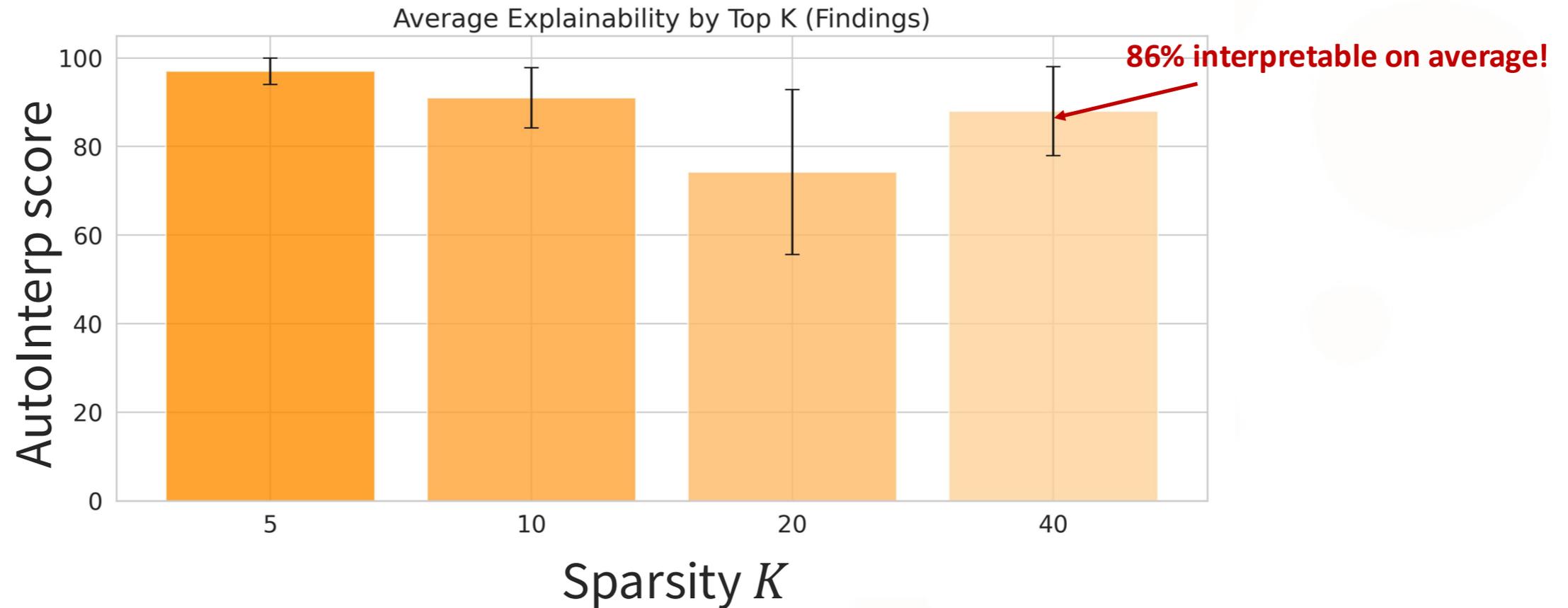
# Running Top-K SAEs

$M = \{128, 512, 2048, 8192\}$

- CT images via Merlin → $N = 8192$ and $K: [5, 10, 20, 40]$ and *Matryoshka* vs *Standard*



$$R^2 = (1 - \frac{MSE}{Var})$$

# Running Top-K SAEs

- How interpretable radiology report concepts in Matryoshka SAEs?



Average Explainability by Top K (Findings)

**86% interpretable on average!**

# $W_{85}$ "*Cirrhotic liver morphology and varices*"

Top 90<sup>th</sup> percentile of $f_{85}^{sparse}$

49 year old male who develops abd. cancer in 1 year, 9 months, 25 days time

cirrhosis with stigmata of portal hypertension… prominent varices including a large recannulized periumbilical vein and large anterior abdominal wall varices

52 year old female who develops abd. cancer in 1 year, 10 months, 23 days time

dysmorphic cirrhotic liver. no suspicious arterial enhancing lesions.patent tips.no biliary dilatation. … vasculature : perigastric varices again noted
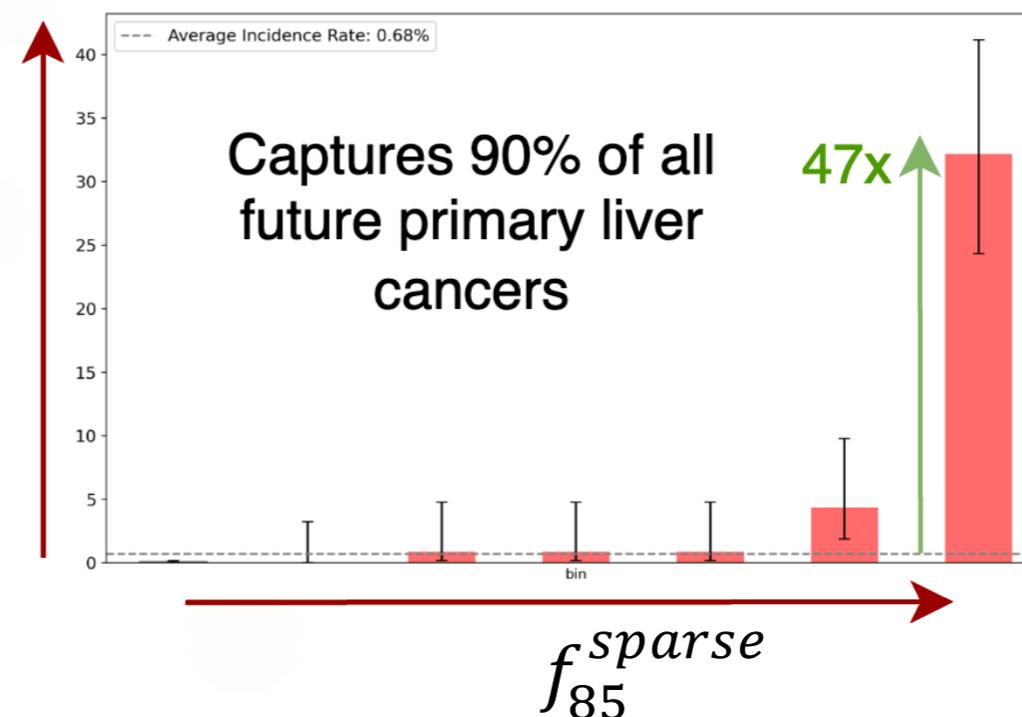
# $W_{85}$ associated with primary liver cancer

- P(developing liver cancer in 4 years | CT was ordered) = 0.68%

- P(developing liver cancer in 4 years | CT was ordered **and** $W_{85} > 0.9$) = 32.7%

"**Cirrhotic liver morphology and varices**"

| Dataset | Stanford (discovery) |
|---|---|
| Rank | #1 / 3315 |
| Odds Ratio | 2.13 |
| Count | 690 (9.3%) |

P(developing liver cancer in under 4 years)



Captures 90% of all future primary liver cancers

47x

$f_{85}^{sparse}$

# $W_{33}$ associated with first heart failure

- P(first recorded heart failure within 4 years | CT was ordered) = 7.6%

- P(first recorded heart failure within 4 years | CT was ordered **and** $W_{33} > 0.9$) = 18.5%