

# Interactive Scalable Interfaces for Machine Learning Interpretability



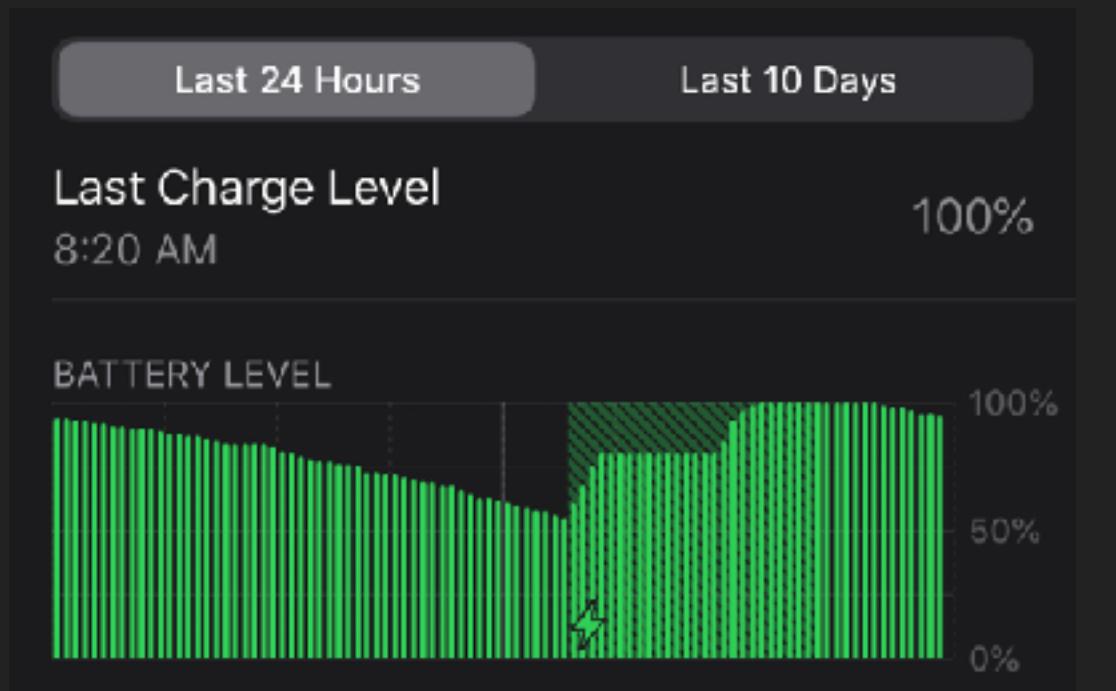
**Fred Hohman**  
@fredhohman



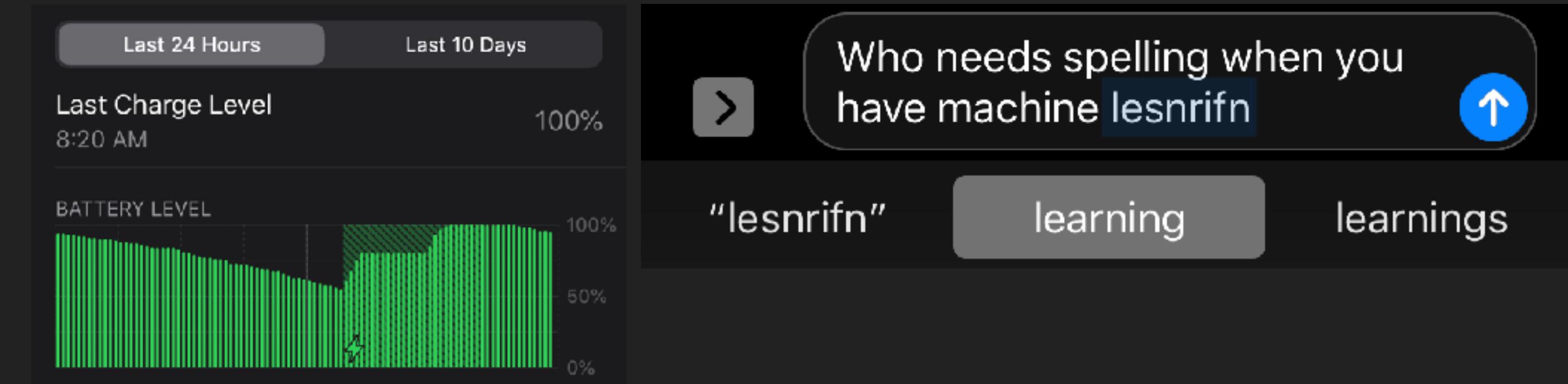
Follow along at [fredhohman.com/talk](http://fredhohman.com/talk)



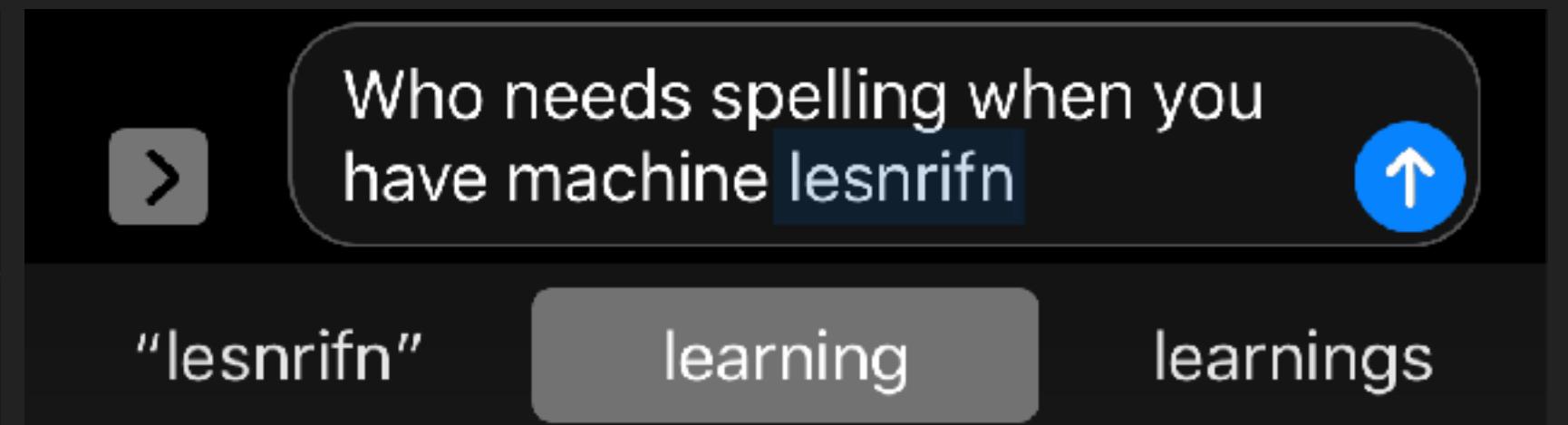
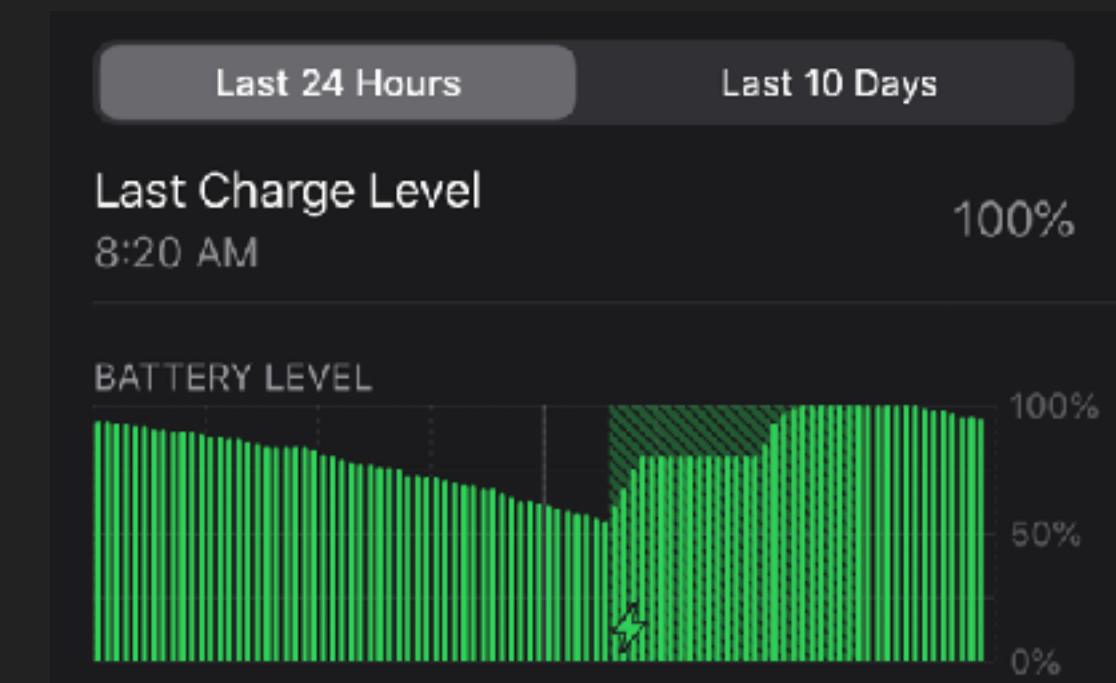
## Optimized phone battery charging



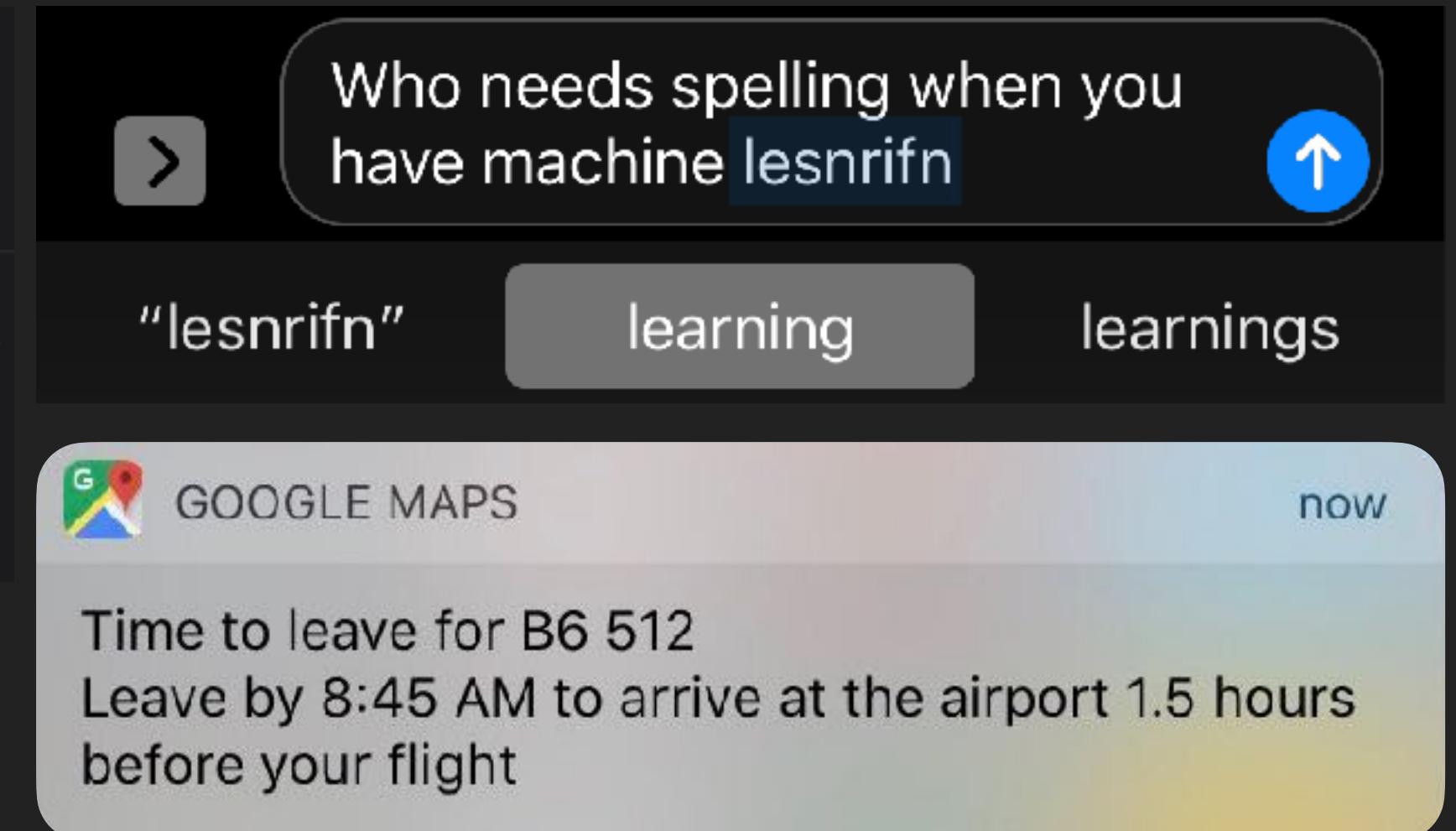
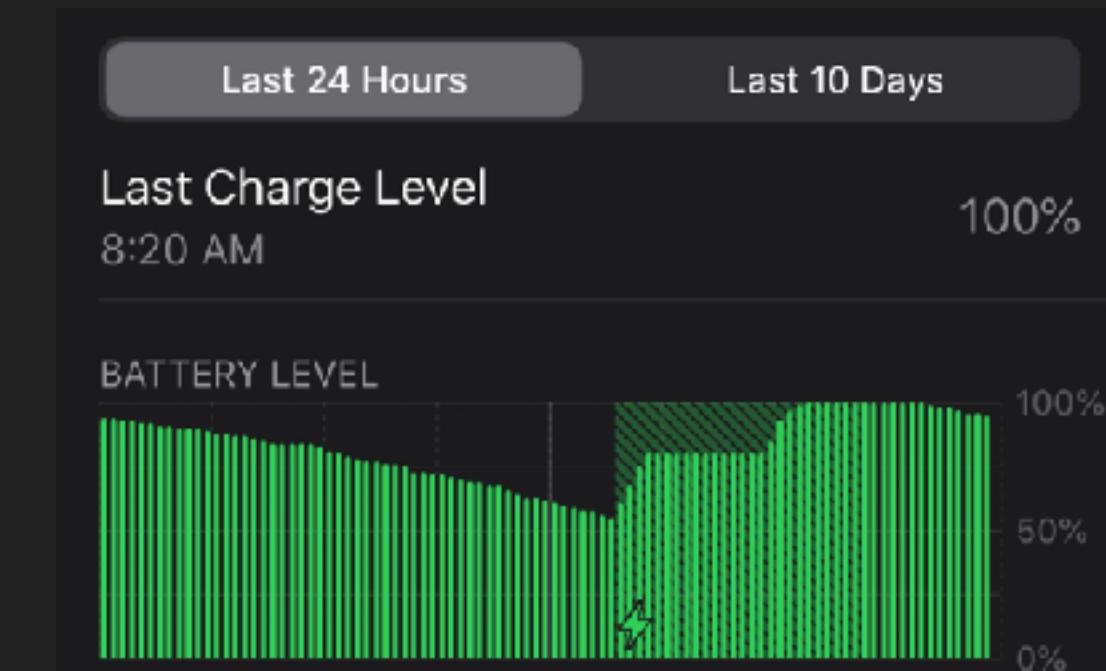
## Optimized phone battery charging Spelling & text prediction



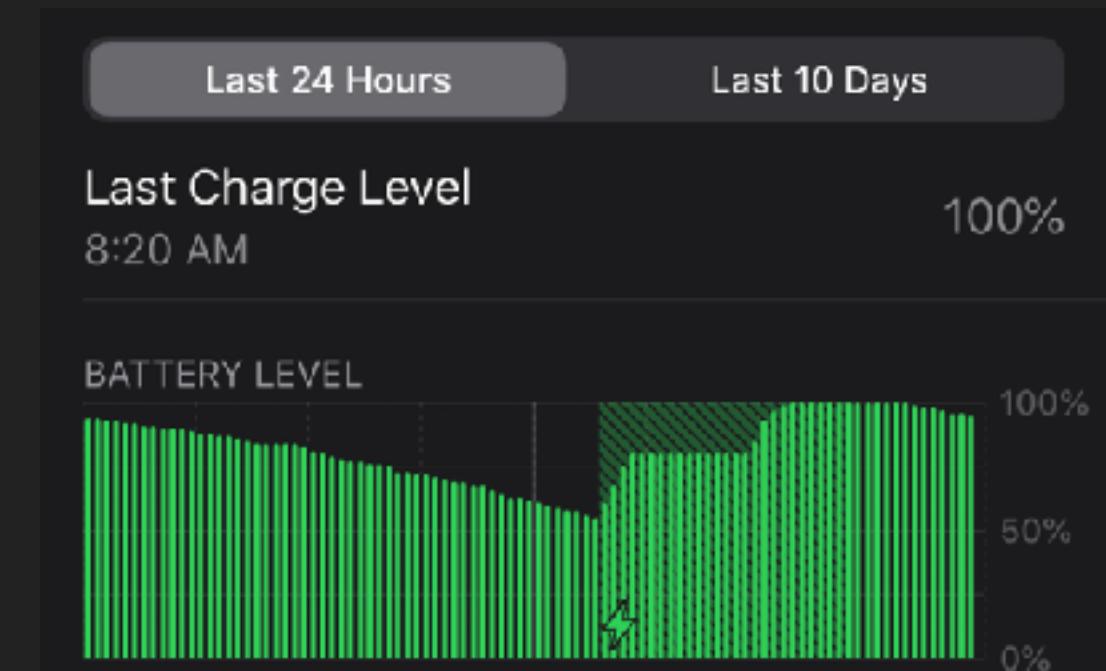
Optimized phone battery charging  
Spelling & text prediction  
Face unlock



Optimized phone battery charging  
Spelling & text prediction  
Face unlock  
Traffic & commute notifications



Optimized phone battery charging  
Spelling & text prediction  
Face unlock  
Traffic & commute notifications  
Email spam filters  
Email "smart" reply  
Search results



Who needs spelling when you have machine lesnrifn ↑

"lesnrifn" learning learnings

GOOGLE MAPS now

Time to leave for B6 512  
Leave by 8:45 AM to arrive at the airport 1.5 hours before your flight

◆ Support Team Nov 14

◆ ?? Fredhohman ?? we need your con... Please confirm your Unsubscribe To co... star

machine learning MICROPHONE



Move your head slowly to complete the circle.

Recipients

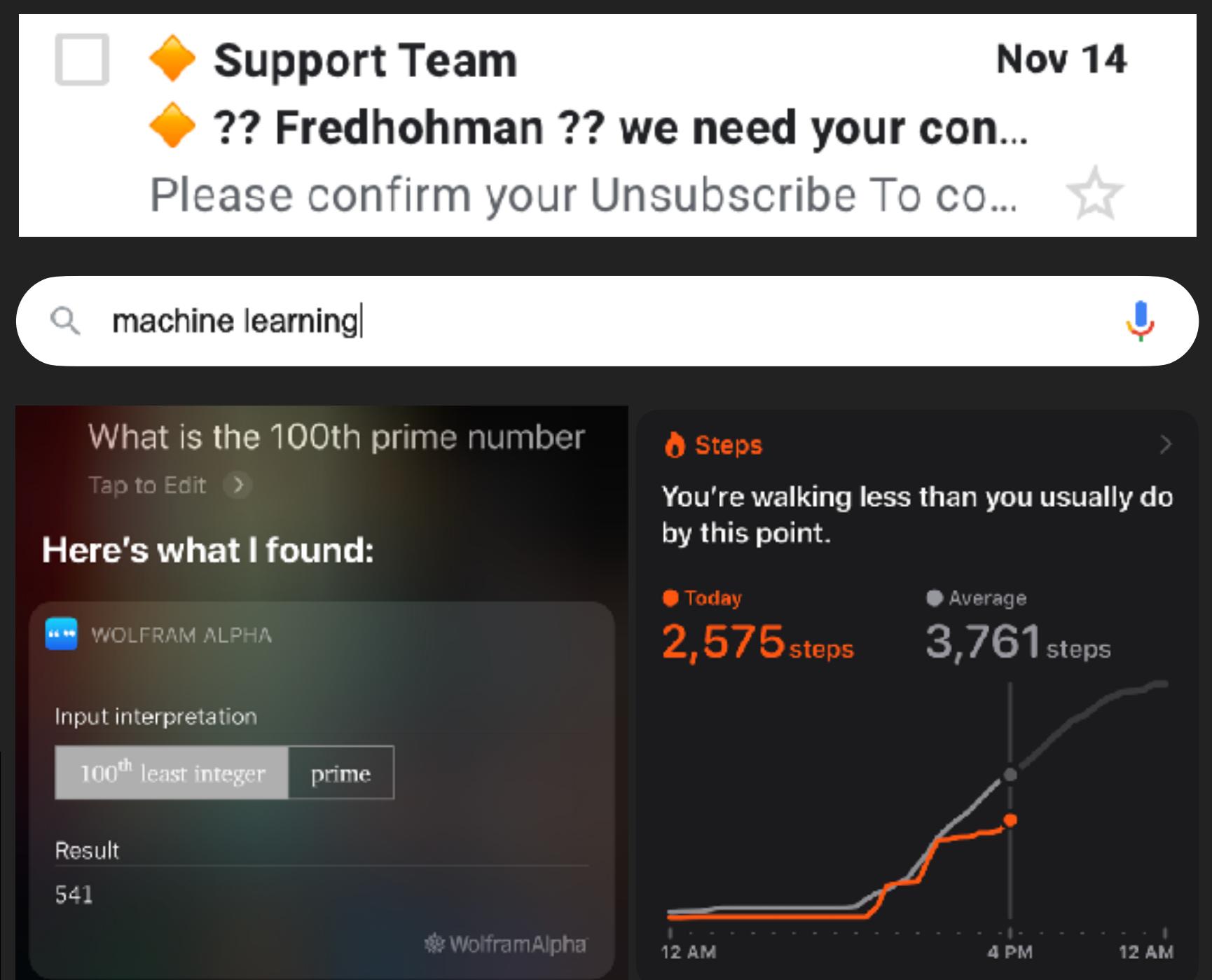
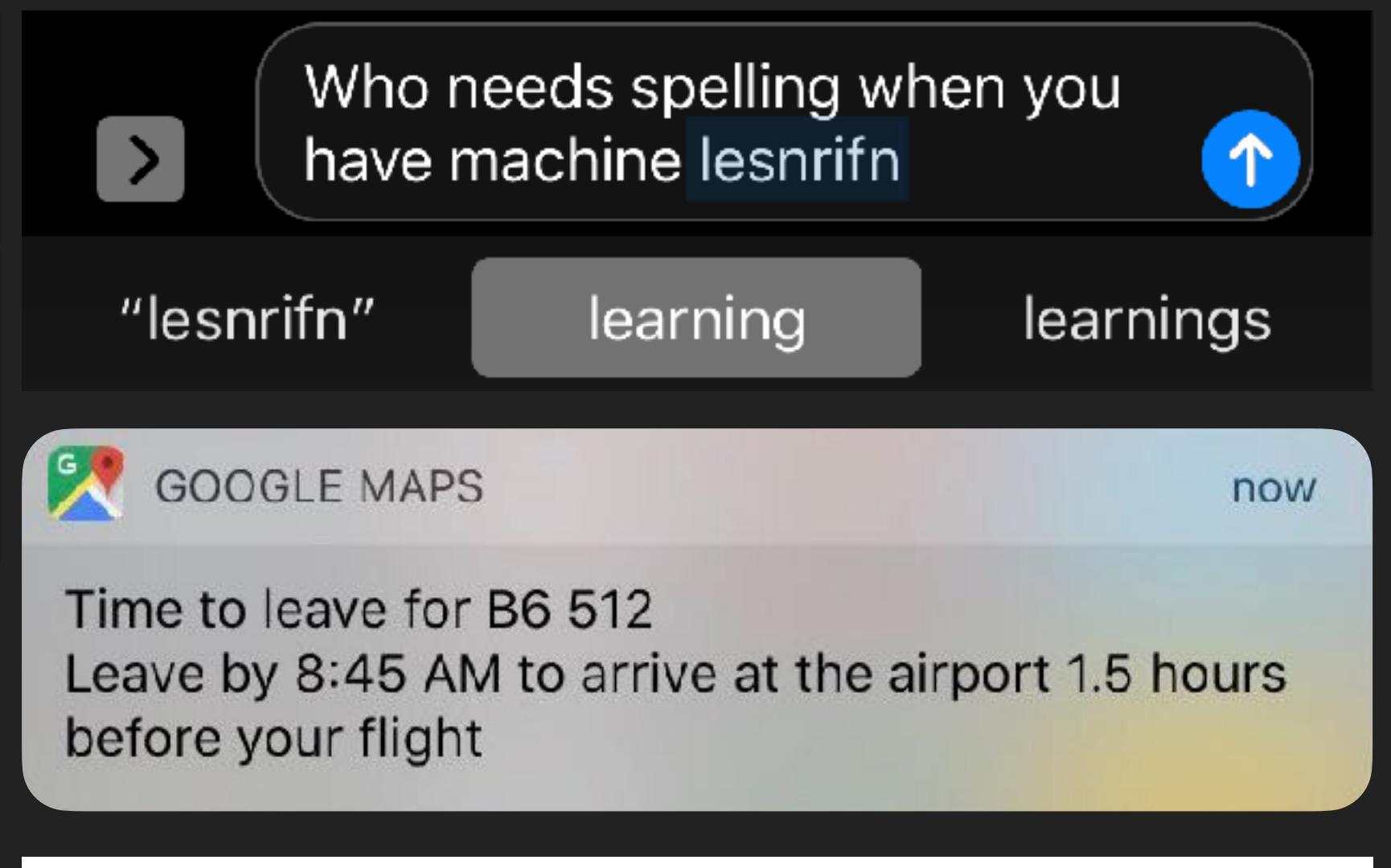
---

Subject

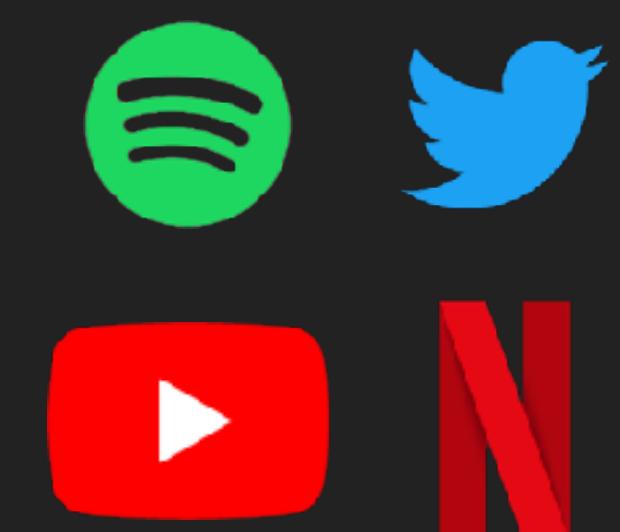
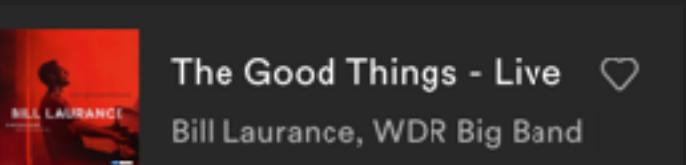
---

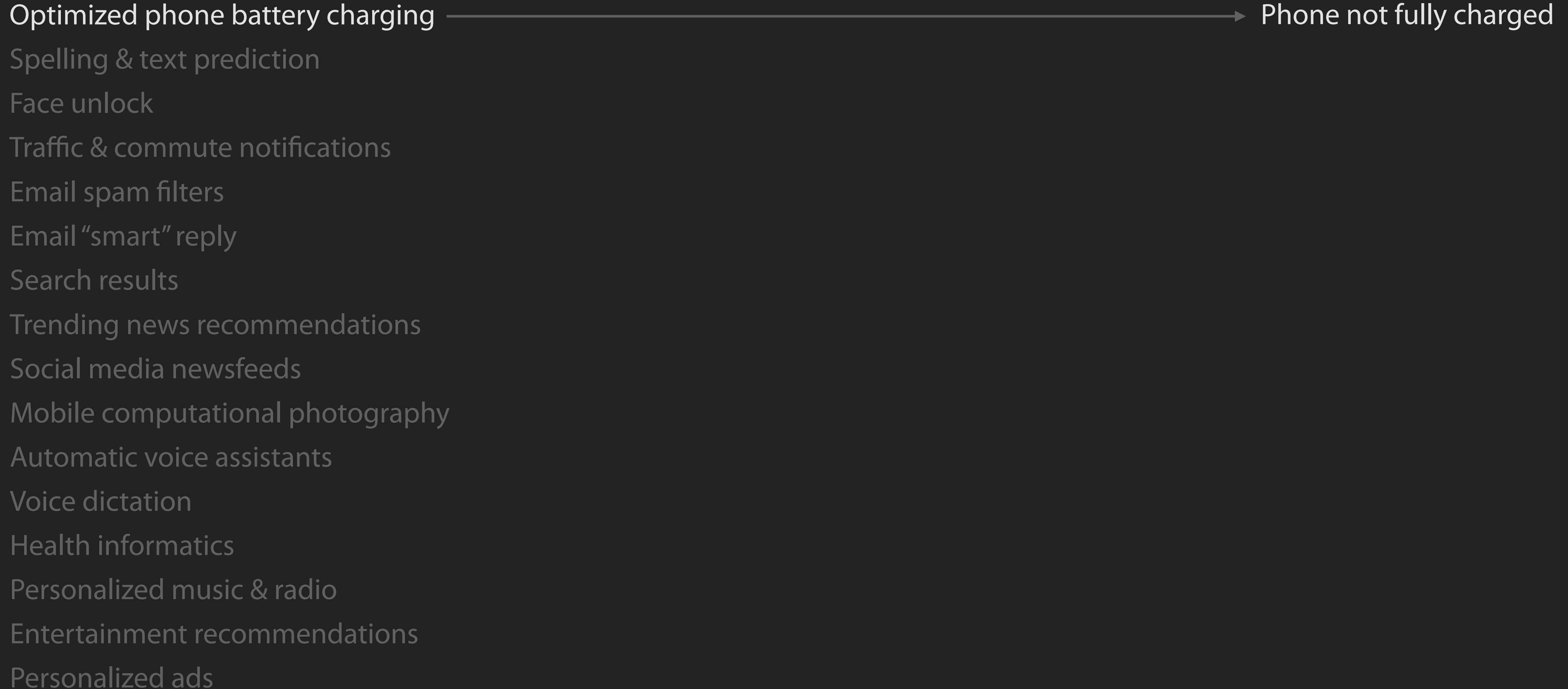
Hope you are well

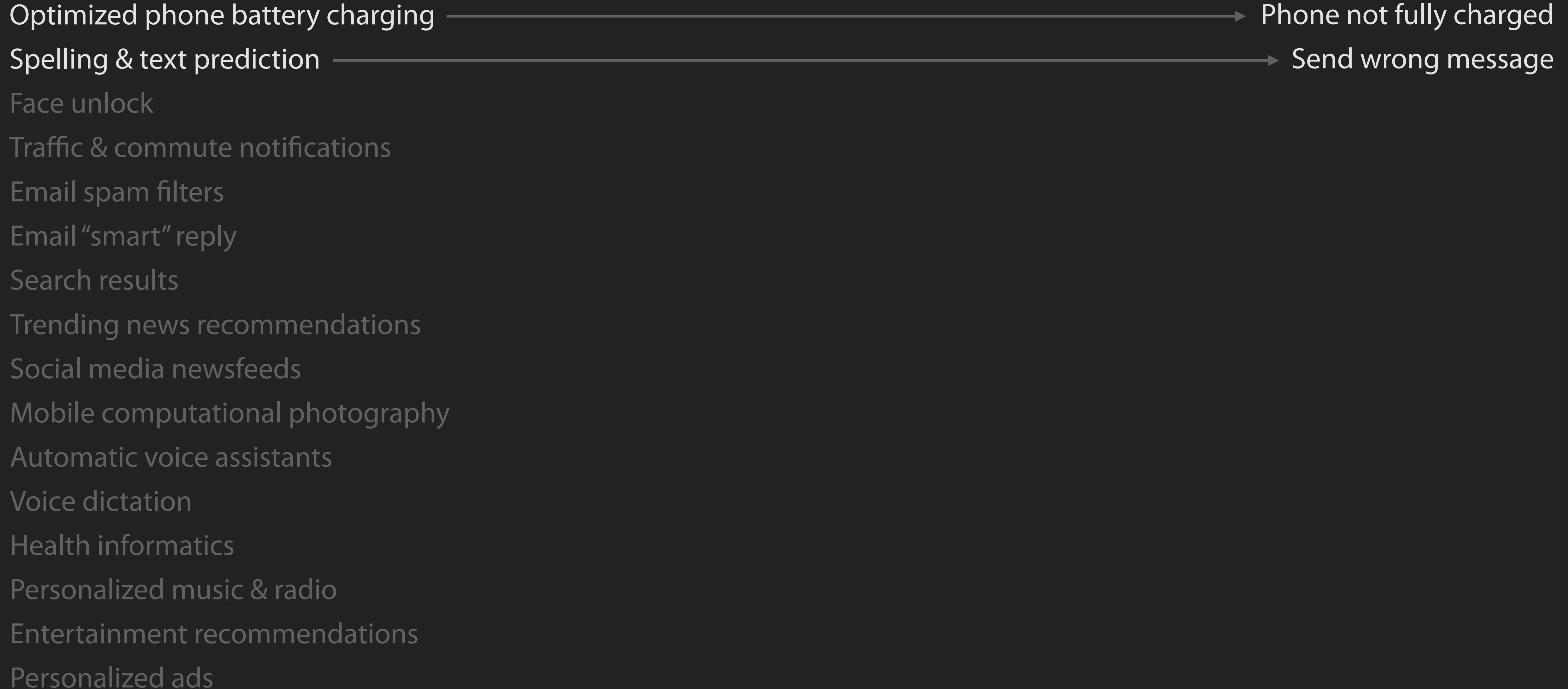
Optimized phone battery charging  
Spelling & text prediction  
Face unlock  
Traffic & commute notifications  
Email spam filters  
Email “smart” reply  
Search results  
Trending news recommendations  
Social media newsfeeds  
Mobile computational photography  
Automatic voice assistants  
Voice dictation  
Health informatics  
Personalized music & radio  
Entertainment recommendations  
Personalized ads

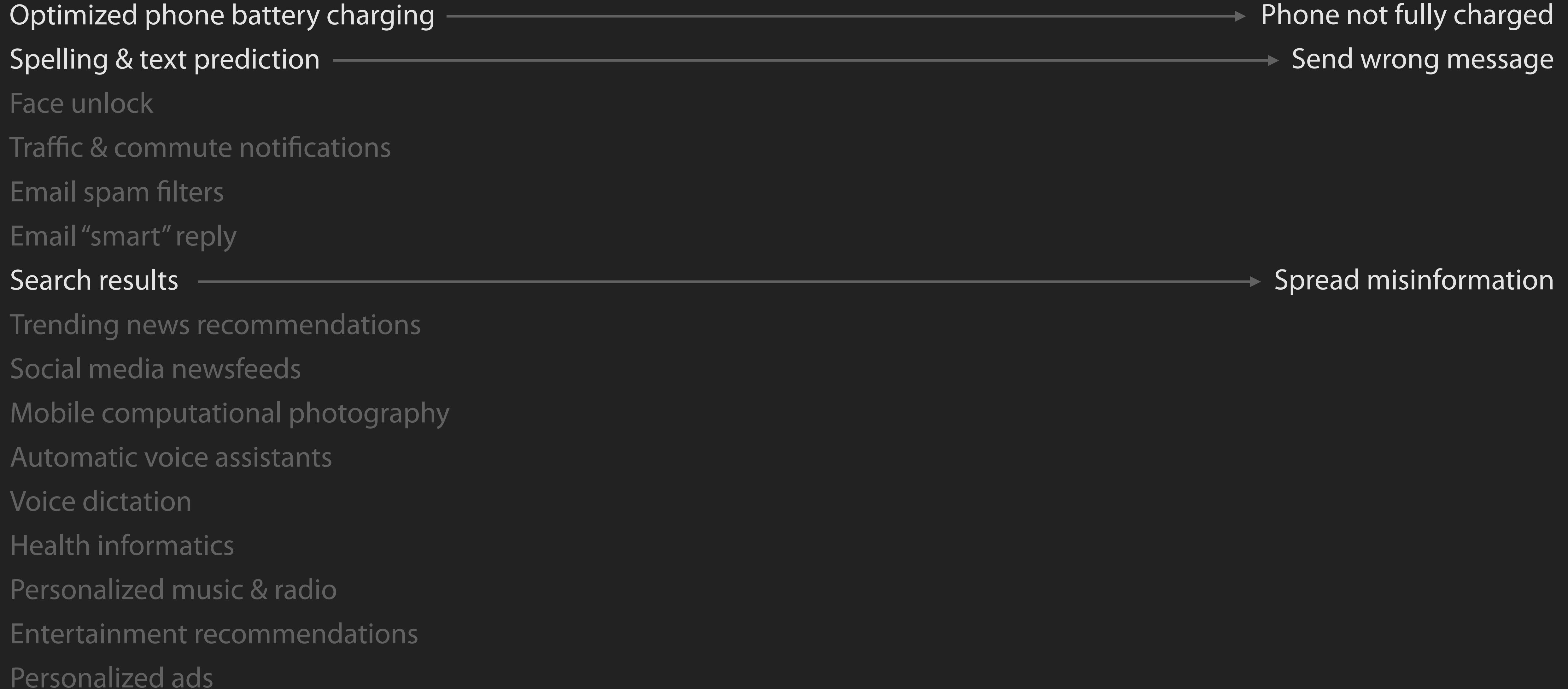


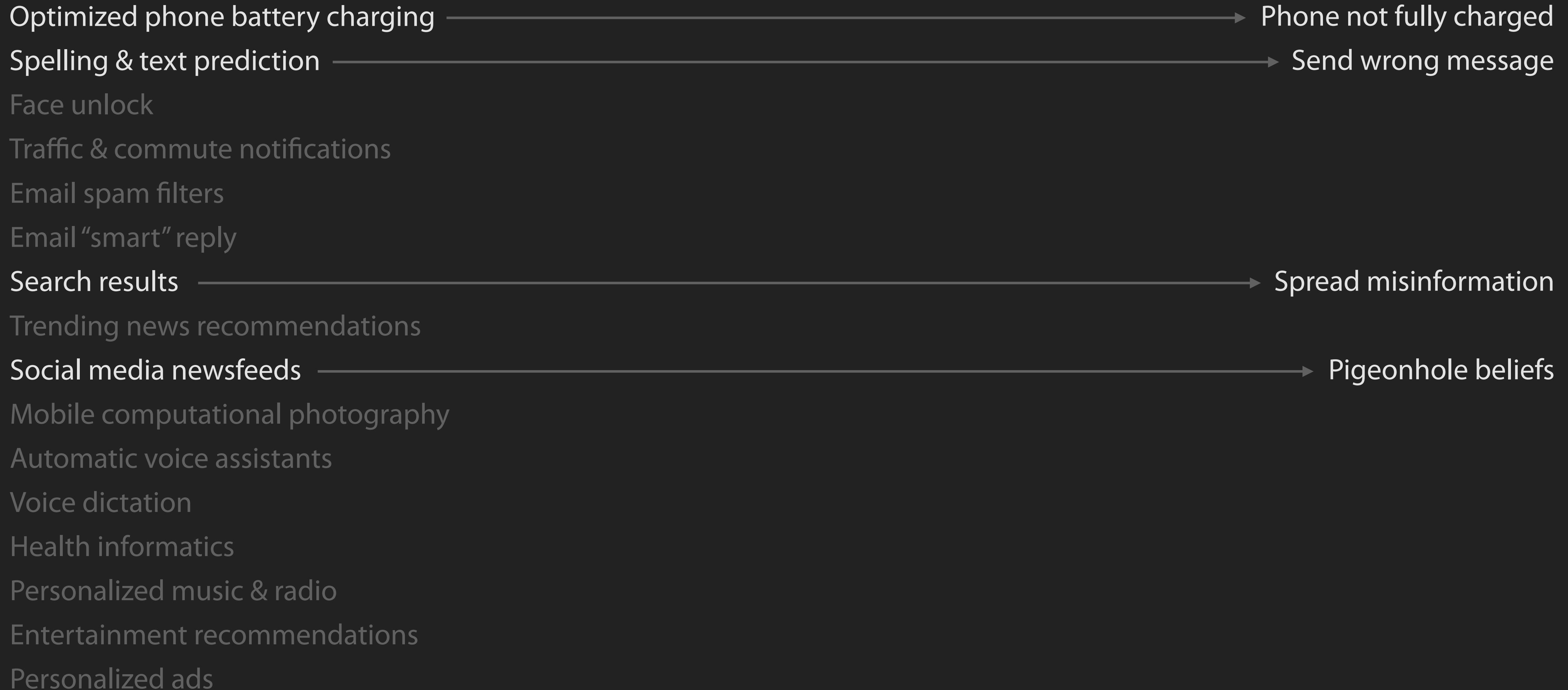
Recipients  
Subject  
Hope you are well

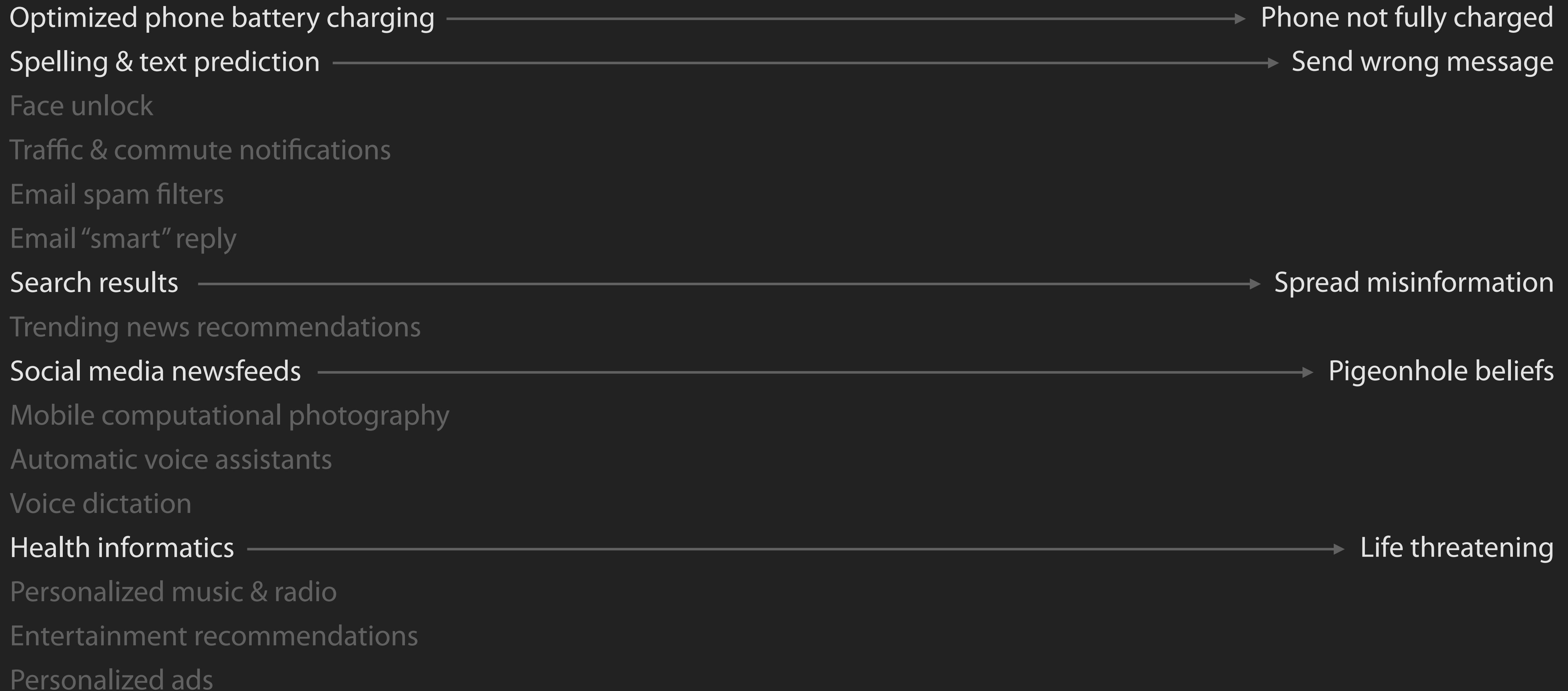


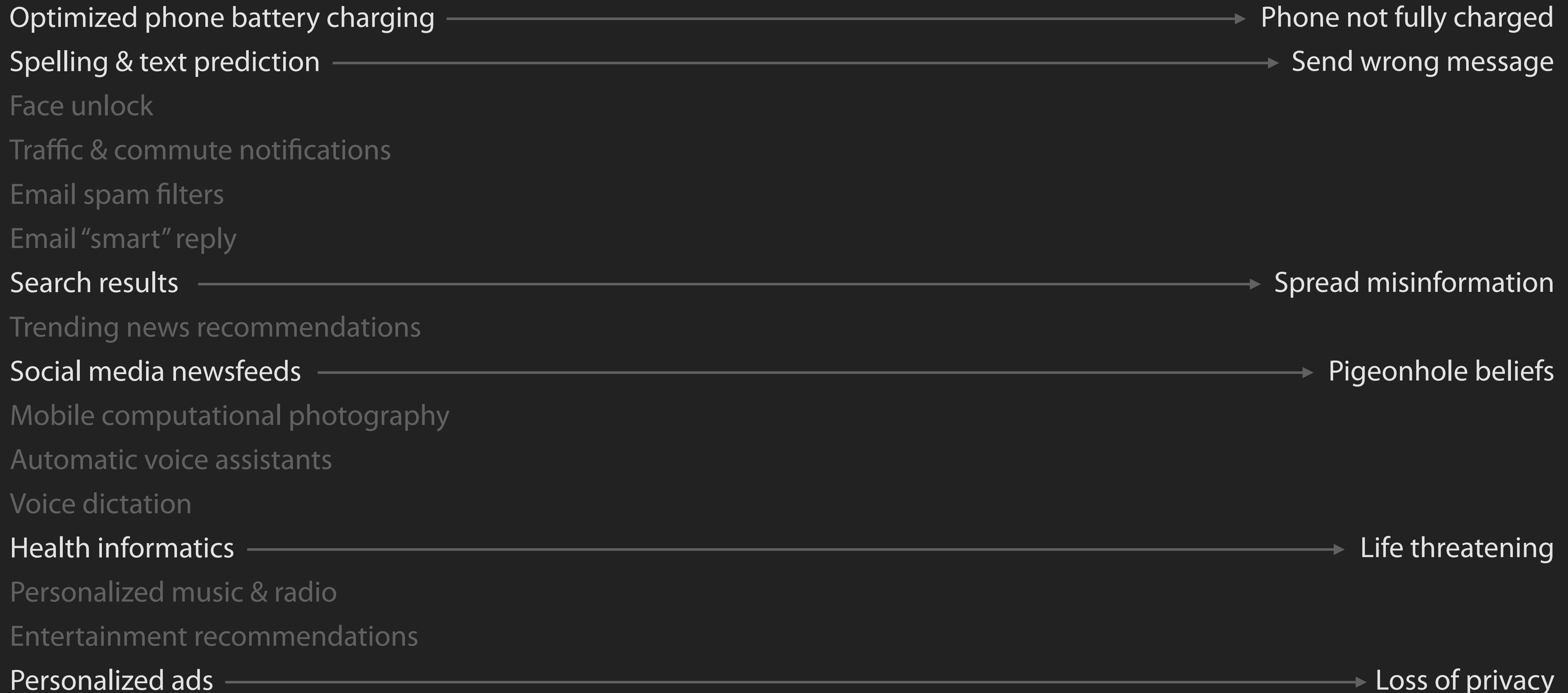








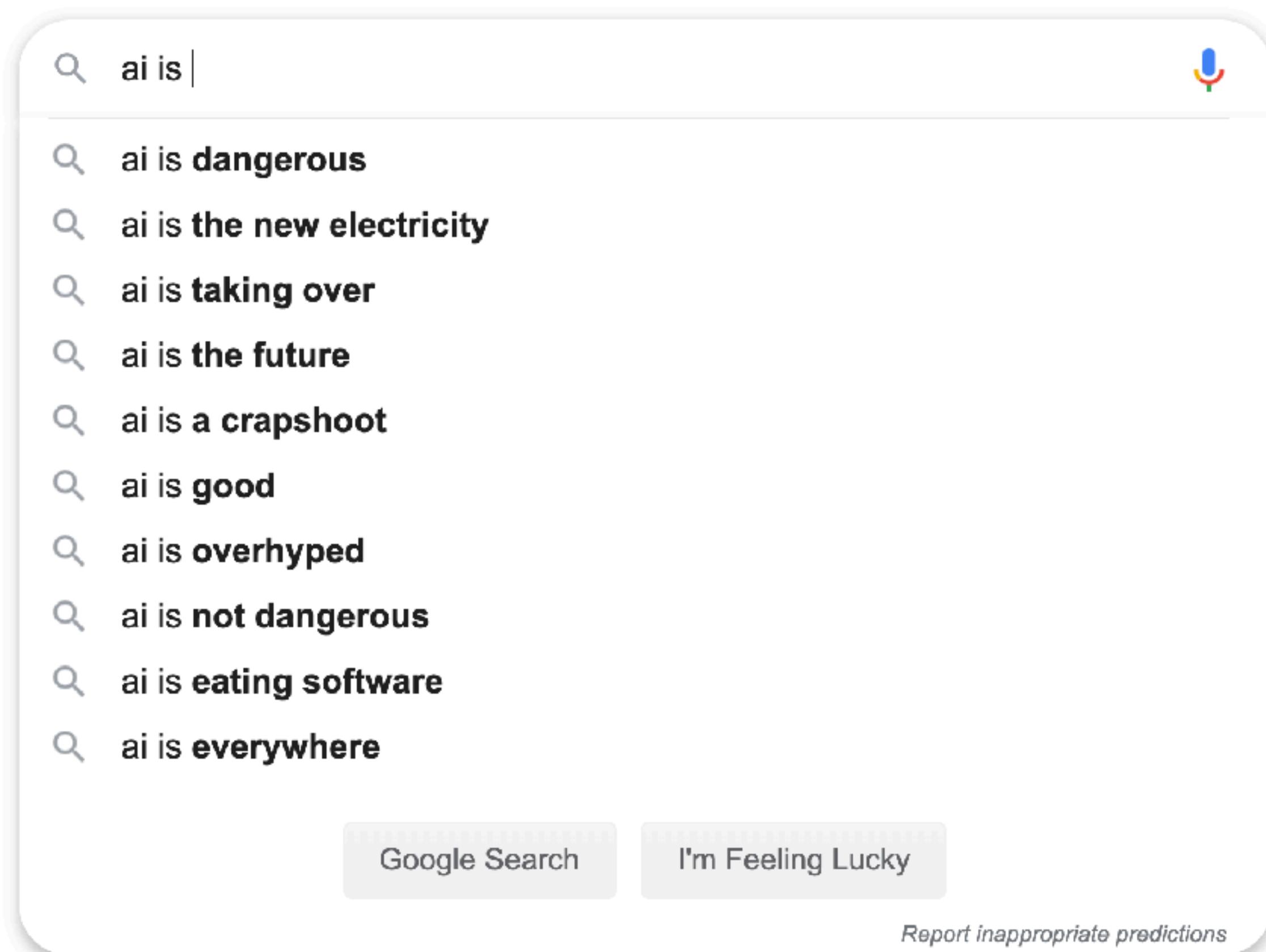




# Machine learning makes mistakes



# Google



A screenshot of a Google search interface showing search suggestions for the query "ai is". The suggestions are listed below the search bar, each preceded by a magnifying glass icon. The suggestions are:

- ai is dangerous
- ai is the new electricity
- ai is taking over
- ai is the future
- ai is a crapshoot
- ai is good
- ai is overhyped
- ai is not dangerous
- ai is eating software
- ai is everywhere

At the bottom of the suggestions box are two buttons: "Google Search" and "I'm Feeling Lucky". Below the suggestions box is a small text link: "Report inappropriate predictions".

Data



Machine Learning  
Model

# Interpretability

*What has a model learned?*

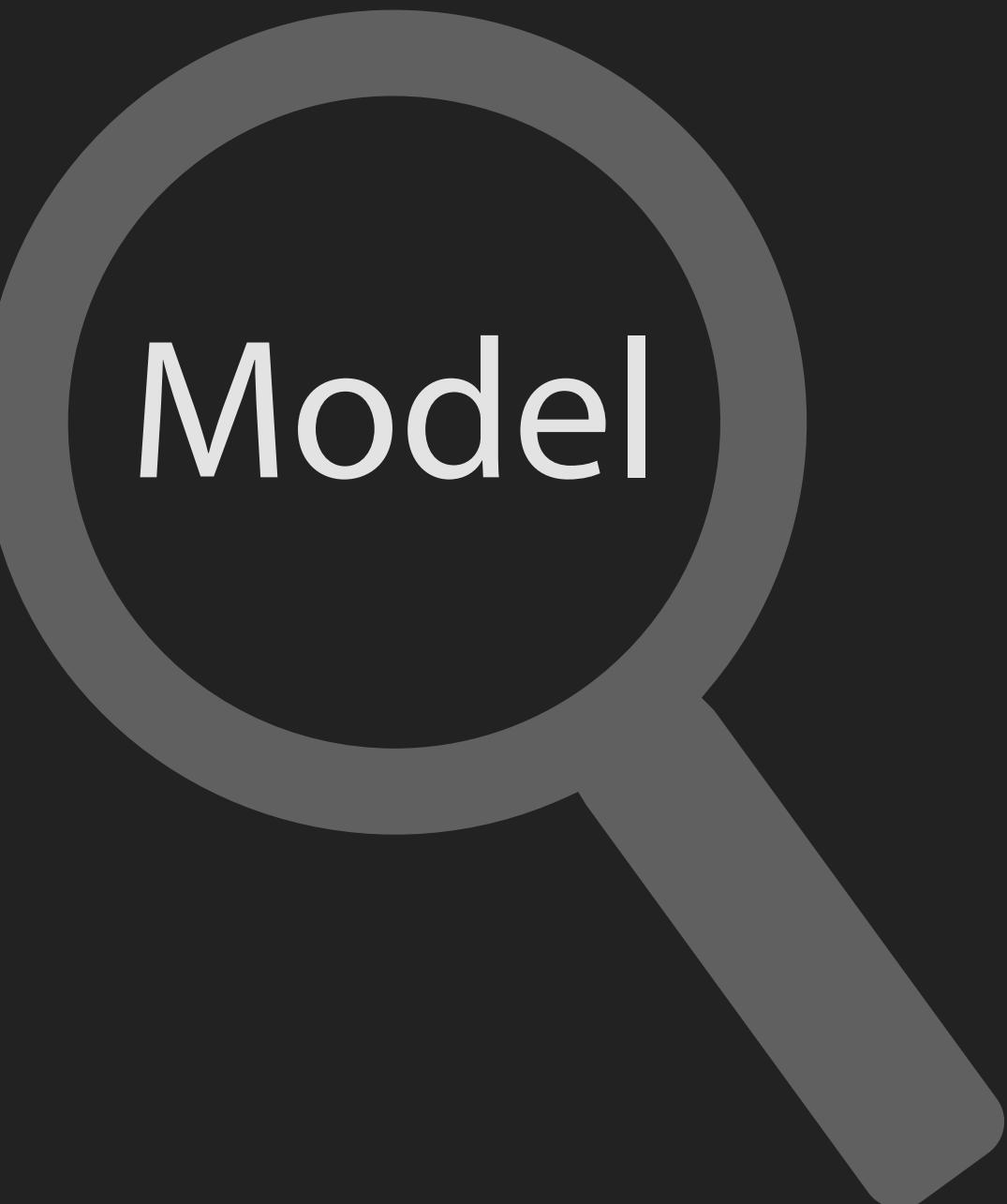
*Can you explain its predictions?*

*Do you understand its behavior?*

Data



Model



# Why Interpretability?

USA TODAY NEWS SPORTS ENTERTAINMENT LIFE MONEY TECH TRAVEL OPINION 72° MORE ▾

REUTERS

Business Markets World Politics TV More

Guardian | Search jobs | Sign in | Search | US edition -

MOTHERBOARD

ion

Media S

Pr

ris

Mac

race

By Lynne Peskoe-Yang

Published 10:00 AM ET, 11/12/2018 | Share | Print

## Algorithms Have Nearly Mastered Human Language. Why Can't They Stop Being Sexist?

To fight gender bias, researchers are training language-processing algorithms to envision a world where it doesn't exist.

# Why Interpretability?

## Fairness & GDPR

REUTERS Business Markets World Politics TV More

Guardian | Search jobs | Sign in | Search | US edition -

MOTHERBOARD

## Algorithms Have Nearly Mastered Human Language. Why Can't They Stop Being Sexist?

To fight gender bias, researchers are training language-processing algorithms to envision a world where it doesn't exist.

By Lynne Peskoe-Yang | Published 10:00 AM ET, Mon, Mar 12, 2018 | Updated 10:00 AM ET, Mon, Mar 12, 2018

# Why Interpretability?

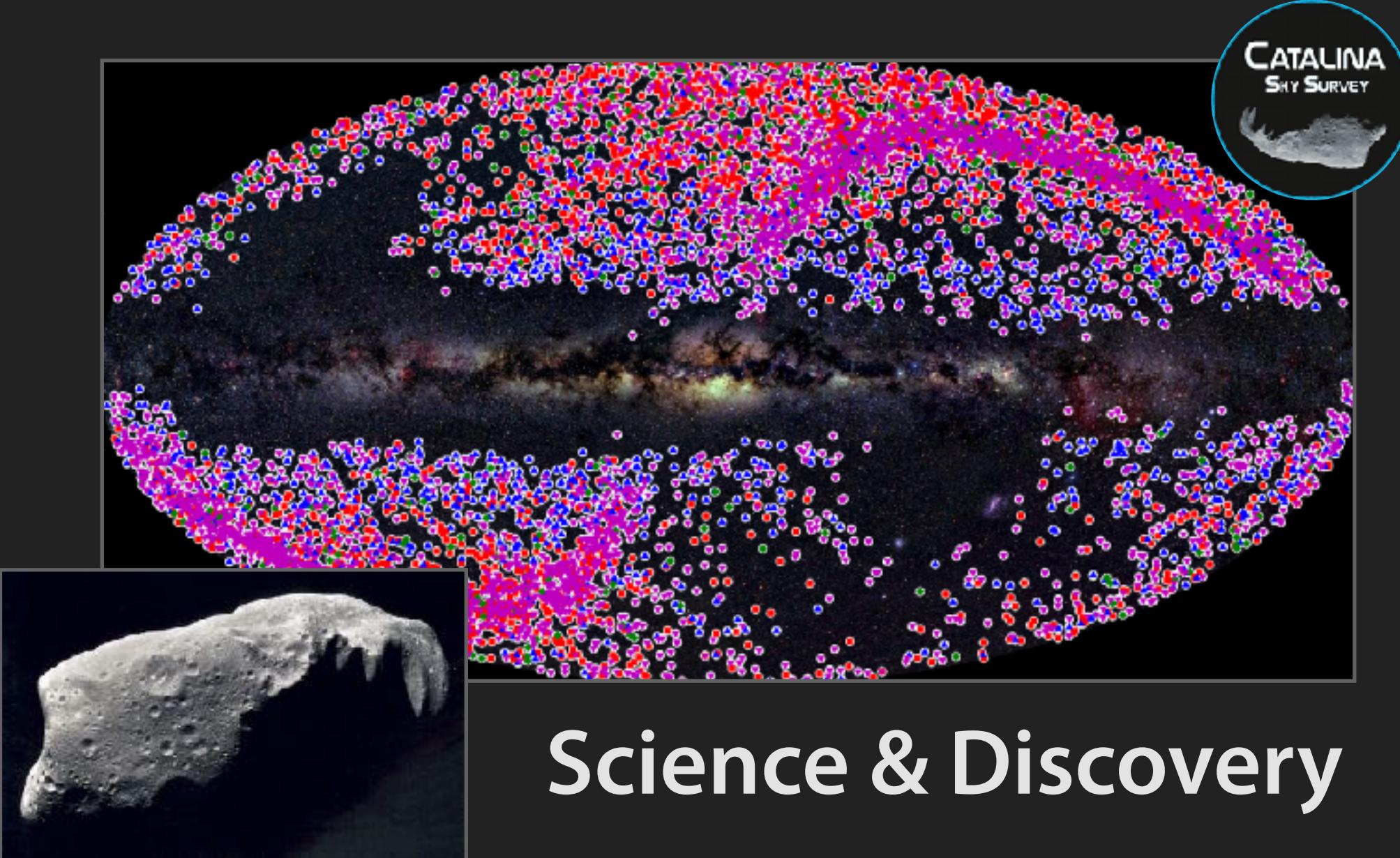
## Fairness & GDPR



## Science & Discovery

A screenshot of a news website featuring a prominent headline: "Algorithms Have Nearly Mastered Human Language. Why Can't They Stop Being Sexist?". Below the headline is a subtext: "To fight gender bias, researchers are training language-processing algorithms to envision a world where it doesn't exist." At the bottom of the article, there is a byline: "By Lynne Peskoe-Yang". The website's navigation bar includes links for Business, Markets, World, Politics, TV, and More.

## Fairness & GDPR



## Science & Discovery

# Why Interpretability?

## Robustness

To data drift & adversarial attacks

## Facilitating User Control

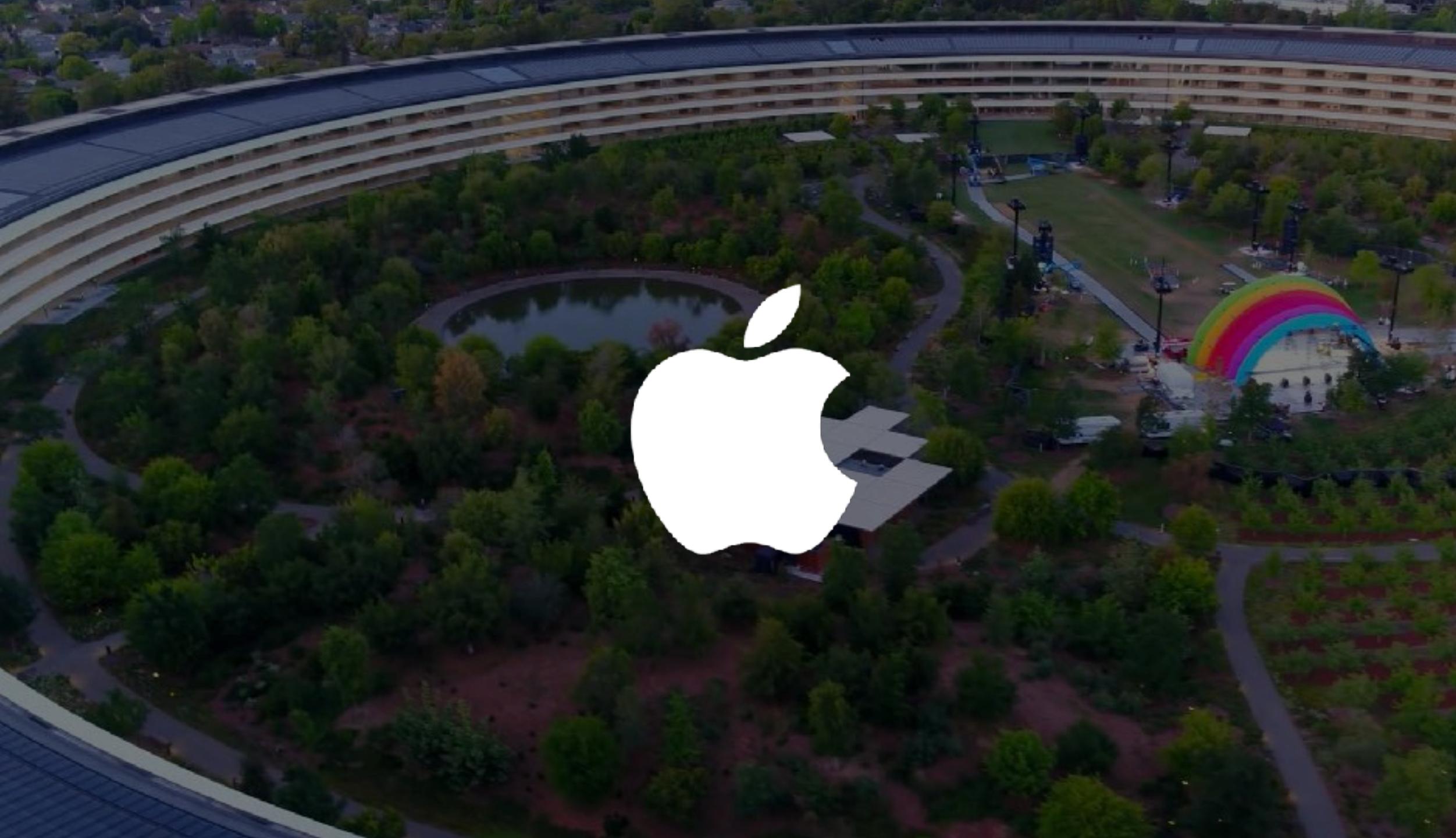
With personalized agents

## User Acceptance

People happier with explanations

## Mismatched Objectives

Is AI right for the right reasons?





Machine learning is a **people** problem.



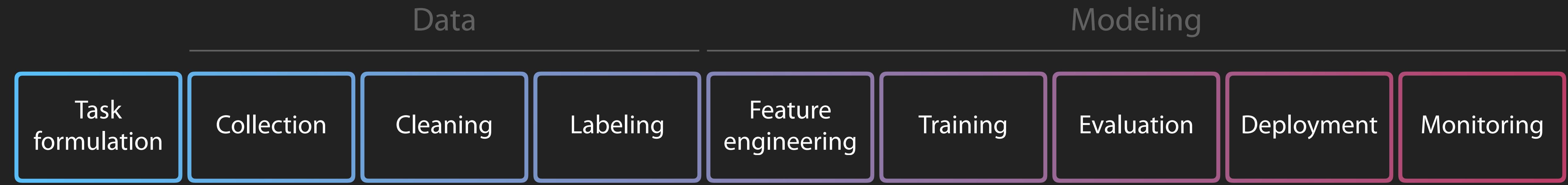
**Jet Propulsion Laboratory**  
California Institute of Technology



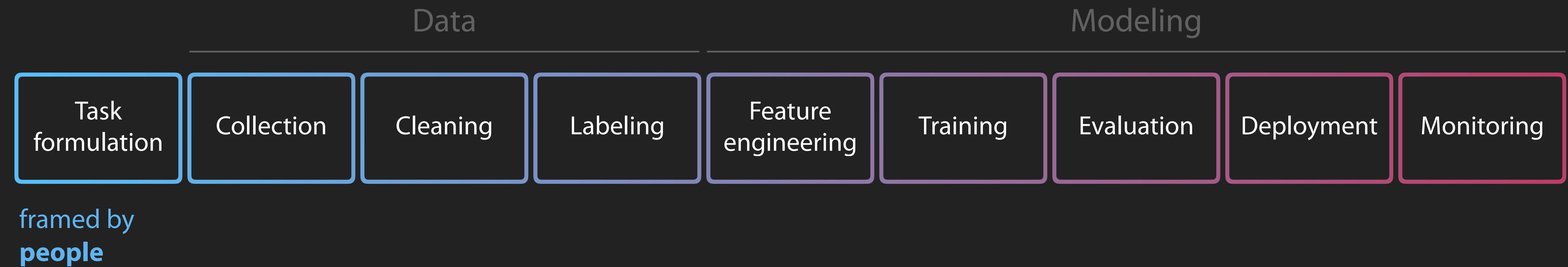
Microsoft Research

**Pacific Northwest**  
NATIONAL LABORATORY

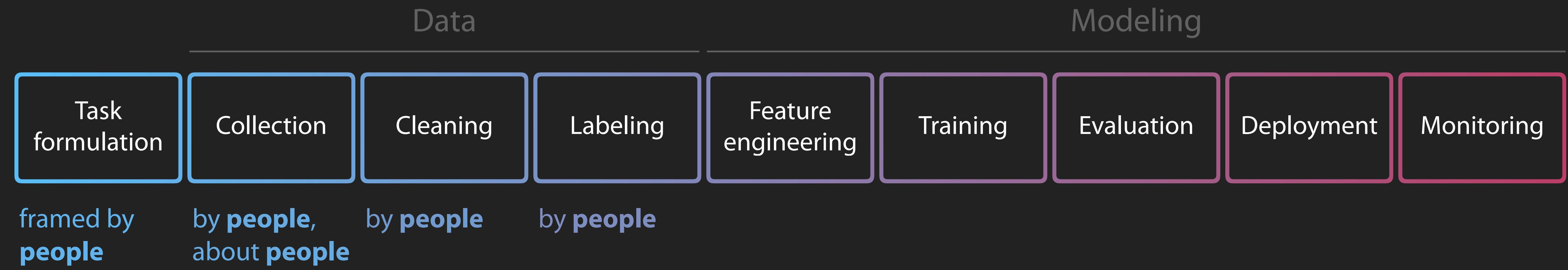
# Machine learning is a **people problem**.



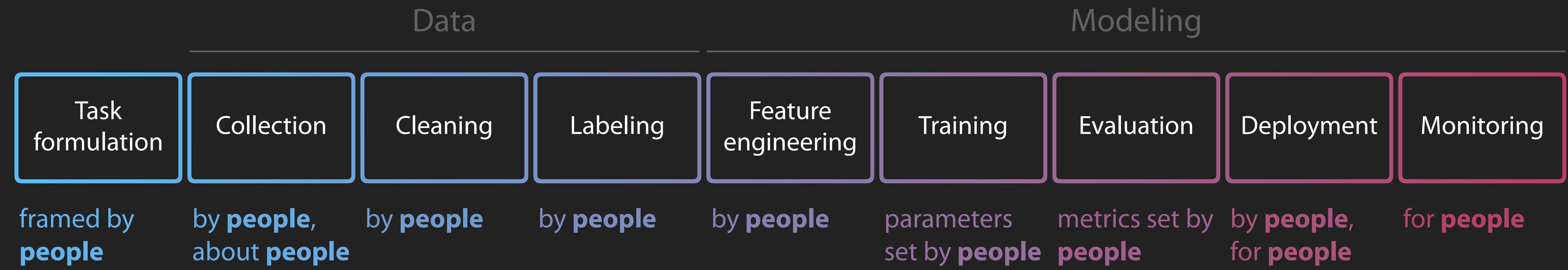
# Machine learning is a **people problem**.



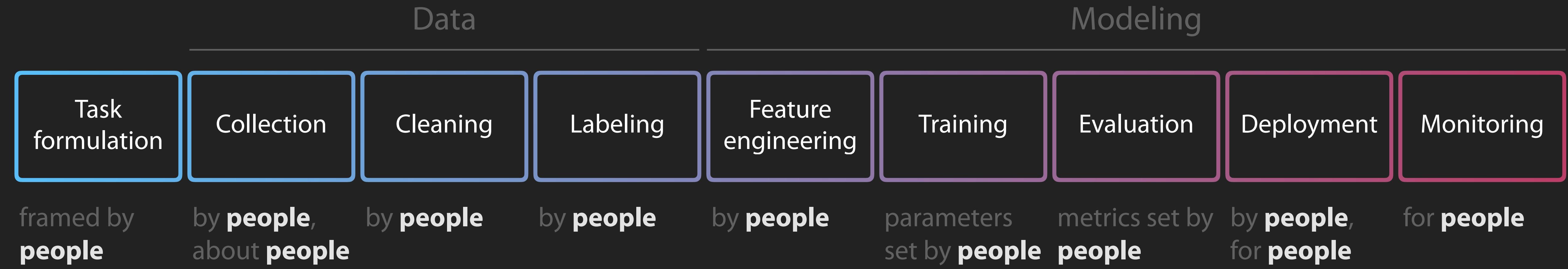
# Machine learning is a **people problem**.



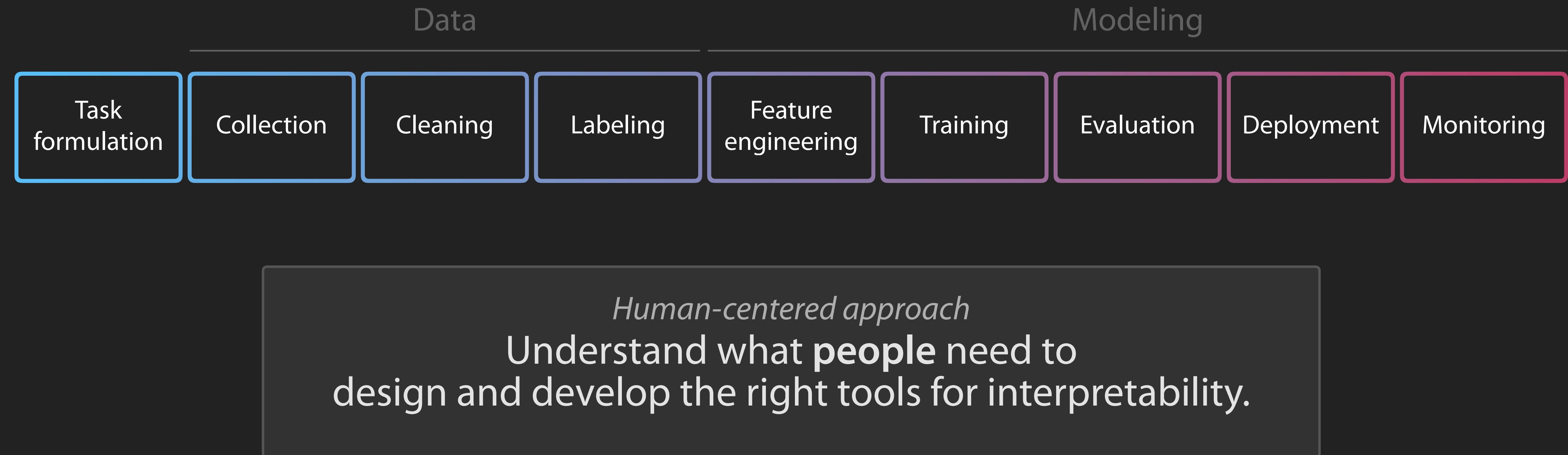
# Machine learning is a **people problem**.



# Machine learning is a **people problem**.



# Machine learning is a **people problem**.



## Research Mission

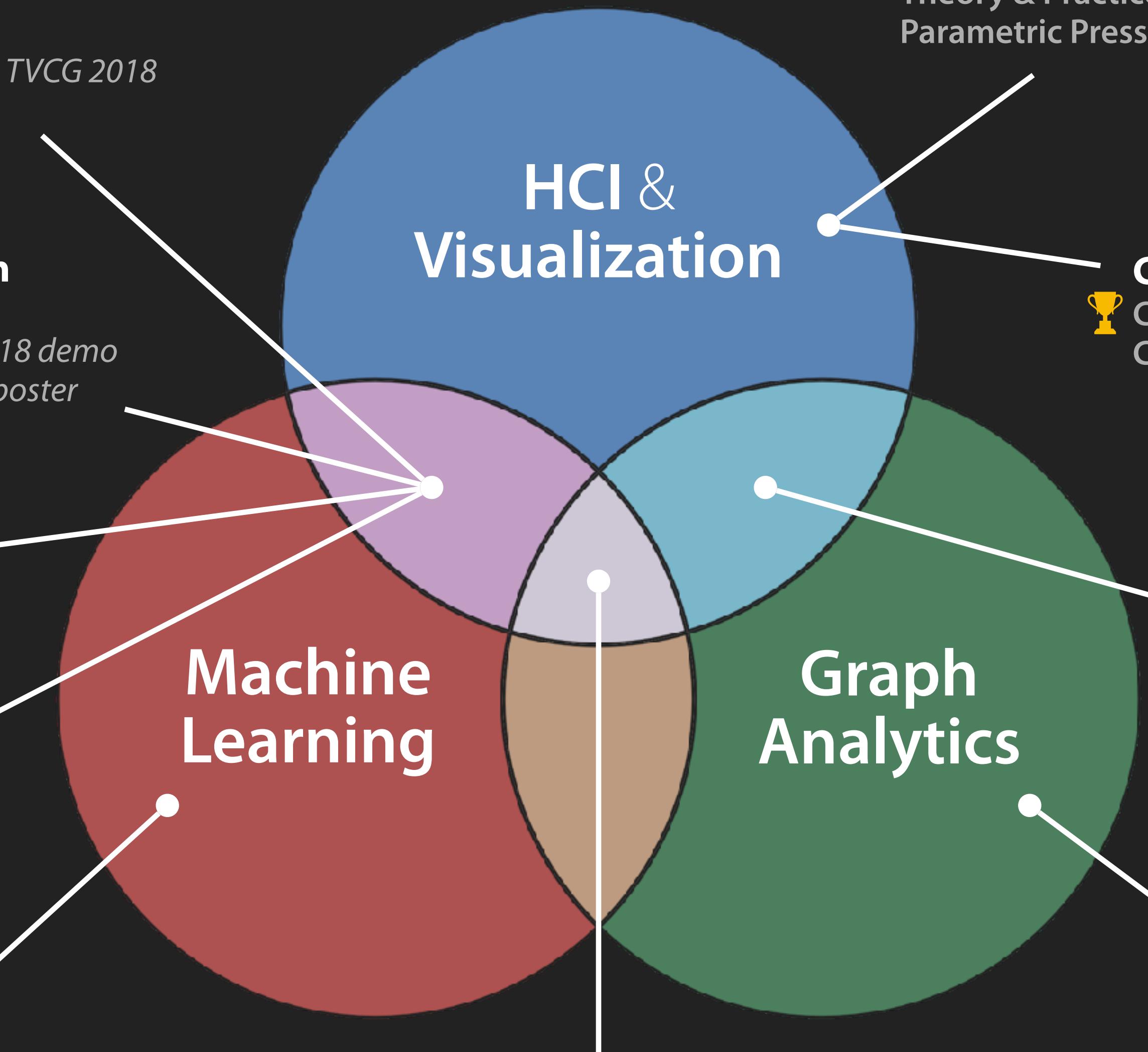
Enable machine learning interpretability at scale and for everyone.

## Research Mission

Enable machine learning interpretability at scale and for everyone.

## Research Mission

Enable machine learning interpretability at scale and for everyone,  
by *designing and developing interactive interfaces* to  
help people confidently understand data-driven systems.



## Visualization Tools for ML

- Gamut *CHI 2019*
- TeleGam *VIS 2019*
- Interrogative Survey *TVCG 2018*
- CHAMELEON *CHI 2020*

## Interactive Experimentation

- SHAPESHOP *CHI 2017 poster*
- Interactive Classification *CVPR 2018 demo*
- NEURALDIVERGENCE *PacificVis 2019 poster*

## ML Education

- Dimensionality Reduction *VISxAI 2018*
- CNN 101 *CHI 2020 poster*

## ML Fairness

- Impartial Machine *Parametric Press 2019*
- FAIRVIS *VAST 2019*
- Intersectional Bias *DebugML 2019*

## ML Robustness

- SHIELD *KDD 2018*
- Compression Defenses *KDD 2018 showcase*
- DeepPop *GeoHum 2018*

## Interactive Articles

- Theory & Practice *Distill 2020*
- Parametric Press *VISCOMM 2019*

## Computational Notebooks

- Code Gathering *CHI 2019*
- Computational Handoff *CHI 2020 poster*

## Visual Graph Sensemaking

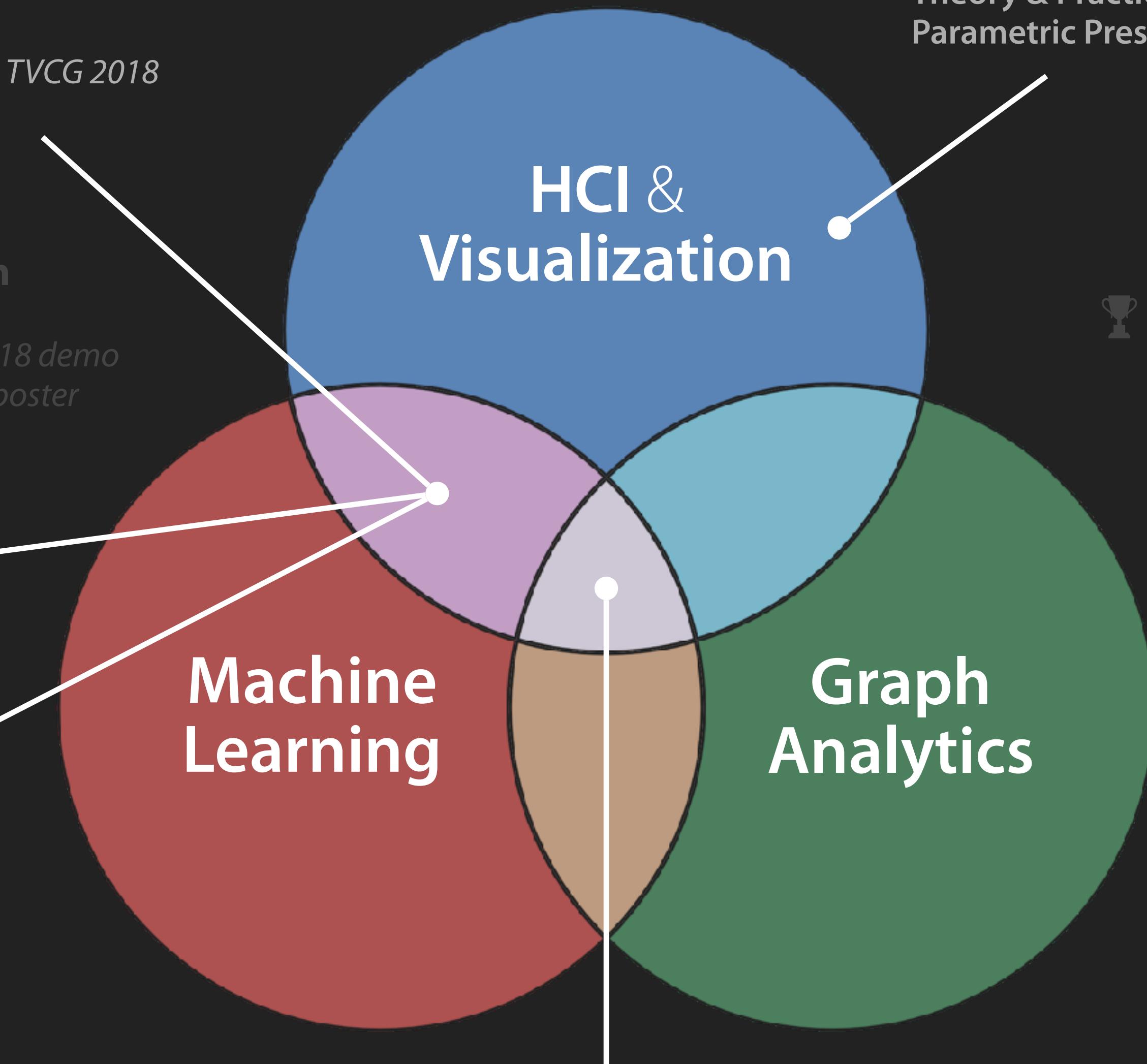
- ATLAS *IUI 2019*
- VIGOR *VAST 2018*
- VISAGE *SIGMOD 2017 demo*
- Graph Layers *VAST 2017 poster*

## Deep Learning Interpretability

- Summit *VAST 2019*
- MASSIF *CHI 2020 poster*

## Graph Mining

- Scalable k-core *Big Data 2019*



## Visualization Tools for ML

Gamut *CHI 2019*  
TeleGam *VIS 2019*  
Interrogative Survey *TVCG 2018*  
CHAMELEON *CHI 2020*

## Interactive Experimentation

SHAPESHOP *CHI 2017 poster*  
Interactive Classification *CVPR 2018 demo*  
NEURALDIVERGENCE *PacificVis 2019 poster*

## ML Education

Dimensionality Reduction *ViSxAI 2018*  
CNN 101 *CHI 2020 poster*

## ML Fairness

Impartial Machine *Parametric Press 2019*  
FAIRVIS *VAST 2019*  
Intersectional Bias *DebugML 2019*

## ML Robustness

SHIELD *KDD 2018*  
Compression Defenses *KDD 2018 showcase*  
DeepPop *GeoHum 2018*

## Interactive Articles

Theory & Practice *Distill 2020*  
Parametric Press *VISCOMM 2019*

## Computational Notebooks

Code Gathering *CHI 2019*  
Computational Handoff *CHI 2020 poster*

## Visual Graph Sensemaking

ATLAS *IUI 2019*  
VIGOR *VAST 2018*  
VISAGE *SIGMOD 2017 demo*  
Graph Layers *VAST 2017 poster*

## Deep Learning Interpretability

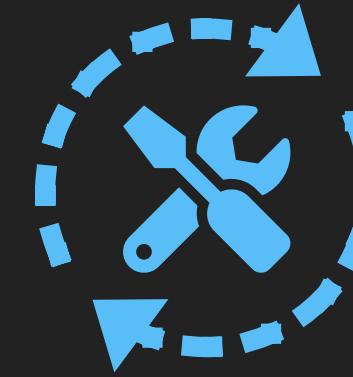
Summit *VAST 2019*  
MASSIF *CHI 2020 poster*

## Graph Mining

Scalable k-core *Big Data 2019*

# Interactive Scalable Interfaces for Machine Learning Interpretability

---



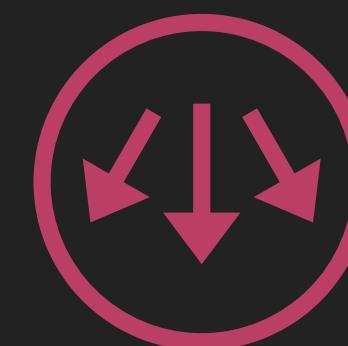
## PART I **Enable** interpretability

**GAMUT** Operationalize interpretability *CHI 2019*  
**TELEGAM** Vis + text for better explanations *VIS 2019*



## PART II **Scale** interpretability

**Interrogative Survey** Summarize interpretability vis *TVCG 2018*  
**SUMMIT** Higher-level explanations for neural networks *VAST 2019*

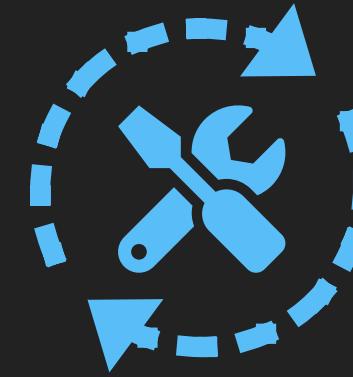


## PART III **Communicate** interpretability

**ML Literacy** Interactive mediums & platforms *VISCOMM 2019, VISxAI 2018*  
**Interactive Articles** Formalizing interactive communication *Distill 2020*

# Interactive Scalable Interfaces for Machine Learning Interpretability

---



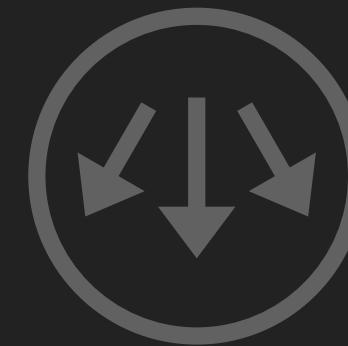
## PART I Enable interpretability

**GAMUT** Operationalize interpretability *CHI 2019*  
**TELEGAM** Vis + text for better explanations *VIS 2019*



## PART II Scale interpretability

**Interrogative Survey** Summarize interpretability vis *TVCG 2018*  
**SUMMIT** Higher-level explanations for neural networks *VAST 2019*



## PART III Communicate interpretability

**ML Literacy** Interactive mediums & platforms *VISCOMM 2019, VISxAI 2018*  
**Interactive Articles** Formalizing interactive communication *Distill 2020*

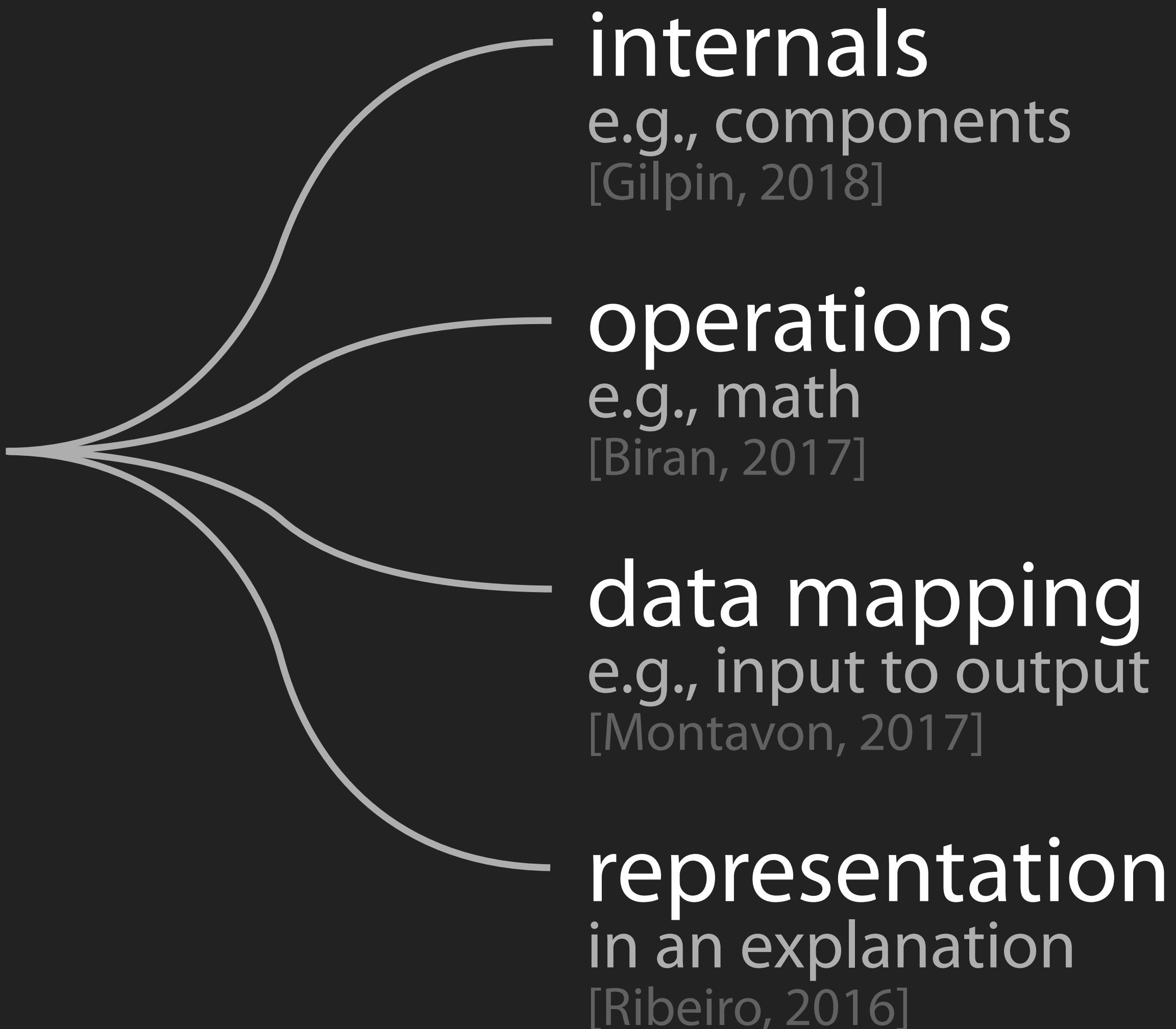
# What is Interpretability?

# What is Interpretability?

*Human understanding  
of a system's...*

# What is Interpretability?

*Human understanding  
of a system's...*



# What is Interpretability?

*Human understanding  
of a system's...*

No formal, agreed upon definition  
[Lipton, 2016]

internals  
e.g., components  
[Gilpin, 2018]

operations  
e.g., math  
[Biran, 2017]

data mapping  
e.g., input to output  
[Montavon, 2017]

representation  
in an explanation  
[Ribeiro, 2016]

# GAMUT

Understanding How Data Scientists  
Understand Machine Learning Models

*CHI 2019*



**Fred Hohman**

Georgia Tech



**Andrew Head**  
UC Berkeley



**Rich Caruana**  
Microsoft Research



**Rob DeLine**  
Microsoft Research



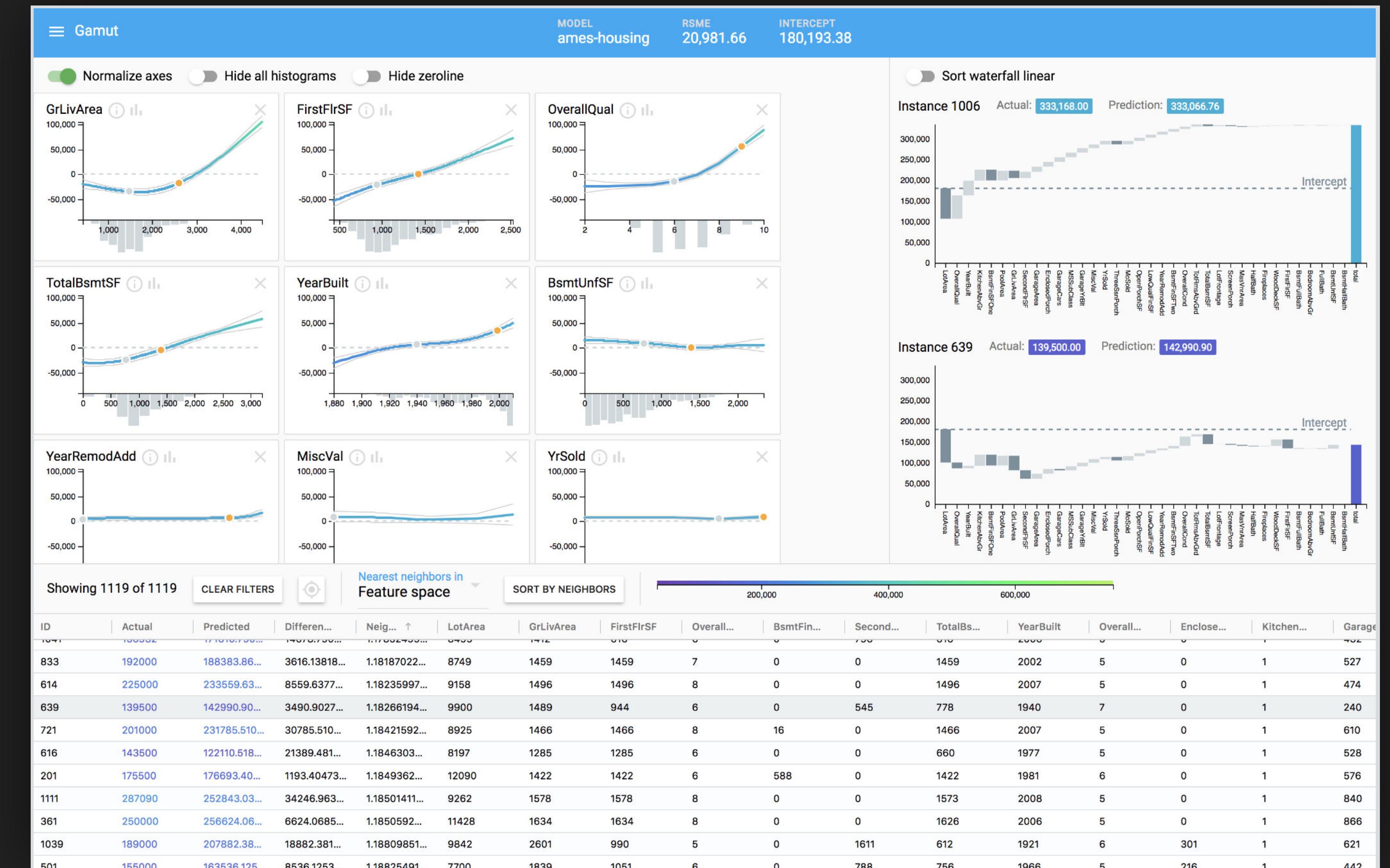
**Steven Drucker**  
Microsoft Research

# GAMUT Contributions

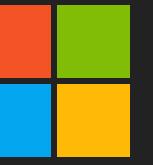
# Operationalize interpretability

# Design Probe embodying operationalization

# Evaluation & Investigation of probe & emerging practice of interpretability with real users



# How to Operationalize Interpretability?

Formative research with professional data scientists @ 

- 4 senior ML researchers
- 5 ML practitioners

*Prompt: In a perfect world, given a machine learning model, what questions would you ask it to help you interpret both the model and its predictions?*

# Explainable ML Interface Questions

Why this prediction?

What is the difference between these two?

What if I added...

What are similar instances?

Where is it wrong?

What is most important?

# Explainable ML Interface Capabilities



Why this prediction?

**Local instance explanations**



What is the difference between these two?

**Instance explanation comparisons**



What if I added...

**Counterfactuals**



What are similar instances?

**Nearest neighbors**



Where is it wrong?

**Regions of error**



What is most important?

**Feature importance**

# Explainable ML Interface Capabilities



Why this prediction?

## Local instance explanations



What is the difference between these two?

## Instance explanation comparisons



What if I added...

## Counterfactuals



What are similar instances?

## Nearest neighbors



Where is it wrong?

## Regions of error



What is most important?

## Feature importance

**GAMUT: A Design Probe to Understand How Data Scientists Understand Machine Learning Models**

Fred Heiman  
Georgia Institute of Technology  
Atlanta, GA, USA  
fred.heiman@gatech.edu

Andrew Head  
UC Berkeley  
Berkeley, CA, USA  
andrewhead@berkeley.edu

Rich Carpano  
Microsoft Research  
Redmond, WA, USA  
rcarpano@microsoft.com

Robert Deline  
Microsoft Research  
Redmond, WA, USA  
rbdeline@microsoft.com

Steven M. Drucker  
Microsoft Research  
Redmond, WA, USA  
sdrucker@microsoft.com

**ABSTRACT**  
Without good models and the right tools to interpret them, data scientists risk making decisions based on hidden biases, spurious correlations, and false generalizations. This has led to a rallying cry for model interpretability. Yet the concept of interpretability remains nebulous, such that researchers and tool designers lack actionable guidelines for how to incorporate interpretability into models and accompanying tools. Through an iterative design process with expert machine learning researchers and practitioners, we designed a visual analytics system, GAMUT, to explore how interactive interfaces could better support model interpretation. Using GAMUT as a probe, we investigated why and how professional data scientists interpret models, and how interface affordances can support data scientists in answering questions about model interpretability. Our investigation showed that interpretability is not a monolithic concept data scientists

**CCS CONCEPTS**  
• Human-centered computing → Empirical studies in visualization; Visualization systems and tools • Computing methodologies → Machine learning

**KEYWORDS**  
Machine learning interpretability, design probe, visual analytics, data visualization, interactive interfaces

**ACM Reference Format:**  
Fred Heiman, Andrew Head, Rich Carpano, Robert Deline, and Steven M. Drucker. 2019. GAMUT: A Design Probe to Understand How Data Scientists Understand Machine Learning Models. In *CHI Conference on Human Factors in Computing Systems Proceedings (CHI 2019)*, May 4–9, 2019, Glasgow, Scotland UK. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3290605.3309809>

*Definitions + examples  
in the paper!*

# How to Test Our Capabilities?

[Hutchinson, 2003]

**Design probe:** “*instrument that is deployed to find out about the unknown—returning with useful or interesting data.*”



Balance of:

- **Design:** inspire reflection on emerging tech
- **Social science:** appreciate needs of users
- **Engineering:** field-testing prototypes



**House 550**  
**\$190,606**

# House 550

\$190,606

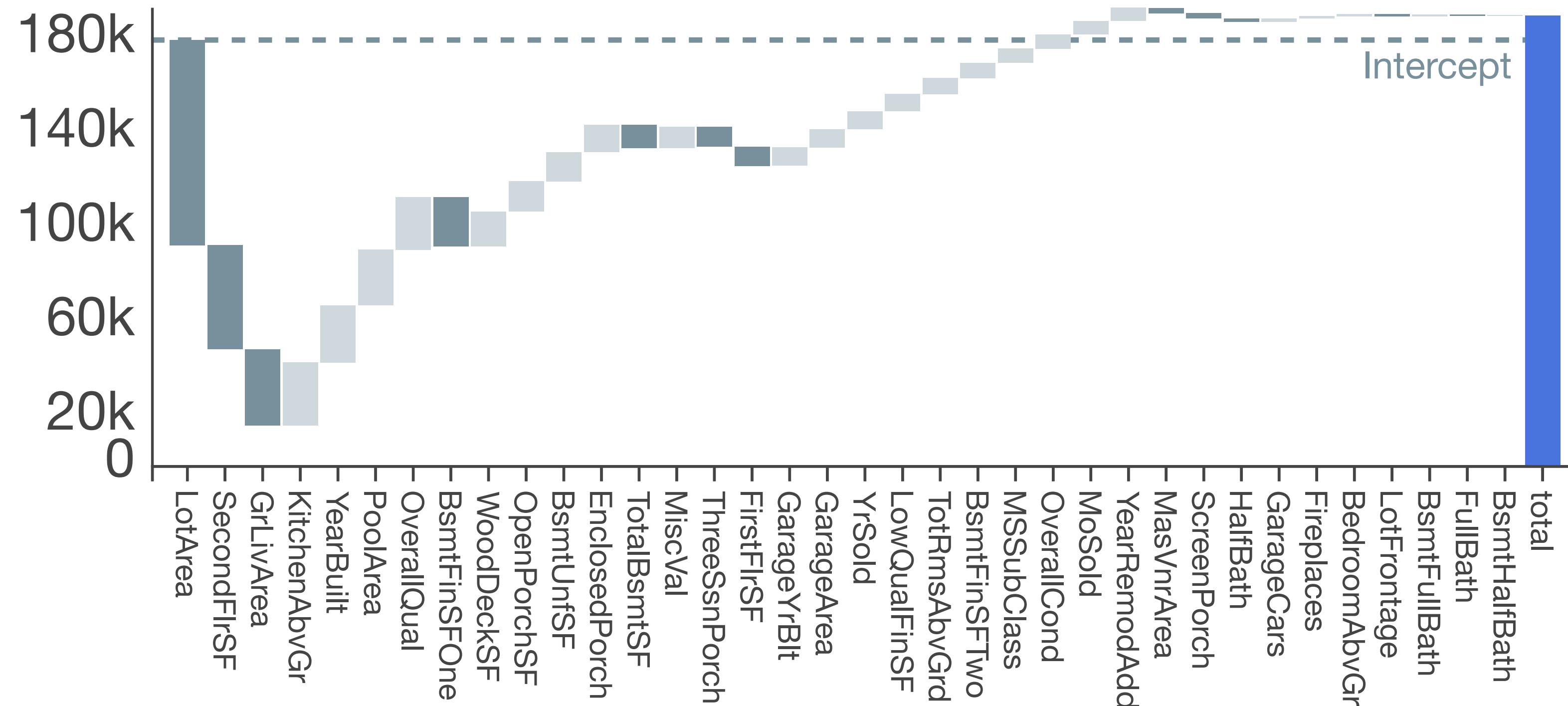
BsmtHalfBath  
FullBath  
BsmtFullBath  
LotFrontage  
BedroomAbvGr  
Fireplaces  
GarageCars  
HalfBath  
ScreenPorch  
MasVnrArea  
YearRemodAdd  
MoSold  
OverallCond  
MSSubClass  
BsmtFinSFTwo  
TotRmsAbvGrd  
LowQualFinSF  
YrSold  
GarageArea  
GarageYrBlt  
FirstFlrSF  
ThreeSsnPorch  
MiscVal  
TotalBsmtSF  
EnclosedPorch  
BsmtUnfSF  
OpenPorchSF  
WoodDeckSF  
BsmtFinSFOne  
OverallQual  
PoolArea  
YearBuilt  
KitchenAbvGr  
GrLivArea  
SecondFlrSF  
LotArea

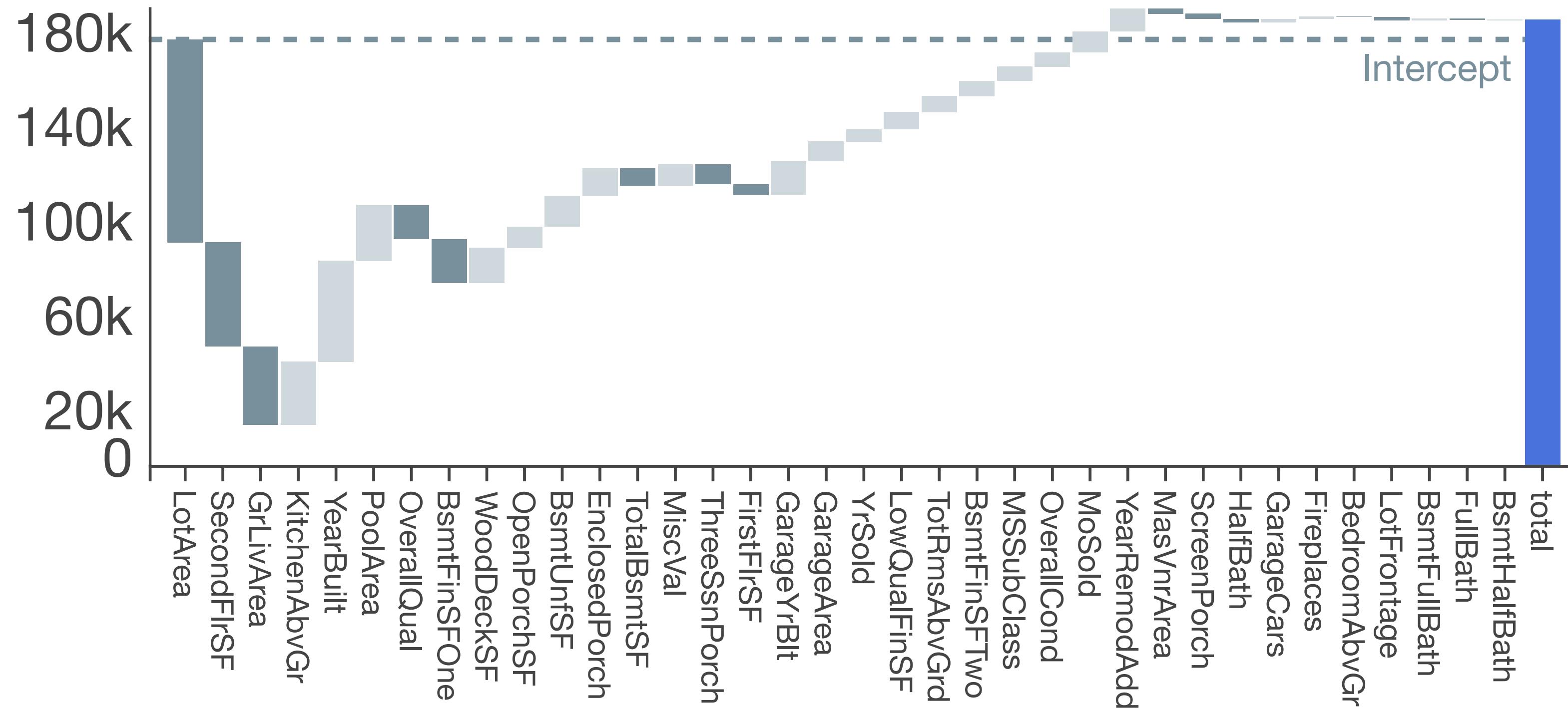
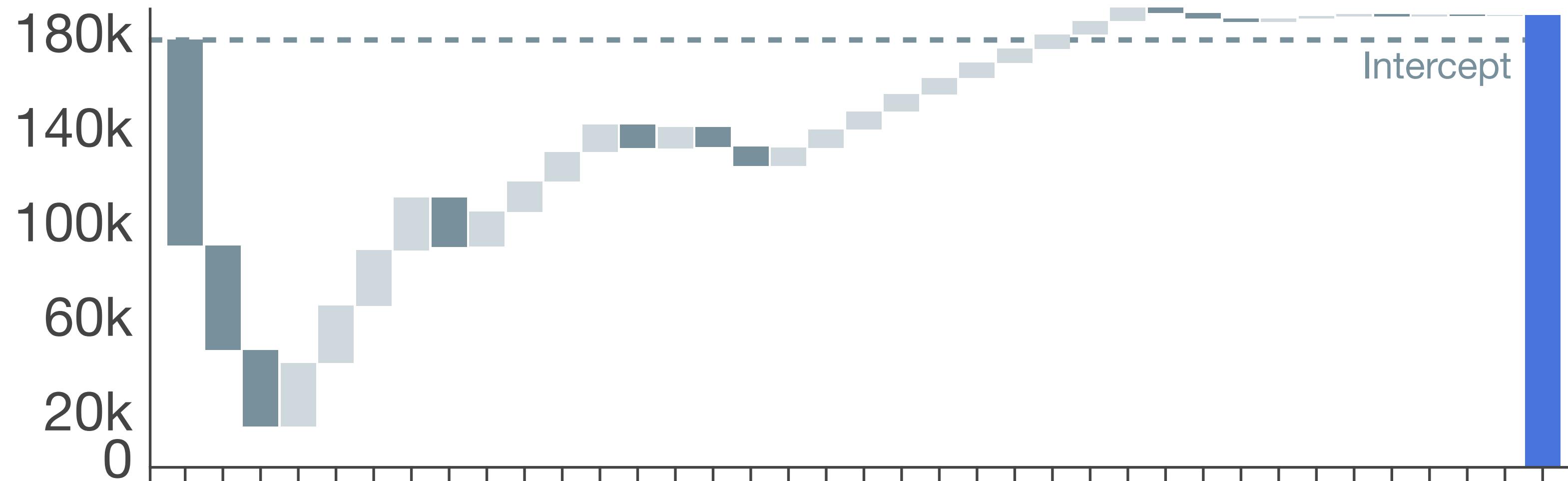
# House 550

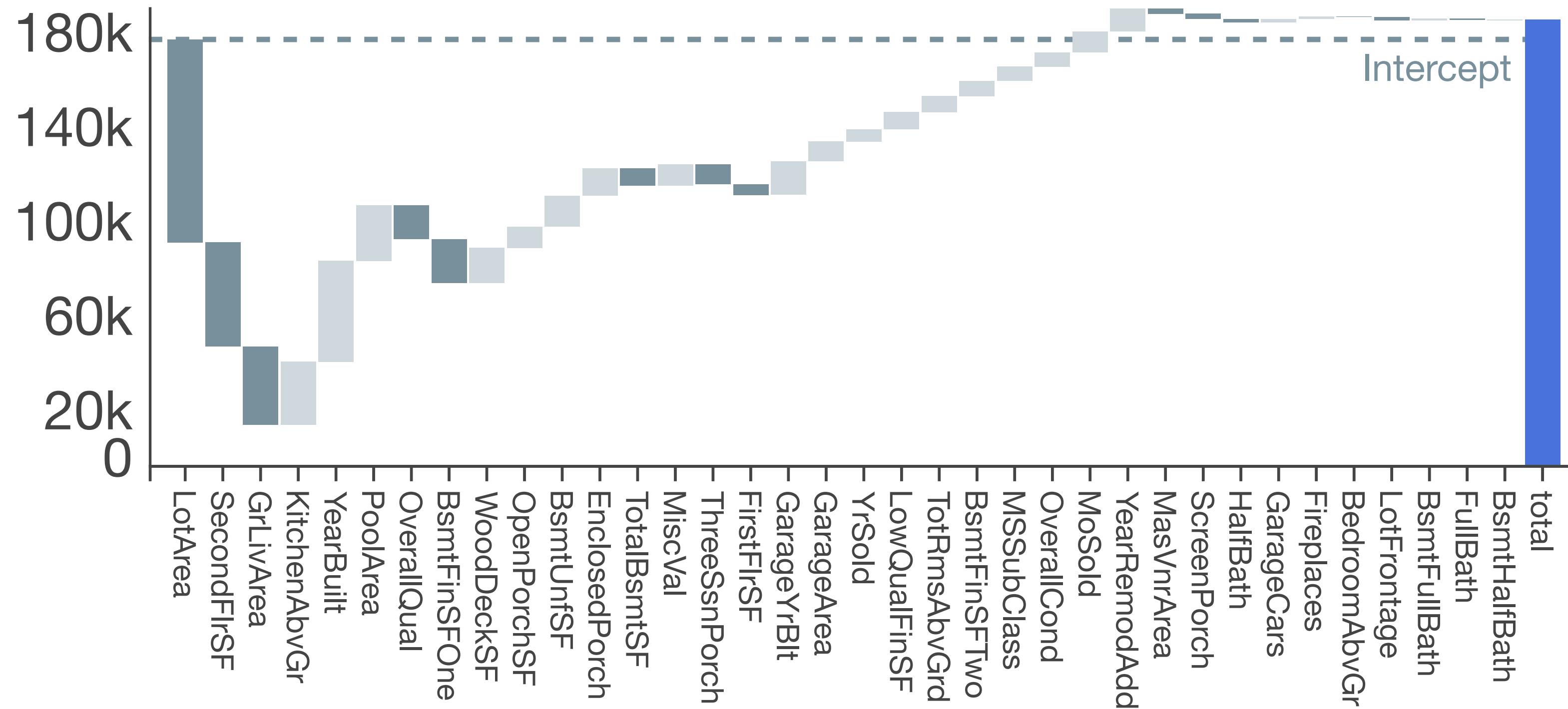
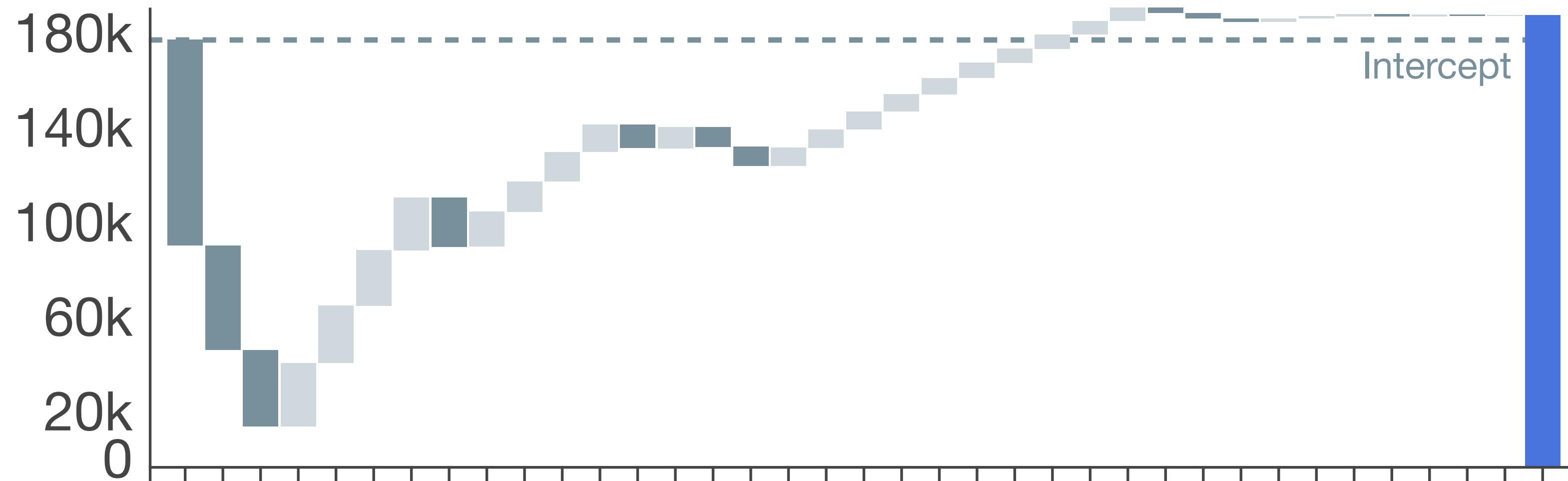
\$190,606

Prediction

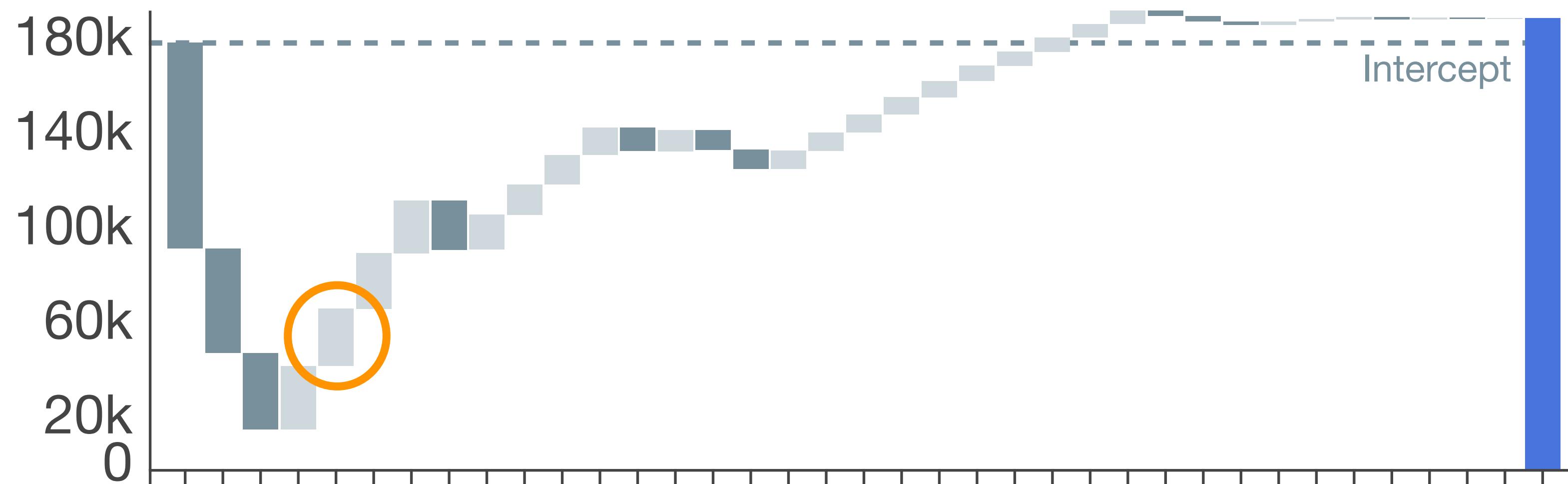
House features



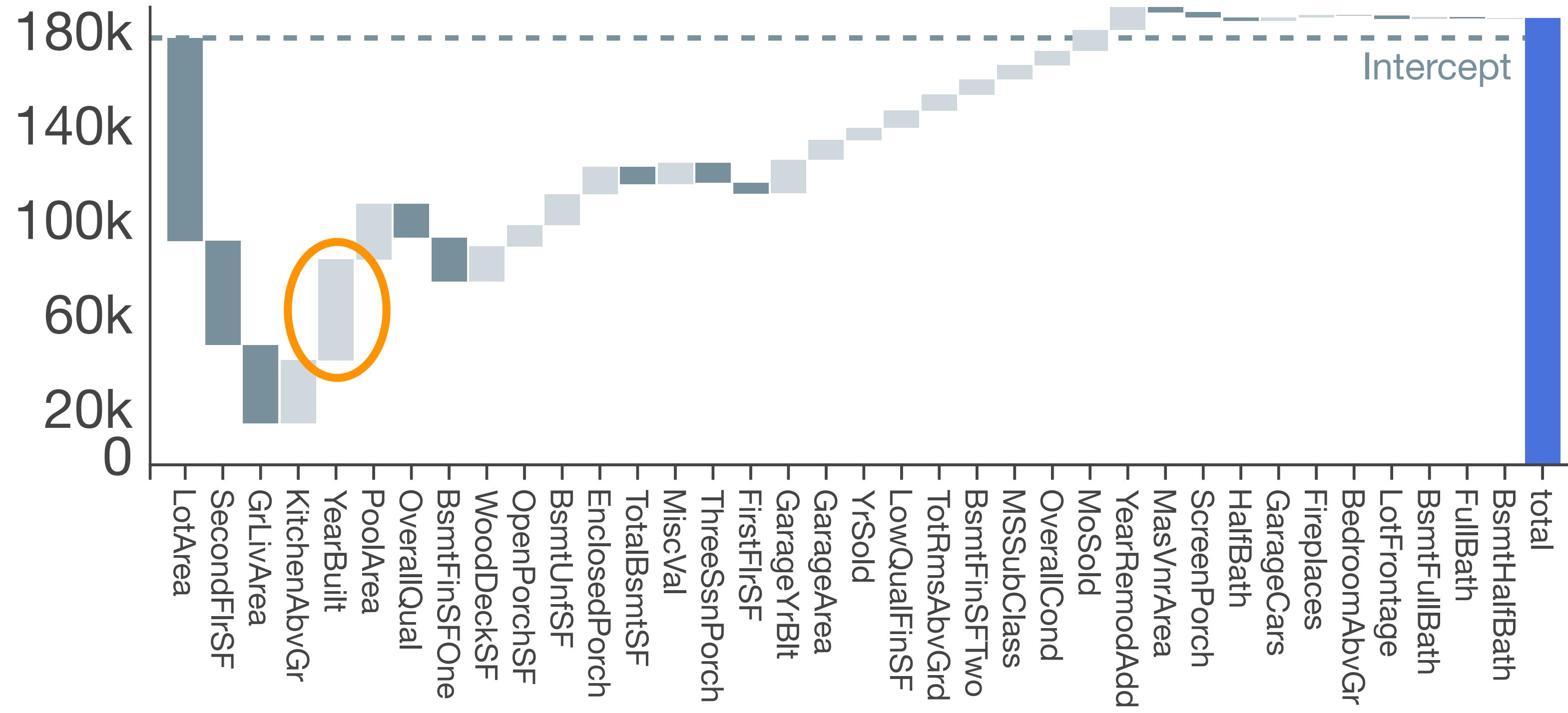


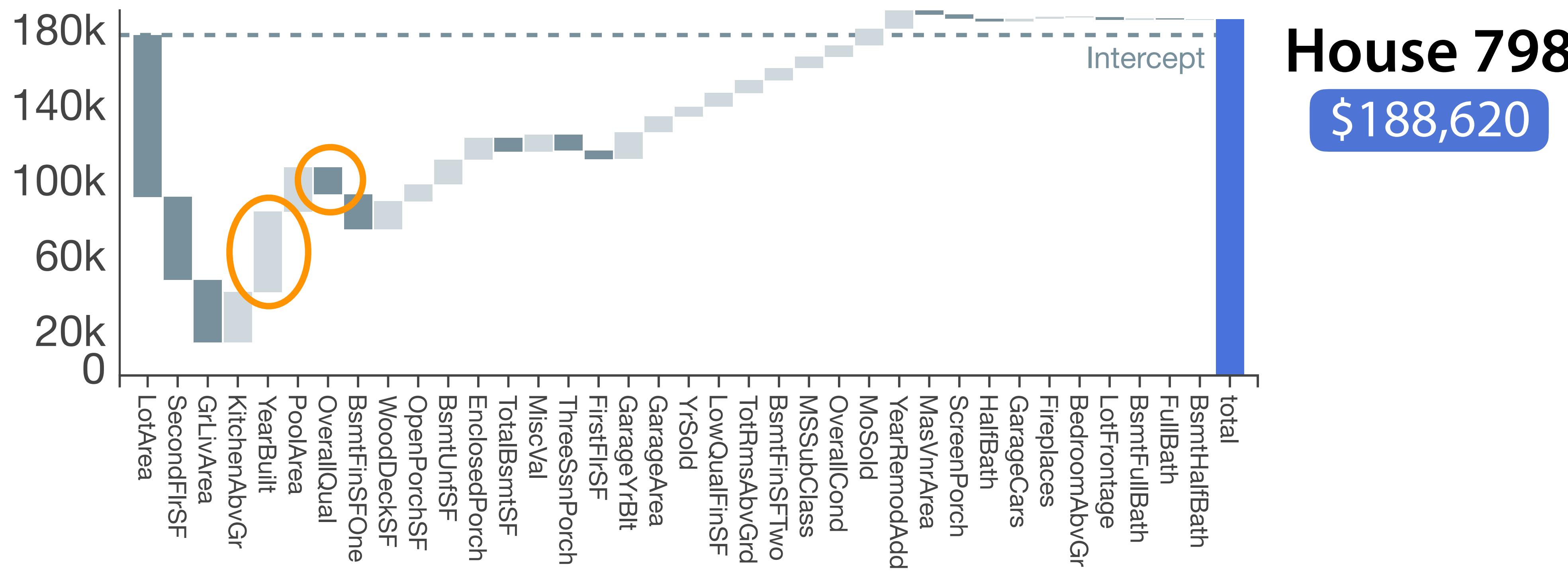
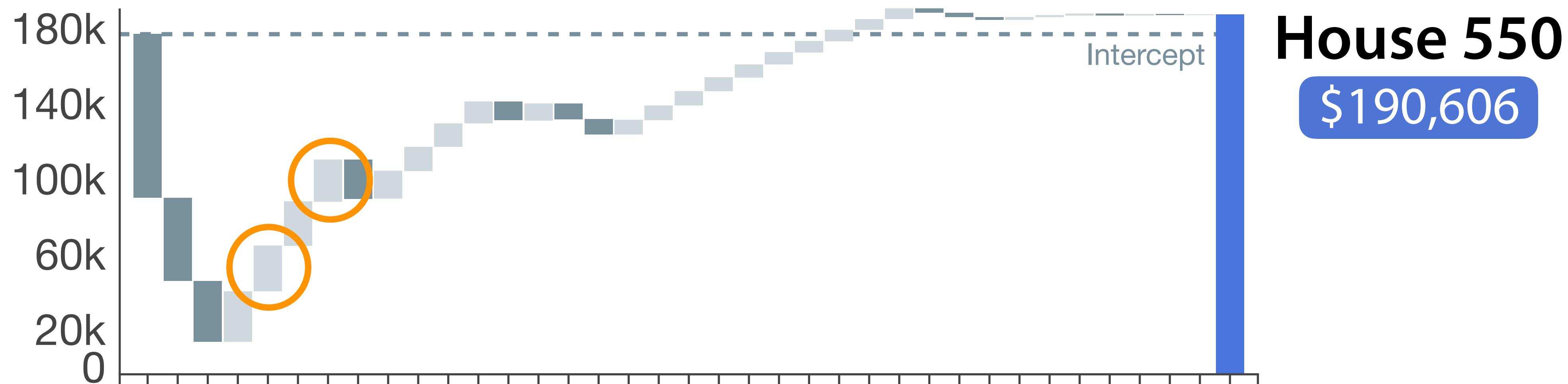


**House 550**  
\$190,606

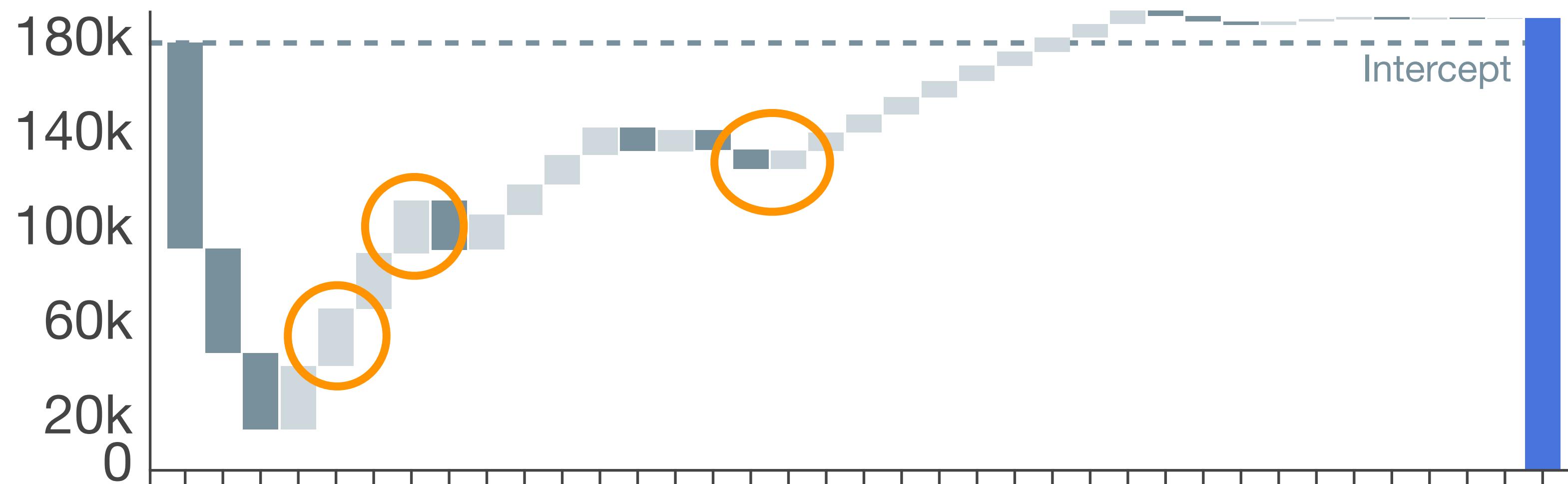


**House 798**  
\$188,620

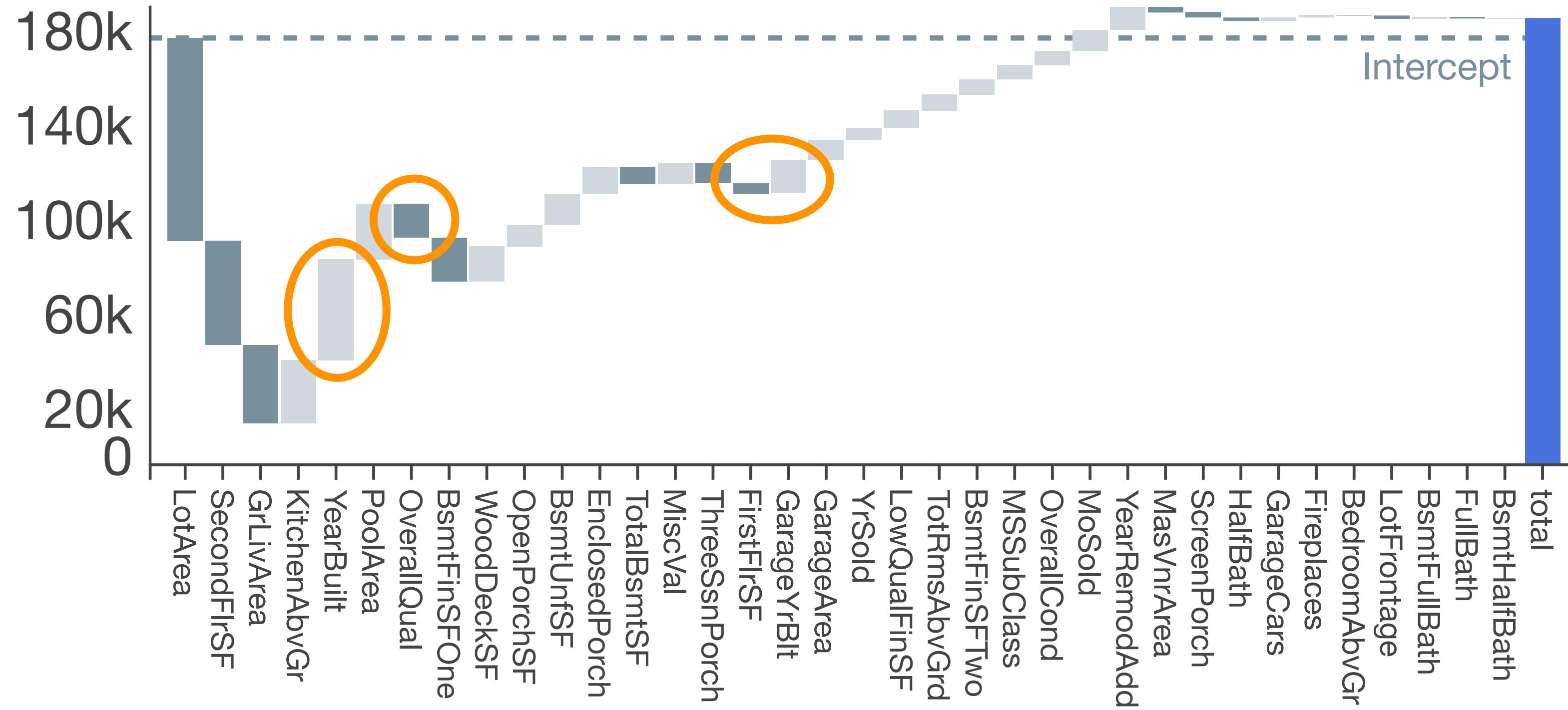


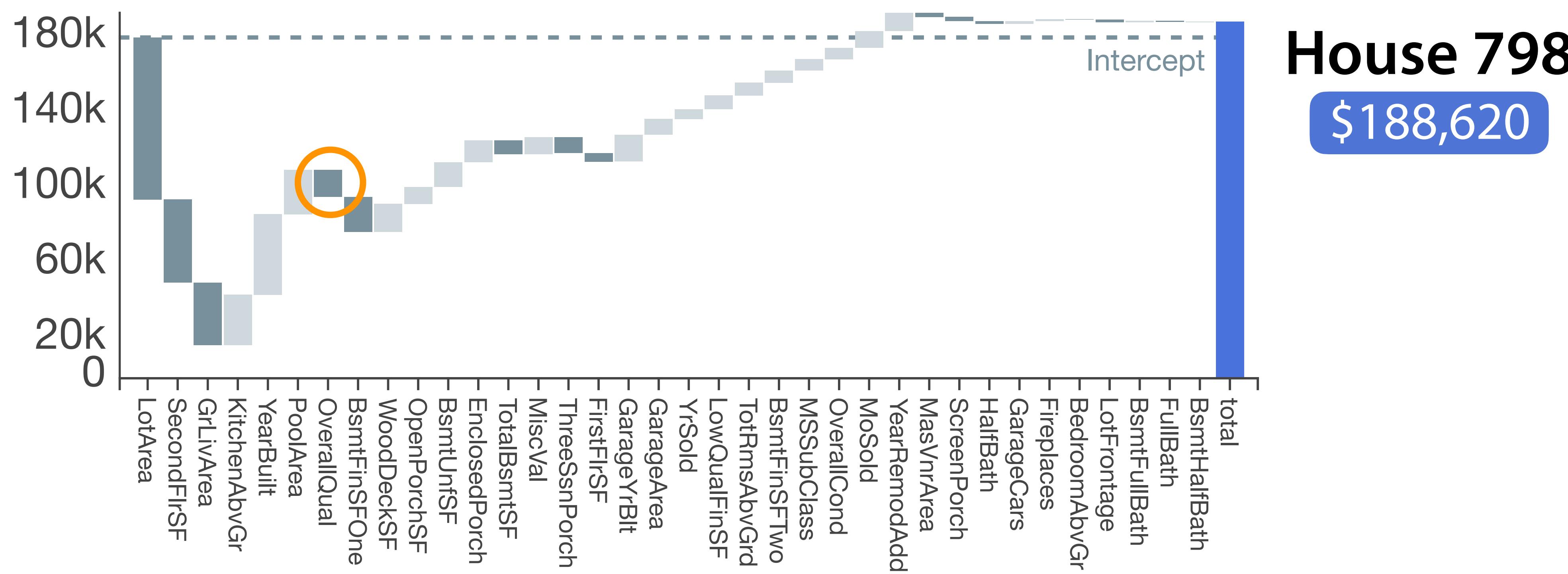
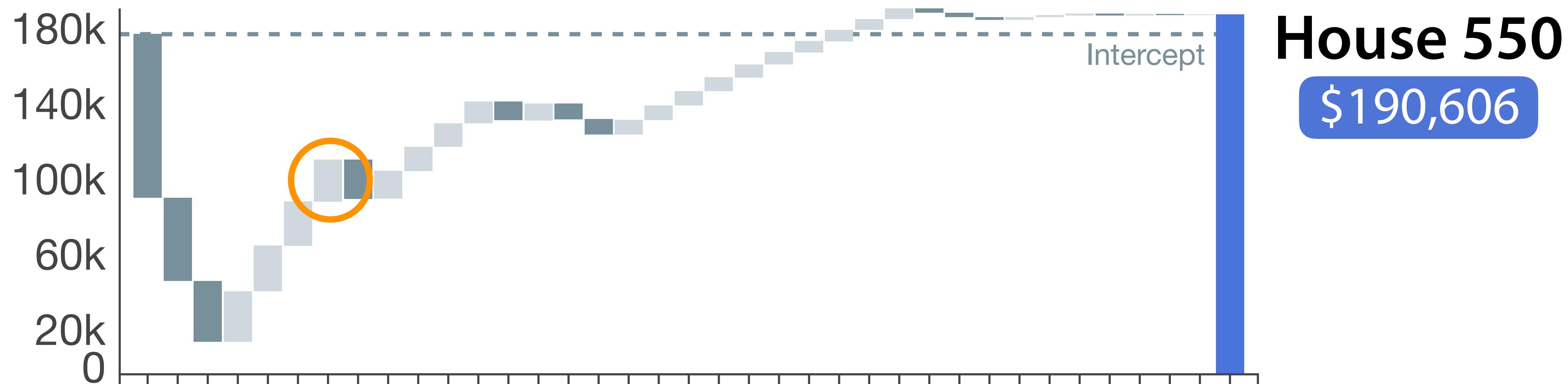


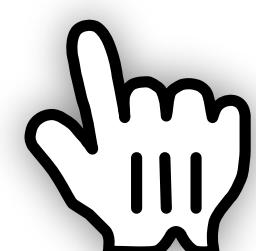
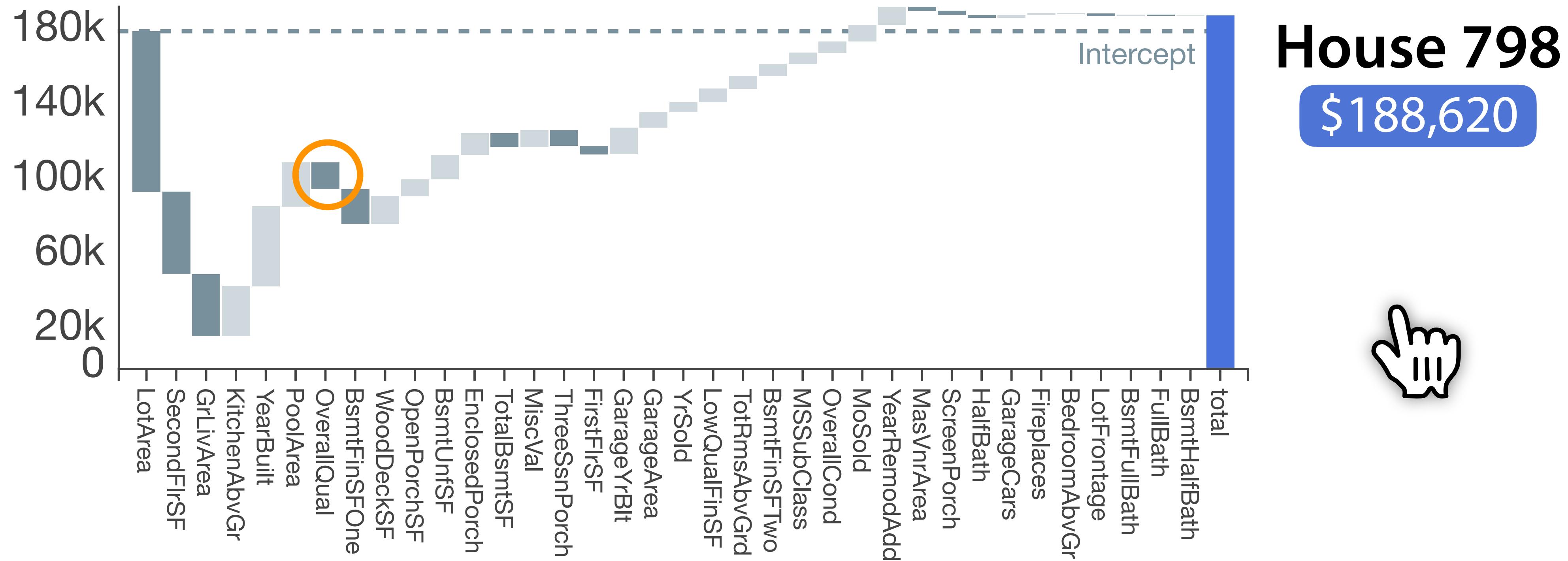
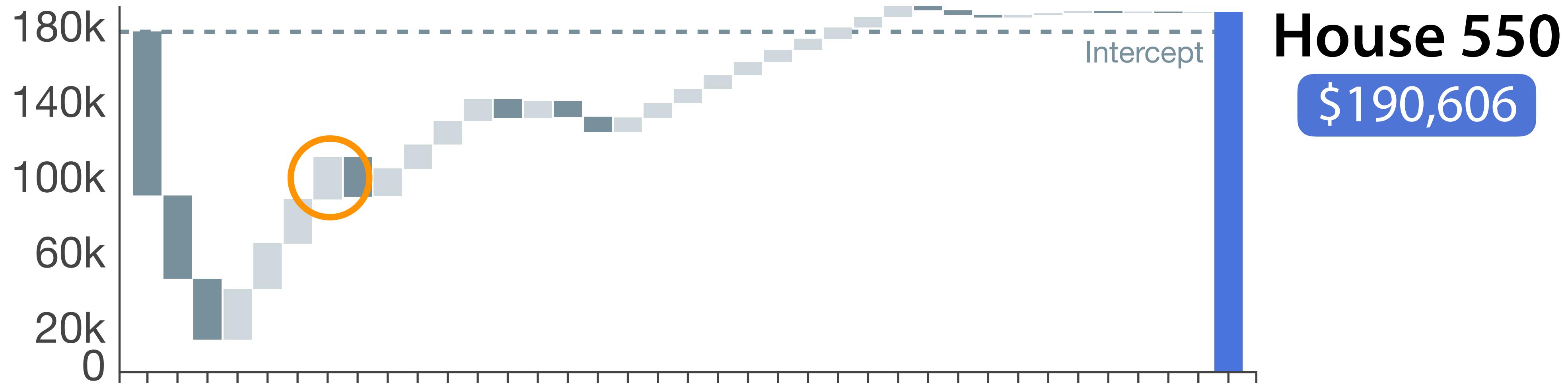
**House 550**  
\$190,606

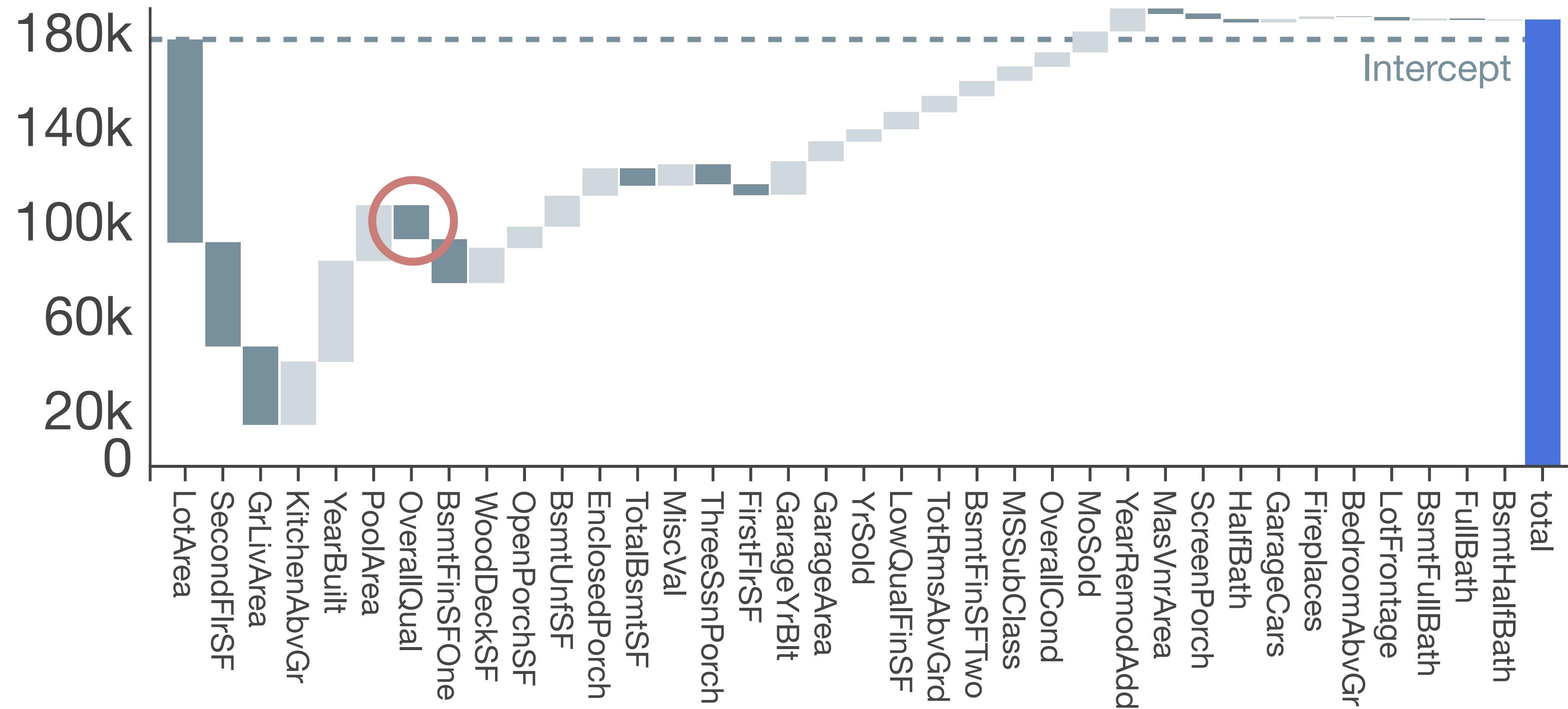
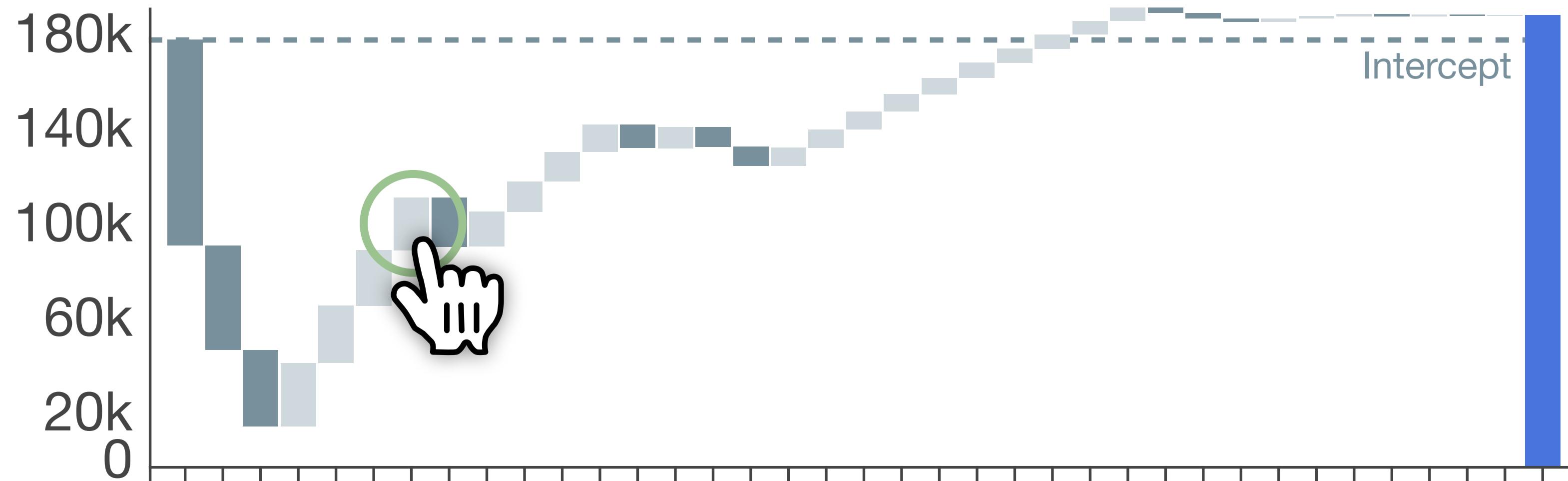


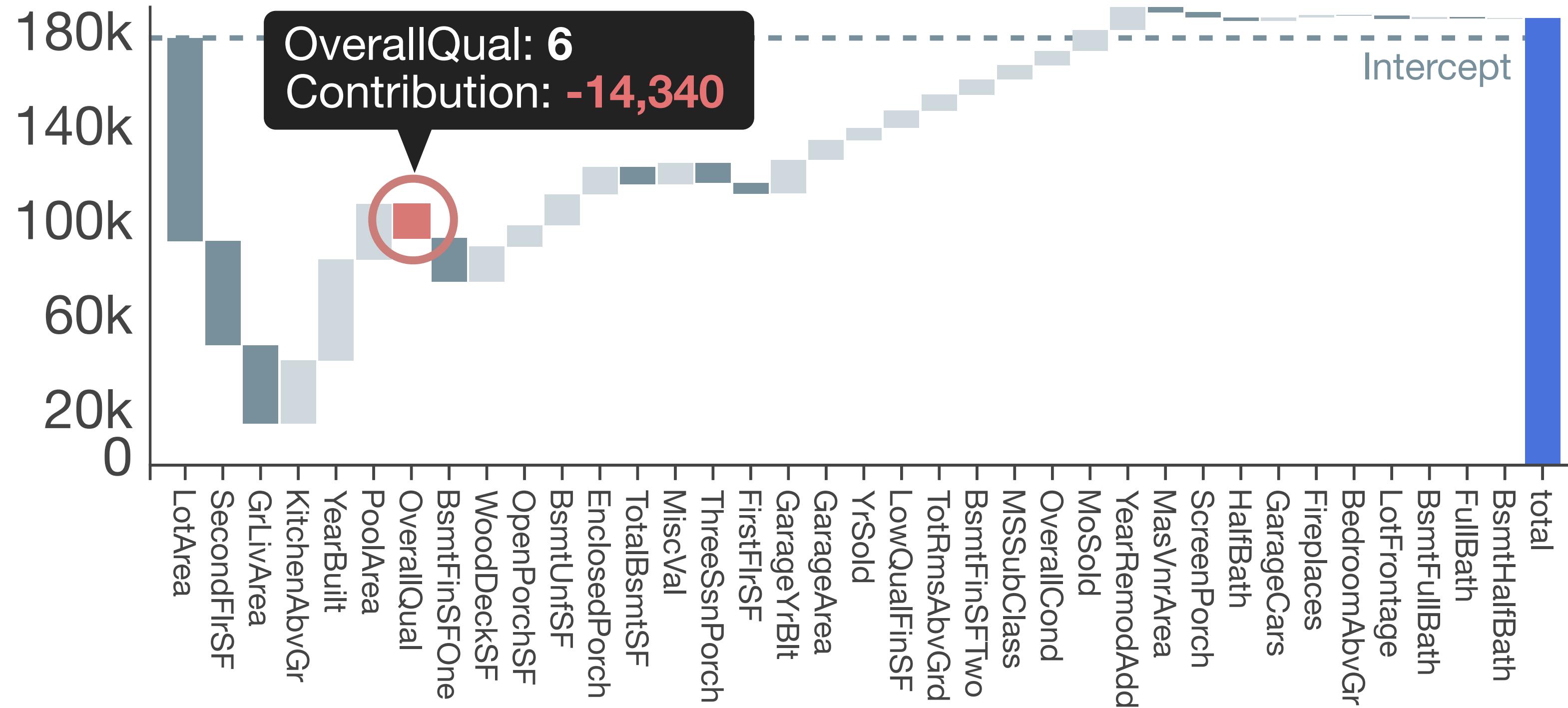
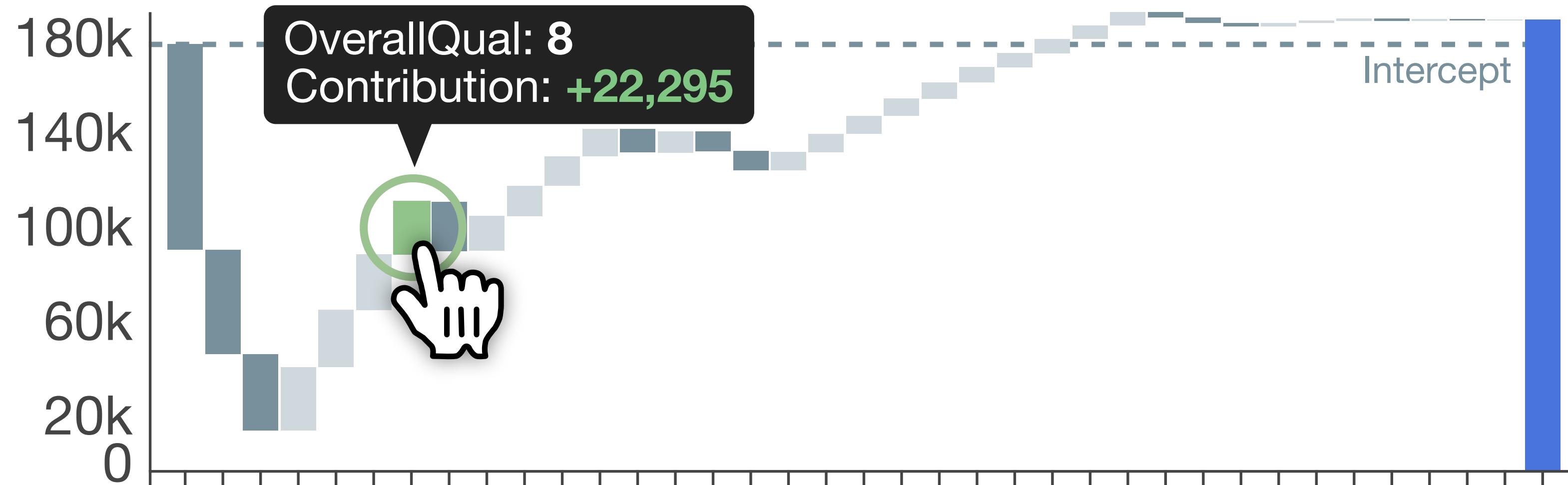
**House 798**  
\$188,620





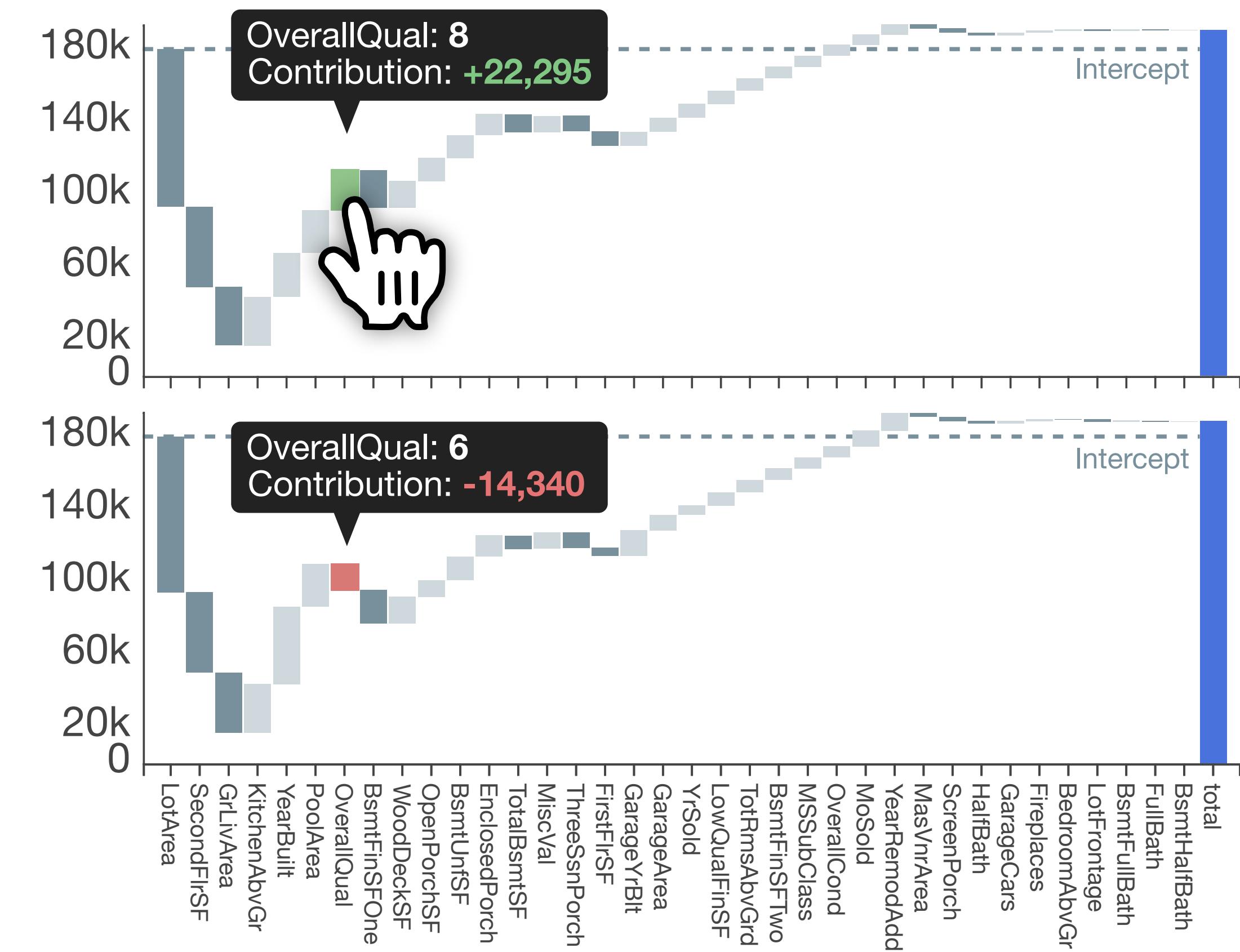
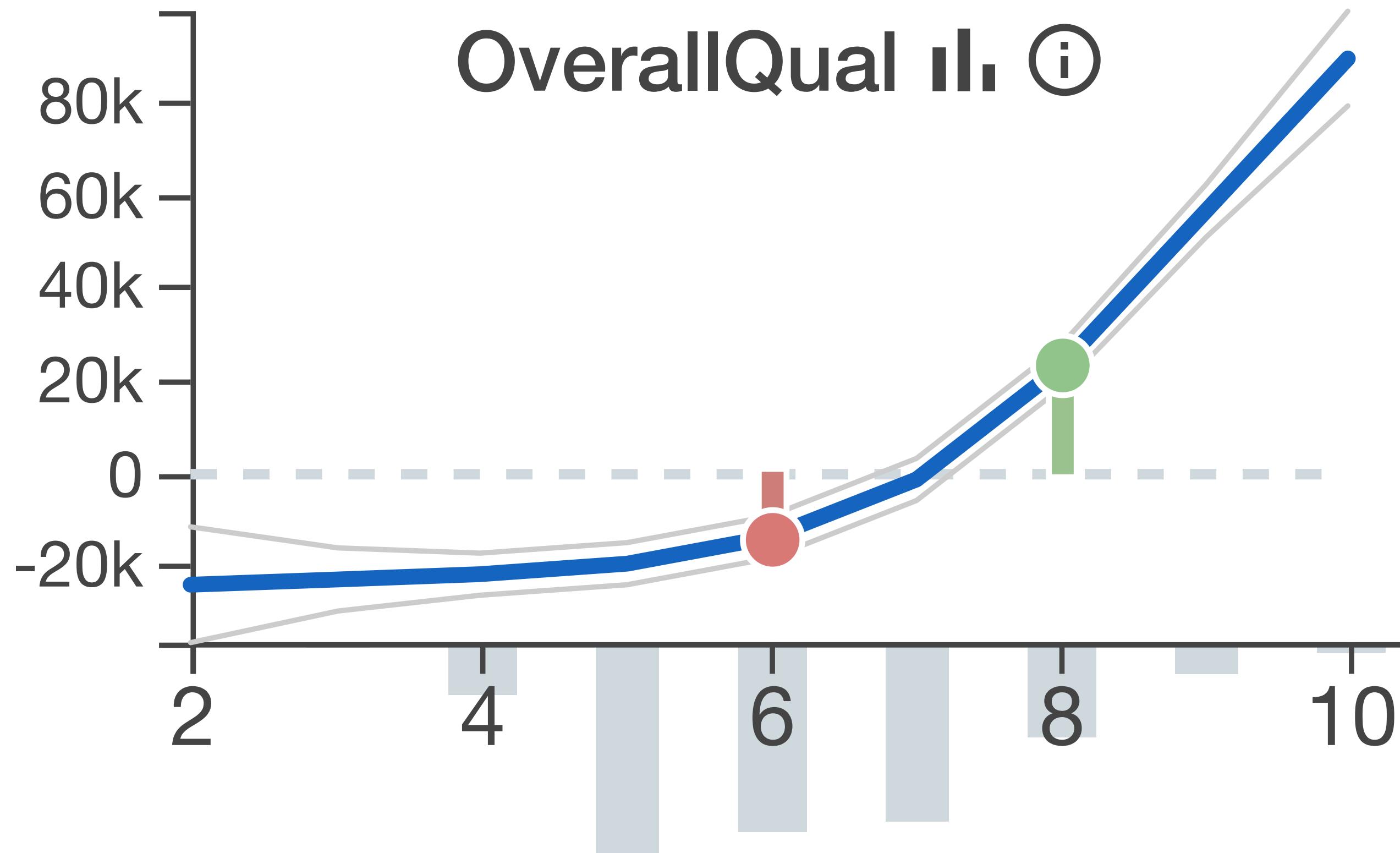






# House 550

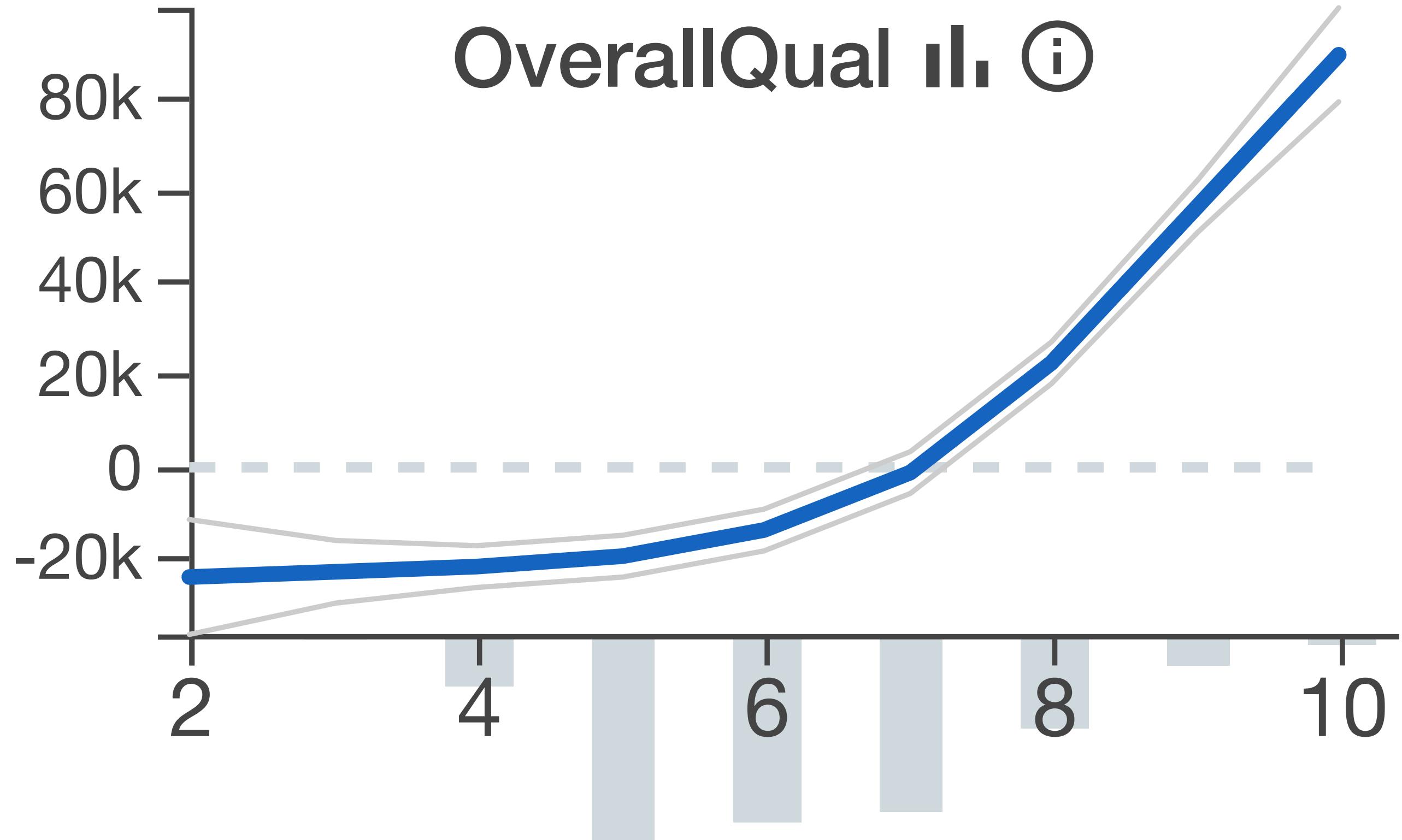
\$190,606



# House 798

\$188,620

# Generalized Additive Model (GAM)



- Global explanation
- Easy to understand:
- Average math skills
- Average graphicacy
- High accuracy, realistic

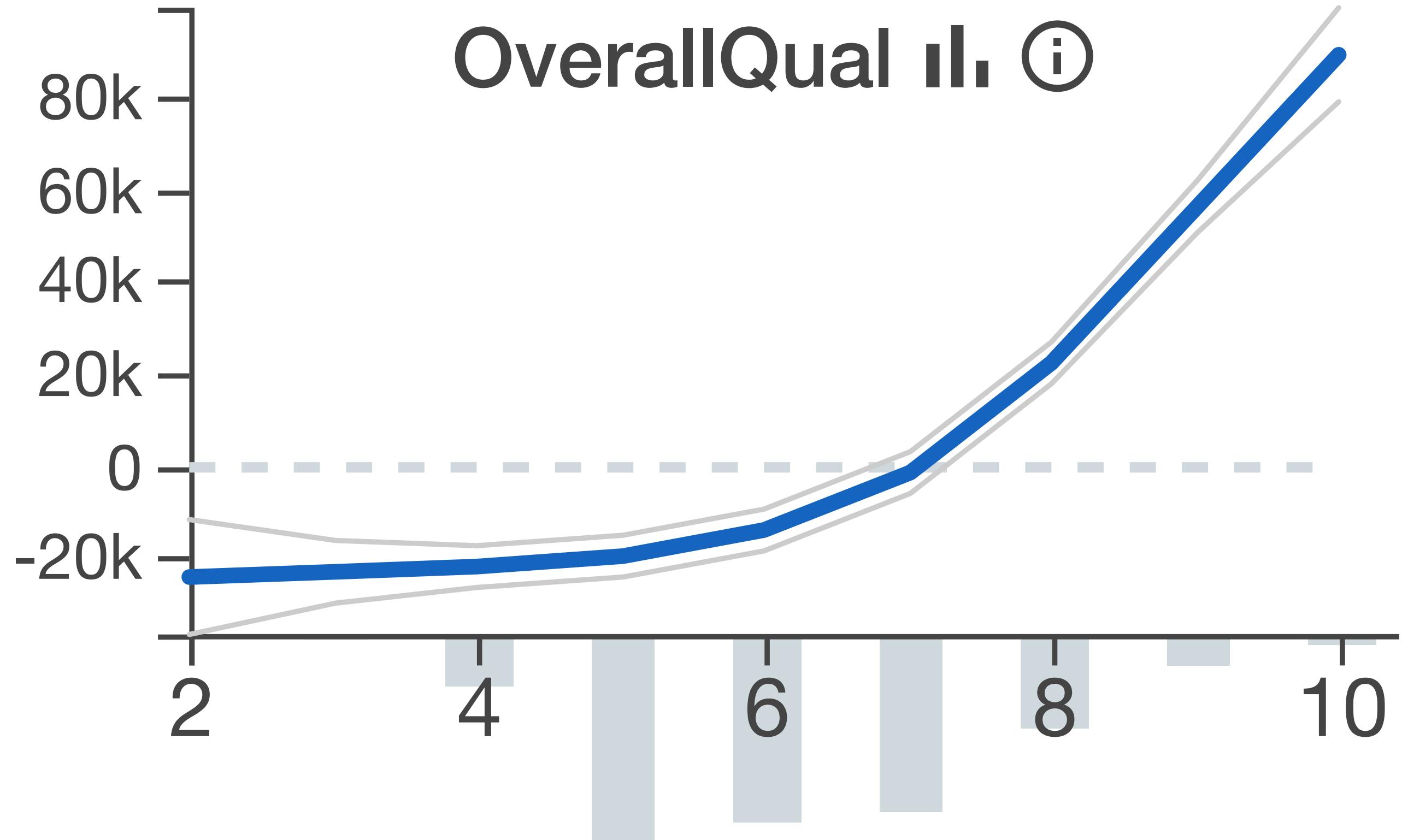
Generalized **linear** model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

Generalized **additive** model

$$y = \beta_0 + f_1(x_1) + f_2(x_2) + \dots + f_n(x_n)$$

# Generalized Additive Model (GAM)



- Global explanation
- Easy to understand:
- Average math skills
- Average graphicacy
- High accuracy, realistic

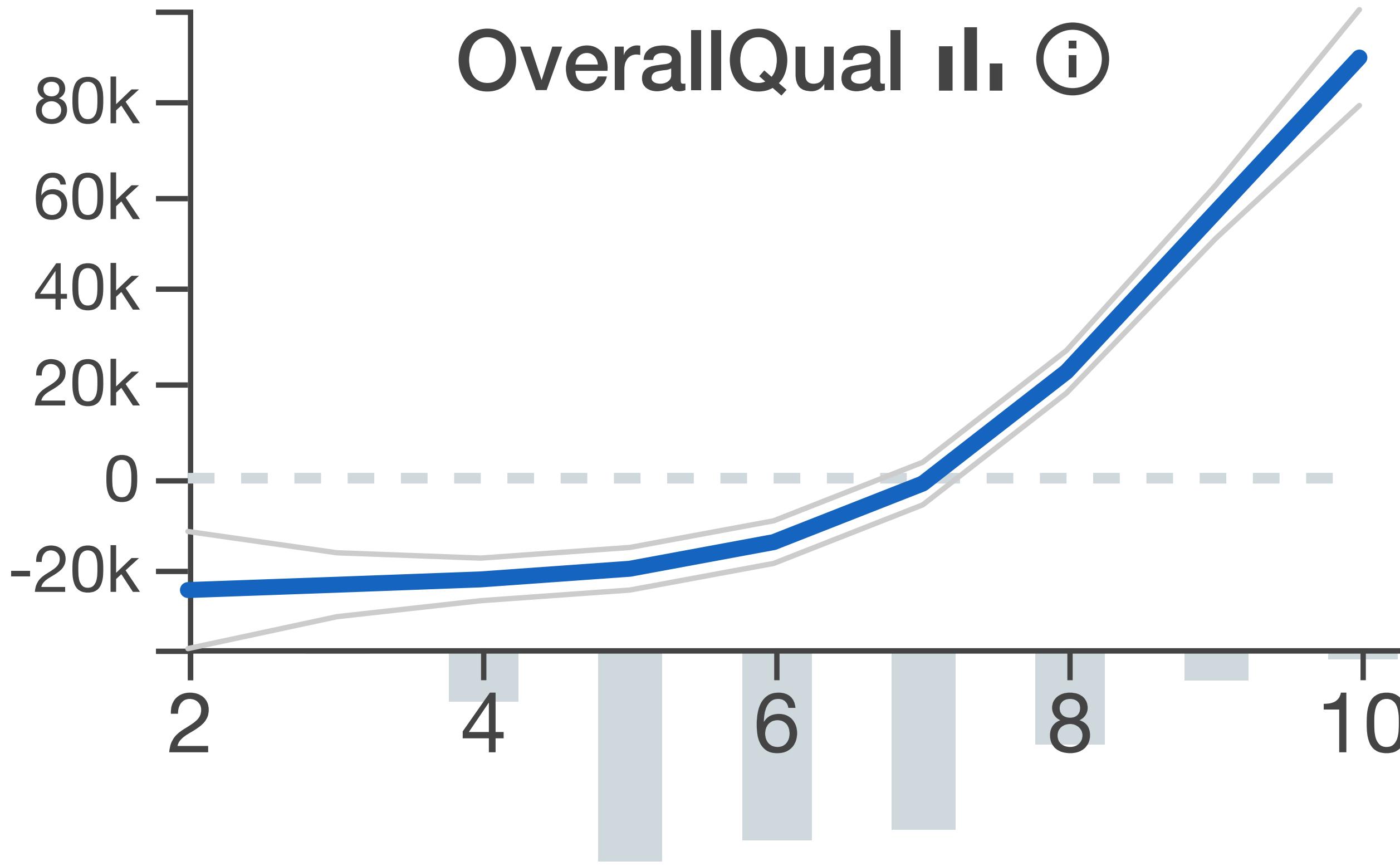
Generalized **linear** model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

Generalized **additive** model

$$y = \beta_0 + f_1(x_1) + f_2(x_2) + \dots + f_n(x_n)$$

## OverallQual II. ⓘ



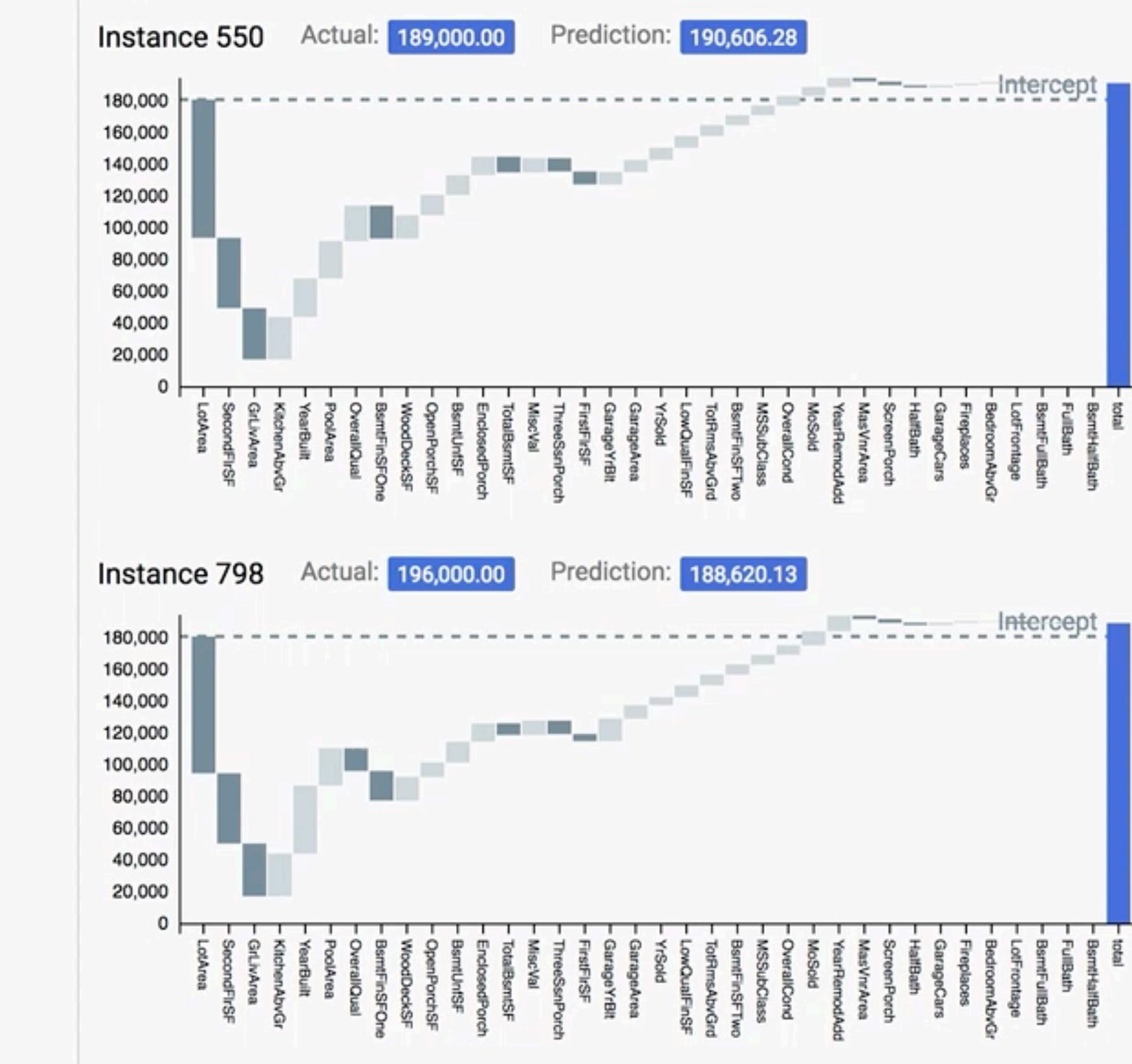
Normalize axes

 Hide all histograms

Hide zero line



## Sort waterfall linear



Showing 1119 of 1119

CLEAR FILTER



## Nearest neighbor Feature space

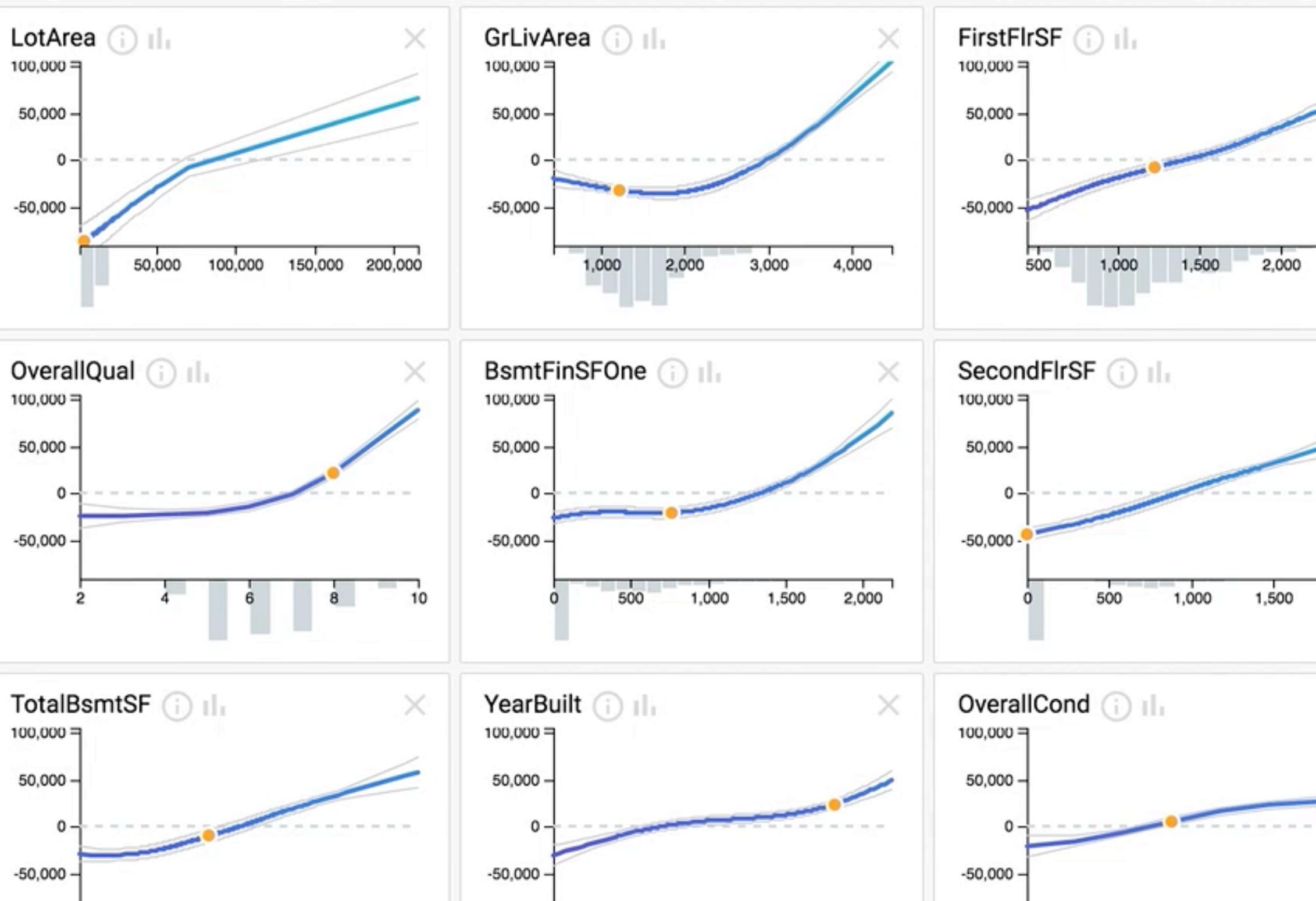
SORT BY NEIGH



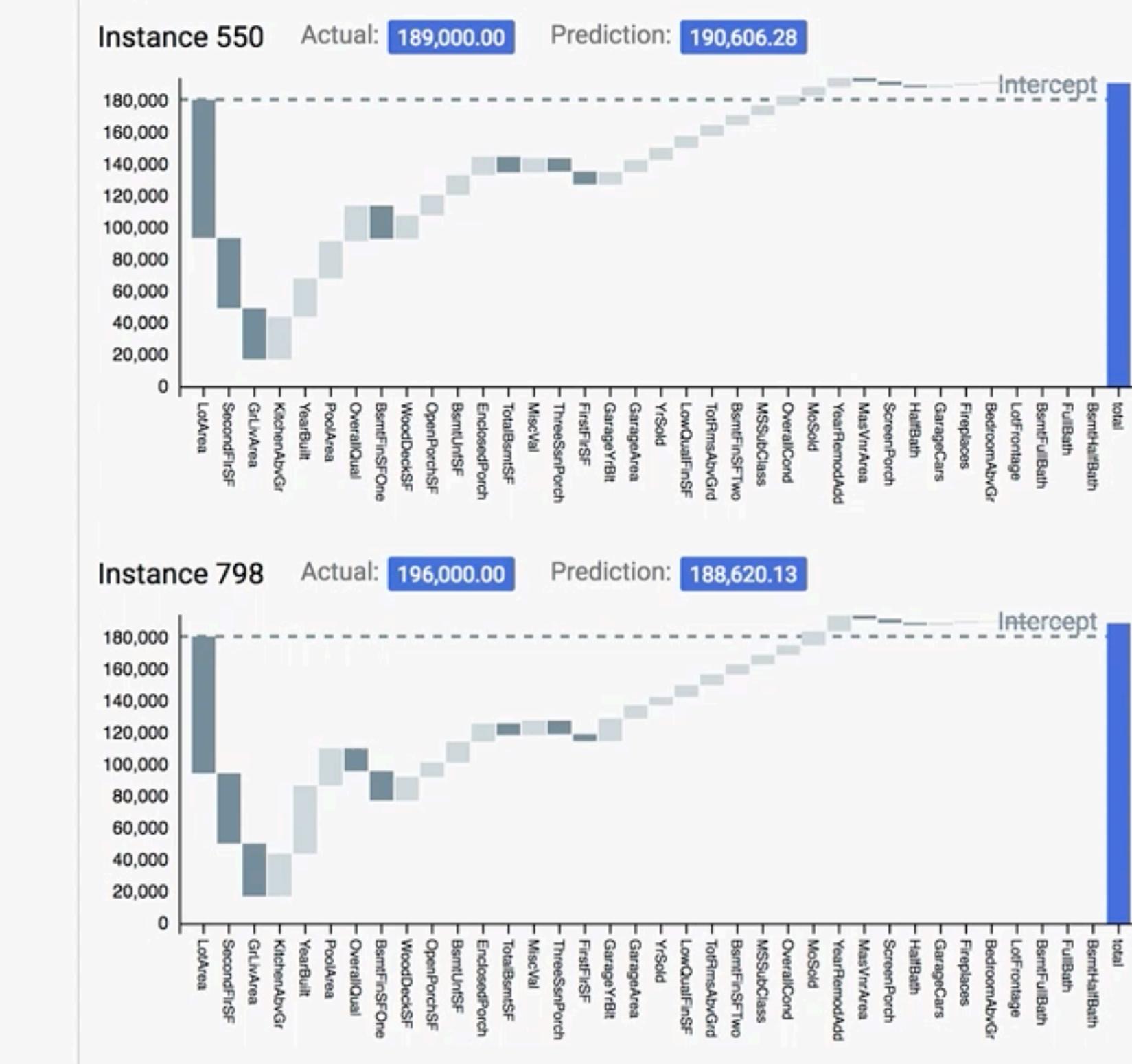
Normalize axes

 Hide all histograms

Hide zero line



## Sort waterfall linear



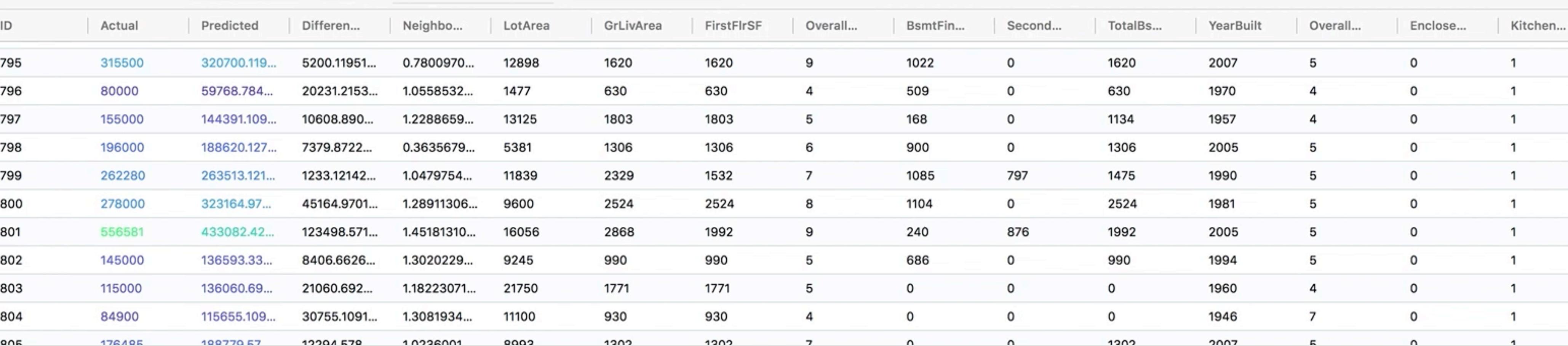
Showing 1119 of 1119

CLEAR FILTER



Nearest neighbor  
Feature space

SORT BY NEIGH



# User Study Findings Summary

## Reasons for Interpretability

*Why do data scientists need interpretability and how do they use it in GAMUT?*

12  data scientists, ~1.5 hours each

Think-aloud + question answering about multiple models  
Tutorial → Study → Interview

## Global v. Local Explanations

*How do data scientists use different explanation paradigms?*

## Interactive Explanations

*How does interactivity play a role in explainable machine learning interfaces?*

# User Study Findings Summary

## Reasons for Interpretability



Model debugging & improvement

Understanding data

Hypothesis generation for building trust

Communication

## Global v. Local Explanations



**Global**  
Familiar  
[3-5 years]

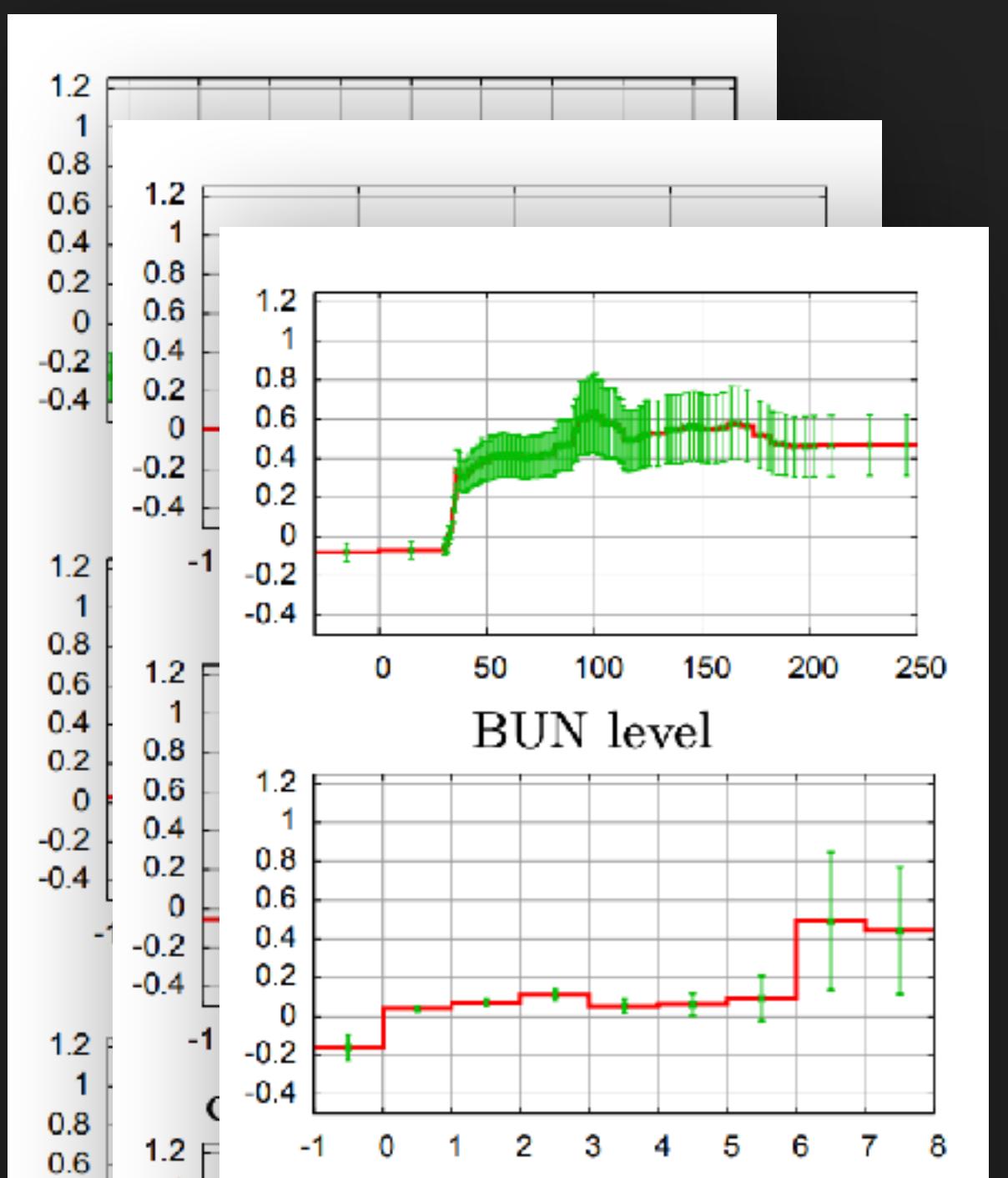
**Both**  
Experts  
[5+ years]

**Local**  
Novices  
[1-2 years]

## Interactive Explanations



Primary mechanism for explaining predictions



# User Study Findings Summary

## Reasons for Interpretability

**Interpretability is not a singular, rigid concept**

Operationalization helps people in practice

## Global v. Local Explanations

**Tailor explanations for specific audiences**

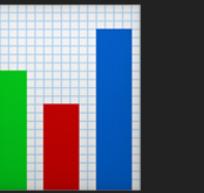
Balance simplicity and completeness

## Interactive Explanations

**Design and integrate effective interaction**

Interaction key to realizing interpretability

# Visualization



## Explanations

- 
- Show model context
  - Interactive analytics
  - Rely on user interpretation

# Visualization



Explanations

- Show model context
- Interactive analytics
- Rely on user interpretation

# Verbalization

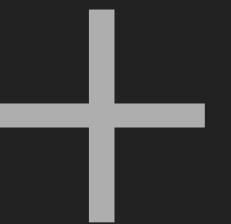
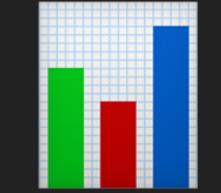


Explanations

- Direct and concise
- Less cognitive load
- No training needed

# Visualization + Verbalization

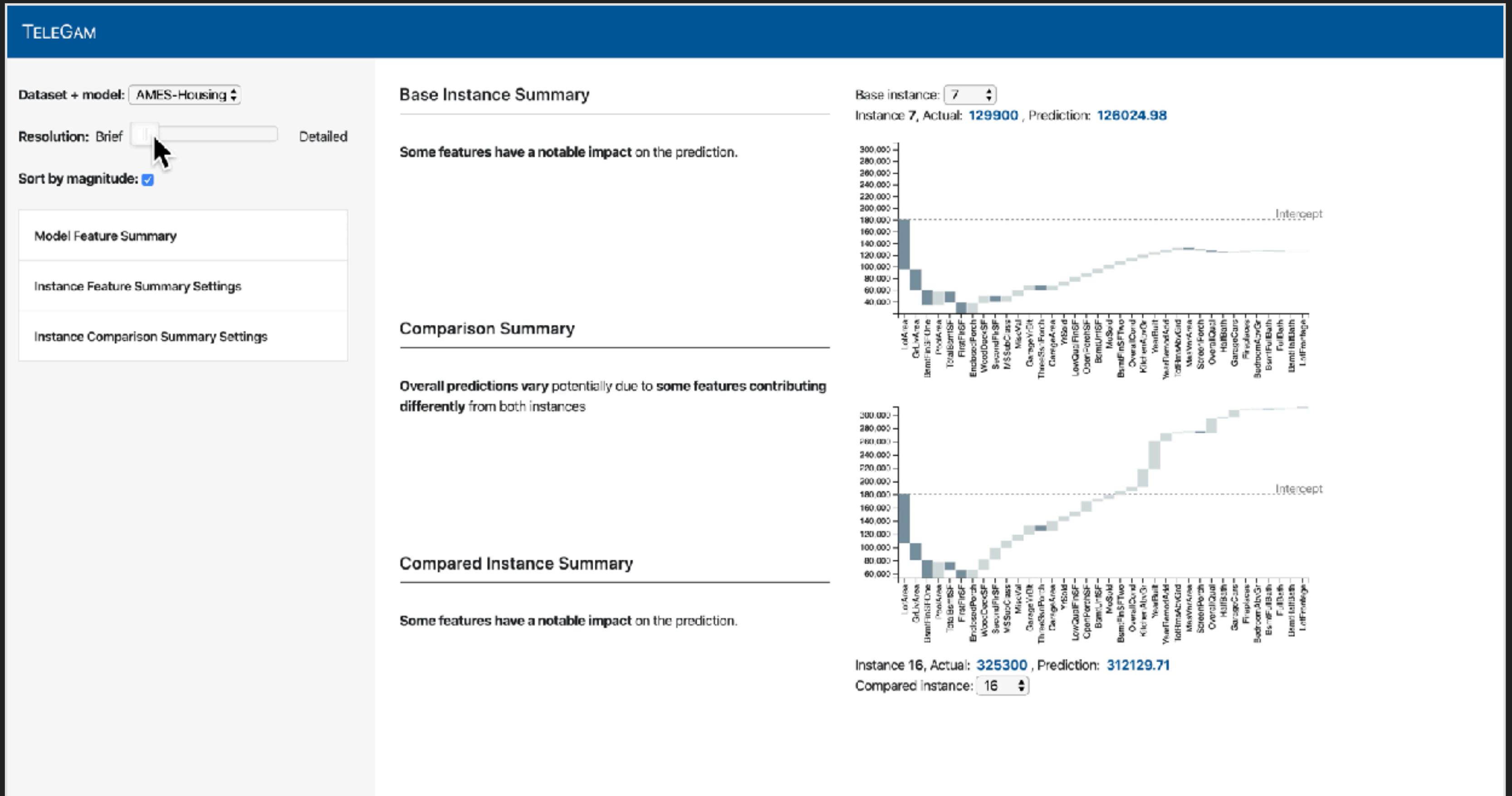
Explanations      Explanations



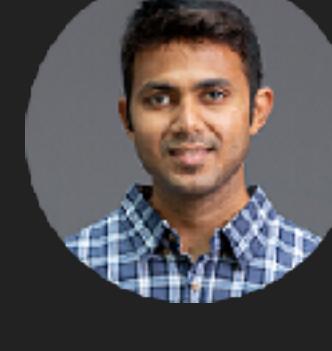
# TELEGAM

VIS 2019

## Combining Visualization & Verbalization for Machine Learning Interpretability



Fred Hohman  
Georgia Tech



Arjun Srinivasan  
Georgia Tech



Steven Drucker  
Microsoft Research

Interactively highlight verbalization in context of the visualization

Users decide resolution, balance simplicity & completeness

VIS 2019

# Combining Visualization & Verbalization for Machine Learning Interpretability

TELEGAM

Dataset + model: AMES-Housing ▾

Resolution: Brief  Detailed

Sort by magnitude:

**Model Feature Summary**

**Instance Feature Summary Settings**

**Instance Comparison Summary Settings**

**Base Instance Summary**

Some features have a notable impact on the prediction.

Base instance: 7 ▾

Instance 7, Actual: 129900 , Prediction: 128024.98

This partial dependence plot shows the relationship between the predicted value (y-axis, 40,000 to 300,000) and various features (x-axis). The features listed include LotArea, GrLivArea, BldgFinSFOne, TotalBsmtSF, TotalFinSF, FirstFlrSF, EnclosedPorch, WoodDeckSF, ScreenPorch, MSSubClass, MiscVal, GarageYrBlt, GarageArea, YrSold, LowQualFinSF, OpenPorchSF, BsmtUnfSF, MuSeu, BsmtFinSF, OverallCond, KitchenGr, YearBuilt, YearRemodAdd, TotalBsmtSF, OverallQual, OverallQual, FullBath, GarageCars, Fireplaces, Bedroms, BsmtFullBath, FullBath, DemRat, LstFrcBeg, and Intercept. The plot shows a significant positive correlation between most features and the prediction, with a notable dip around the 'TotalBsmtSF' feature.

**Comparison Summary**

Overall predictions vary potentially due to some features contributing differently from both instances

**Compared Instance Summary**

Some features have a notable impact on the prediction.

Instance 16, Actual: 325300 , Prediction: 312129.71

Compared instance: 16 ▾

This partial dependence plot shows the relationship between the predicted value (y-axis, 60,000 to 300,000) and various features (x-axis). The features listed are identical to the first plot: LotArea, GrLivArea, BldgFinSFOne, TotalBsmtSF, TotalFinSF, FirstFlrSF, EnclosedPorch, WoodDeckSF, ScreenPorch, MSSubClass, MiscVal, GarageYrBlt, GarageArea, YrSold, LowQualFinSF, OpenPorchSF, BsmtUnfSF, MuSeu, BsmtFinSF, OverallCond, KitchenGr, YearBuilt, YearRemodAdd, TotalBsmtSF, OverallQual, OverallQual, FullBath, GarageCars, Fireplaces, Bedroms, BsmtFullBath, FullBath, DemRat, LstFrcBeg, and Intercept. The plot shows a similar trend to the first one, with a significant positive correlation between most features and the prediction.



# Fred Hohman

## Georgia Tech



# Arjun Srinivasan

## Georgia Tech



# Steven Drucker

## Microsoft Research

# Interactively highlight verbalization in context of the visualization

Users decide resolution,  
balance **simplicity** &  
**completeness**

# Impact



**GAMUT**  
Deployed at Microsoft Research  
Demoed at for executive leadership

[msrgamut.microsoft.com](http://msrgamut.microsoft.com)

**INTERPRETML** 2.4k+ stars  
Influenced visualization design for popular interpretability toolkit  
[github.com/interpretml/interpret](https://github.com/interpretml/interpret)

**SANDDANCE** 4.2k+ stars  
GAMUT linked to visualization and data exploration tools  
[sanddance.js.org](https://sanddance.js.org)

# Interactive Scalable Interfaces for Machine Learning Interpretability

---



## PART I Enable interpretability

**GAMUT** Operationalize interpretability *CHI 2019*

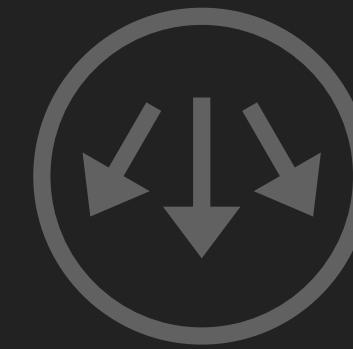
**TELEGAM** Vis + text for better explanations *VIS 2019*



## PART II Scale interpretability

**Interrogative Survey** Summarize interpretability vis *TVCG 2018*

**SUMMIT** Higher-level explanations for neural networks *VAST 2019*

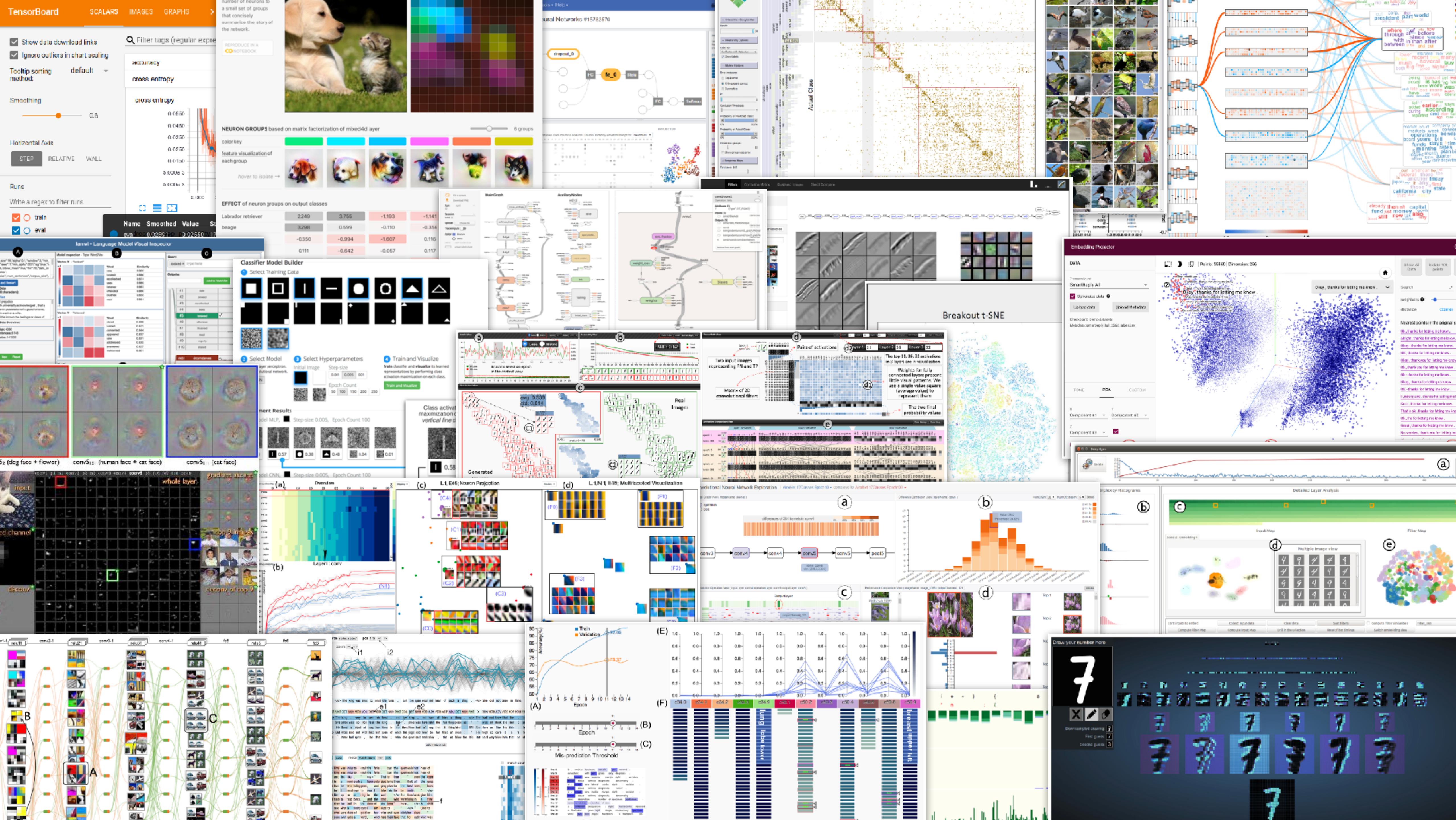


## PART III Communicate interpretability

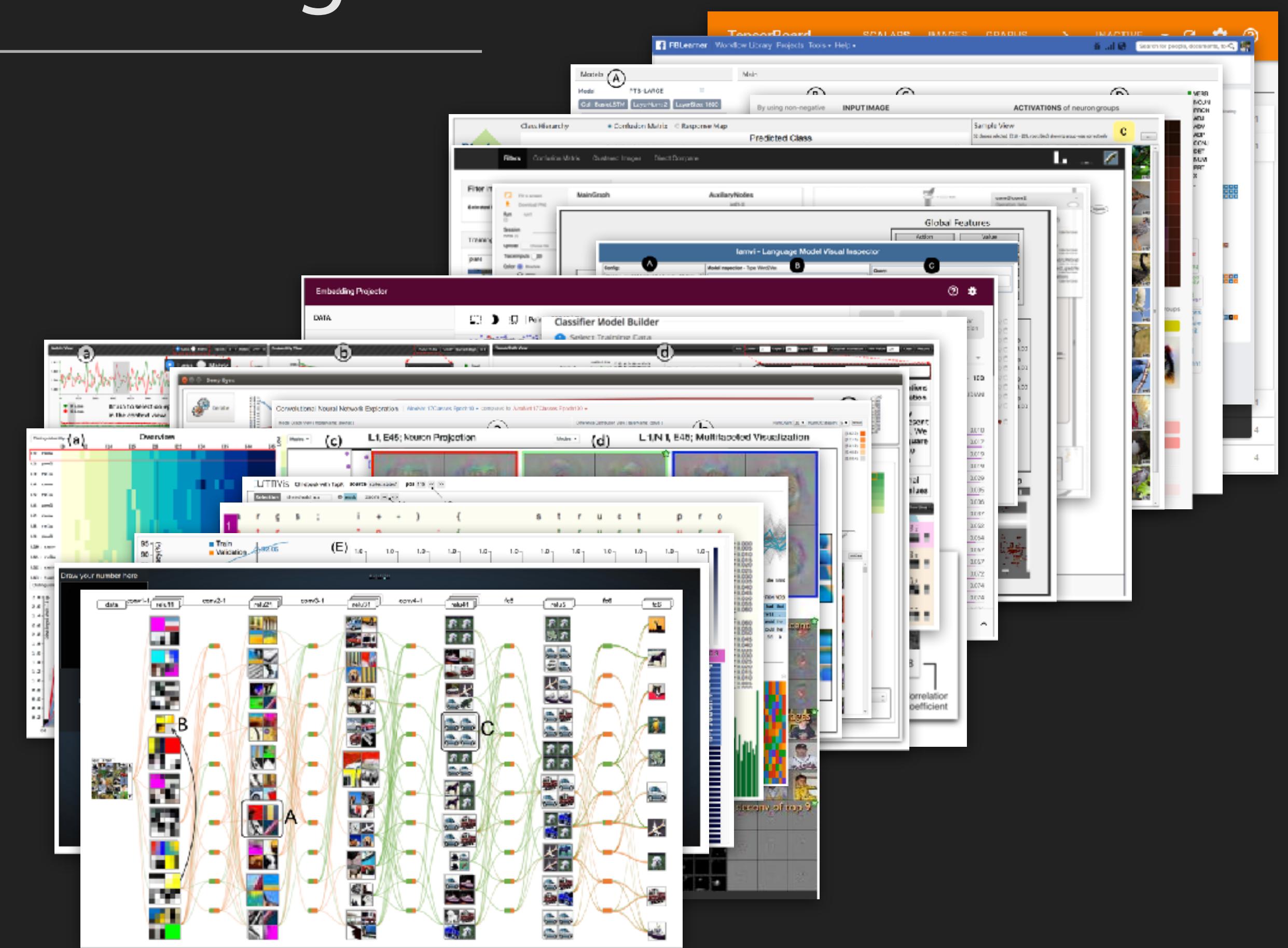
**ML Literacy** Interactive mediums & platforms *VISCOMM 2019, VISxAI 2018*

**Interactive Articles** Formalizing interactive communication *Distill 2020*





# Organize & Summarize Research for Visual Analytics in Deep Learning



# INTERROGATIVE SURVEY

For Visual Analytics in Deep Learning

*TVCG 2018*



**Fred Hohman**

Georgia Tech



**Minsuk Kahng**

Georgia Tech



**Robert Pienta**

Georgia Tech



**Polo Chau**

Georgia Tech

# Survey Highlights

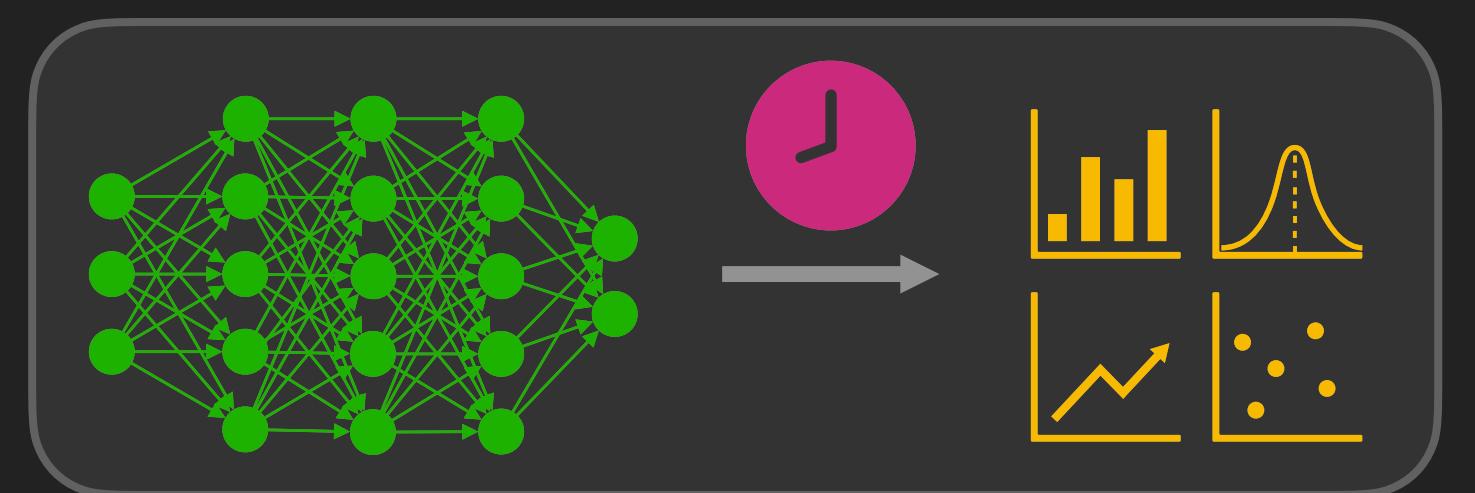
## WHY

*Why would one want to use visualization in deep learning?*



## WHAT

*What data, features, and relationships in deep learning can be visualized?*



## WHO

*Who would use and benefit from visualizing deep learning?*

## HOW

*How can we visualize deep learning data, features, and relationships?*

## WHEN

*When in the deep learning process is visualization used?*



## WHERE

*Where has deep learning visualization been used?*

Author	Year	WHY			WHO			WHAT			HOW			WHEN		WHERE					
		Interpretability & Explainability	Debugging & Improving Models	Comparing & Selecting Models	Education	Model Developers & Builders	Model Users	Non-experts	Learned Model Parameters	Individual Computational Units	Neurons in High-dimensional Space	Aggregated Information	Node-link Diagrams for Network Architecture	Dimensionality Reduction & Scatter Plots	Line Charts for Temporal Metrics	Instance-based Analysis & Exploration	Interactive Experimentation	Algorithms for Attribution & Feature Visualization	During Training	After Training	Publication Venue
Abadi, et al.	2016	■	■	■		■	■											■	■	■	arXiv
Bau, et al.	2017	■		■	■	■		■		■	■	■				■	■	■	■	■	CVPR
Bilal, et al.	2017	■	■	■		■	■	■		■	■	■				■	■	■	■	■	TVCg
Bojarski, et al.	2016	■	■	■		■	■	■													arXiv
Bruckner	2014	■	■			■	■	■		■	■	■				■	■	■	■	■	MS Thesis
Carter, et al.	2016	■			■	■	■	■	■	■	■	■				■	■	■	■	■	Distill
Cashman, et al.	2017	■	■			■	■	■		■	■	■				■	■	■	■	■	VADL
Chae, et al.	2017	■	■			■	■	■		■	■	■				■	■	■	■	■	VADL
Chung, et al.	2016	■	■			■	■	■		■	■	■				■	■	■	■	■	FILM
Goyal, et al.	2016	■				■	■	■		■	■	■				■	■	■	■	■	arXiv
Harley	2015	■				■	■	■		■	■	■				■	■	■	■	■	ISVC
Hohman, et al.	2017	■	■	■	■	■	■	■	■	■	■	■				■	■	■	■	■	CHI
Kahng, et al.	2018	■				■	■	■		■	■	■				■	■	■	■	■	TVCg
Karpathy, et al.	2015	■				■	■	■		■	■	■				■	■	■	■	■	arXiv
Li, et al.	2015	■				■	■	■		■	■	■				■	■	■	■	■	arXiv
Liu, et al.	2017	■	■			■	■	■		■	■	■				■	■	■	■	■	TVCg
Liu, et al.	2018	■	■			■	■	■		■	■	■				■	■	■	■	■	TVCg
Ming, et al.	2017	■				■	■	■		■	■	■				■	■	■	■	■	VAST
Norton & Qi	2017	■	■			■	■	■		■	■	■				■	■	■	■	■	VizSec
Olah	2014	■				■	■	■		■	■	■				■	■	■	■	■	Web
Olah, et al.	2018	■	■	■		■	■	■		■	■	■				■	■	■	■	■	Distill
Pezzotti, et al.	2017	■	■			■	■	■		■	■	■				■	■	■	■	■	TVCg
Rauber, et al.	2017	■	■			■	■	■		■	■	■				■	■	■	■	■	TVCg
Robinson, et al.	2017	■				■	■	■		■	■	■				■	■	■	■	■	GeoHum.
Rong, et al.	2016	■				■	■	■		■	■	■				■	■	■	■	■	ICML VIS
Smilkov, et al.	2016	■				■	■	■		■	■	■				■	■	■	■	■	NIPS Workshop
Smilkov, et al.	2017	■	■			■	■	■		■	■	■				■	■	■	■	■	ICML VIS
Strobell, et al.	2017	■	■			■	■	■		■	■	■				■	■	■	■	■	TVCg
Tzeng & Ma	2005	■				■	■	■		■	■	■				■	■	■	■	■	VIS
Wang, et al.	2018	■	■	■		■	■	■		■	■	■				■	■	■	■	■	TVCg
Webster, et al.	2017					■	■	■		■	■	■				■	■	■	■	■	Web
Wongsuphasawat, et al.	2018	■				■	■	■		■	■	■				■	■	■	■	■	TVCg
Yosinski, et al.	2015	■				■	■	■		■	■	■				■	■	■	■	■	ICML DL
Zahavy, et al.	2016	■	■			■	■	■		■	■	■				■	■	■	■	■	ICML
Zeiler, et al.	2014	■	■			■	■	■		■	■	■				■	■	■	■	■	ECCV
Zeng, et al.	2017	■	■			■	■	■		■	■	■				■	■	■	■	■	VADL
Zhong, et al.	2017	■	■			■	■	■		■	■	■				■	■	■	■	■	ICML VIS
Zhu, et al.	2016	■				■	■	■		■	■	■				■	■	■	■	■	ECCV

# Survey Highlights

Author	Year	WHY			WHO			WHAT			HOW			WHEN		WHERE						
		Interpretability & Explainability	Debugging & Improving Models	Comparing & Selecting Models	Education	Model Developers & Builders	Model Users	Non-experts	Computational Graph & Network Architecture	Learned Model Parameters	Individual Computational Units	Neurons in High-dimensional Space	Aggregated Information	Node-link Diagrams for Network Architecture	Dimensionality Reduction & Scatter Plots	Line Charts for Temporal Metrics	Instance-based Analysis & Exploration	Interactive Experimentation	Algorithms for Attribution & Feature Visualization	During Training	After Training	Publication Venue
Abadi, et al.	2016	■	■	■		■	■												■	■	■	arXiv
Bau, et al.	2017	■		■		■				■								■	■	■	CVPR	
Bilal, et al.	2017	■	■			■				■								■	■	■	TVCN	
Bojarski, et al.	2016	■	■			■												■	■	■	arXiv	
Bruckner	2014	■	■			■												■	■	■	MS Thesis	
Carter, et al.	2016	■		■	■	■	■	■										■	■	■	Distill	
Cashman, et al.	2017	■	■			■	■											■	■	■	VADL	
Chae, et al.	2017	■	■			■												■	■	■	VADL	
Chung, et al.	2016	■	■			■												■	■	■	FILM	
Goyal, et al.	2016	■				■												■	■	■	arXiv	
Harey	2015	■				■												■	■	■	ISVC	
Hohman, et al.	2017	■	■	■		■	■											■	■	■	CHI	
Kahng, et al.	2018	■	■			■	■											■	■	■	TVCN	
Karpathy, et al.	2015	■				■	■											■	■	■	arXiv	
Li, et al.	2015	■				■	■											■	■	■	arXiv	
Liu, et al.	2017	■	■			■												■	■	■	TVCN	
Liu, et al.	2018	■	■			■												■	■	■	TVCN	
Ming, et al.	2017	■		■		■												■	■	■	VAST	
Norton & Qi	2017	■	■	■		■	■	■										■	■	■	VizSec	
Olah	2014	■				■												■	■	■	Web	
Olah, et al.	2018	■				■	■	■	■									■	■	■	Distill	
Pezzotti, et al.	2017	■	■			■												■	■	■	TVCN	
Rauber, et al.	2017	■	■	■		■												■	■	■	TVCN	
Robinson, et al.	2017	■				■												■	■	■	GeoHum.	
Rong, et al.	2016	■	■			■	■											■	■	■	ICML VIS	
Smilkov, et al.	2016	■				■												■	■	■	NIPS Workshop	
Smilkov, et al.	2017	■	■	■		■												■	■	■	ICML VIS	
Strobell, et al.	2017	■	■			■	■											■	■	■	TVCN	
Tzeng & Ma	2005	■				■												■	■	■	VIS	
Wang, et al.	2018	■	■	■		■												■	■	■	TVCN	
Webster, et al.	2017					■												■	■	■	Web	
Wongsuphasawat, et al.	2018	■				■												■	■	■	TVCN	
Yosinski, et al.	2015	■				■												■	■	■	ICML DL	
Zahavy, et al.	2016	■	■			■												■	■	■	ICML	
Zeiler, et al.	2014	■	■			■												■	■	■	ECCV	
Zeng, et al.	2017	■	■			■												■	■	■	VADL	
Zhong, et al.	2017	■	■			■												■	■	■	ICML VIS	
Zhu, et al.	2016	■				■	■	■										■	■	■	ECCV	

**WHO**

30 / 38

# Survey Highlights

works designed for  
**model developers**

Author	Year	WHY			WHO		WHAT		HOW		WHEN		WHERE									
		Interpretability & Explainability	Debugging & Improving Models	Comparing & Selecting Models	Education	Model Developers & Builders	Model Users	Non-experts	Computational Graph & Network Architecture	Learned Model Parameters	Individual Computational Units	Neurons in High-dimensional Space	Aggregated Information	Node-link Diagrams for Network Architecture	Dimensionality Reduction & Scatter Plots	Line Charts for Temporal Metrics	Instance-based Analysis & Exploration	Interactive Experimentation	Algorithms for Attribution & Feature Visualization	During Training	After Training	Publication Venue
Abadi, et al.	2016	■	■	■																■	■	arXiv
Bau, et al.	2017	■		■															■	■	CVPR	
Bilal, et al.	2017	■	■																	■	■	TVCg
Bojarski, et al.	2016	■	■																	■	■	arXiv
Bruckner	2014	■	■																	■	■	MS Thesis
Carter, et al.	2016	■		■															■	■	Distill	
Cashman, et al.	2017	■	■																■	■	VADL	
Chae, et al.	2017	■	■																■	■	VADL	
Chung, et al.	2016	■	■																■	■	FILM	
Goyal, et al.	2016	■																	■	■	arXiv	
Harley	2015	■																	■	■	ISVC	
Hohman, et al.	2017	■	■	■															■	■	CHI	
Kahng, et al.	2018	■	■																■	■	TVCg	
Karpathy, et al.	2015	■																	■	■	arXiv	
Li, et al.	2015	■																	■	■	arXiv	
Liu, et al.	2017	■	■																■	■	TVCg	
Liu, et al.	2018	■	■																■	■	TVCg	
Ming, et al.	2017	■		■															■	■	VAST	
Norton & Qi	2017	■	■	■															■	■	VizSec	
Olah	2014	■		■														■	■	Web		
Olah, et al.	2018	■	■	■														■	■	Distill		
Pezzotti, et al.	2017	■	■															■	■	TVCG		
Rauber, et al.	2017	■	■	■														■	■	TVCG		
Robinson, et al.	2017	■																■	■	GeoHum.		
Rong, et al.	2016	■	■															■	■	ICML VIS		
Smilkov, et al.	2016	■																■	■	NIPS Workshop		
Smilkov, et al.	2017	■	■	■														■	■	ICML VIS		
Strobell, et al.	2017	■	■															■	■	TVCG		
Tzeng & Ma	2005	■																■	■	VIS		
Wang, et al.	2018	■	■	■														■	■	TVCG		
Webster, et al.	2017		■															■	■	Web		
Wongsuphasawat, et al.	2018		■															■	■	TVCG		
Yosinski, et al.	2015	■		■														■	■	ICML DL		
Zahavy, et al.	2016	■	■															■	■	ICML		
Zeiler, et al.	2014	■	■															■	■	ECCV		
Zeng, et al.	2017	■	■															■	■	VADL		
Zhong, et al.	2017	■	■															■	■	ICML VIS		
Zhu, et al.	2016	■																■	■	ECCV		

# Survey Highlights

**WHO**

30 / 38

works designed for  
**model developers**

11 / 38

works designed for  
**non-experts**

Author	Year	WHY			WHO		WHAT		HOW		WHEN		WHERE									
		Interpretability & Explainability	Debugging & Improving Models	Comparing & Selecting Models	Education	Model Developers & Builders	Model Users	Non-experts	Computational Graph & Network Architecture	Learned Model Parameters	Individual Computational Units	Neurons in High-dimensional Space	Aggregated Information	Node-link Diagrams for Network Architecture	Dimensionality Reduction & Scatter Plots	Line Charts for Temporal Metrics	Instance-based Analysis & Exploration	Interactive Experimentation	Algorithms for Attribution & Feature Visualization	During Training	After Training	Publication Venue
Abadi, et al.	2016	■	■	■																■	■	arXiv
Bau, et al.	2017	■		■															■	■	CVPR	
Bilal, et al.	2017	■	■	■															■	■	TVCN	
Bojarski, et al.	2016	■	■	■															■	■	arXiv	
Bruckner	2014	■	■																■	■	MS Thesis	
Carter, et al.	2016	■		■														■	■	Distill		
Cashman, et al.	2017	■	■															■	■	VADL		
Chae, et al.	2017	■	■	■														■	■	VADL		
Chung, et al.	2016	■	■															■	■	FILM		
Goyal, et al.	2016	■																■	■	arXiv		
Harley	2015	■																■	■	ISVC		
Hohman, et al.	2017	■	■	■														■	■	CHI		
Kahng, et al.	2018	■	■															■	■	TVCN		
Karpathy, et al.	2015	■																■	■	arXiv		
Li, et al.	2015	■																■	■	arXiv		
Liu, et al.	2017	■	■															■	■	TVCN		
Liu, et al.	2018	■	■															■	■	TVCN		
Ming, et al.	2017	■																■	■	VAST		
Norton & Qi	2017	■	■	■														■	■	VizSec		
Olah	2014	■																■	■	Web		
Olah, et al.	2018	■	■	■														■	■	Distill		
Pezzotti, et al.	2017	■	■															■	■	TVCN		
Rauber, et al.	2017	■	■	■														■	■	TVCN		
Robinson, et al.	2017	■																■	■	GeoHum.		
Rong, et al.	2016	■	■															■	■	ICML VIS		
Smilkov, et al.	2016	■																■	■	NIPS Workshop		
Smilkov, et al.	2017	■	■	■														■	■	ICML VIS		
Strobell, et al.	2017	■	■															■	■	TVCN		
Tzeng & Ma	2005	■																■	■	VIS		
Wang, et al.	2018	■	■	■														■	■	TVCN		
Webster, et al.	2017					■												■	■	Web		
Wongsuphasawat, et al.	2018		■															■	■	TVCN		
Yosinski, et al.	2015	■																■	■	ICML DL		
Zahavy, et al.	2016	■	■															■	■	ICML		
Zeiler, et al.	2014	■	■															■	■	ECCV		
Zeng, et al.	2017	■	■															■	■	VADL		
Zhong, et al.	2017	■	■															■	■	ICML VIS		
Zhu, et al.	2016	■																■	■	ECCV		

# Survey Highlights

**WHO**

30 / 38

works designed for  
**model developers**

**HOW**

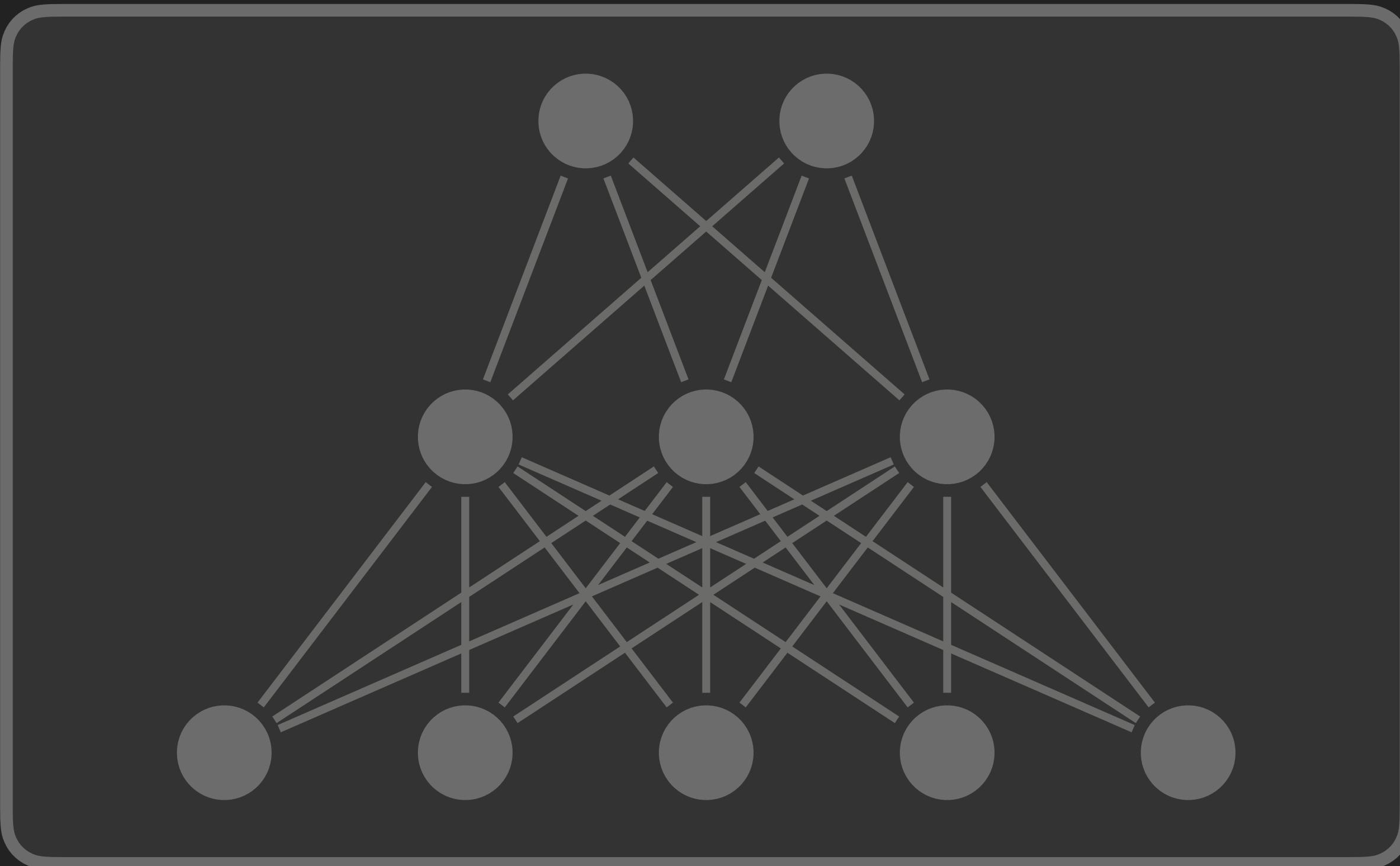
11 / 38

works designed for  
**non-experts**

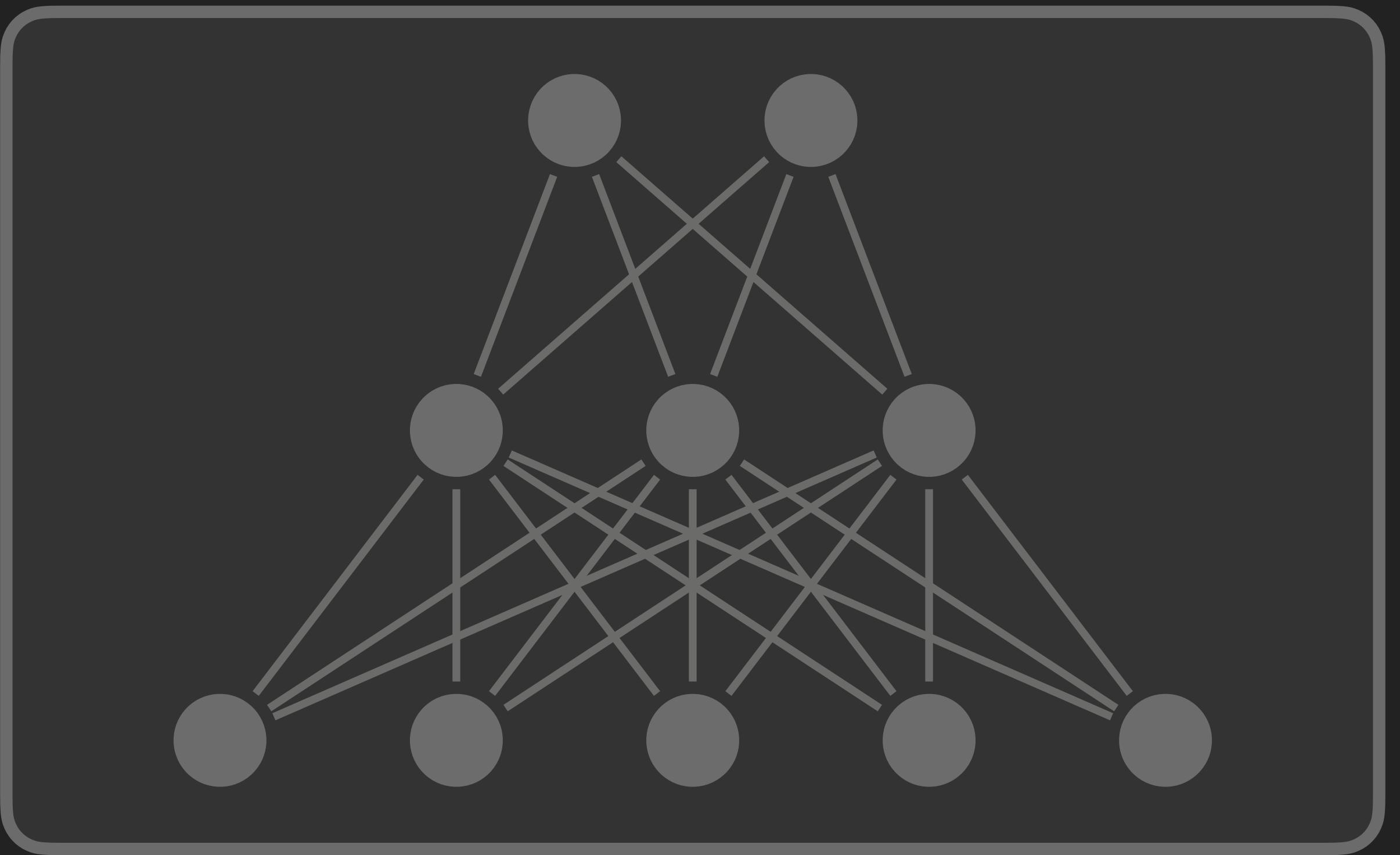
33 / 38

works use  
**instance-based analysis**

Author	Year	WHY		WHO		WHAT		HOW		WHEN		WHERE								
		Interpretability & Explainability	Debugging & Improving Models	Comparing & Selecting Models	Education	Model Developers & Builders	Model Users	Learned Model Parameters	Individual Computational Units	Neurons in High-dimensional Space	Aggregated Information	Line Charts for Temporal Metrics	Node-link Diagrams for Network Architecture	Dimensionality Reduction & Scatter Plots	Instance-based Analysis & Exploration	Interactive Experimentation	Algorithms for Attribution & Feature Visualization	During Training	After Training	Publication Venue
Abadi, et al.	2016	■	■	■		■											■	■	■	arXiv
Bau, et al.	2017	■		■													■	■	■	CVPR
Bilal, et al.	2017	■	■	■													■	■	■	TVCG
Bojarski, et al.	2016	■	■	■													■	■	■	arXiv
Bruckner	2014	■	■														■	■	■	MS Thesis
Carter, et al.	2016	■		■													■	■	■	Distill
Cashman, et al.	2017	■	■														■	■	■	VADL
Chae, et al.	2017	■	■	■													■	■	■	VADL
Chung, et al.	2016	■	■														■	■	■	FILM
Goyal, et al.	2016	■															■	■	■	arXiv
Harley	2015	■															■	■	■	ISVC
Hohman, et al.	2017	■	■	■													■	■	■	CHI
Kahng, et al.	2018	■	■														■	■	■	TVCG
Karpathy, et al.	2015	■															■	■	■	arXiv
Li, et al.	2015	■															■	■	■	arXiv
Liu, et al.	2017	■	■														■	■	■	TVCG
Liu, et al.	2018	■	■														■	■	■	TVCG
Ming, et al.	2017	■															■	■	■	VAST
Norton & Qi	2017	■	■	■													■	■	■	VizSec
Olah	2014	■															■	■	■	Web
Olah, et al.	2018	■	■	■													■	■	■	Distill
Pezzotti, et al.	2017	■	■														■	■	■	TVCG
Rauber, et al.	2017	■	■														■	■	■	TVCG
Robinson, et al.	2017	■															■	■	■	GeoHum.
Rong, et al.	2016	■	■														■	■	■	ICML VIS
Smilkov, et al.	2016	■															■	■	■	NIPS Workshop
Smilkov, et al.	2017	■	■	■													■	■	■	ICML VIS
Strobell, et al.	2017	■	■														■	■	■	TVCG
Tzeng & Ma	2005	■															■	■	■	VIS
Wang, et al.	2018	■	■	■													■	■	■	TVCG
Webster, et al.	2017					■											■	■	■	Web
Wongsuphasawat, et al.	2018	■				■											■	■	■	TVCG
Yosinski, et al.	2015	■				■											■	■	■	ICML DL
Zahavy, et al.	2016	■	■			■											■	■	■	ICML
Zeiler, et al.	2014	■	■			■											■	■	■	ECCV
Zeng, et al.	2017	■	■			■											■	■	■	VADL
Zhong, et al.	2017	■	■			■											■	■	■	ICML VIS
Zhu, et al.	2016	■				■											■	■	■	ECCV

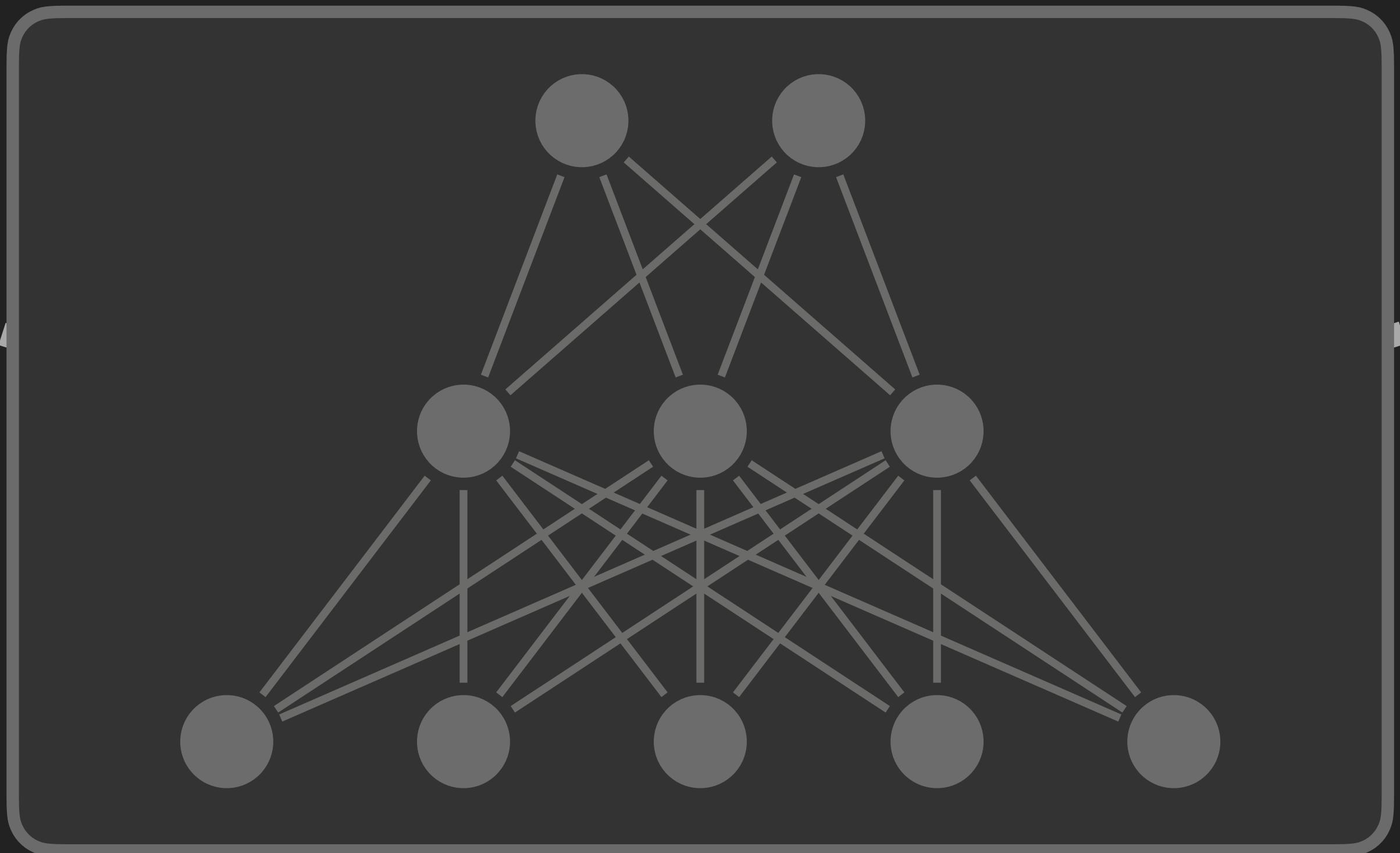


Neural Network



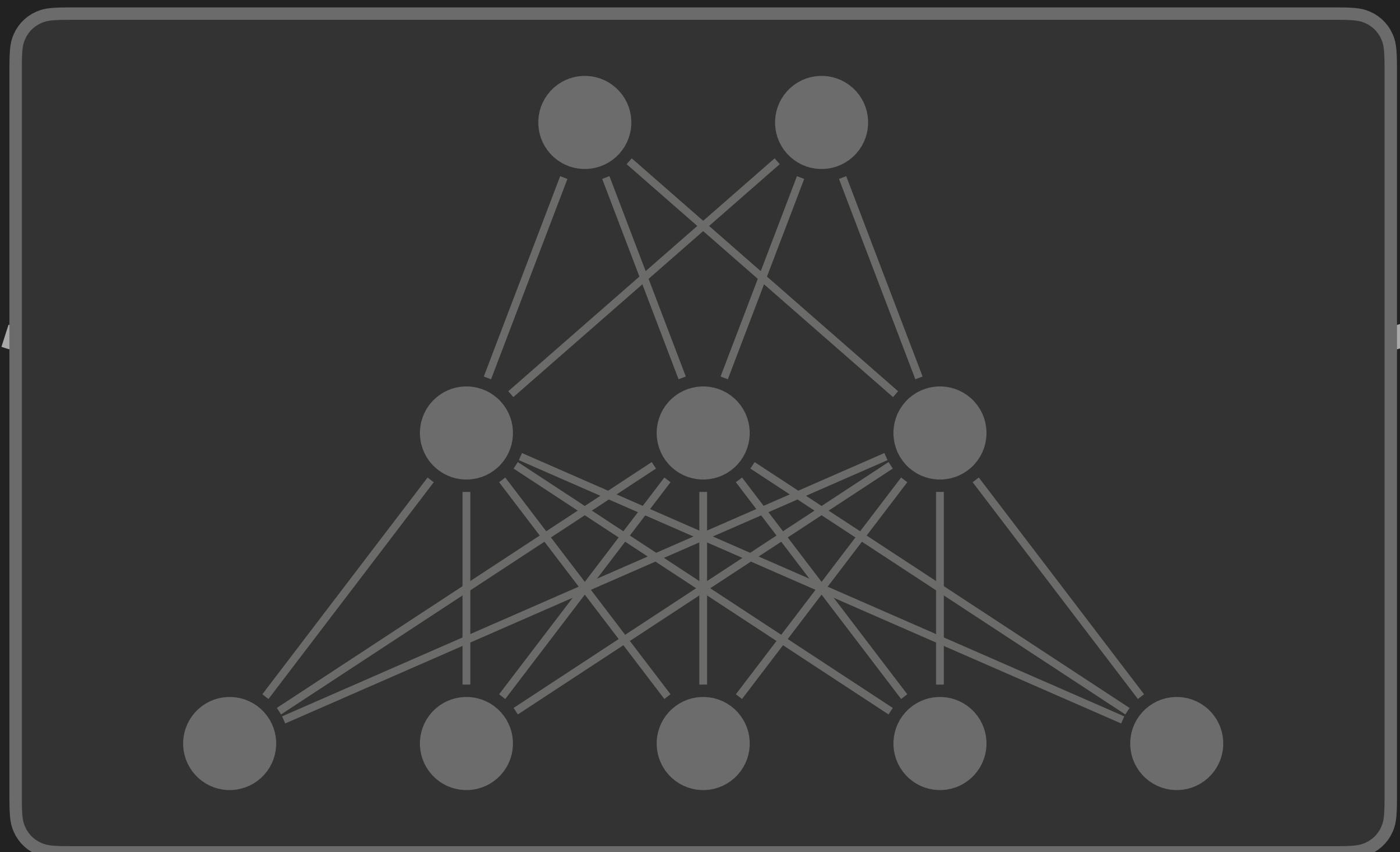
Neural Network

bike ✓



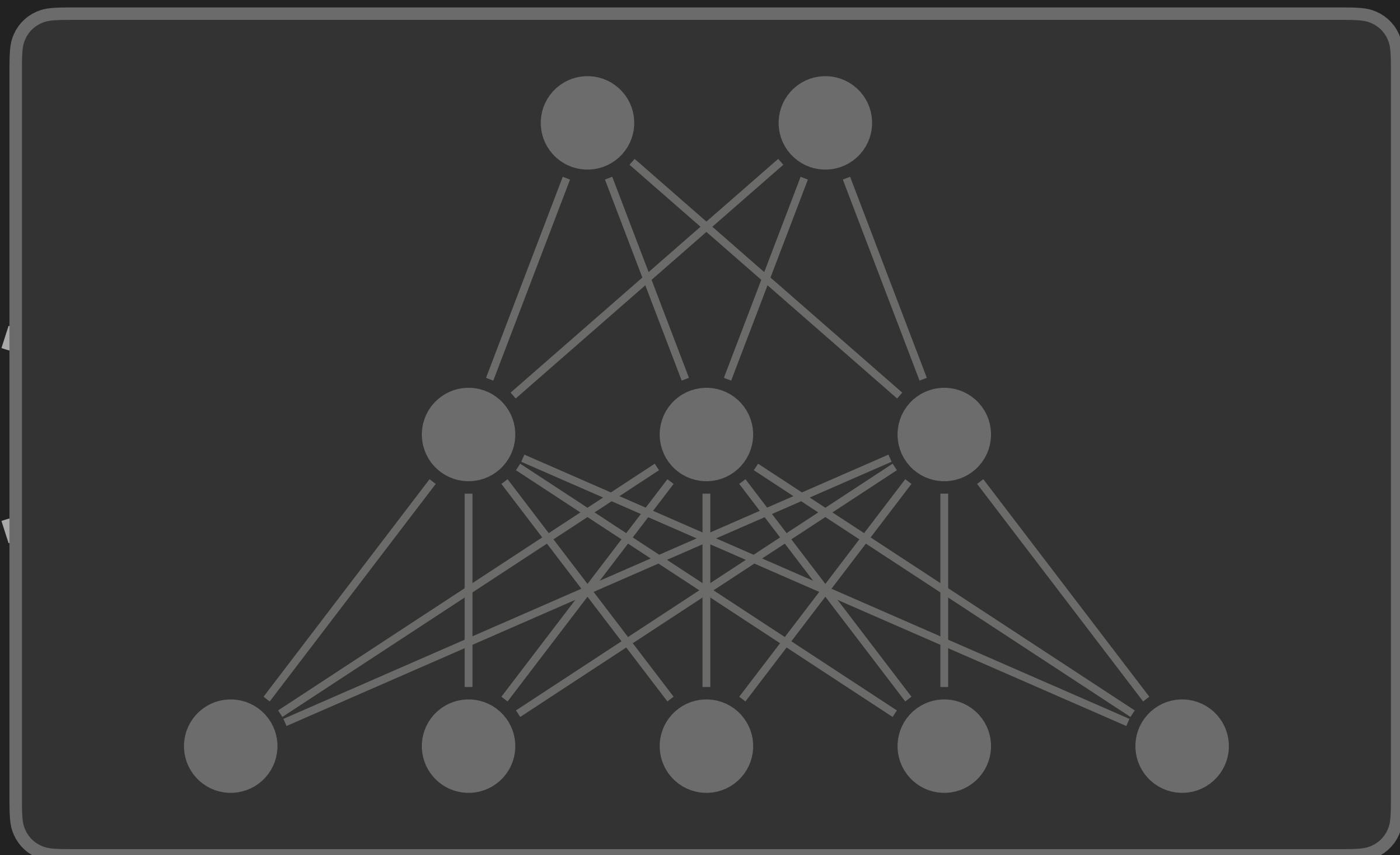
Neural Network

bike ✓



Neural Network

bike ✓



Neural Network

truck ✗

?

truck X



?

truck X

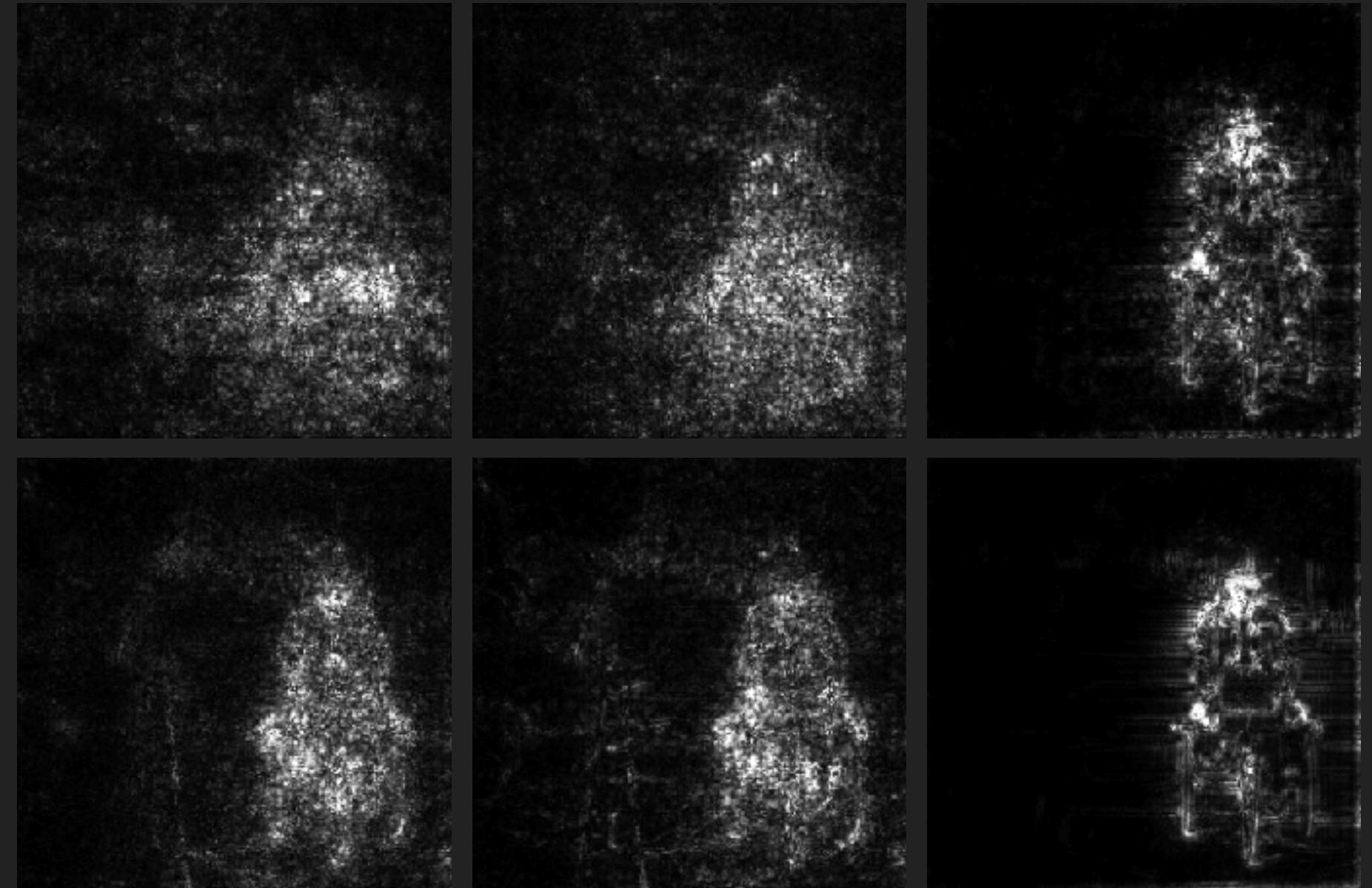


Attention

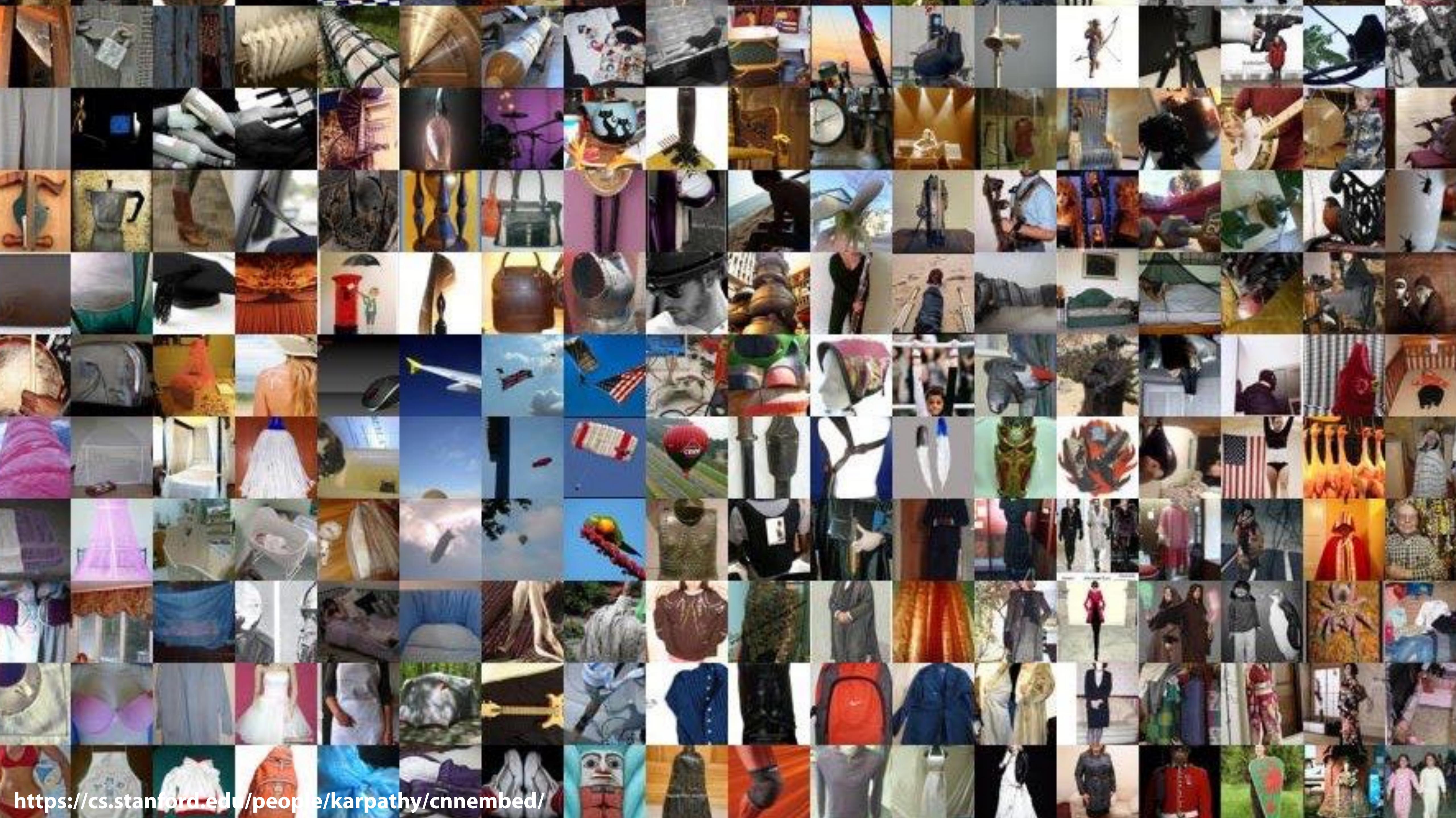


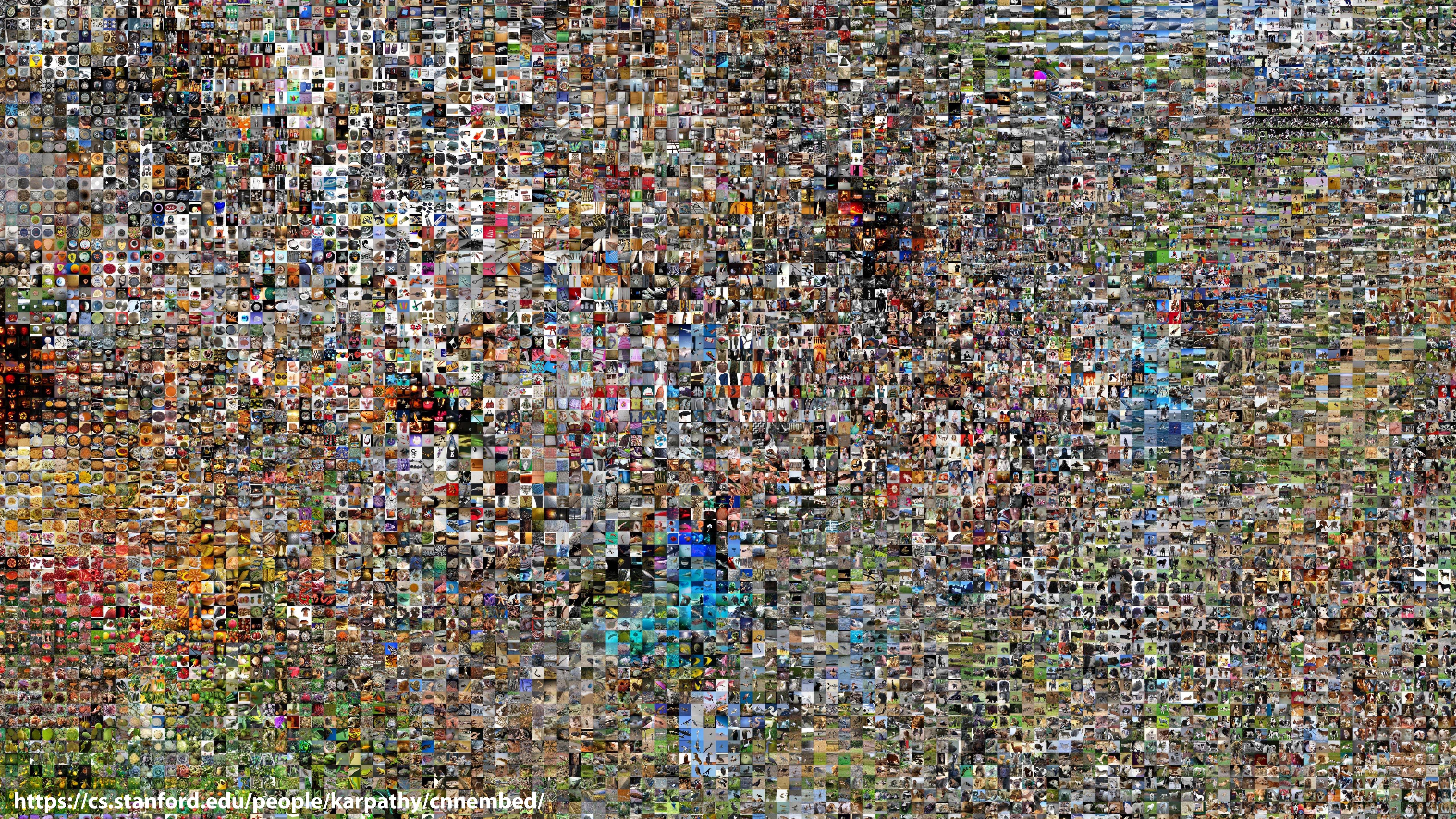
[Selvaraju, et al., ICCV, 2017]

Sensitivity



[Smilkov, et al., arXiv, 2017]





<https://cs.stanford.edu/people/karpathy/cnnembed/>

# SUMMIT

Scaling Deep Learning Interpretability by  
Visualizing Activation & Attribution Summarizations

*VAST 2019*



**Fred Hohman**

Georgia Tech



**Haekyu Park**

Georgia Tech



**Caleb Robinson**

Georgia Tech



**Polo Chau**

Georgia Tech

# SUMMIT

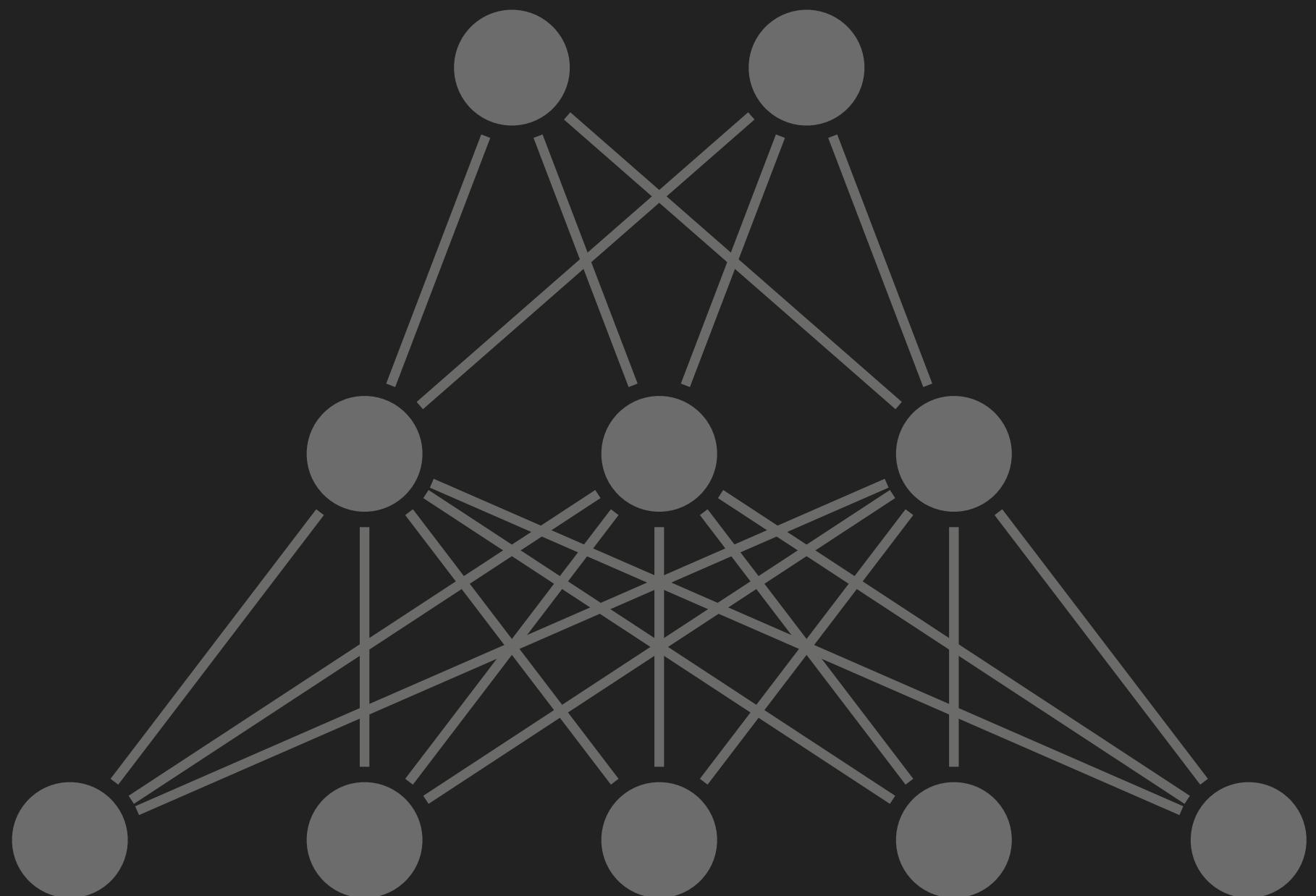
**Scalably summarize and interactively visualize**  
neural network feature representations  
for millions of images

# SUMMIT

**Scalably summarize and interactively visualize**  
neural network feature representations  
for millions of images



*white wolf*

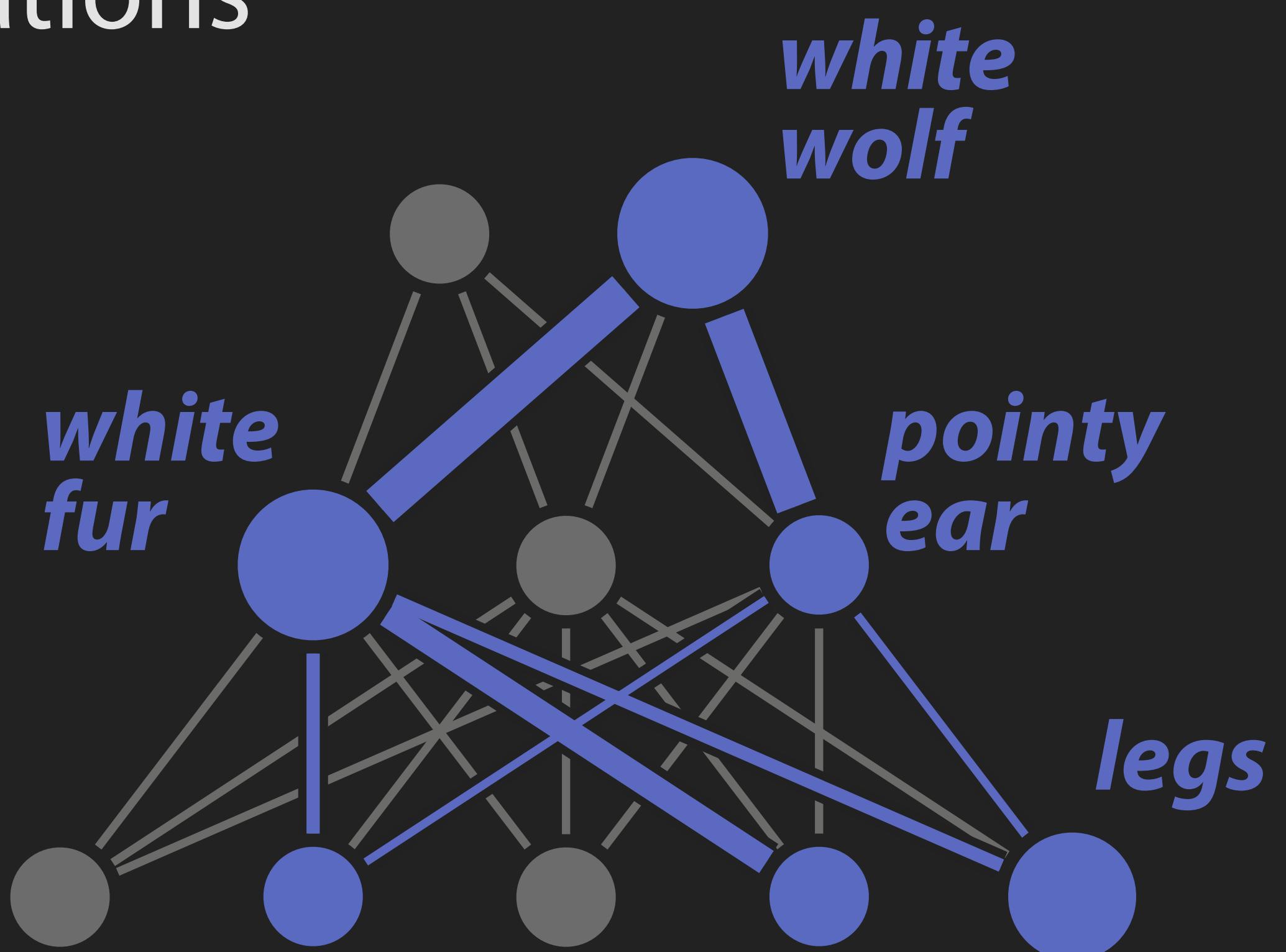


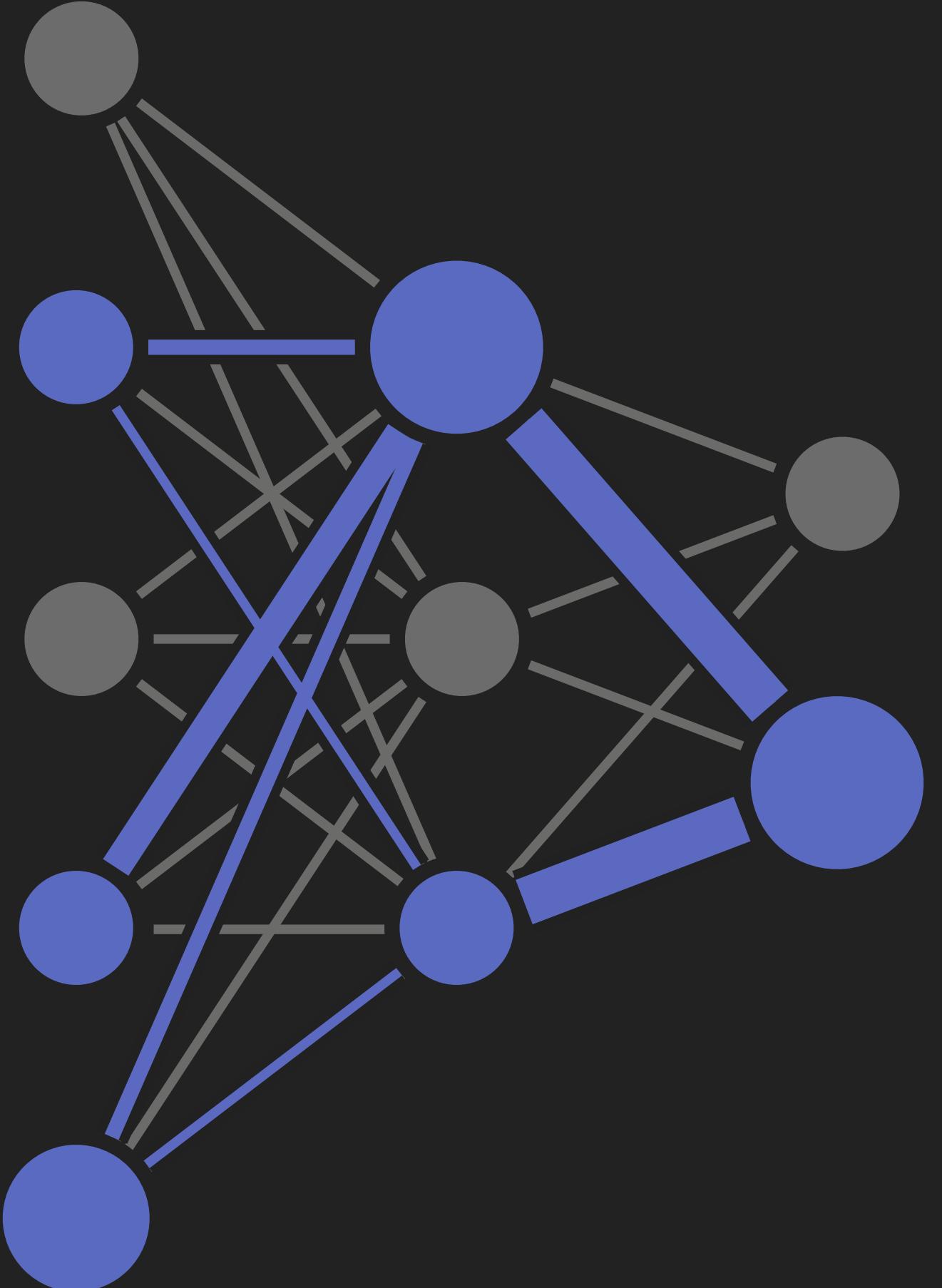
# SUMMIT

**Scalably summarize and interactively visualize**  
neural network feature representations  
for millions of images

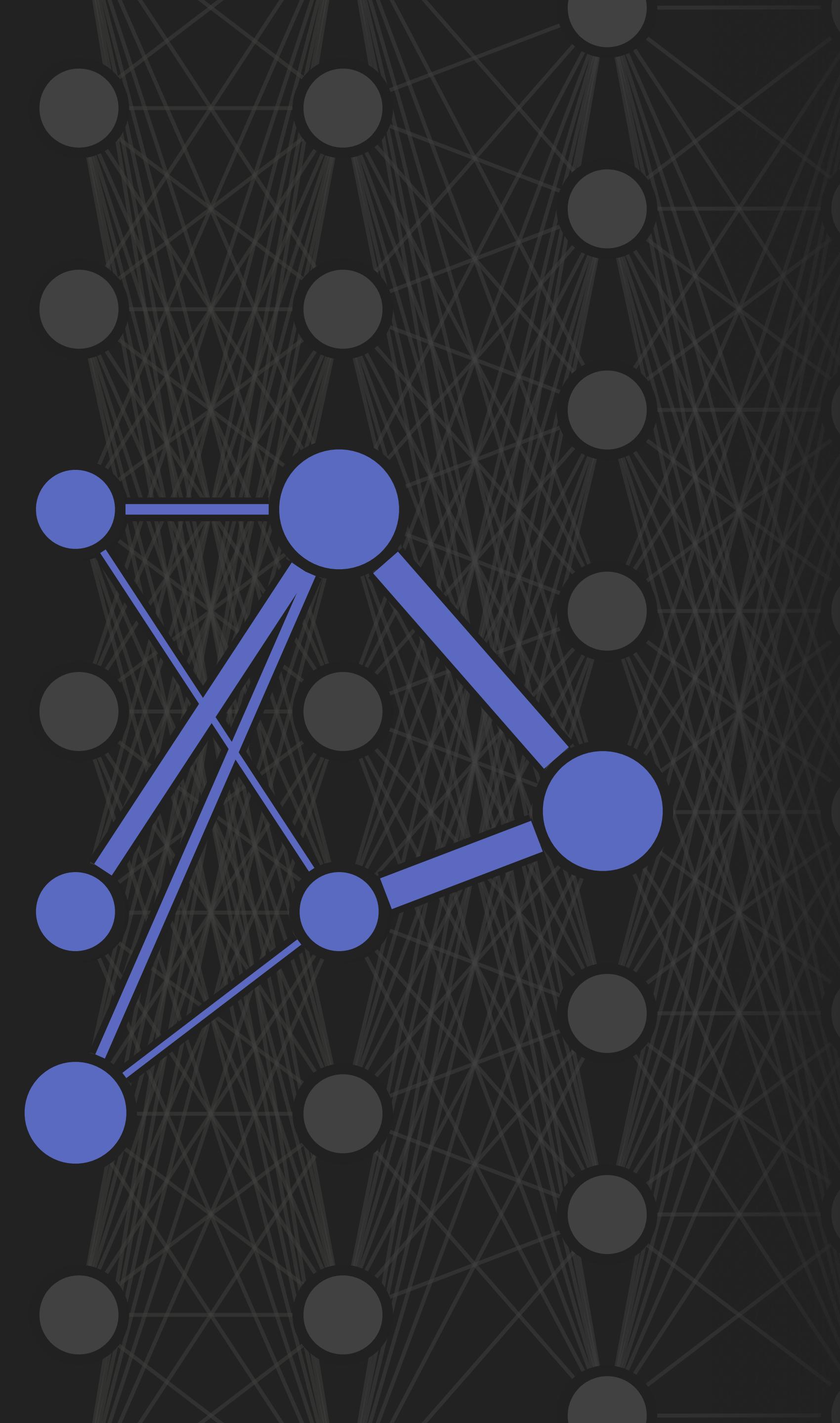


***white wolf***





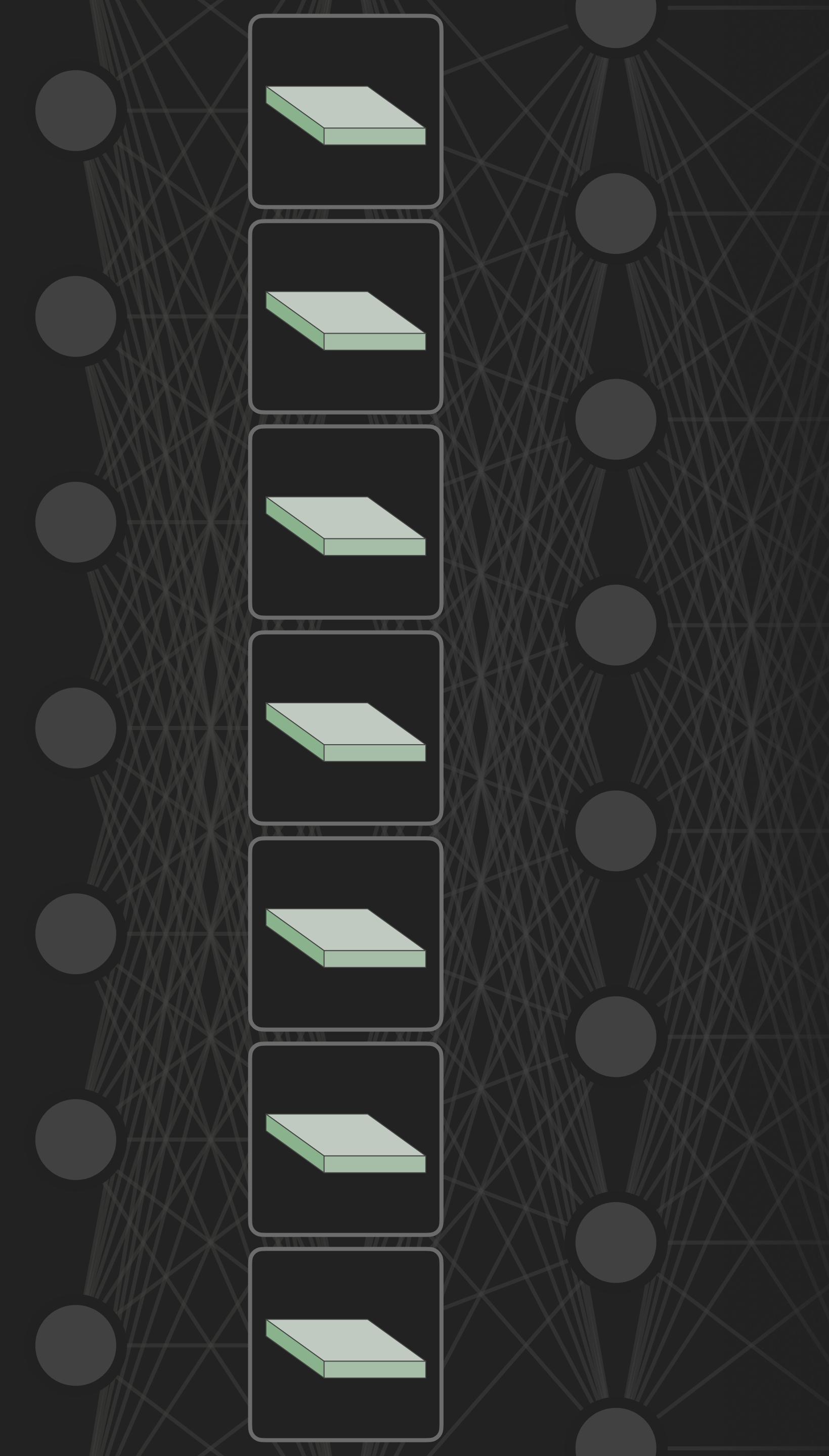
How do we make  
**attribution graphs?**



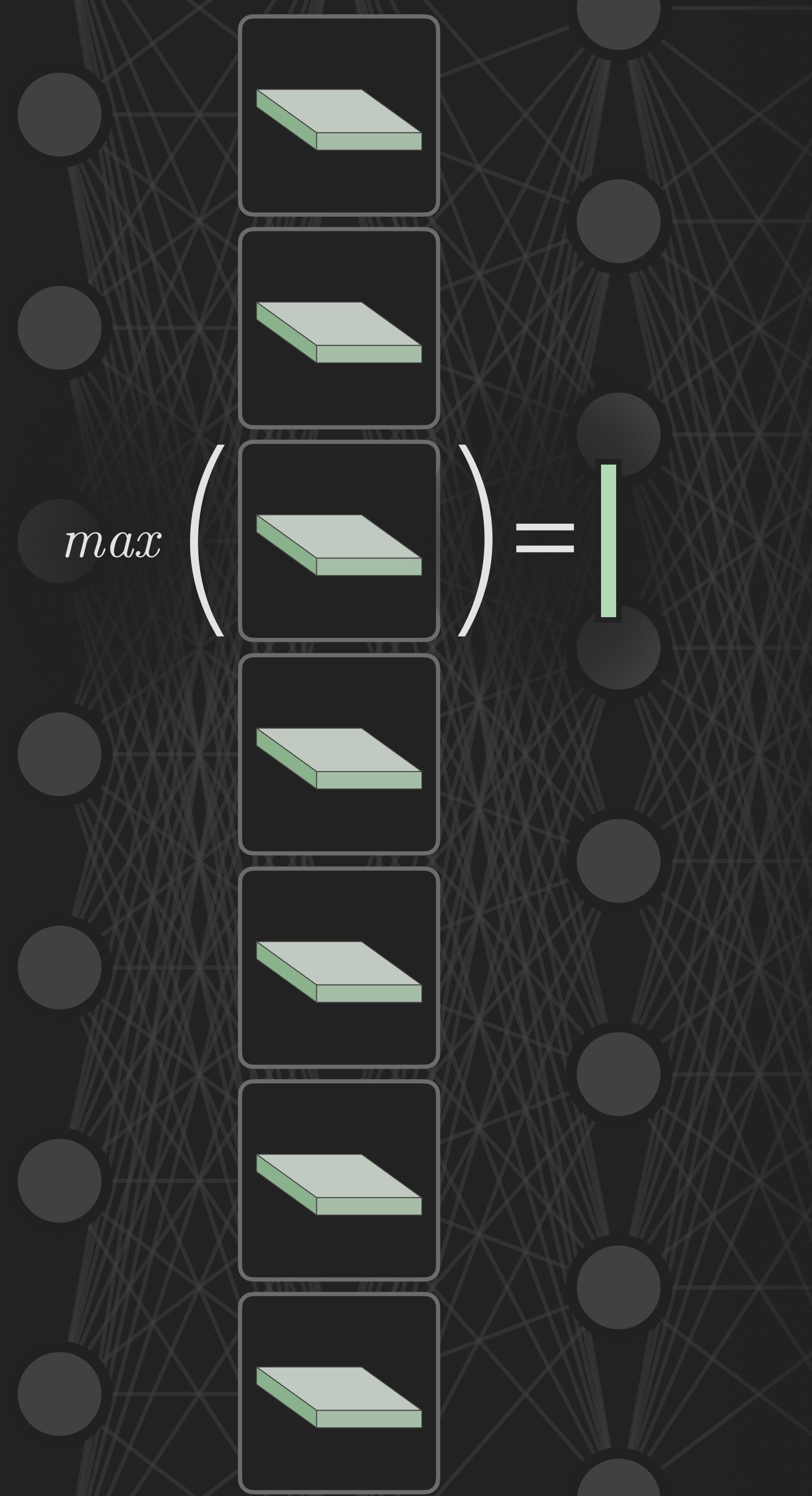
How do we make  
**attribution graphs?**



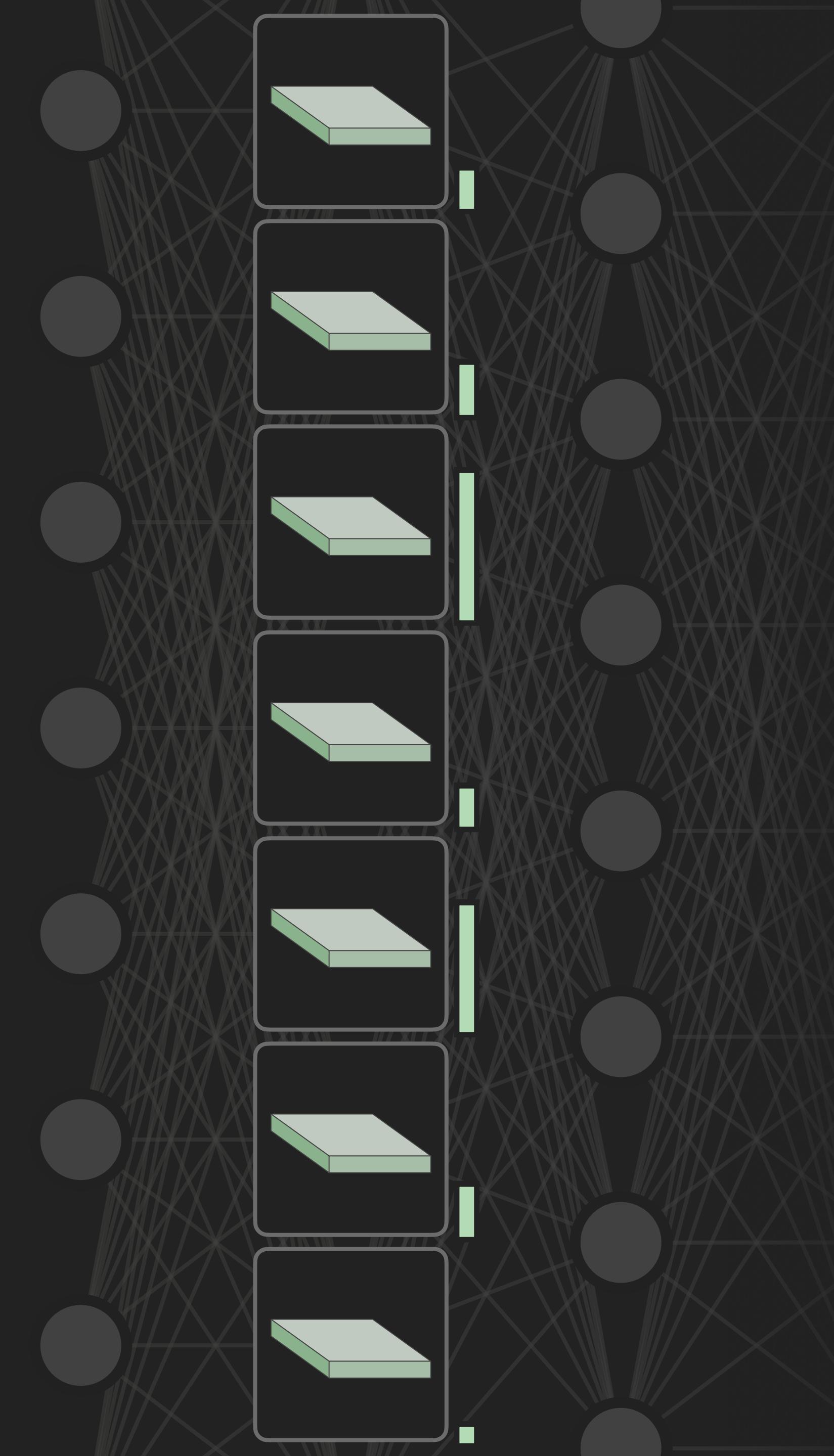
How do we make  
**attribution graphs?**



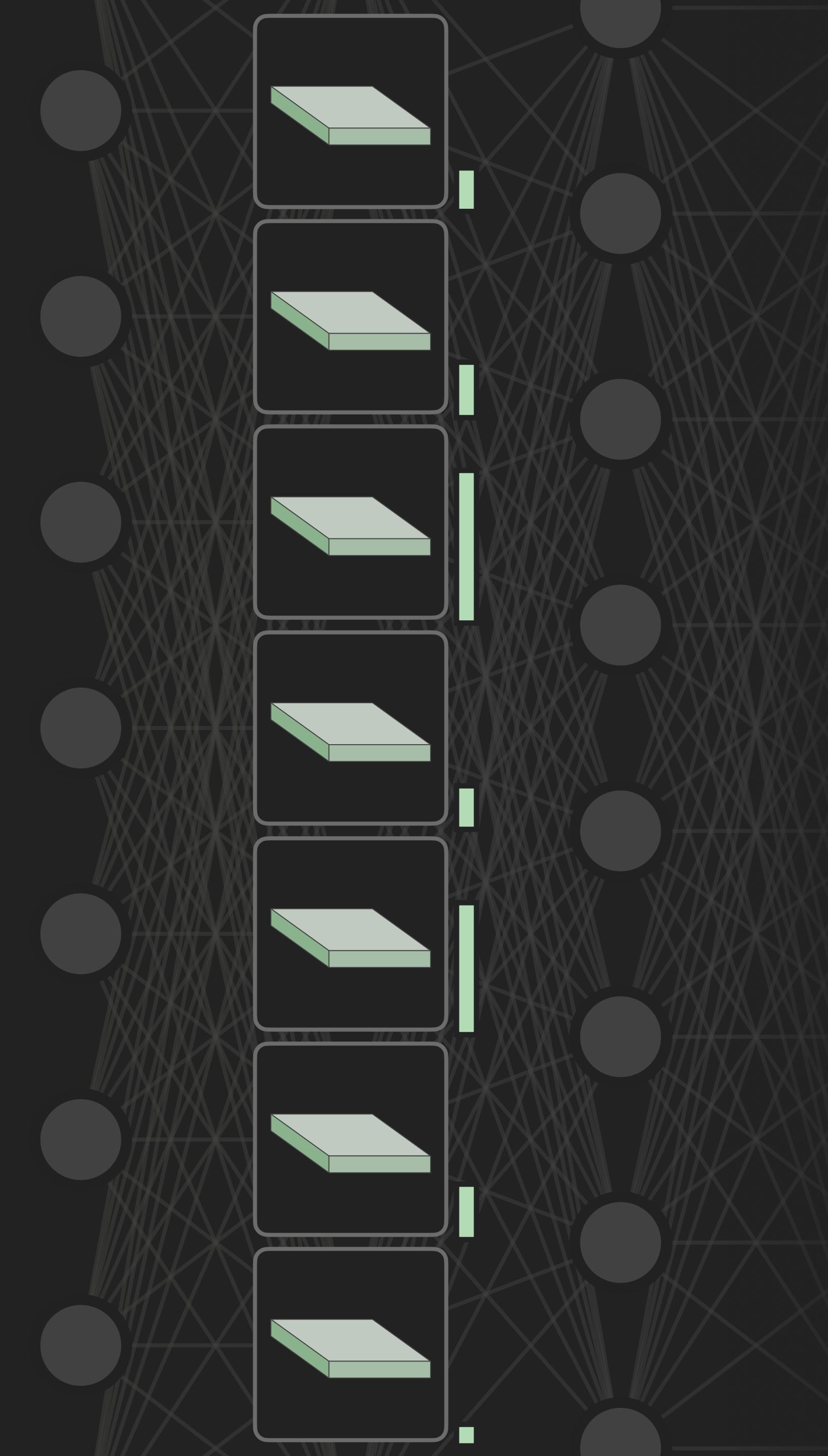
Aggregate network  
**activations** (nodes)



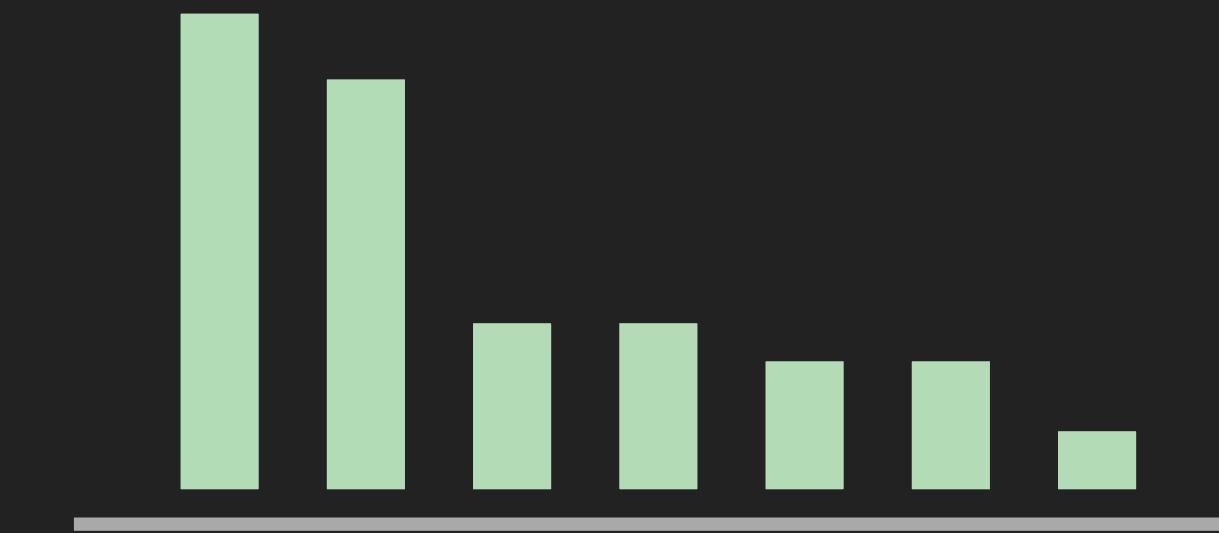
Aggregate network  
**activations** (nodes)

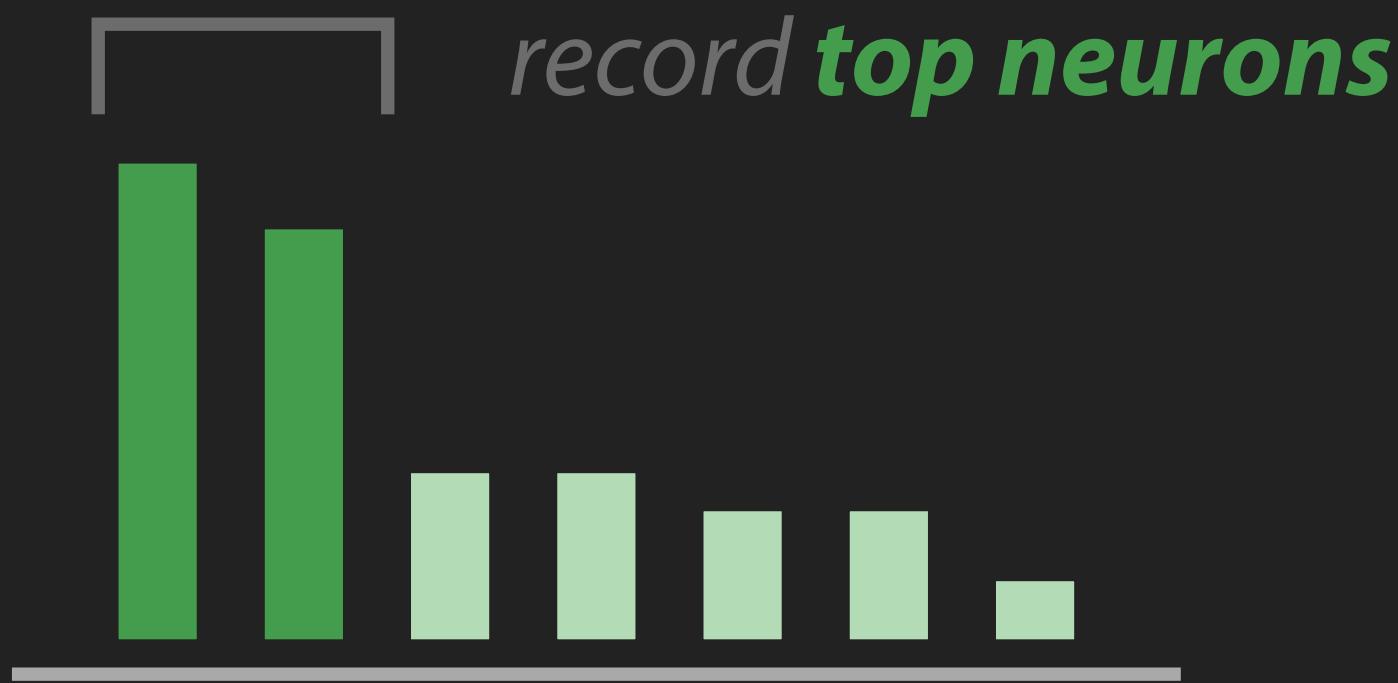
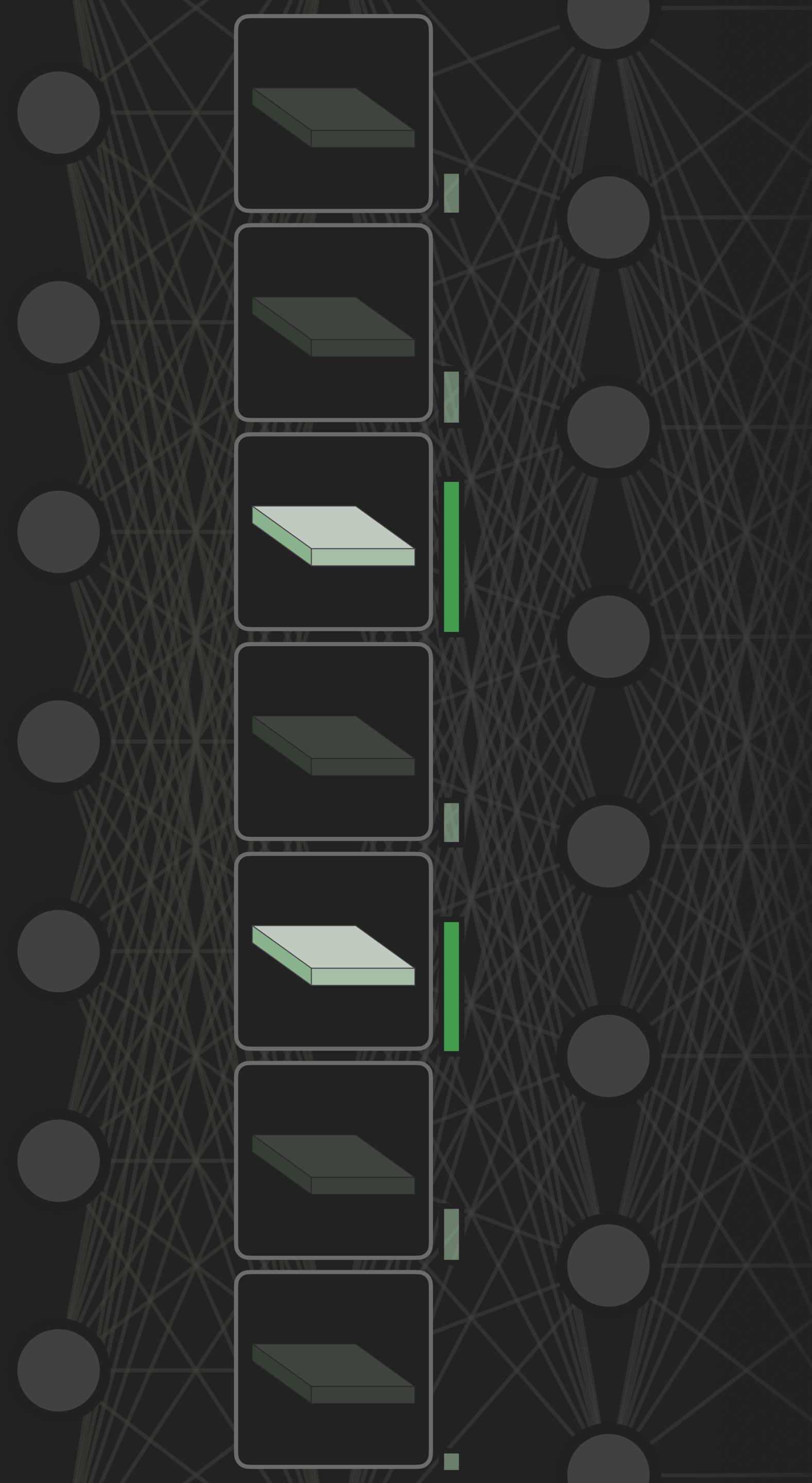


Aggregate network  
**activations** (nodes)

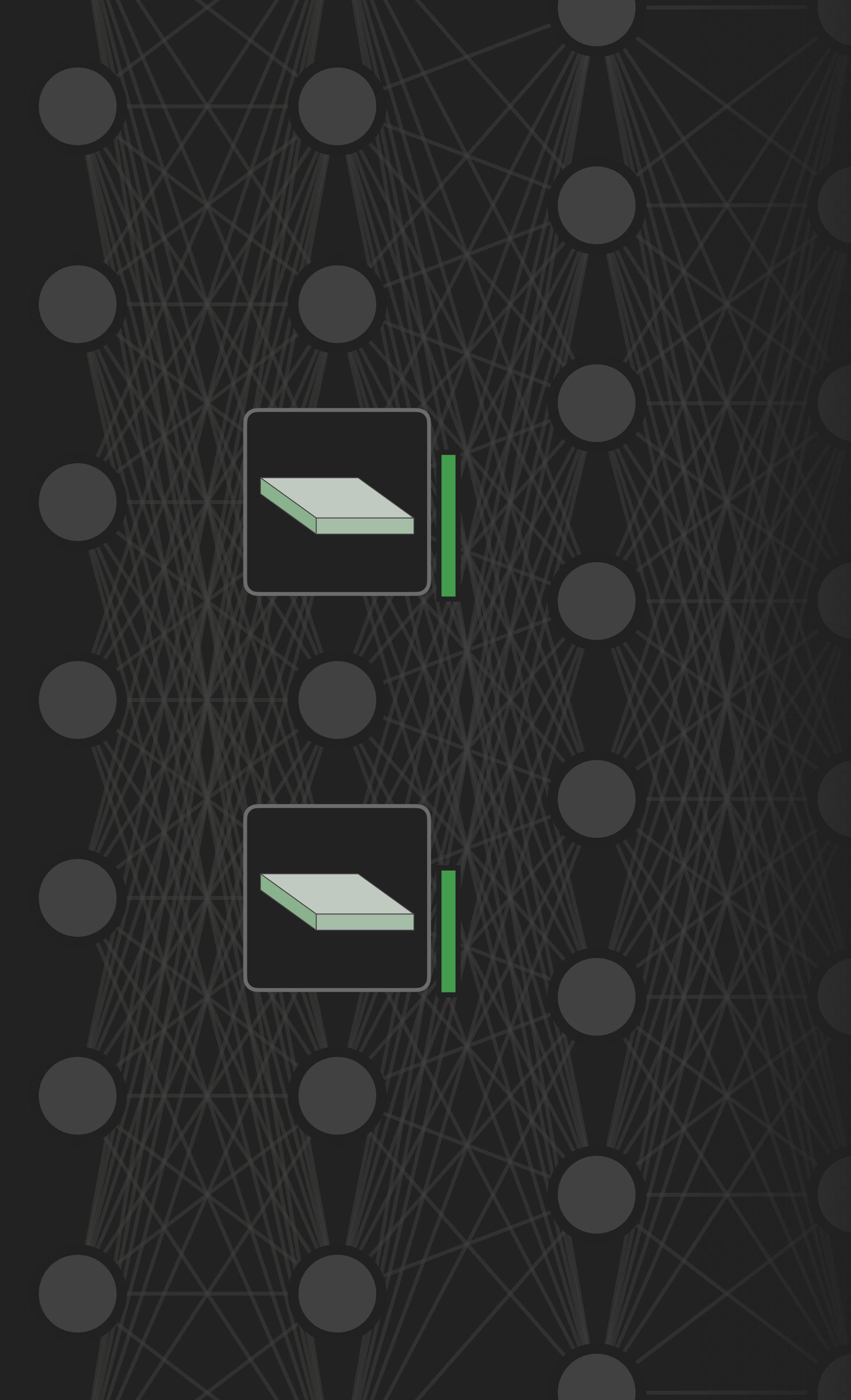


Aggregate network  
**activations** (nodes)

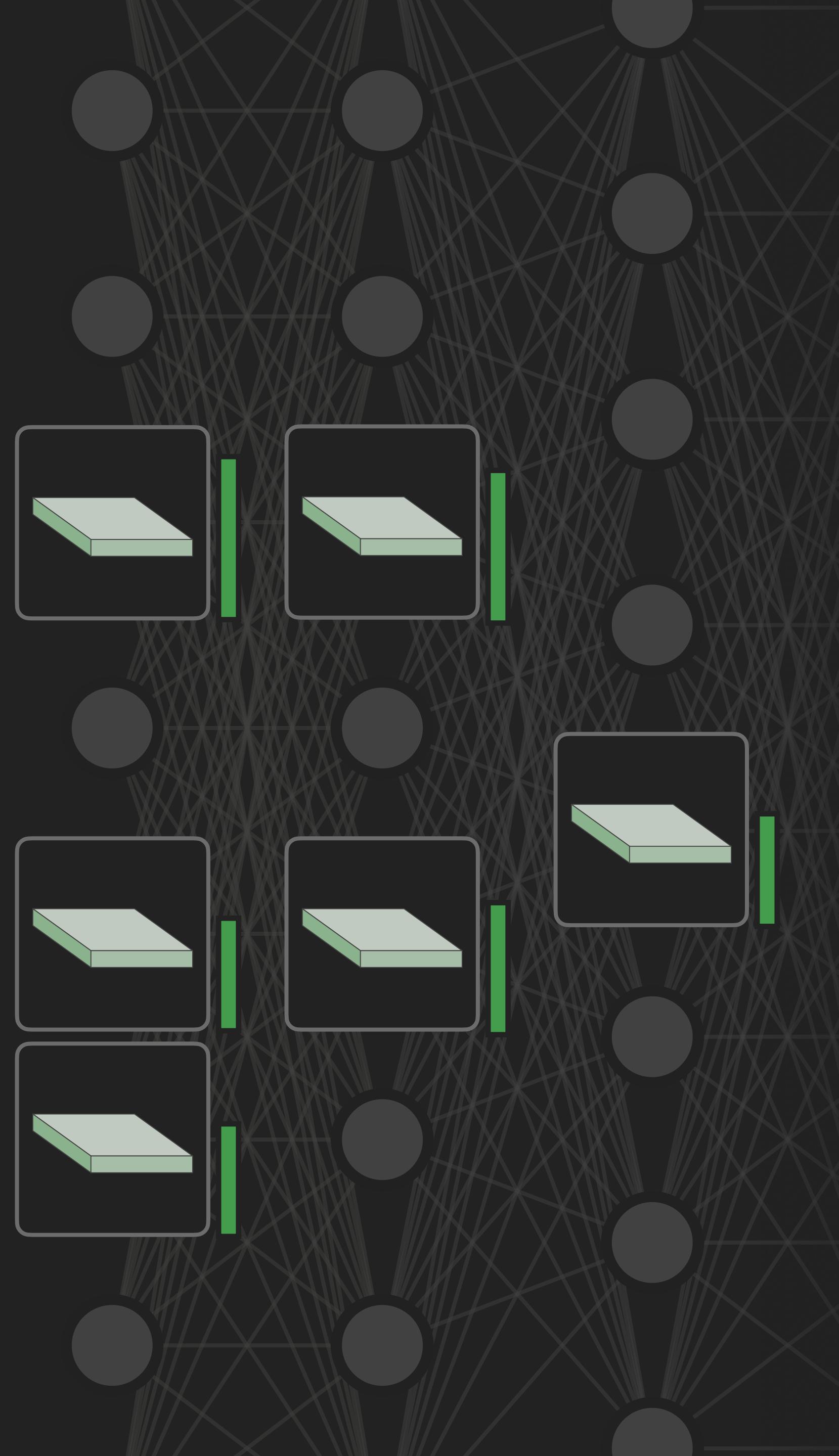




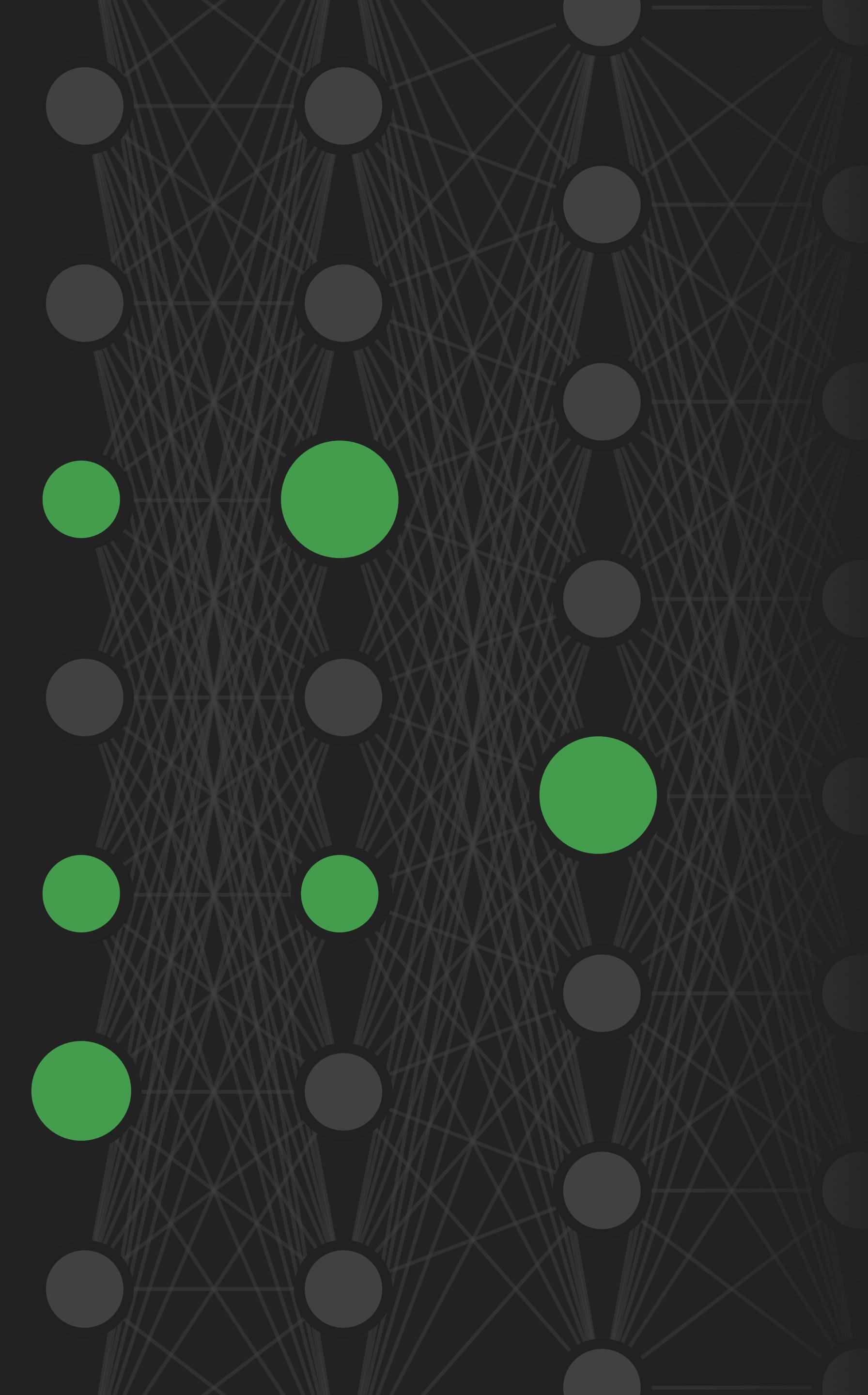
Aggregate network  
**activations** (nodes)



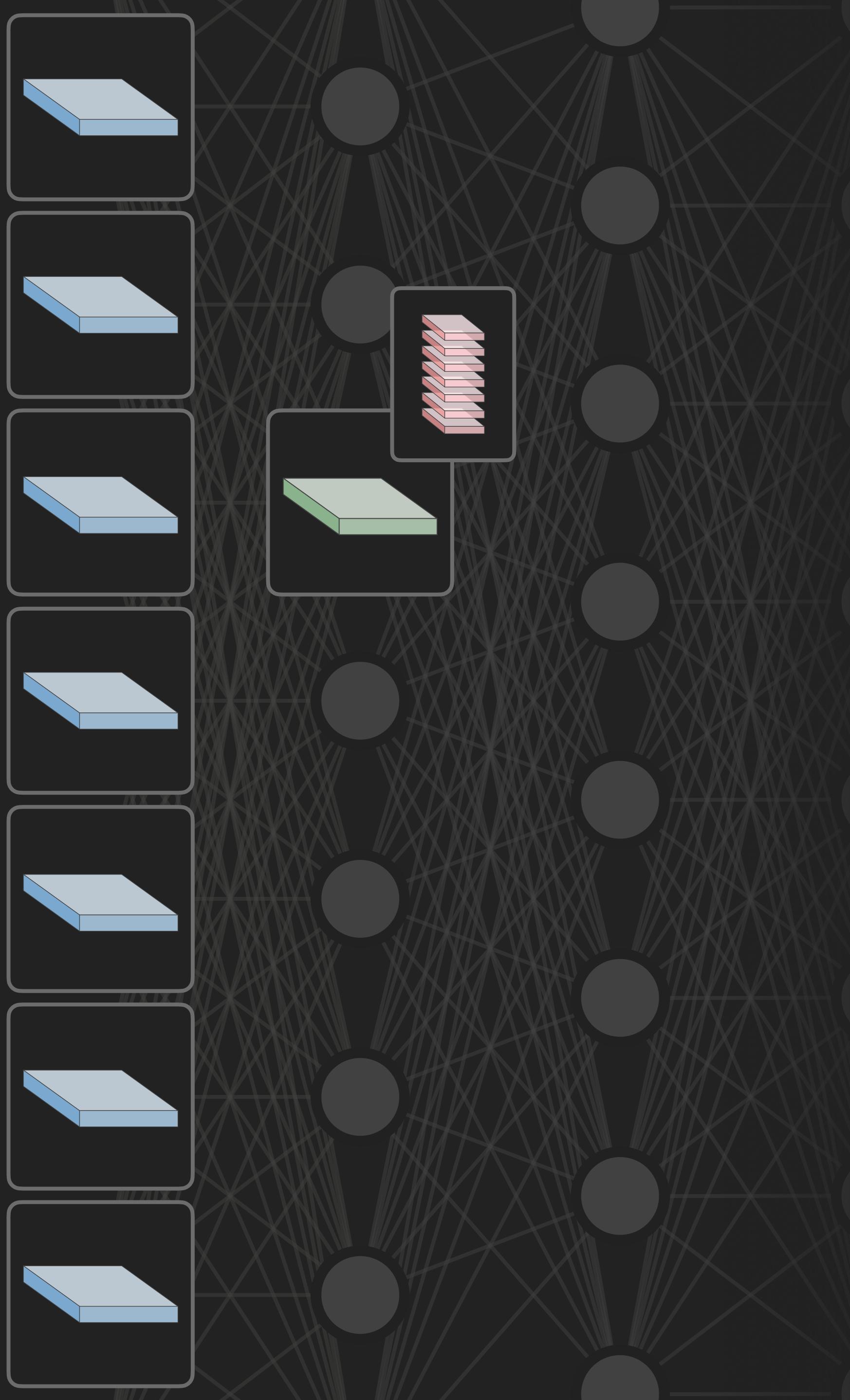
Aggregate network  
**activations** (nodes)



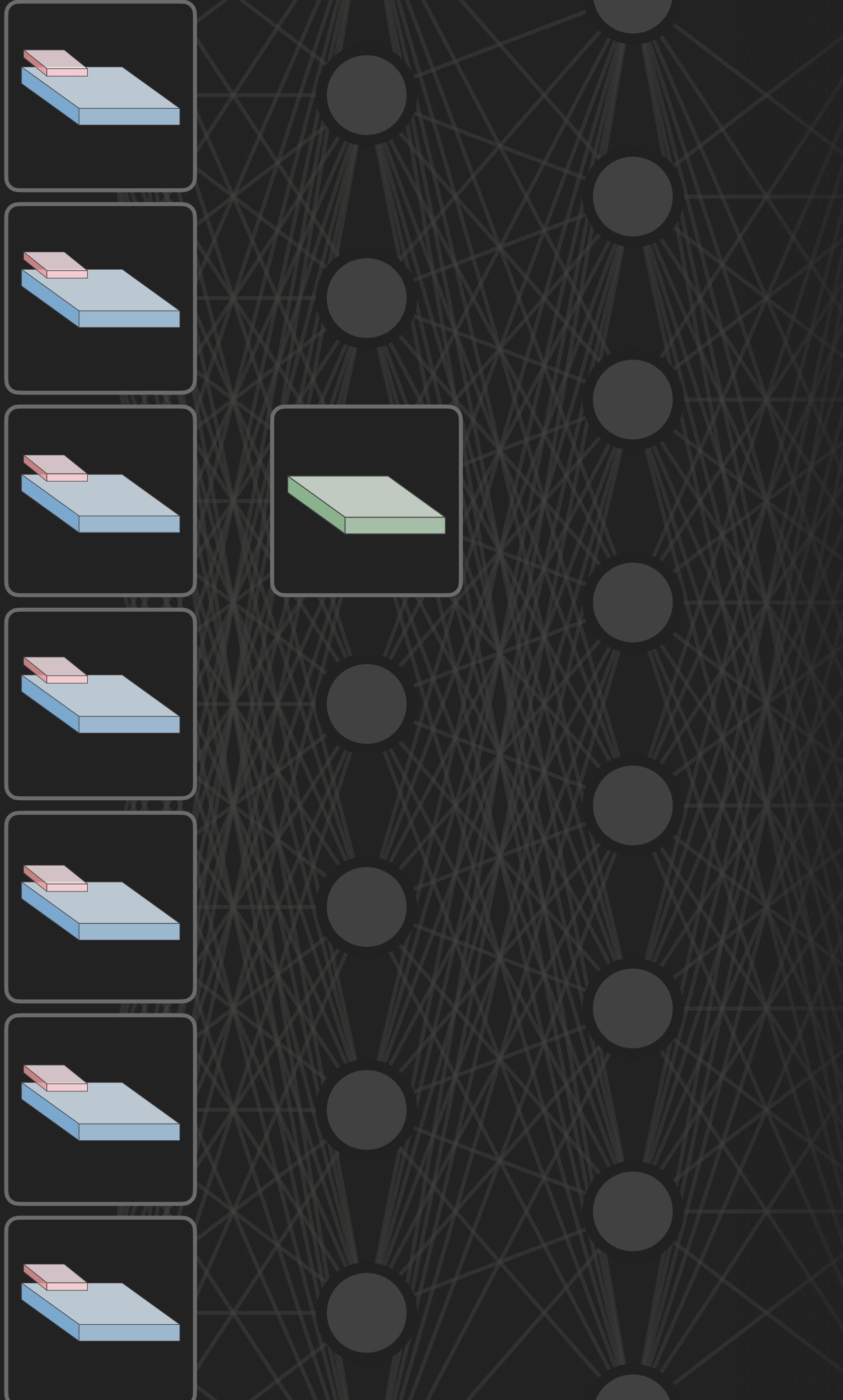
Aggregate network  
**activations** (nodes)



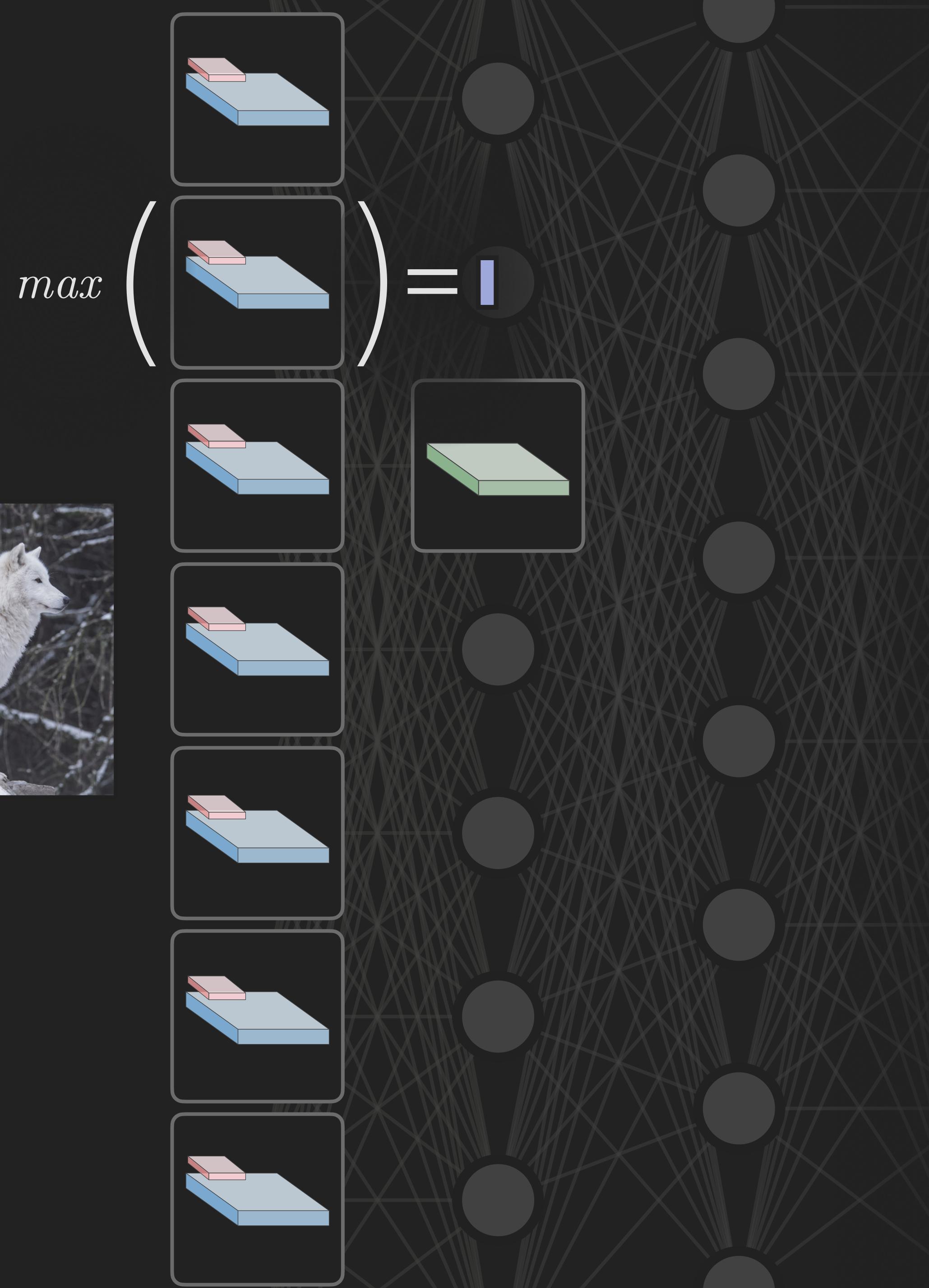
Aggregate network  
**activations** (nodes)



Aggregate network  
**influences** (edges)



Aggregate network  
**influences** (edges)



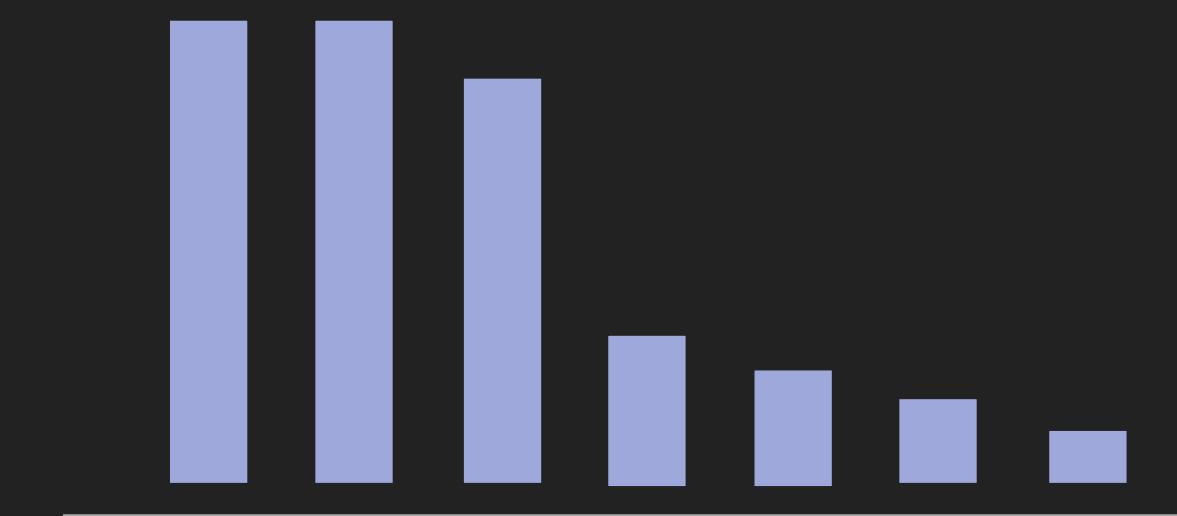
Aggregate network  
**influences** (edges)

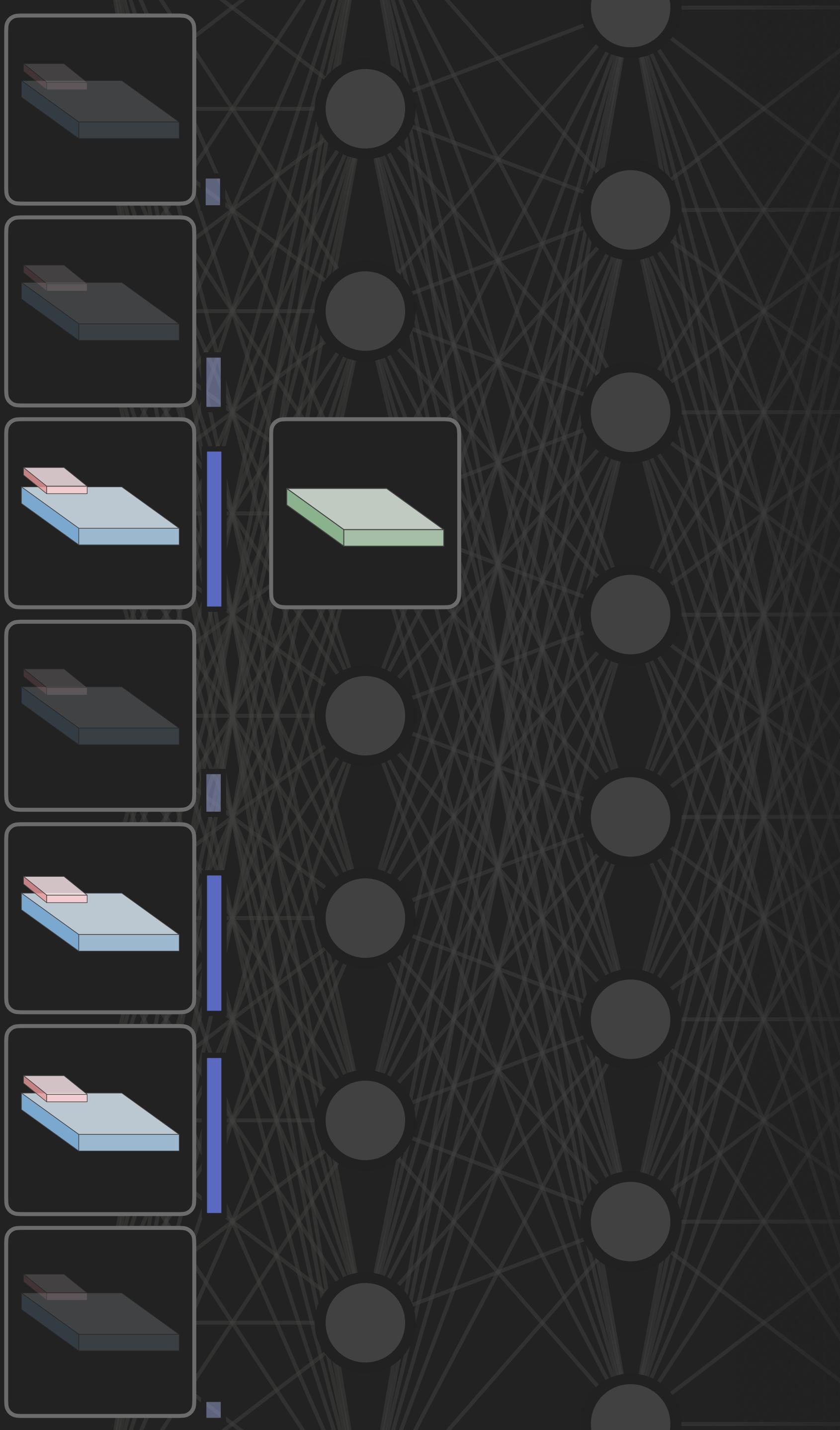


Aggregate network  
**influences** (edges)

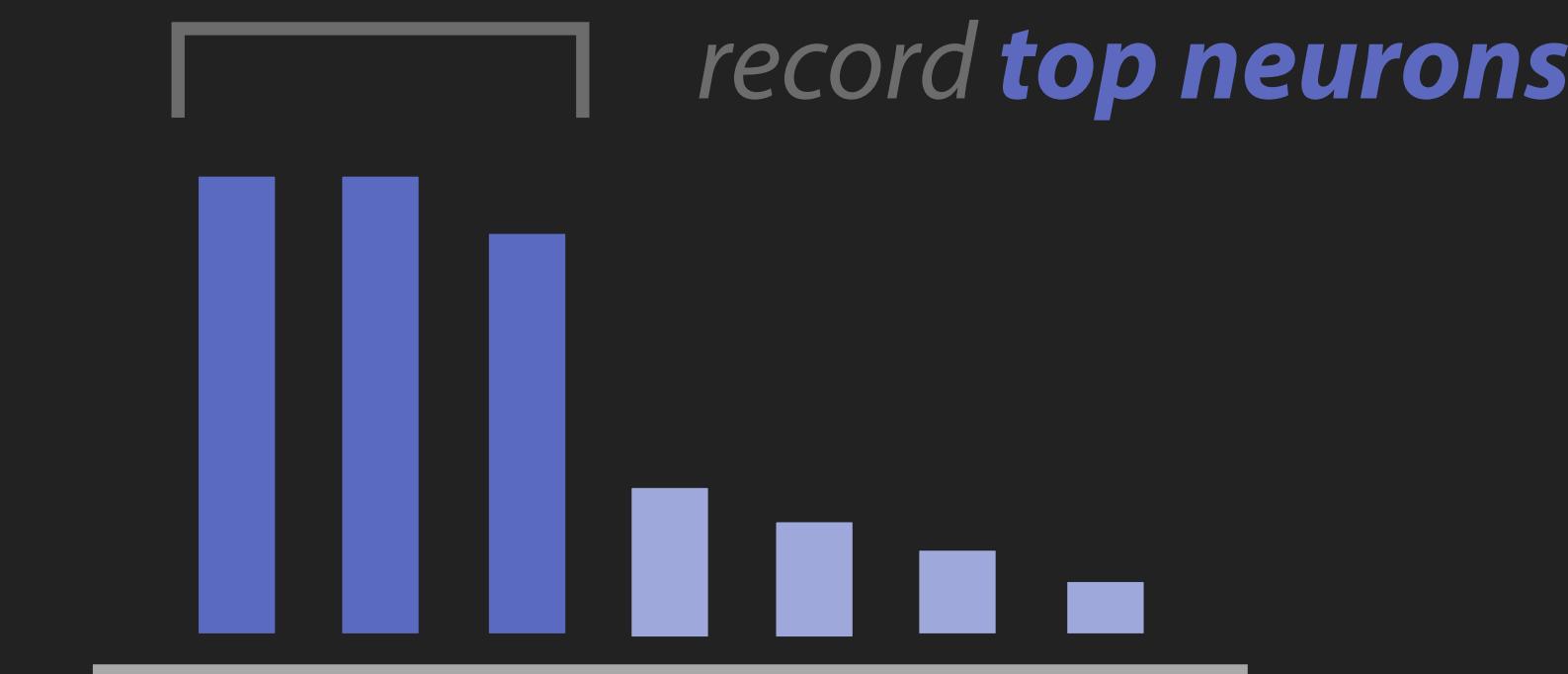


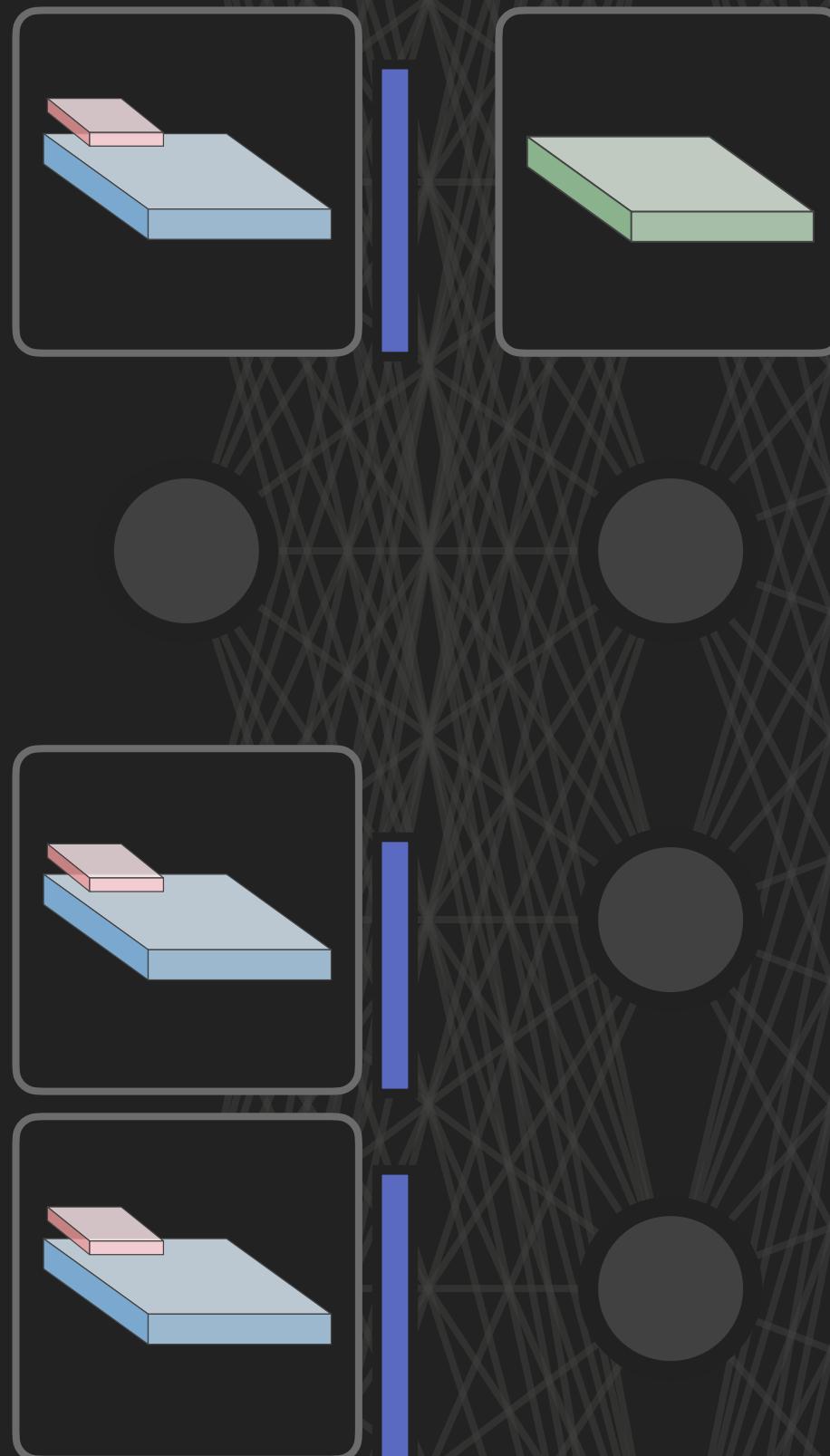
Aggregate network  
**influences** (edges)



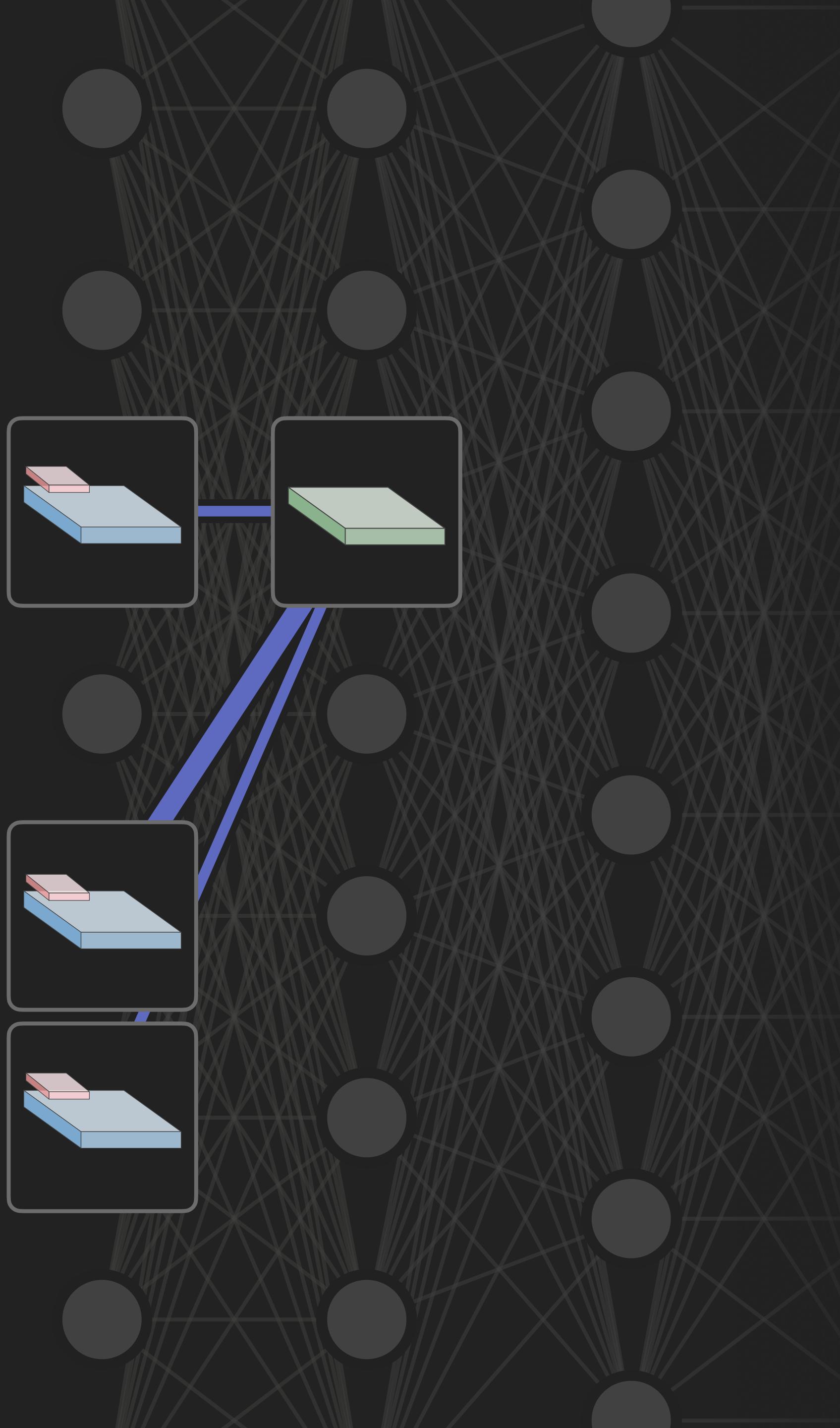


## Aggregate network influences (edges)

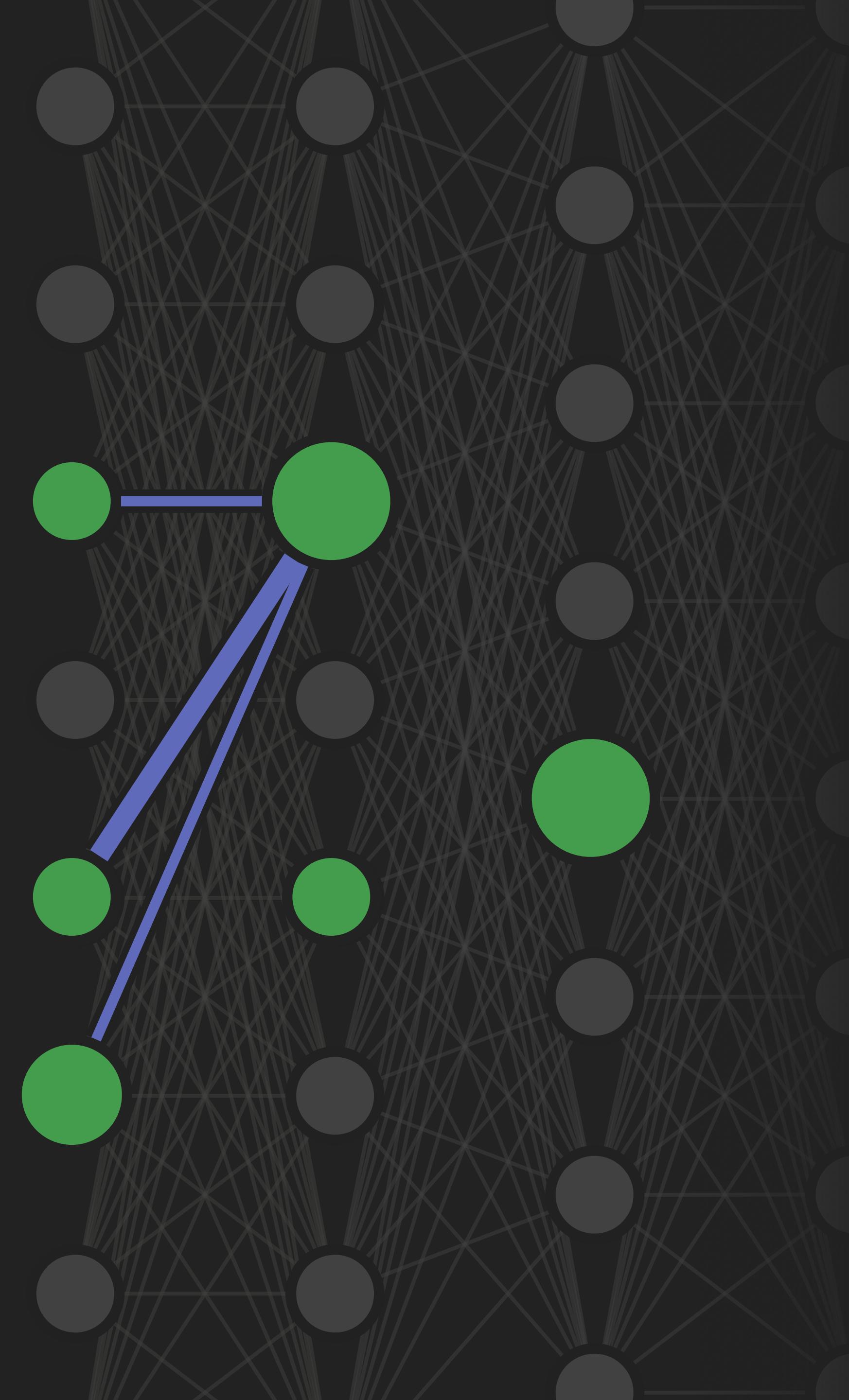




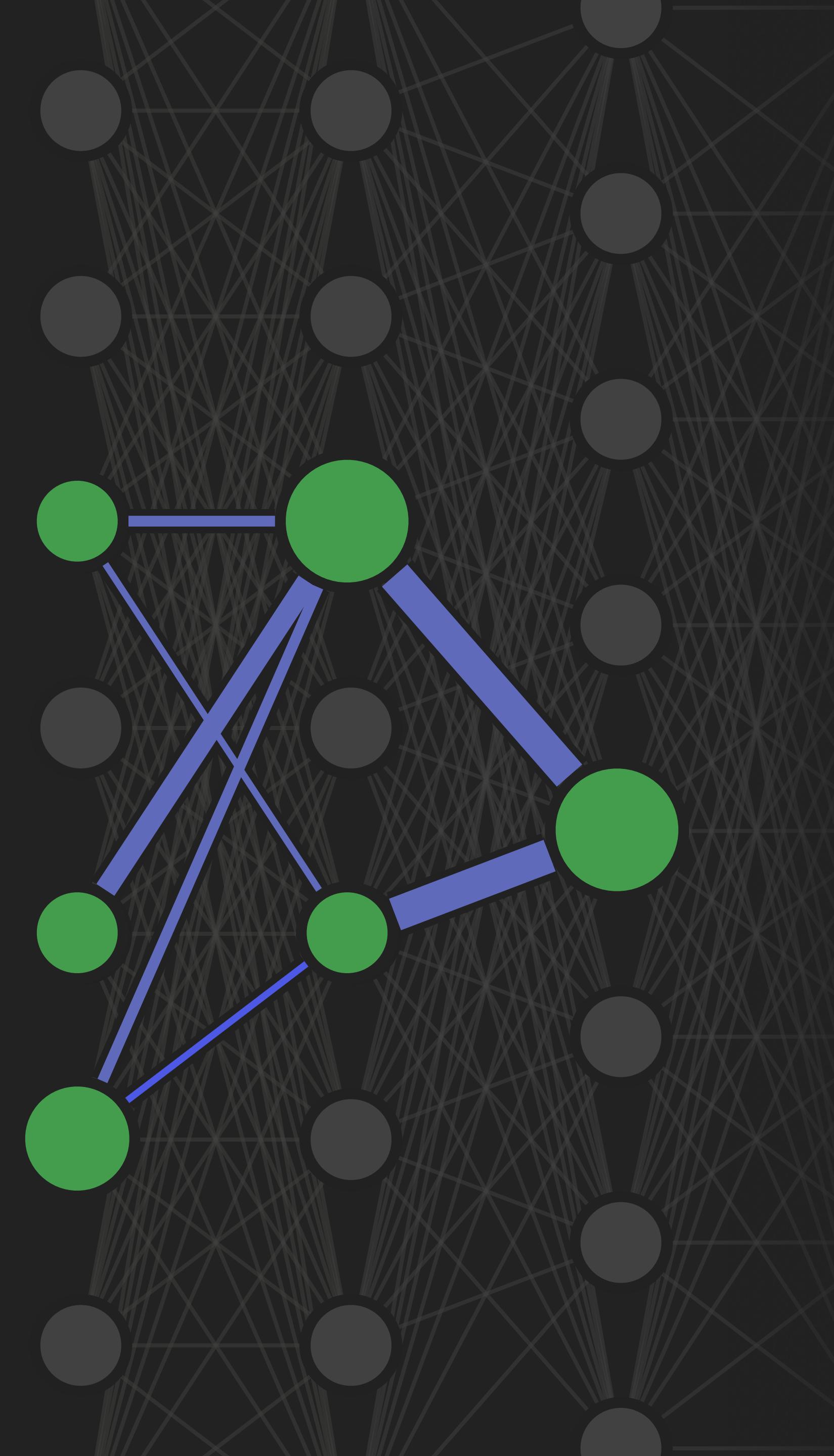
Aggregate network  
**influences** (edges)



Aggregate network  
**influences** (edges)



Combine **activations**  
and **influences**



Combine **activations**  
and **influences**

Further summarize graph  
**personalized PageRank**

# Feature Visualization

What kind of input would cause a neuron to maximally activate?

# Feature Visualization

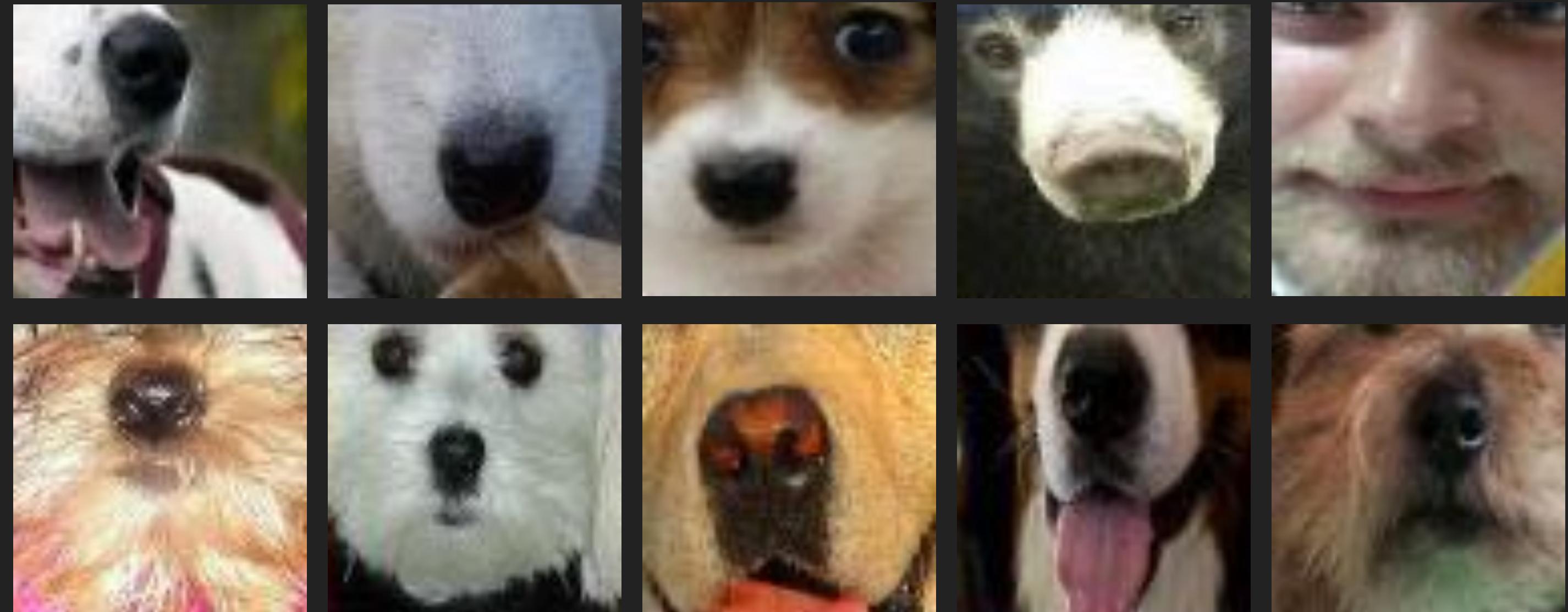
What kind of input would cause a neuron to maximally activate?

***Generate examples: starting from random noise, optimize an image to activate a particular neuron***

# Feature Visualization

What kind of input would cause a neuron to maximally activate?

***Generate examples: starting from random noise, optimize an image to activate a particular neuron***



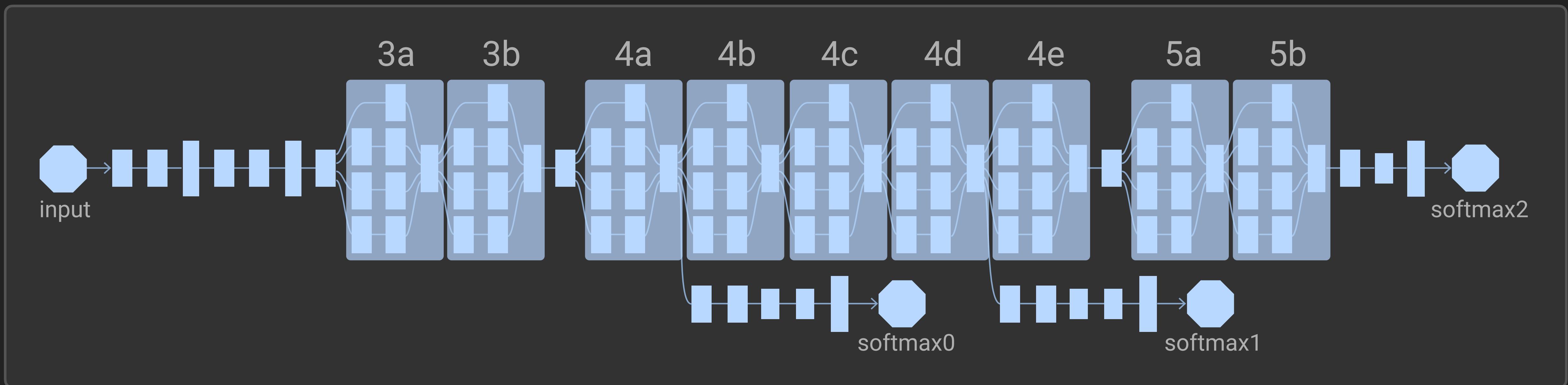
*mixed4b, 409*

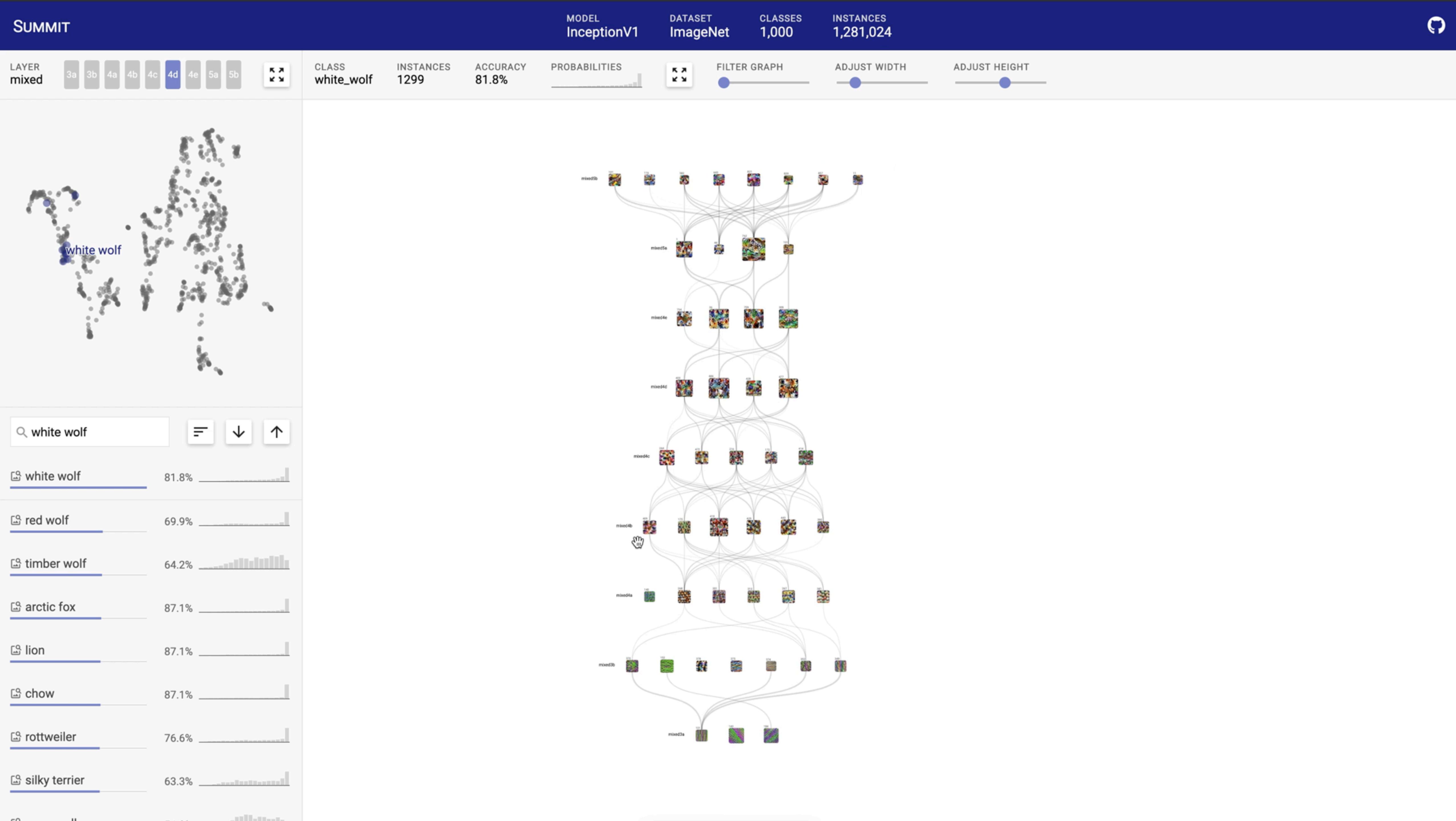
[Olah, et al., Distill, 2017]

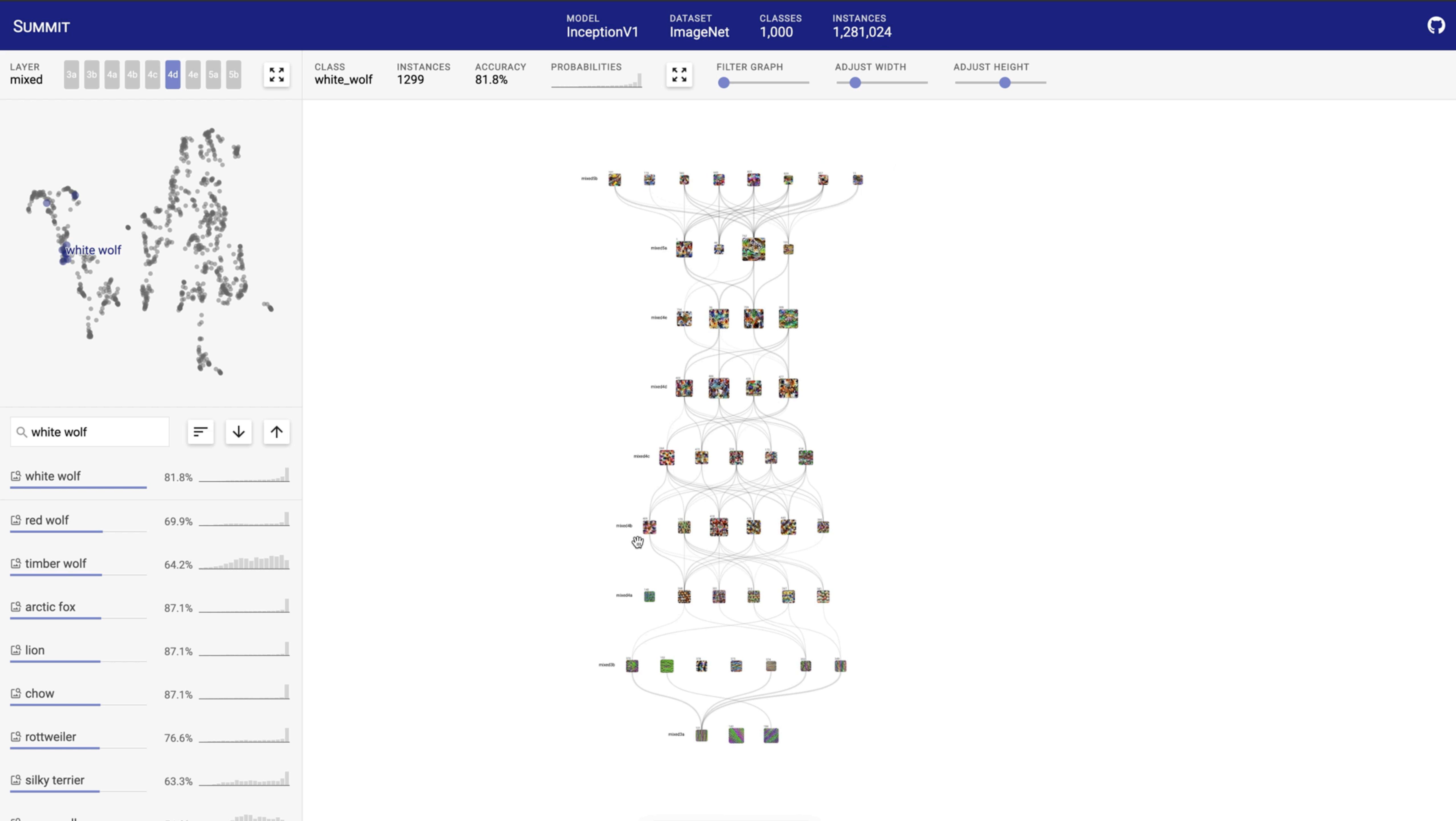
# Demo

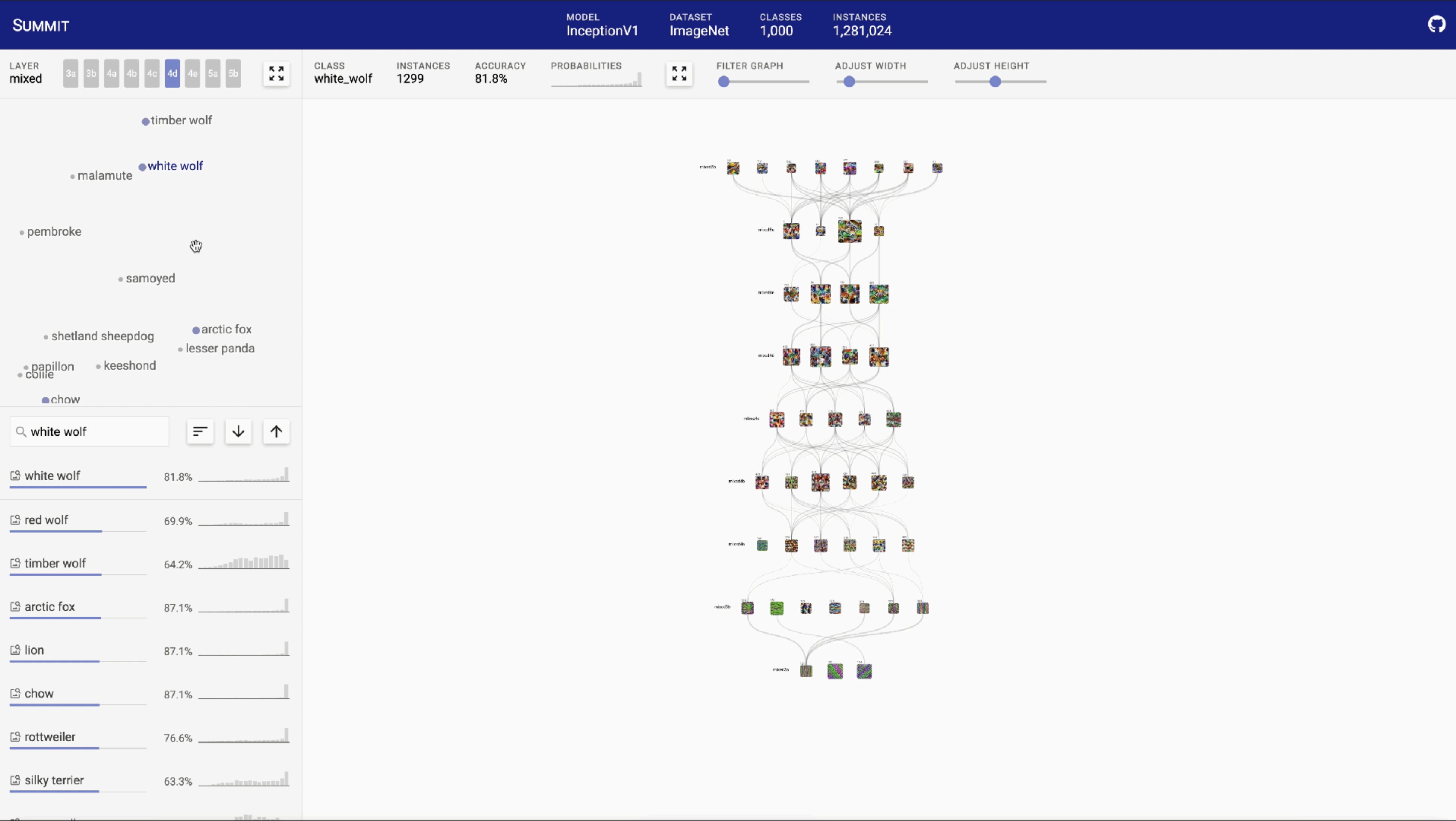
InceptionV1  
Large-scale,  
prevalent CNN

ImageNet (ILSVRC)  
~1.3M images  
1,000 classes







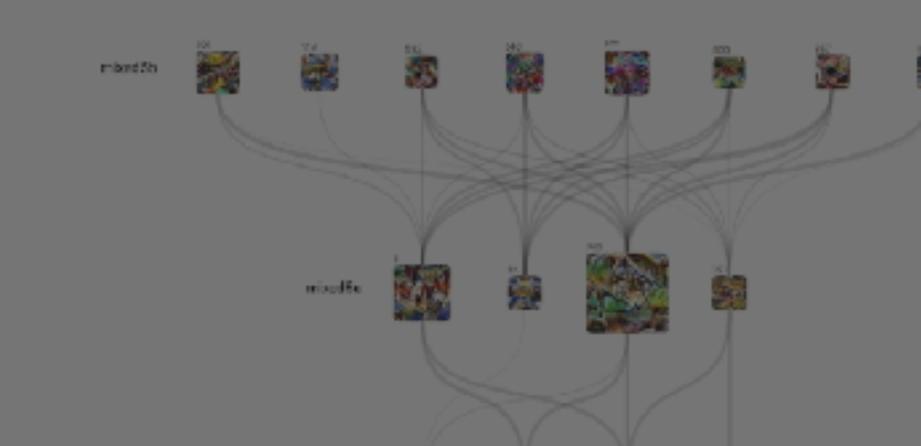
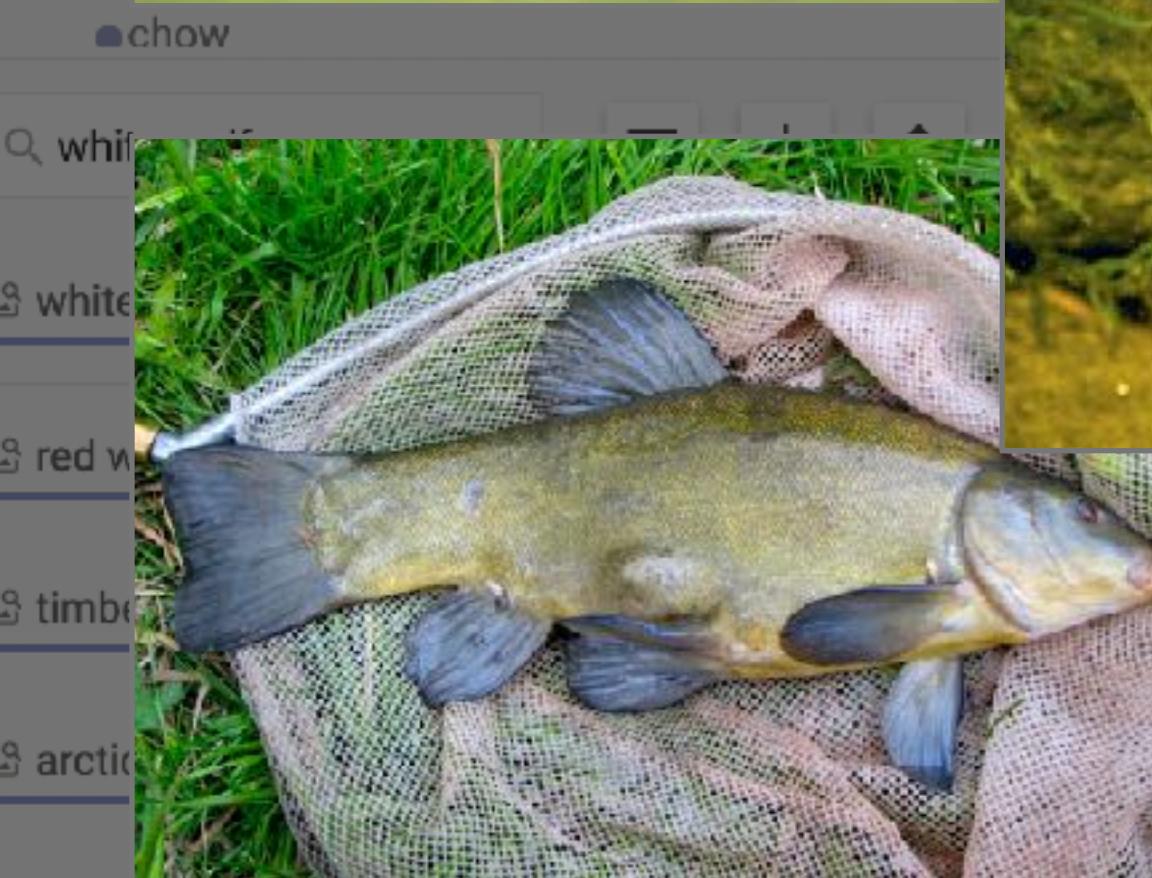
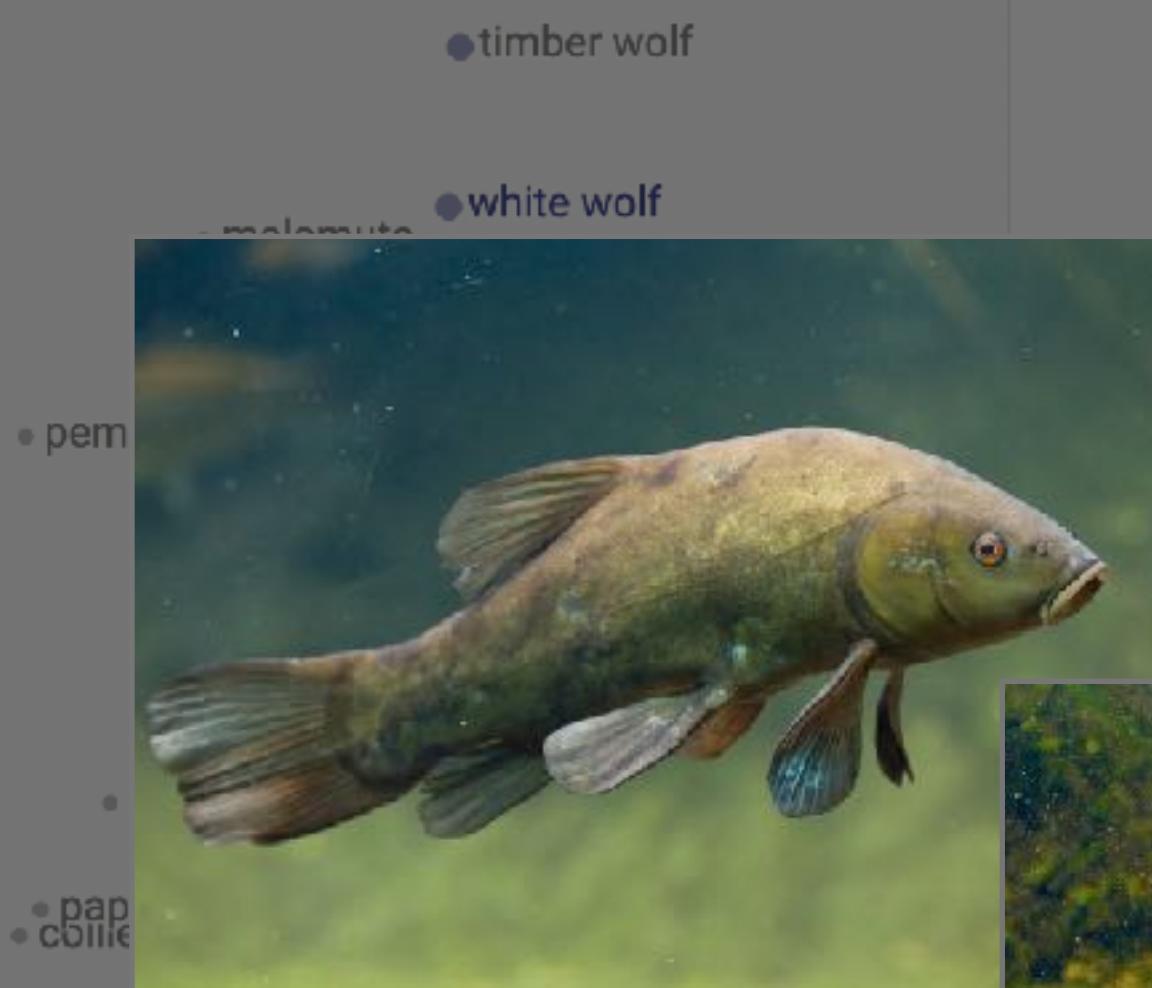


LAYER  
mixed    3a 3b 4a 4b 4c 4d 4e 5a 5bCLASS  
white\_wolf    INSTANCES  
1299    ACCURACY  
81.8%    PROBABILITIES

FILTER GRAPH

ADJUST WIDTH

ADJUST HEIGHT



What features has a neural network learned for **tench**?  
How are those features related?



LAYER  
mixed

3a 3b 4a 4b 4c 4d 4e 5a 5b

CLASS  
white\_wolfINSTANCES  
1299ACCURACY  
81.8%

PROBABILITIES



FILTER GRAPH



ADJUST WIDTH



ADJUST HEIGHT

- timber wolf
- malamute ● white wolf
- pembroke
- samoyed
- shetland sheepdog ● arctic fox
- papillon ● lesser panda
- collie ● keeshond
- chow

white wolf 81.8%

red wolf 69.9%

timber wolf 64.2%

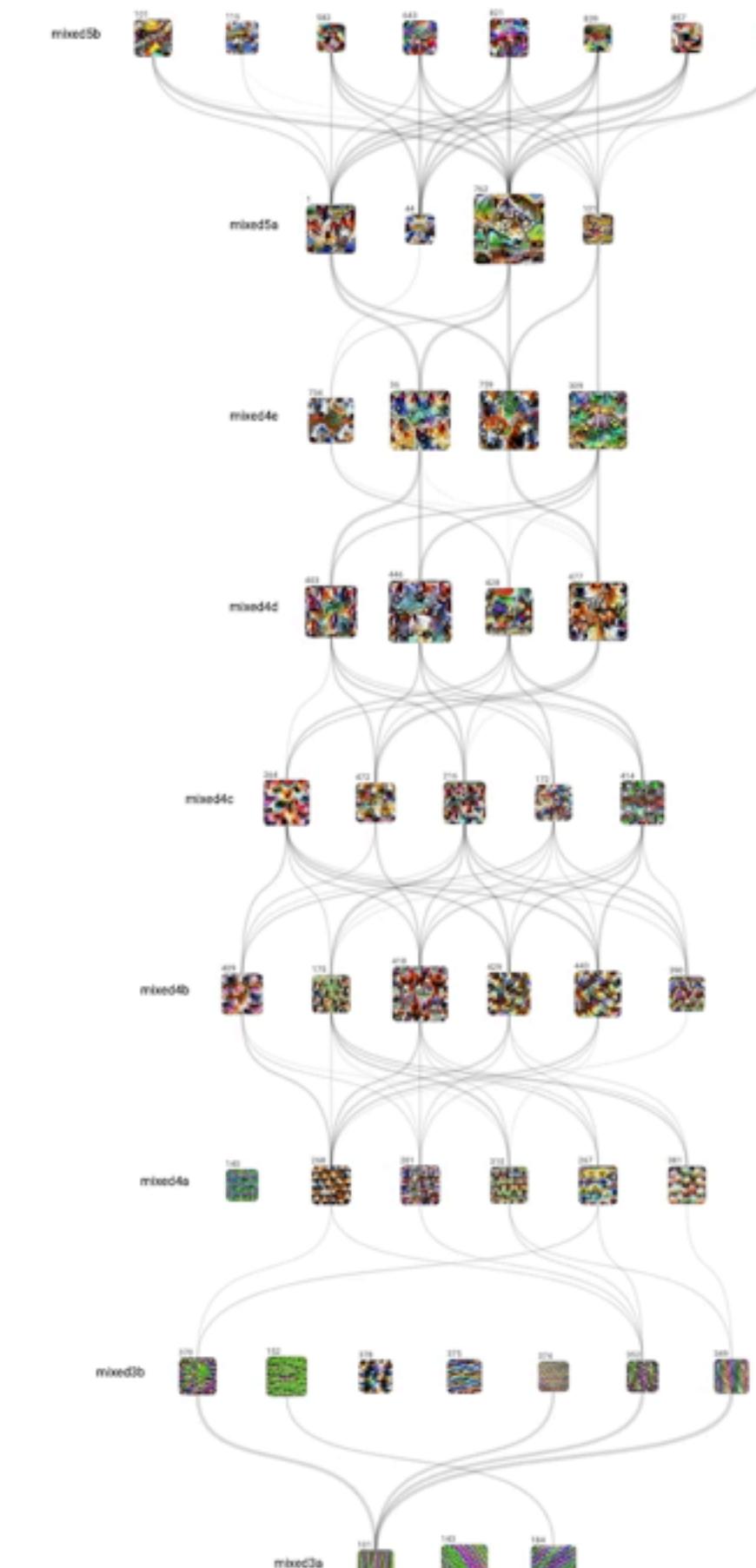
arctic fox 87.1%

lion 87.1%

chow 87.1%

rottweiler 76.6%

silky terrier 63.3%

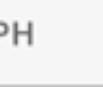


LAYER  
mixed    3a 3b 4a 4b 4c 4d 4e 5a 5bCLASS  
white\_wolfINSTANCES  
1299ACCURACY  
81.8%

PROBABILITIES



FILTER GRAPH



ADJUST WIDTH

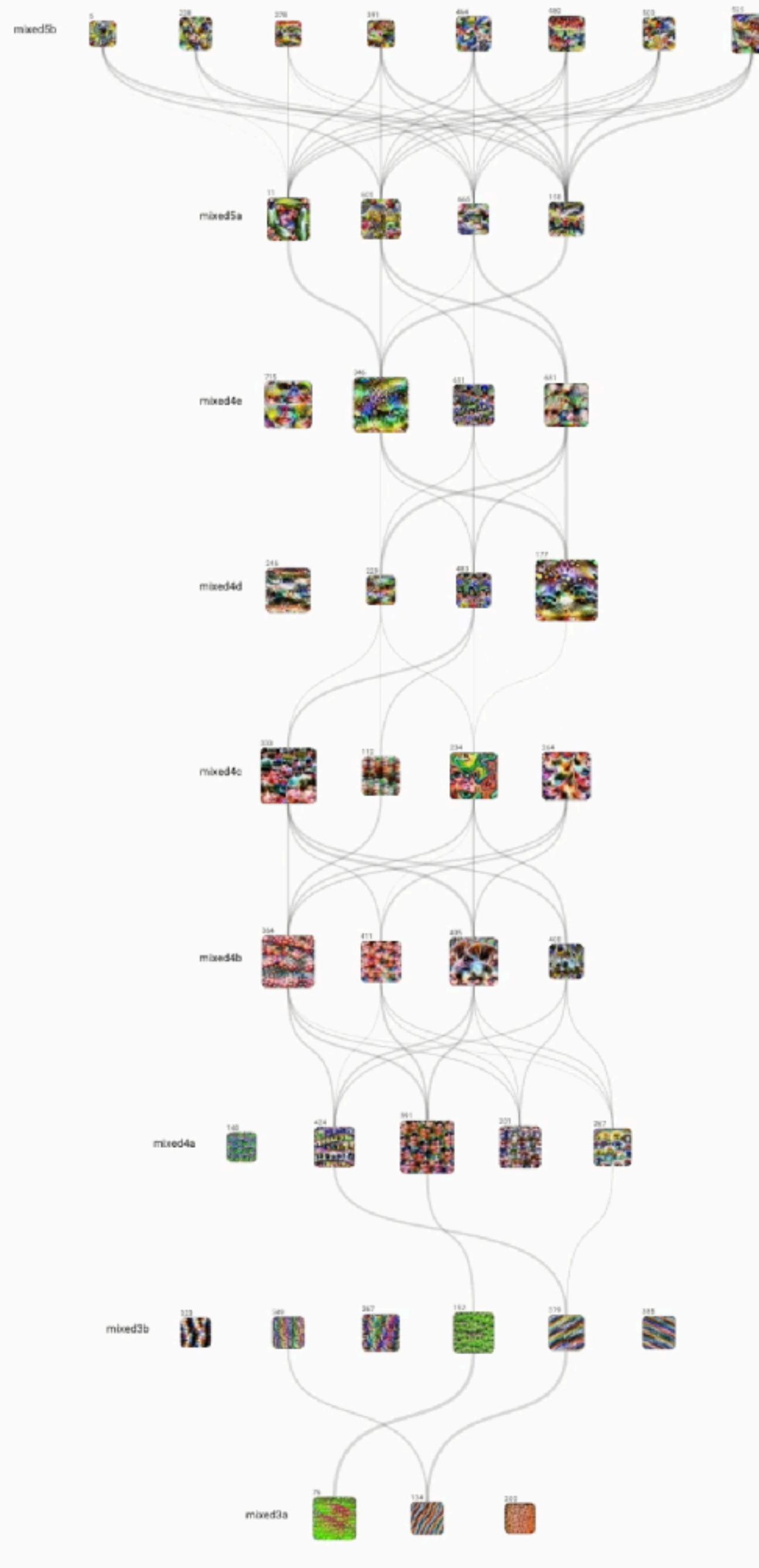


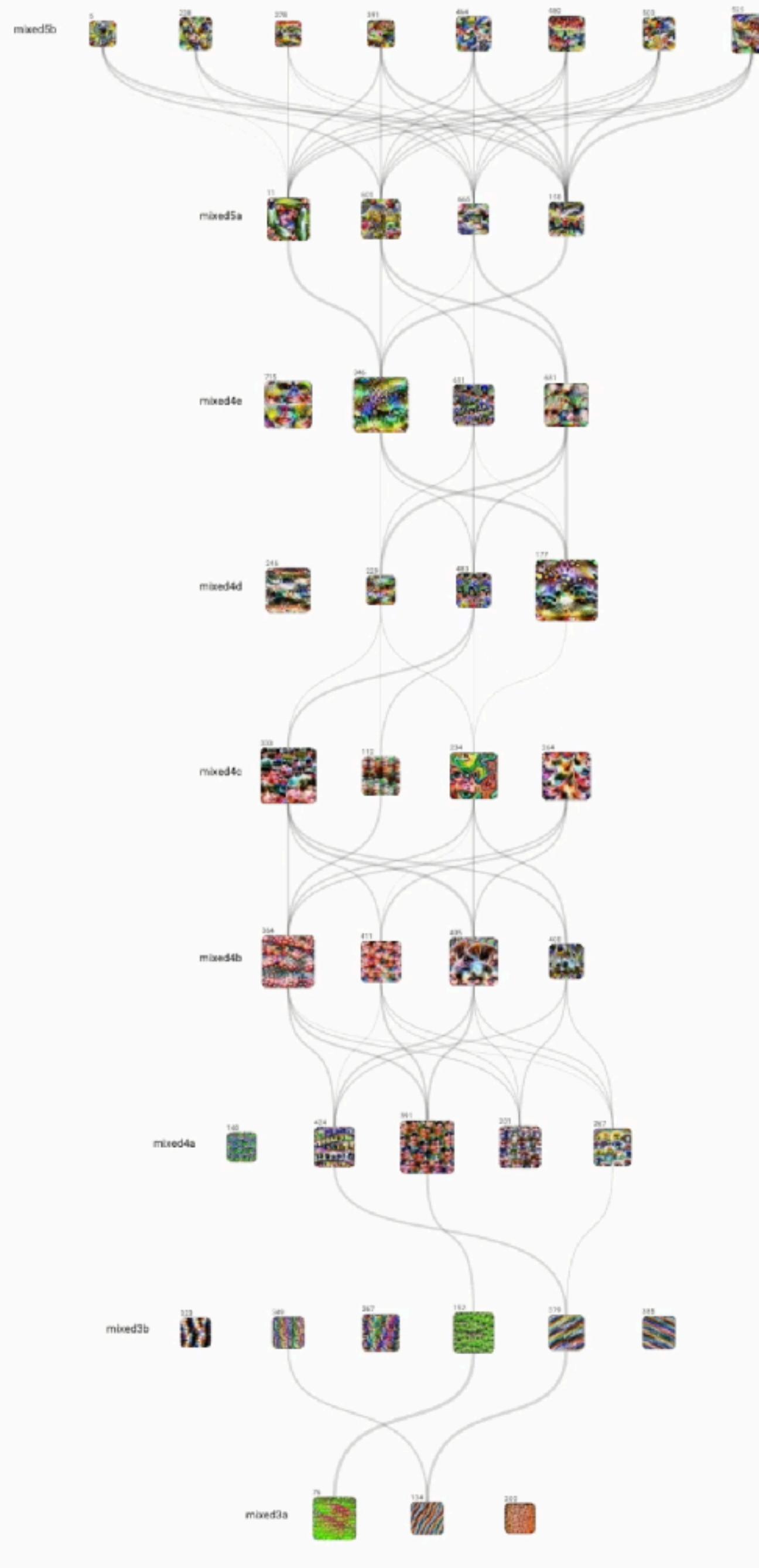
ADJUST HEIGHT

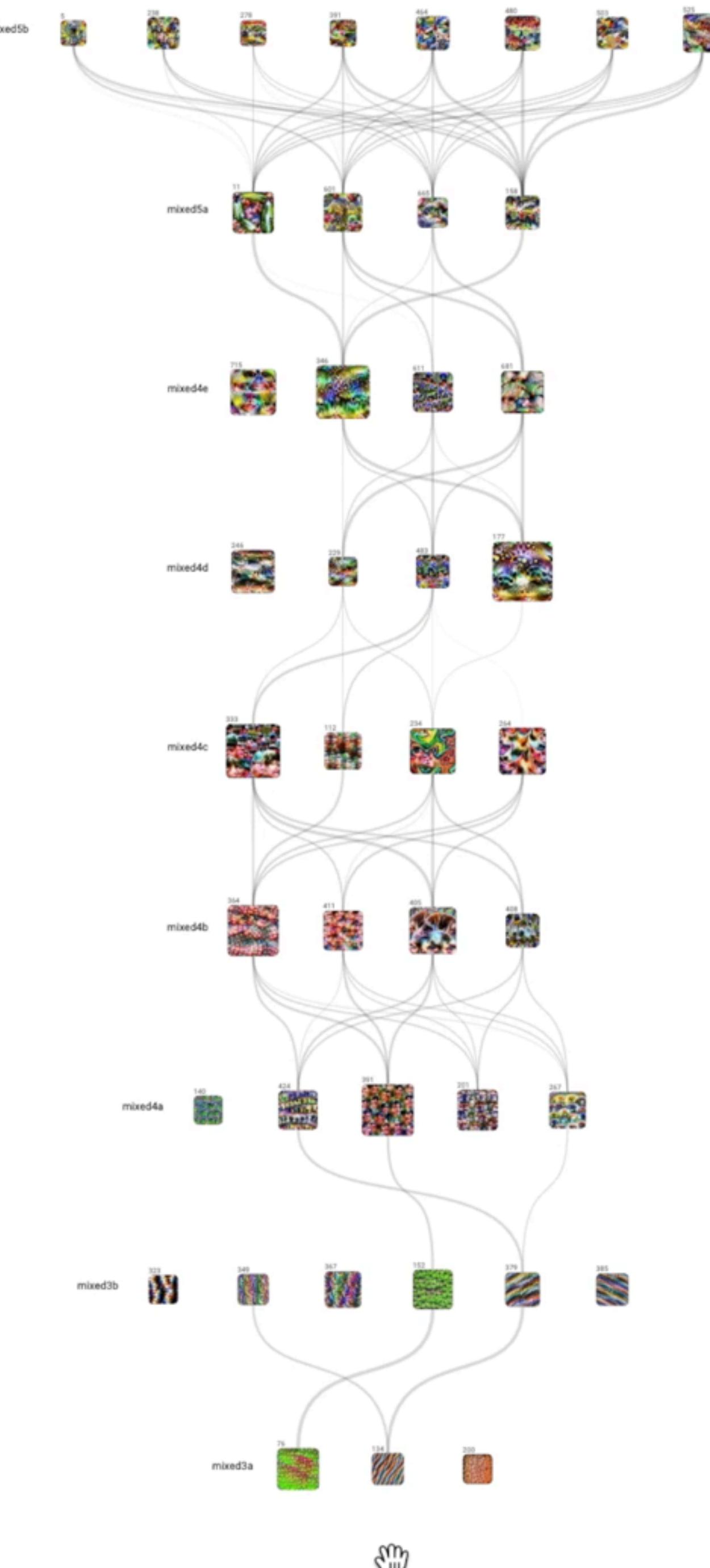
- timber wolf
- malamute • white wolf
- pembroke
- samoyed
- shetland sheepdog • arctic fox
- papillon • lesser panda
- collie • keeshond
- chow

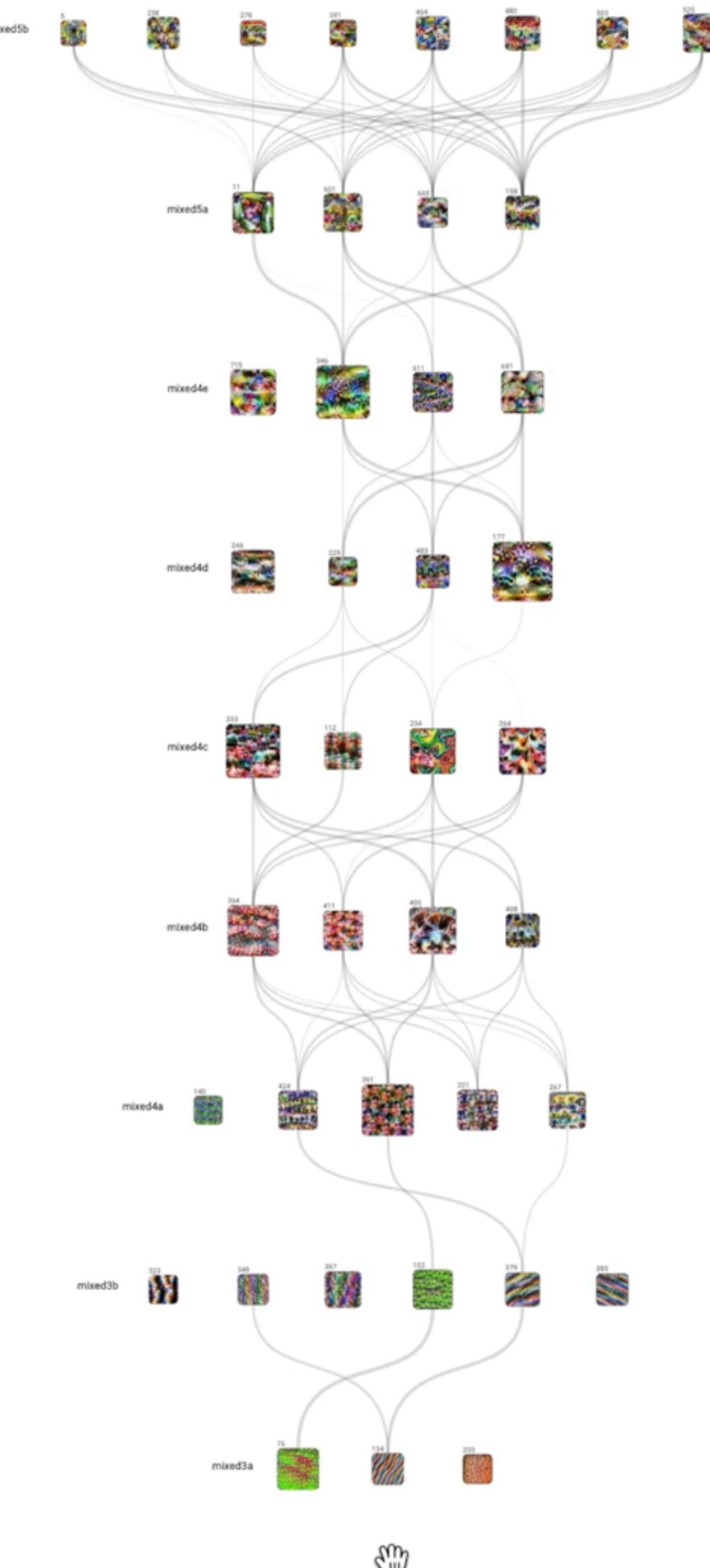
white wolf	81.8%	
red wolf	69.9%	
timber wolf	64.2%	
arctic fox	87.1%	
lion	87.1%	
chow	87.1%	
rottweiler	76.6%	
silky terrier	63.3%	





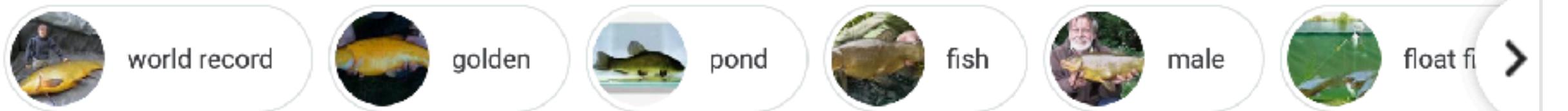










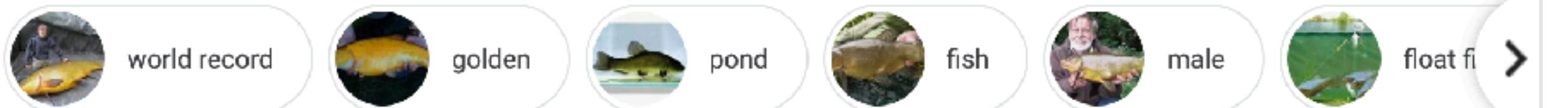
[All](#) [\*\*Images\*\*](#) [News](#) [Shopping](#) [Videos](#) [More](#) [Settings](#) [Tools](#)[Tench - Wikipedia](#)[en.wikipedia.org](https://en.wikipedia.org)[Top Tench Fishing Baits & Tactics...](#)[dynamitebaits.com](https://dynamitebaits.com)[Tench Fishing Guide - What Is Tench ...](#)[badangling.com](https://badangling.com)[Early season tench fishing tips ...](#)[dynamitebaits.com](https://dynamitebaits.com)[SPRING SPECIMENS Article | Korum ...](#)[korum.co.uk](https://korum.co.uk)[Boilie Approach For Tench | Drennan ...](#)[drennattackle.com](https://drennattackle.com)

Google

tench



All Images News Shopping Videos More Settings Tools



Tench - Wikipedia  
en.wikipedia.org



Top Tench Fishing Baits & Tactics...  
dynamitebaits.com



Tench Fishing Guide - What Is Tench ...  
badangling.com



Early season tench fishing tips ...  
dynamitebaits.com



SPRING SPECIMENS Article | Korum ...  
korum.co.uk



Boilie Approach For Tench | Drennan ...  
drennattackle.com

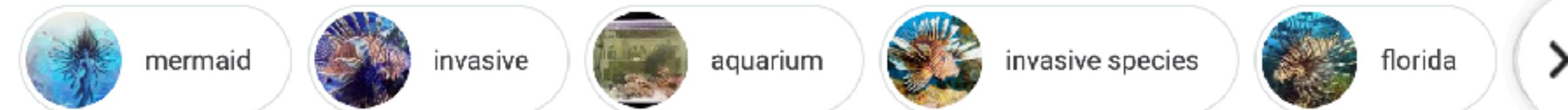


Google

lionfish



All Images News Shopping Videos More Settings Tools



Pterois - Wikipedia  
en.wikipedia.org



Invasive lionfish are delicious — but ...  
oceana.org



Invasive lionfish are delicious — but ...  
oceana.org



Lionfish: The Beautiful and Dange...  
livescience.com



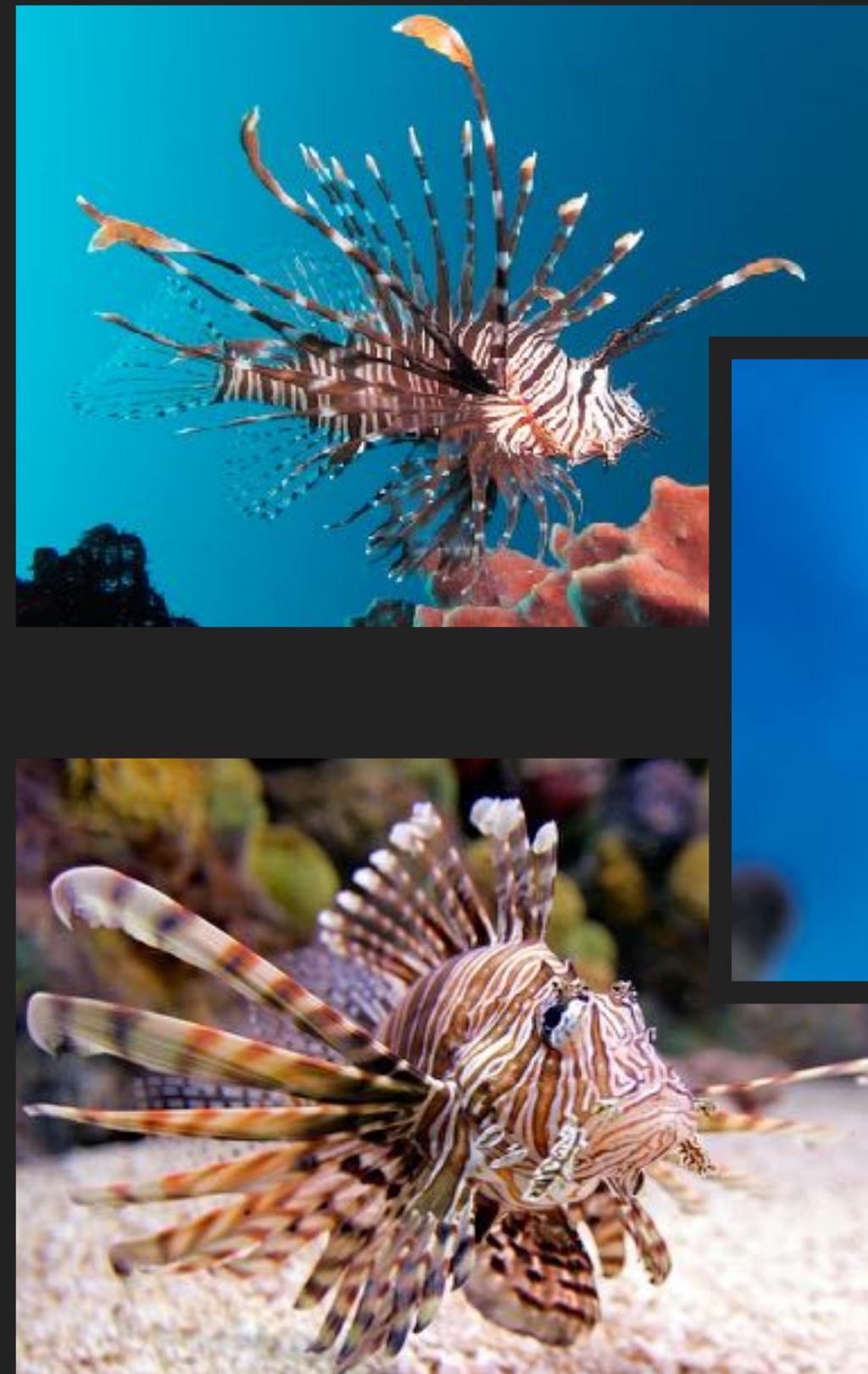
What is a lionfish?  
oceanservice.noaa.gov



lionfish | Invasive Species, Sting ...  
britannica.com

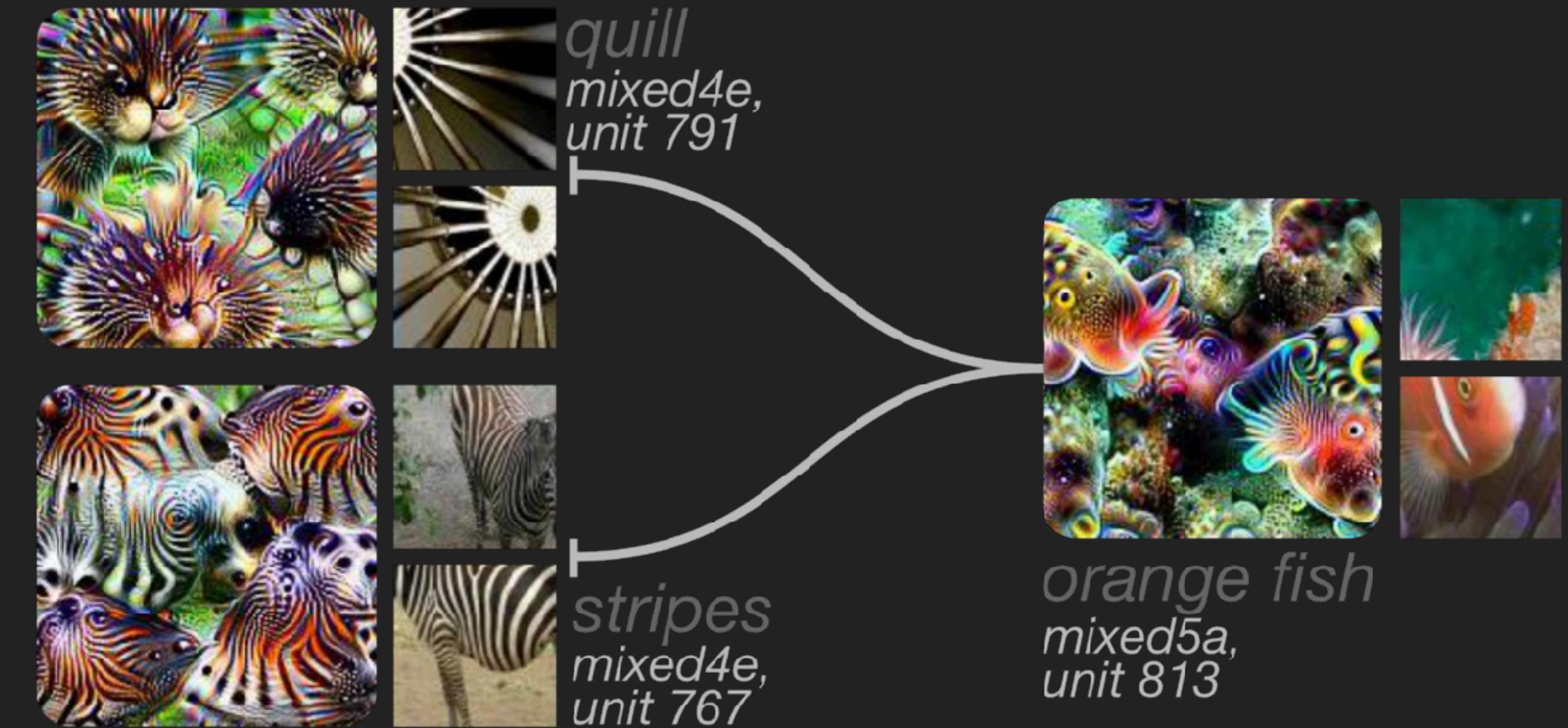


# Unexpected Features

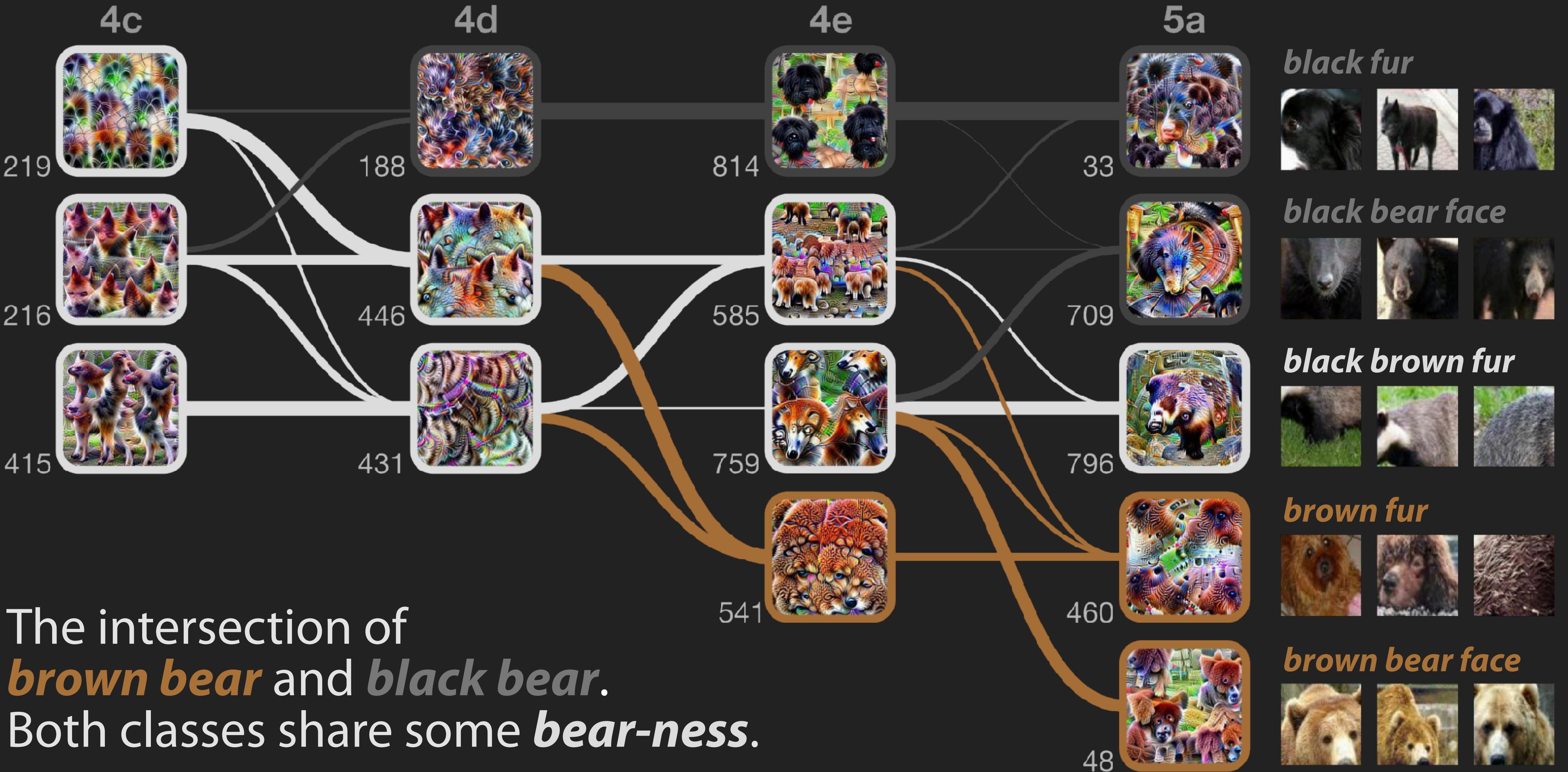


**lionfish**

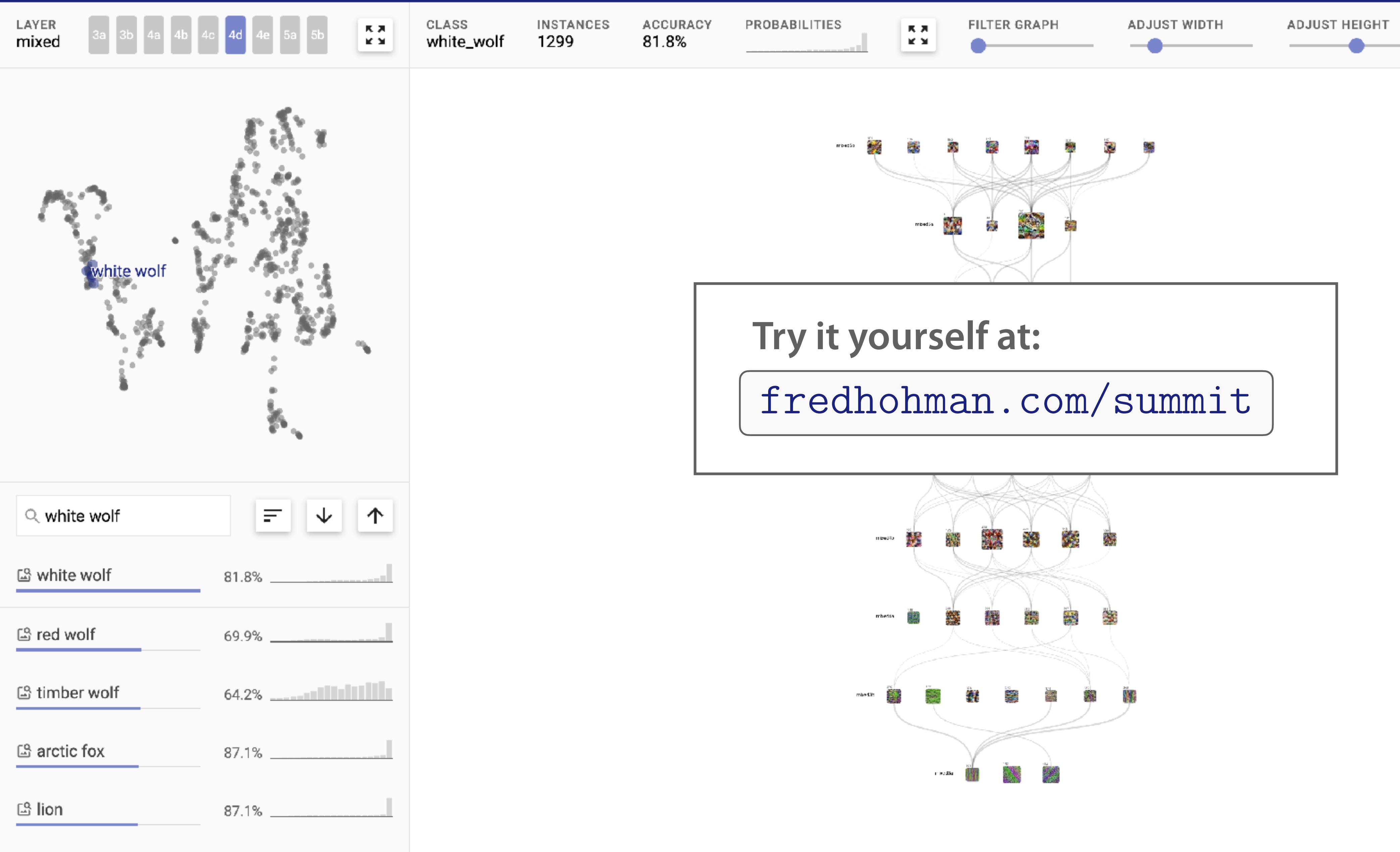
No more “people” features.  
But few “fish” features! Mostly textures.



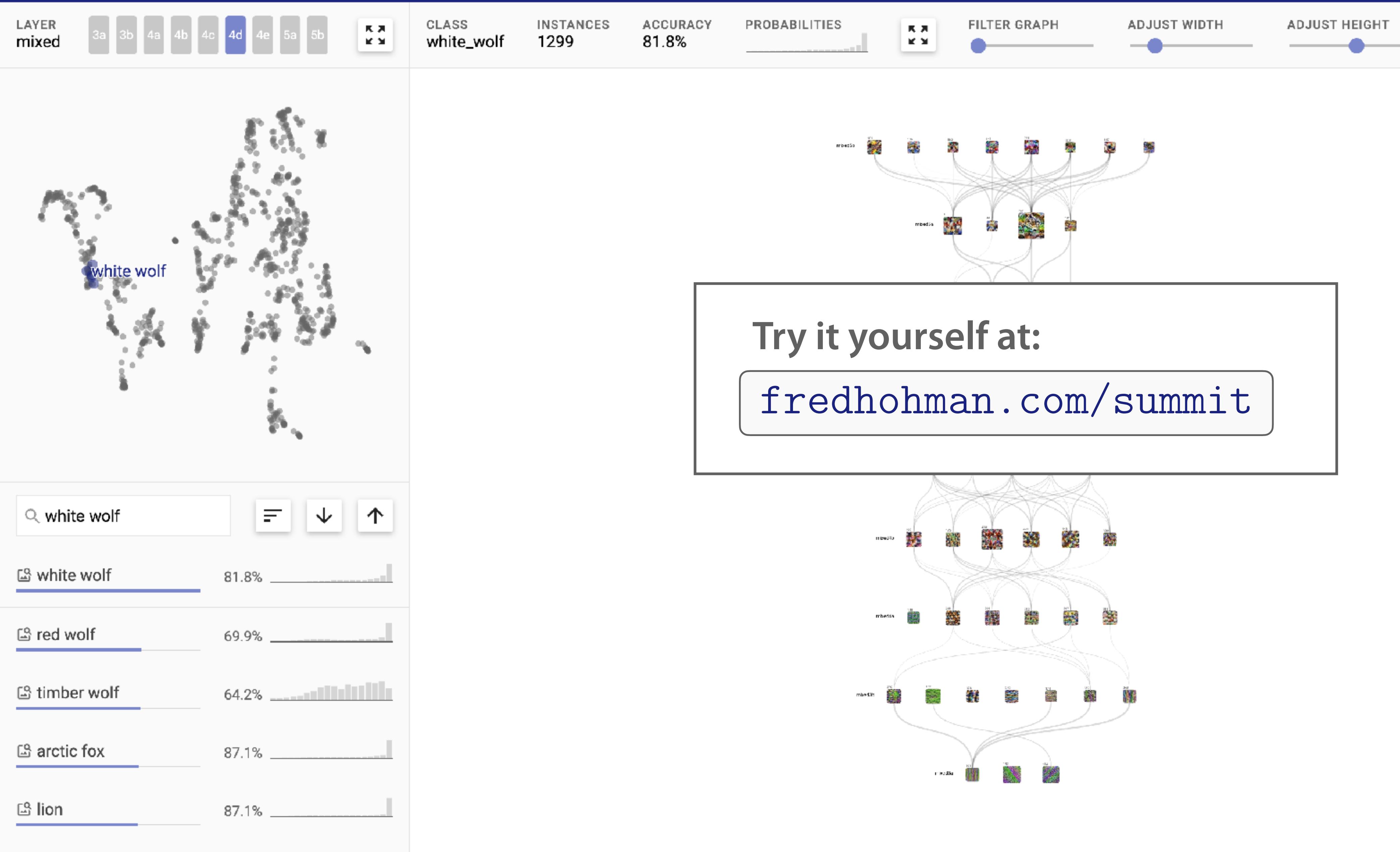
Attribution graph substructure from **lionfish** class.

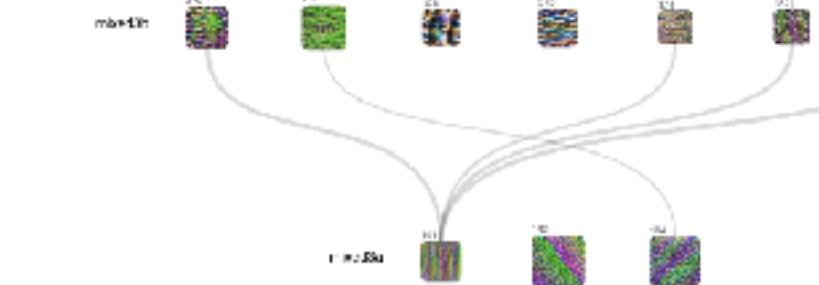
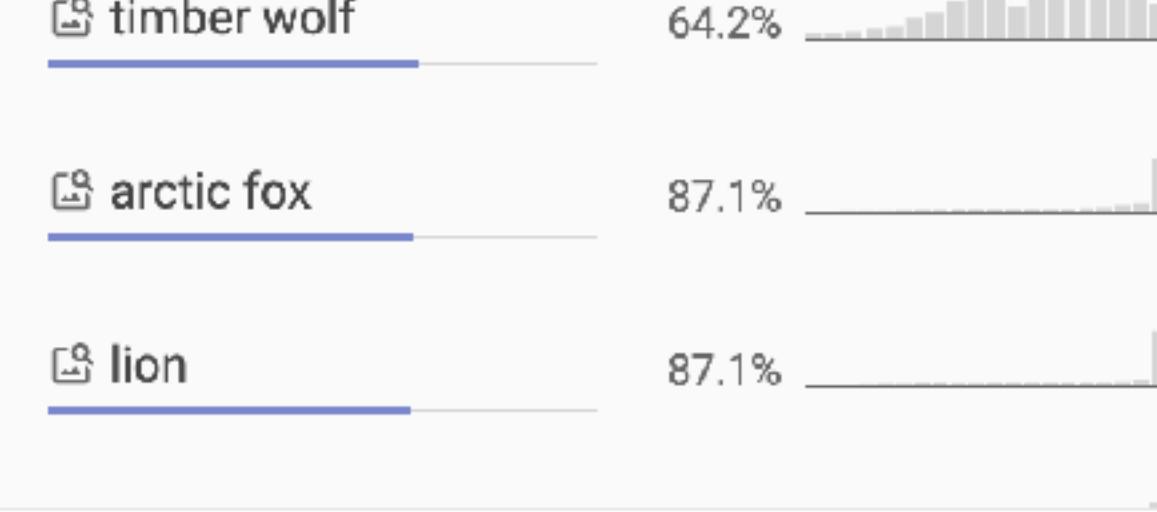


The intersection of  
**brown bear** and **black bear**.  
Both classes share some **bear-ness**.









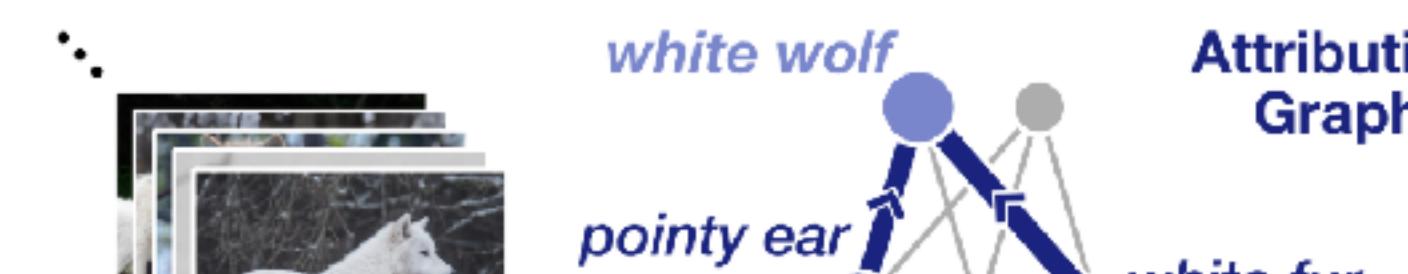
## What is **SUMMIT**?

Understanding how neural networks make predictions remains a fundamental challenge. Existing work on interpreting neural network predictions for images often focuses on explaining predictions for single images or neurons, yet predictions are computed from millions of weights optimized over millions of images—such explanations can easily miss a bigger picture.

We present **SUMMIT**, an interactive visualization that scalably summarizes what features a deep learning model has learned and how those features interact to make predictions.

## How does it work?

**SUMMIT** introduces two new scalable summarization techniques that aggregate activations and neuron-influences to create *attribution graphs*: a class-specific visualization that simultaneously highlights *what* features a neural network detects and *how* they are related.



Our work joins a growing body of open-access research that aims to use interactive visualization to explain complex inner workings of modern machine learning techniques. We believe our summarization approach that builds entire class representations is an important step for developing higher-level explanations for neural networks. We hope our work will inspire deeper engagement from both the information visualization and machine learning communities to further develop human-centered tools for artificial intelligence.

## Credits

**SUMMIT** was created by [Fred Hohman](#), [Haekyu Park](#), [Caleb Robinson](#), and [Polo Chau](#) at Georgia Tech. We also thank Nilaksh Das and the Georgia Tech Visualization Lab for their support and constructive feedback. This work is supported by a NASA Space Technology Research Fellowship and NSF grants IIS-1563816, CNS-1704701, and TWC-1526254.



**Summit: Scaling Deep Learning Interpretability by Visualizing Activation and Attribution Summarizations**  
[Fred Hohman](#), [Haekyu Park](#), [Caleb Robinson](#), and [Duen Horng \(Polo\) Chau](#).  
*IEEE Transactions on Visualization and Computer Graphics (TVCG, Proc. VAST'19)*, 2020.

⚠ **Live demo:** [fredhohman.com/summit](http://fredhohman.com/summit)

📘 **Paper:** <https://fredhohman.com/papers/19-summit-vast.pdf>

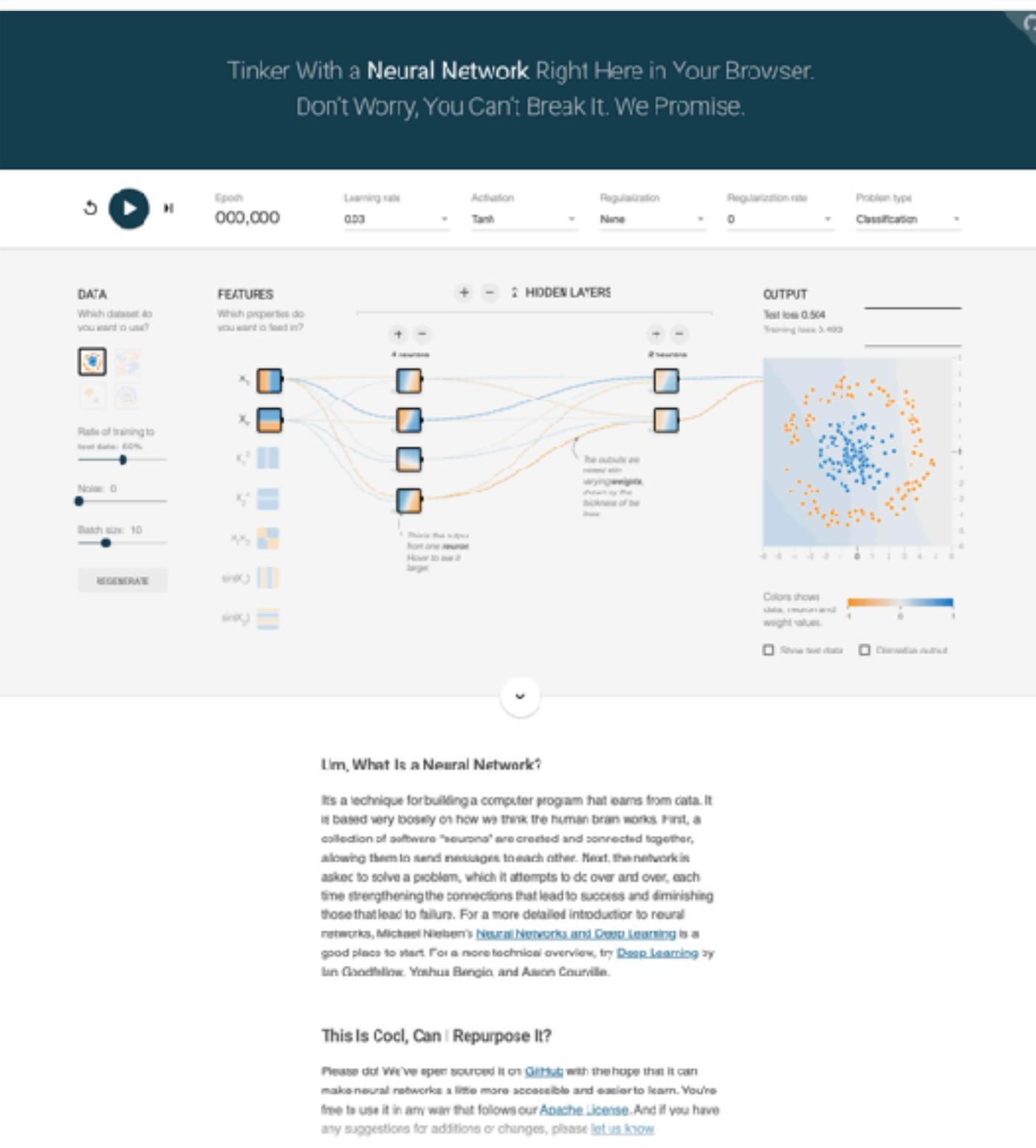
🎥 **Video:** <https://youtu.be/J4GMLvoH1ZU>

💻 **Code:** <https://github.com/fredhohman/summit>

🕒 **Slides:** coming October 2019!

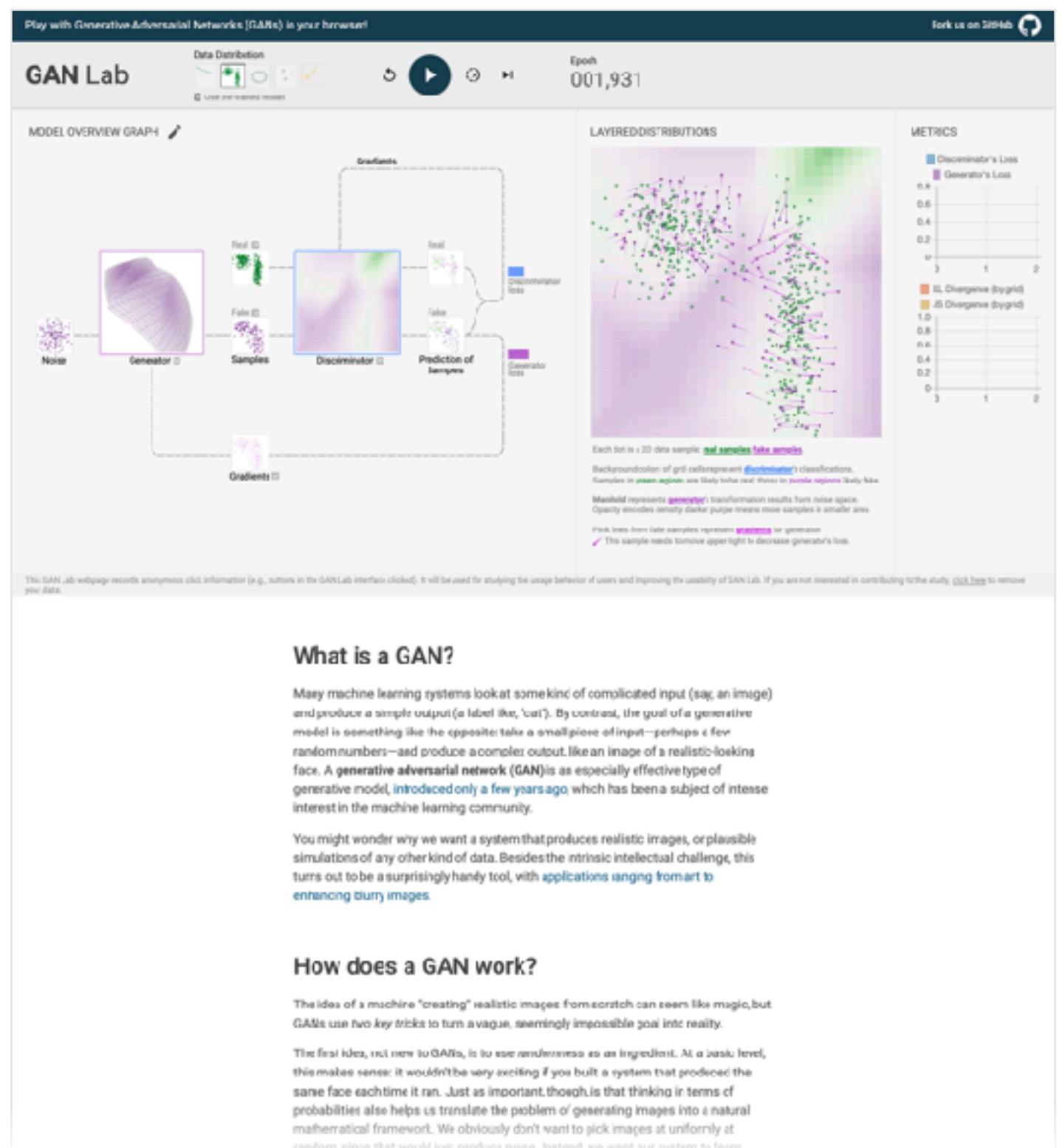
# TensorFlow Playground

2016



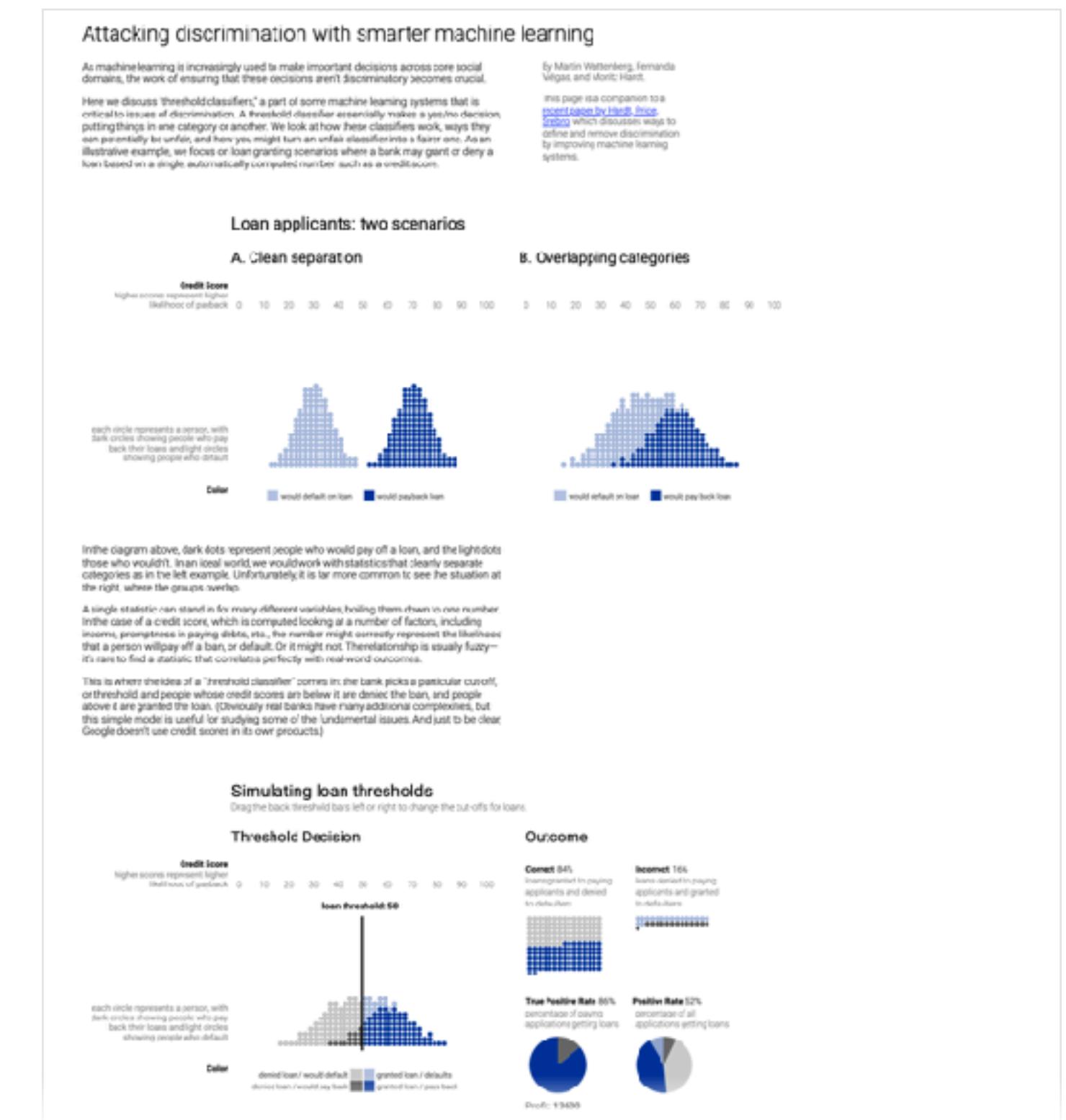
# GAN Lab

2018



# Attacking Discrimination with Smarter Machine Learning

2016



# Interactive Scalable Interfaces for Machine Learning Interpretability

---



## PART I Enable interpretability

**GAMUT** Operationalize interpretability *CHI 2019*

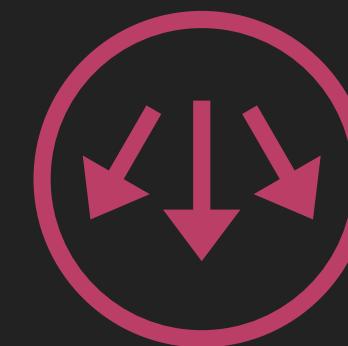
**TELEGAM** Vis + text for better explanations *VIS 2019*



## PART II Scale interpretability

**Interrogative Survey** Summarize interpretability vis *TVCG 2018*

**SUMMIT** Higher-level explanations for neural networks *VAST 2019*



## PART III **Communicate** interpretability

**ML Literacy** Interactive mediums & platforms *VISCOMM 2019, VISxAI 2018*

**Interactive Articles** Formalizing interactive communication *Distill 2020*

# Interactive Articles for Politics, Pop Culture, ...

## The Pudding

### The Largest Vocabulary In Hip Hop

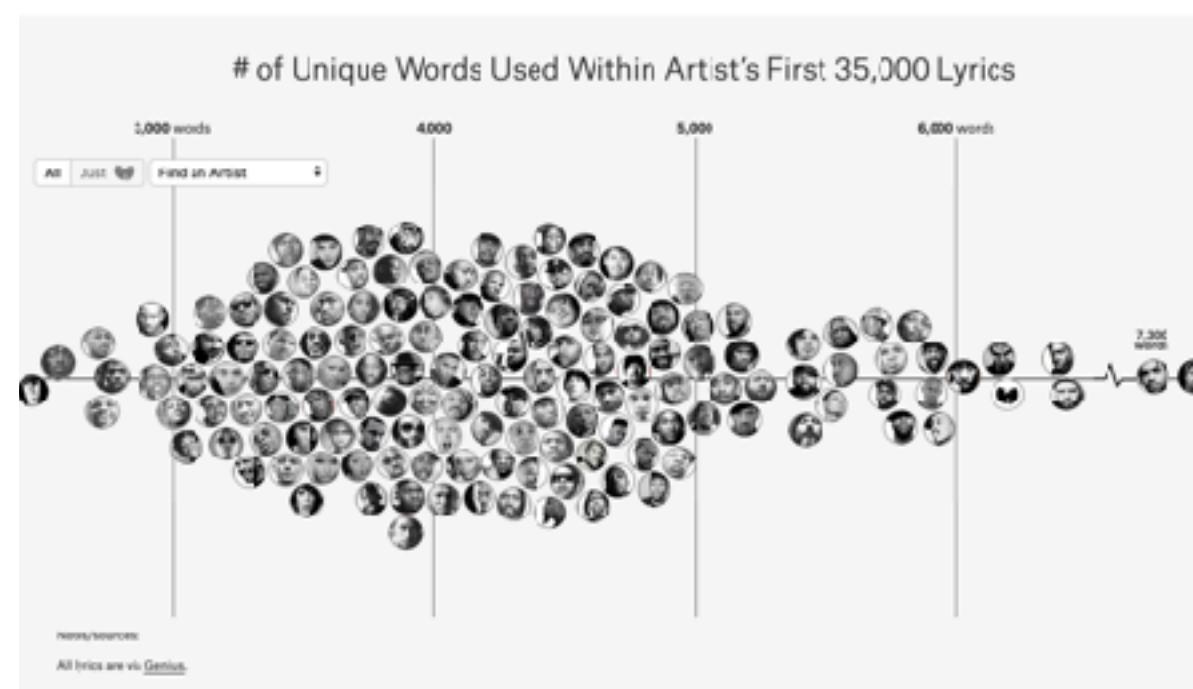
Rappers, ranked by the number of unique words used in their lyrics

By Matt Daniels

Updated on January 21, 2019  
with 70 new rappers including Brockhampton, Deafheaven, Drizzy, Lil Uzi Vert, Travis Scott, and Migos

This project was originally published in 2014 and recently updated in January 2019 with never lyrics data and 75 additional artists, including Lil Uzi Vert, Lil Yachty, Migos, and 21 Savage.

It compares the number of unique words used by some of the most famous artists in hip hop that is, an example of a quantitative view of lyricalism, once proposed by Tahir Hemphill. I used each artist's first 35,000 lyrics. This way, prolific artists, such as Jay-Z, can be compared to lesser artists, such as Drake.



35,000 words covers 2 to 5 studio albums and EPs. I included mixtapes if the artist was short of the 35,000 words. Quite a few rappers don't have enough official material to be included (for example, Biggie, Chance the Rapper, Queen Latifah, and Erykah Badu).

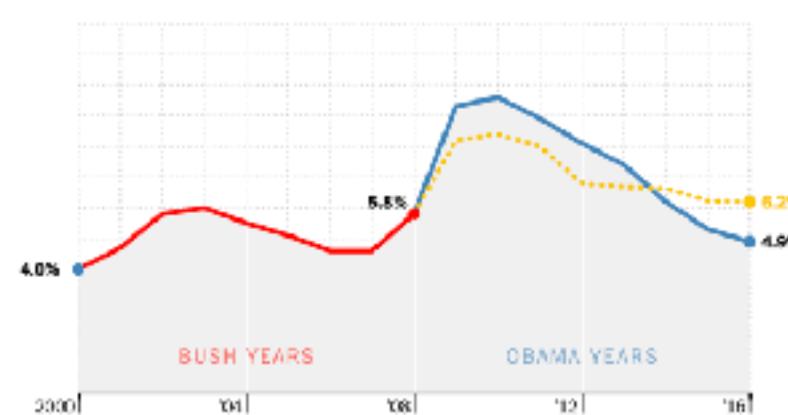
## The New York Times

### You Draw It: What Got Better or Worse During Obama's Presidency

By LINDY RICHARDSON, HANNAH PHILLIPS and ASHLEY PHARRELL JUNE 13, 2017

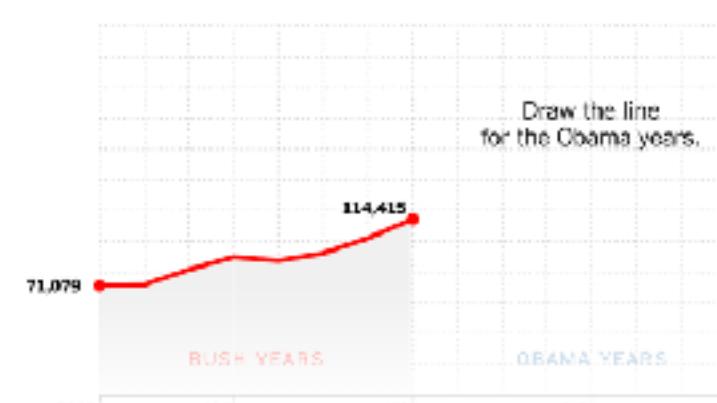
Draw your guesses on the charts below to see if you're as smart as you think you are.

Under President Obama, the **unemployment rate** ...



... reached its **lowest level since 2007**. The current rate is a stunning decline from the 9.3 percent in 2009, the year Mr. Obama took office.

Under Mr. Obama, the **number of immigrants convicted of crimes who were deported** ...



## FiveThirtyEight

### Forecasting the race for the House

Updated Nov. 3, 2018, at 11:00 PM

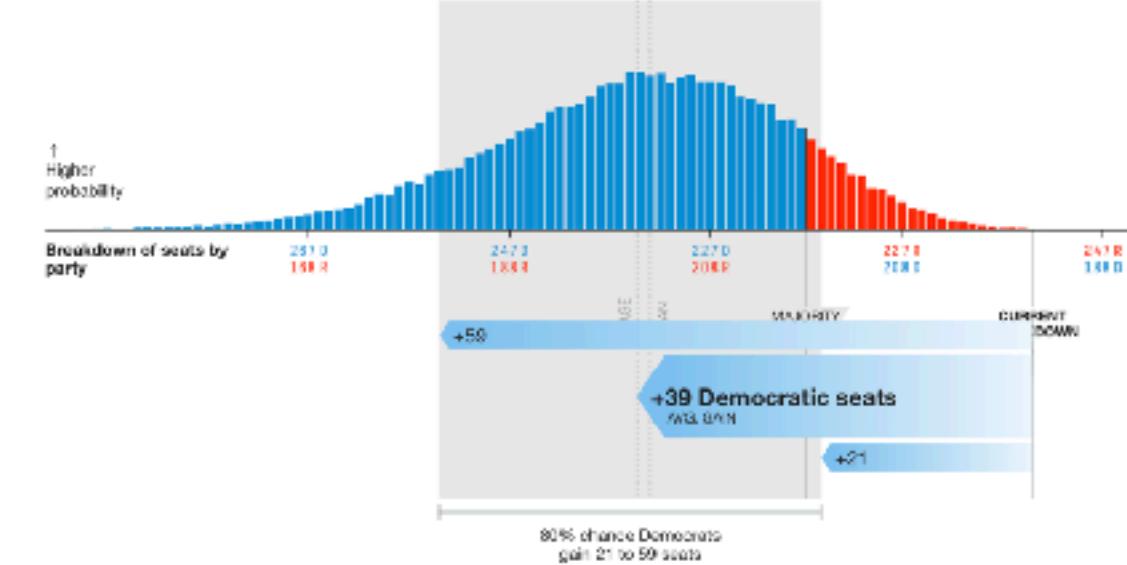
**7 in 8**

Chance Democrats win control (87.9%)

f v

**1 in 8**

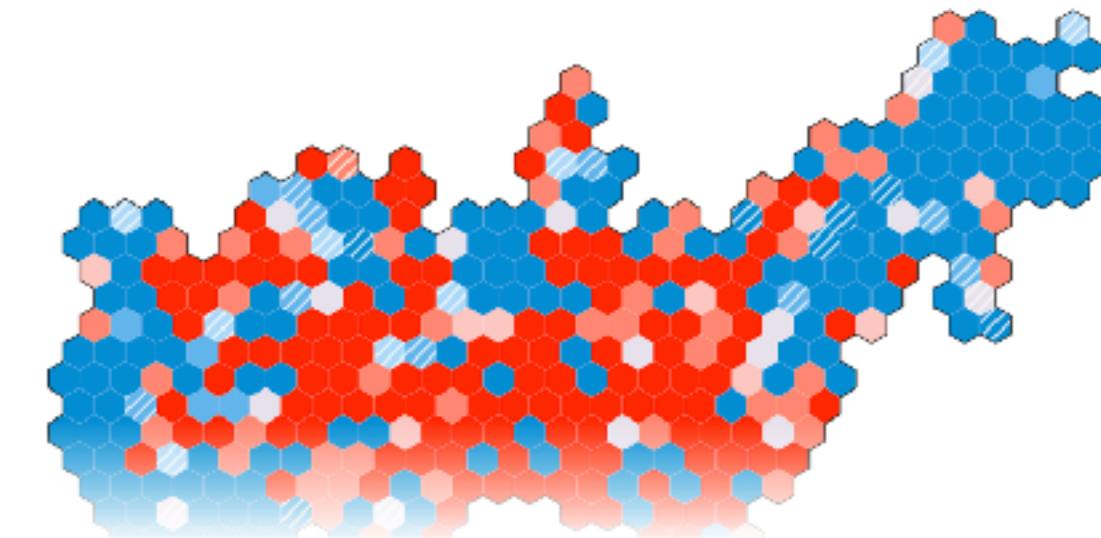
Chance Republicans keep control (12.1%)



### Our forecast for every district

The chance of each candidate winning, with all 435 House districts shown at the same size

Cartogram Map

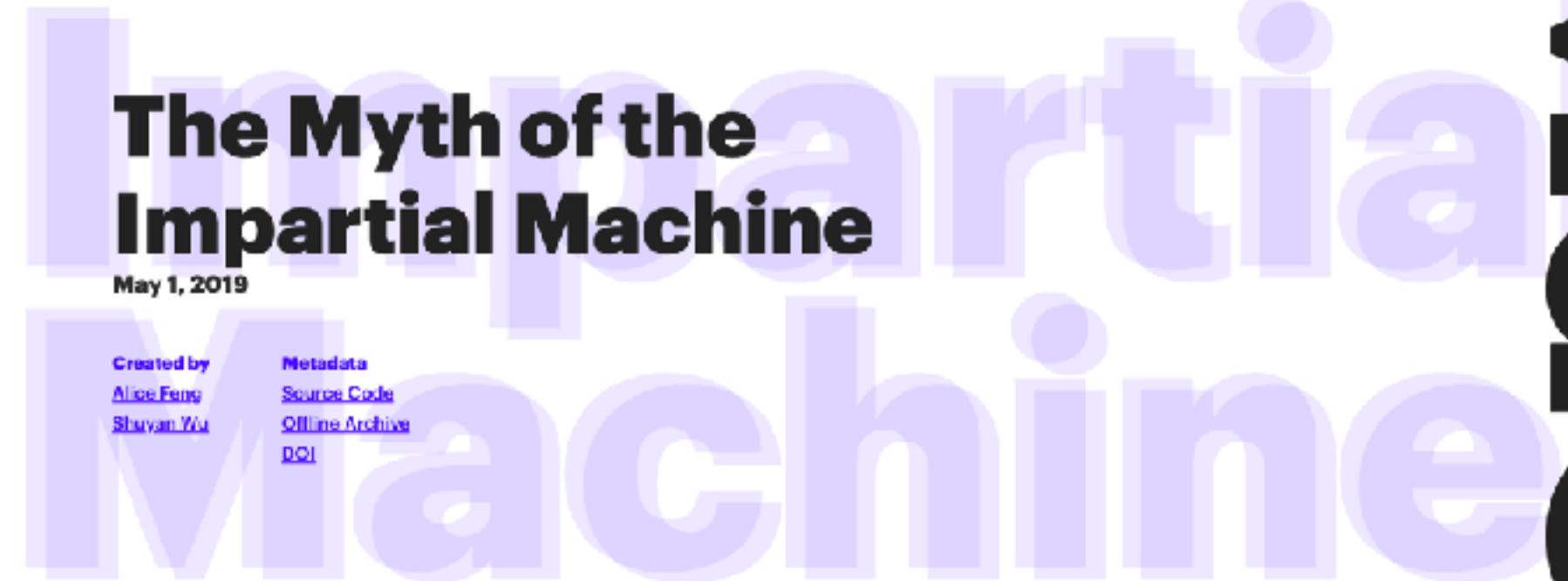


# Interactive Articles now for Machine Learning!

Parametric Press

Issue 01  
Science + Society

Table of Contents



Wide-ranging applications of data science bring utopian proposals of a world free from bias, but in reality, machine learning models reproduce the inequalities that shape the data they're fed. Can programmers free their models from prejudice?

From voice assistants to image recognition, fraud detection to social media feeds, machine learning (ML) and artificial intelligence (AI) are becoming an increasingly important part of society. The two fields have made enormous strides in recent years thanks to gains in computing power and the so-called "information explosion." Such algorithms are being used in fields as varied as medicine, agriculture, insurance, transportation and art, and the number of companies rushing to embrace what ML and AI can offer has increased rapidly in recent years.

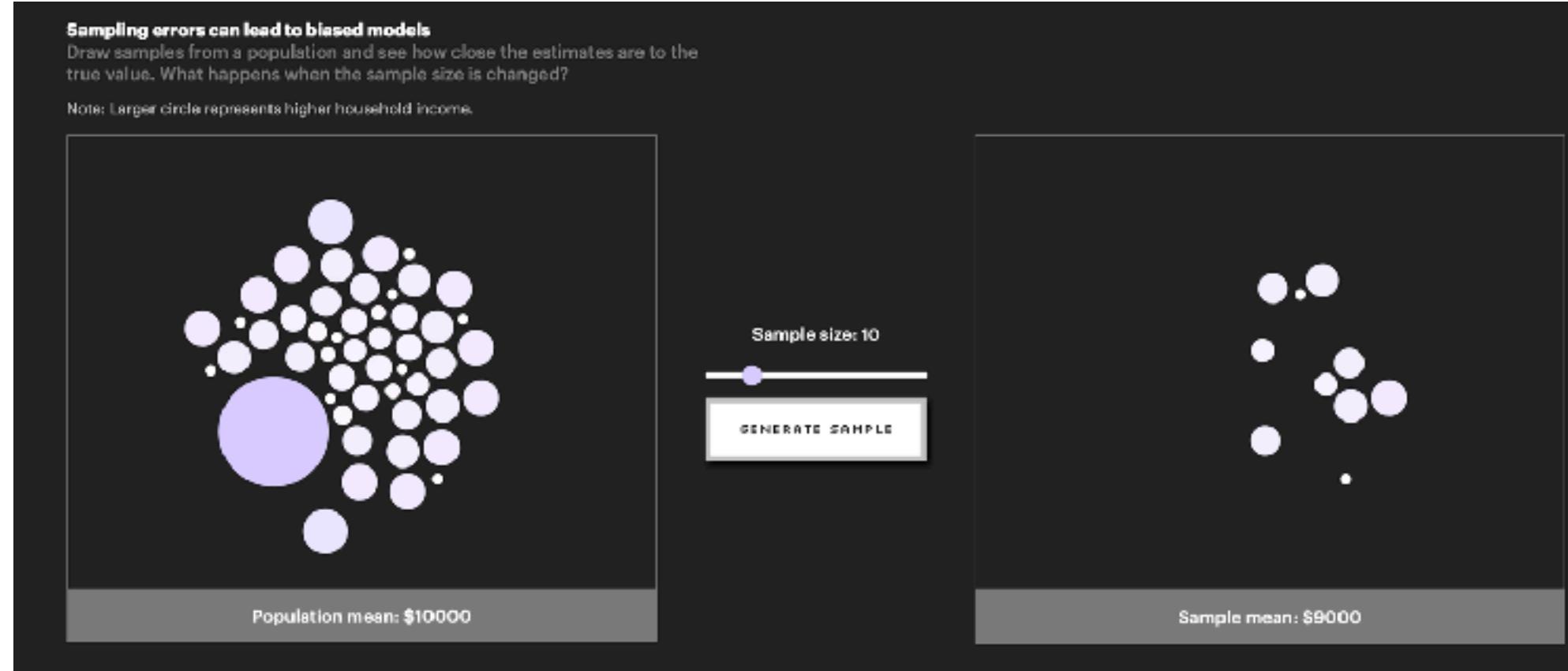
According to a [survey](#) conducted by Teradata in July 2017, 80% of enterprises have already begun investing in AI technologies and 30% plan to increase their spending over the next 36 months. Investment in such models is also [forecasted to grow](#) from \$12 billion in 2017 to over \$50 billion by 2021. Billed as being more

Impartial  
Machine  
Science  
Society  
Machine  
Learning  
Bias  
Prejudice  
Diversity  
Inclusion  
Equity  
Ethics  
Bias  
Prejudice  
Diversity  
Inclusion  
Equity  
Ethics

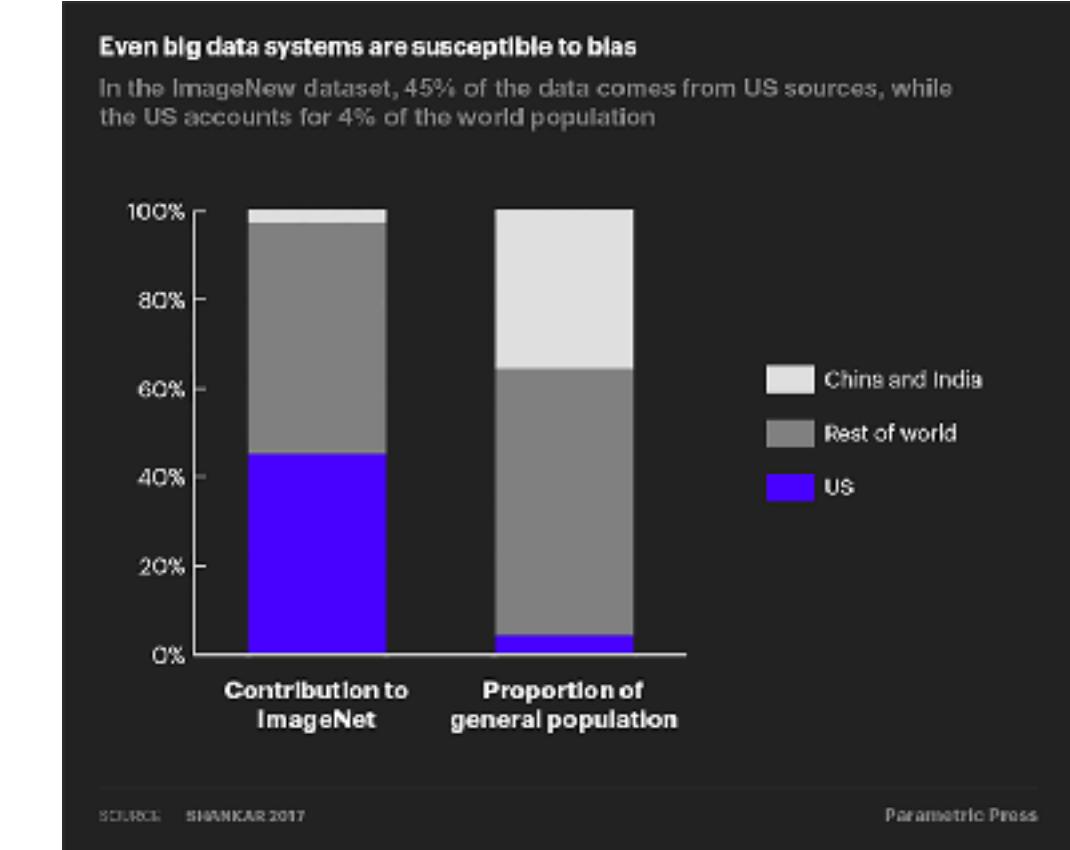




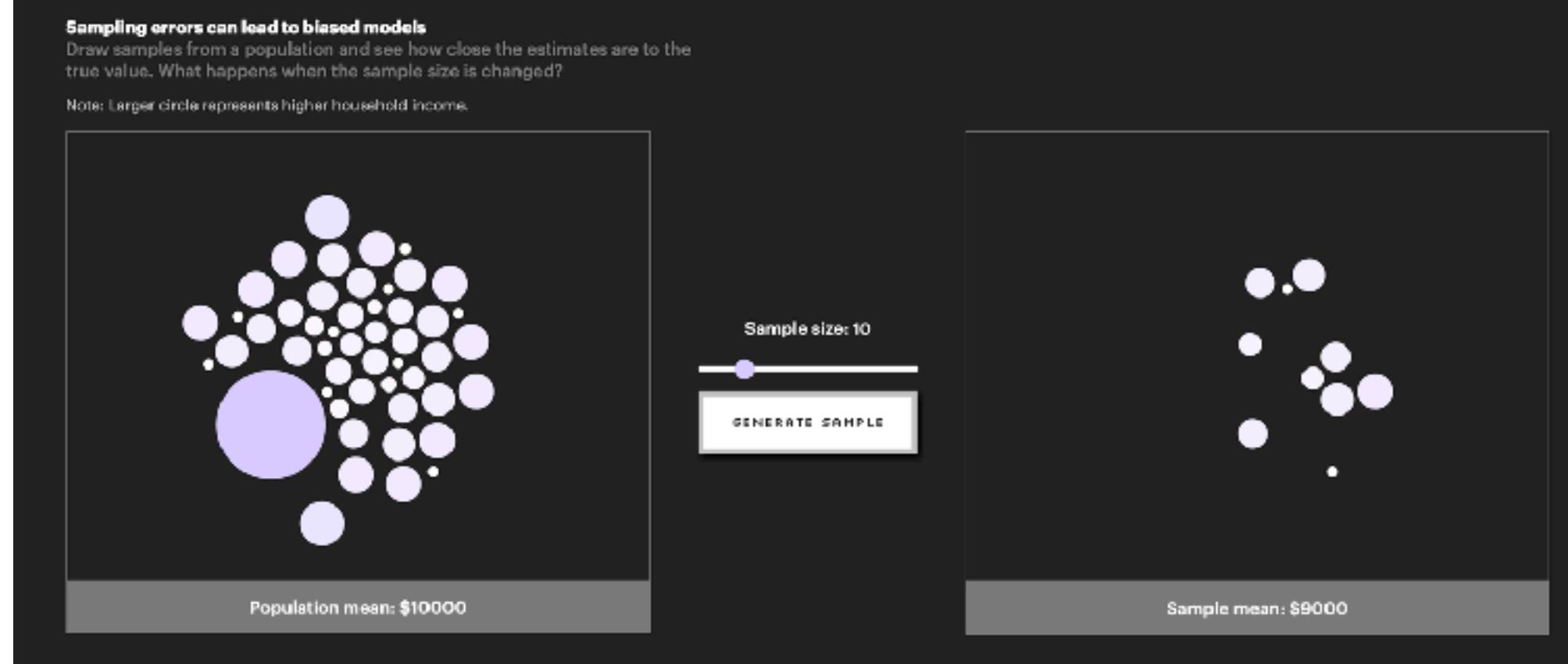
# Interactive Graphics



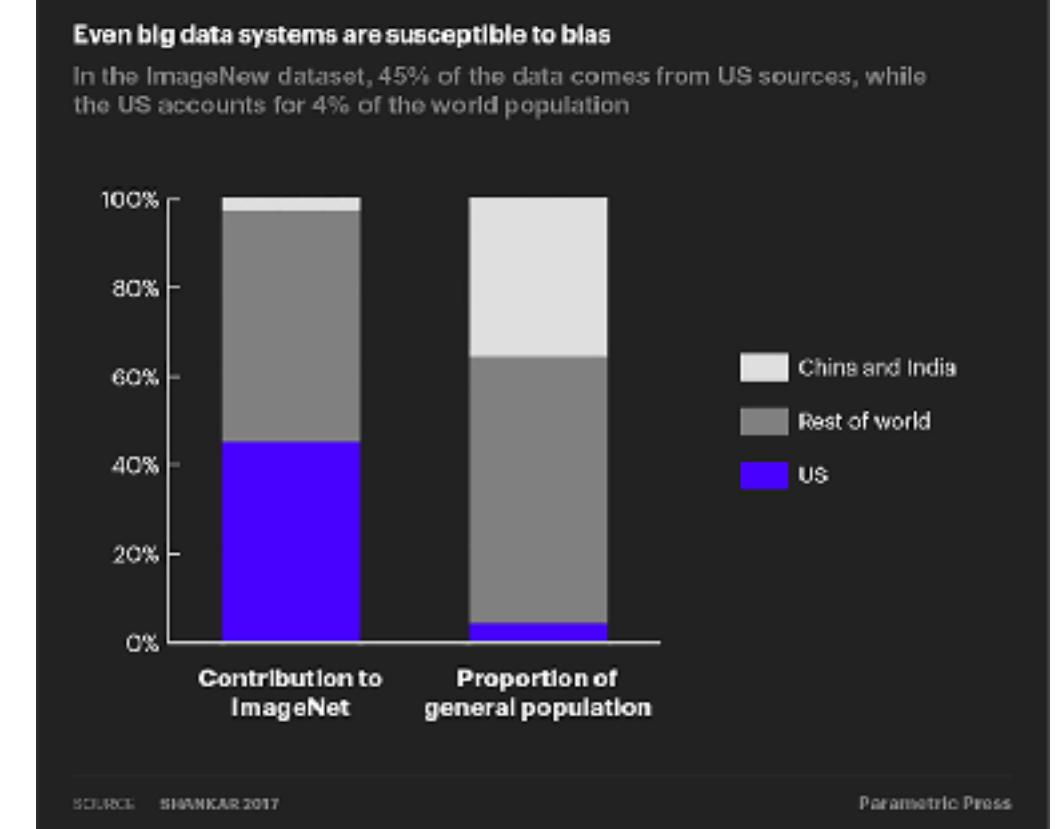
# Visualization



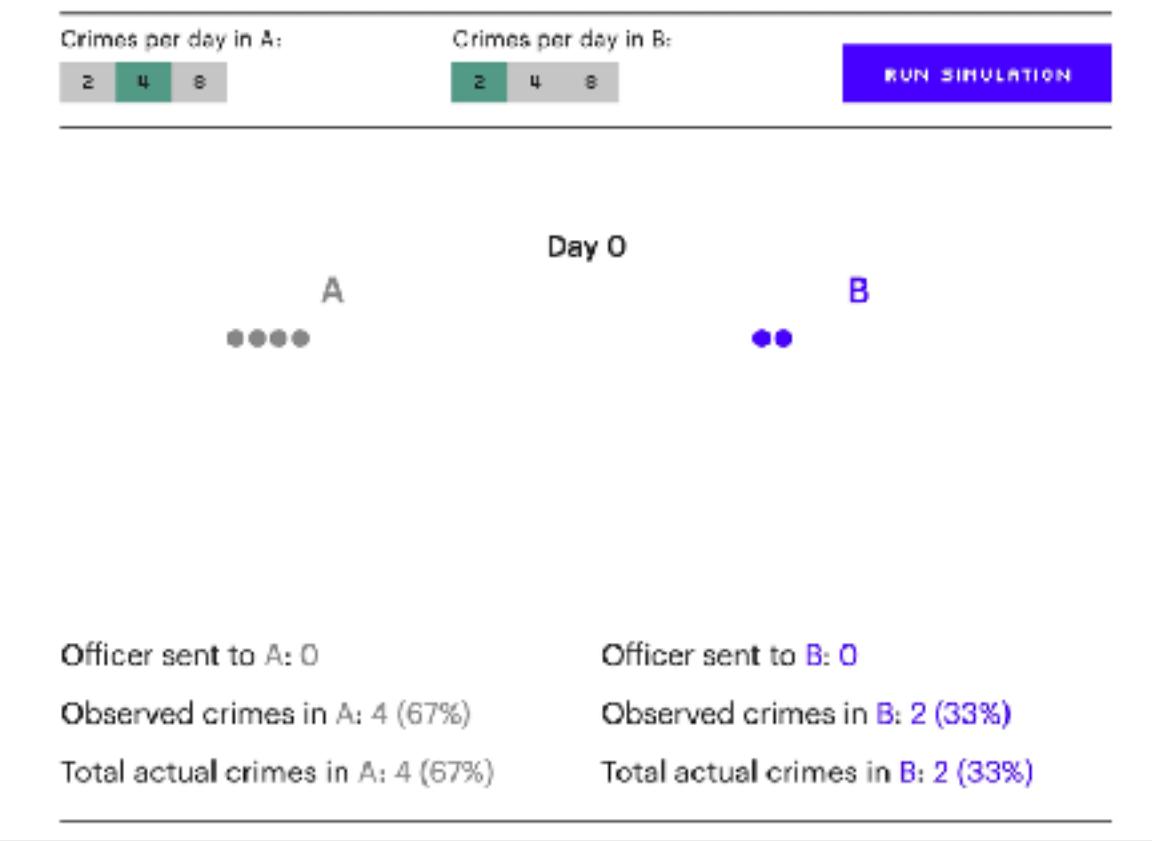
# Interactive Graphics



# Visualization



# Simulation



Parametric Press  
Interactive Data Visualization

**The Myth of the Impartial Machine**

Table of Contents

This page contains a collection of interactive data visualizations from the book "The Myth of the Impartial Machine". Each visualization is designed to demonstrate a specific concept related to machine learning and data bias.

Sampling errors can lead to biased models

Even big data systems are susceptible to bias

Crimes per day in A: 2 4 8

Crimes per day in B: 2 4 8

RUN SIMULATION

Day 0

A

B

Officer sent to A: 0

Officer sent to B: 0

Observed crimes in A: 4 (67%)

Observed crimes in B: 2 (33%)

Total actual crimes in A: 4 (67%)

Total actual crimes in B: 2 (33%)

Population mean: \$10000

Sample size: 10

GENERATE SAMPLE

Sample mean: \$9000

Contribution to ImageNet

Proportion of general population

Chin and India

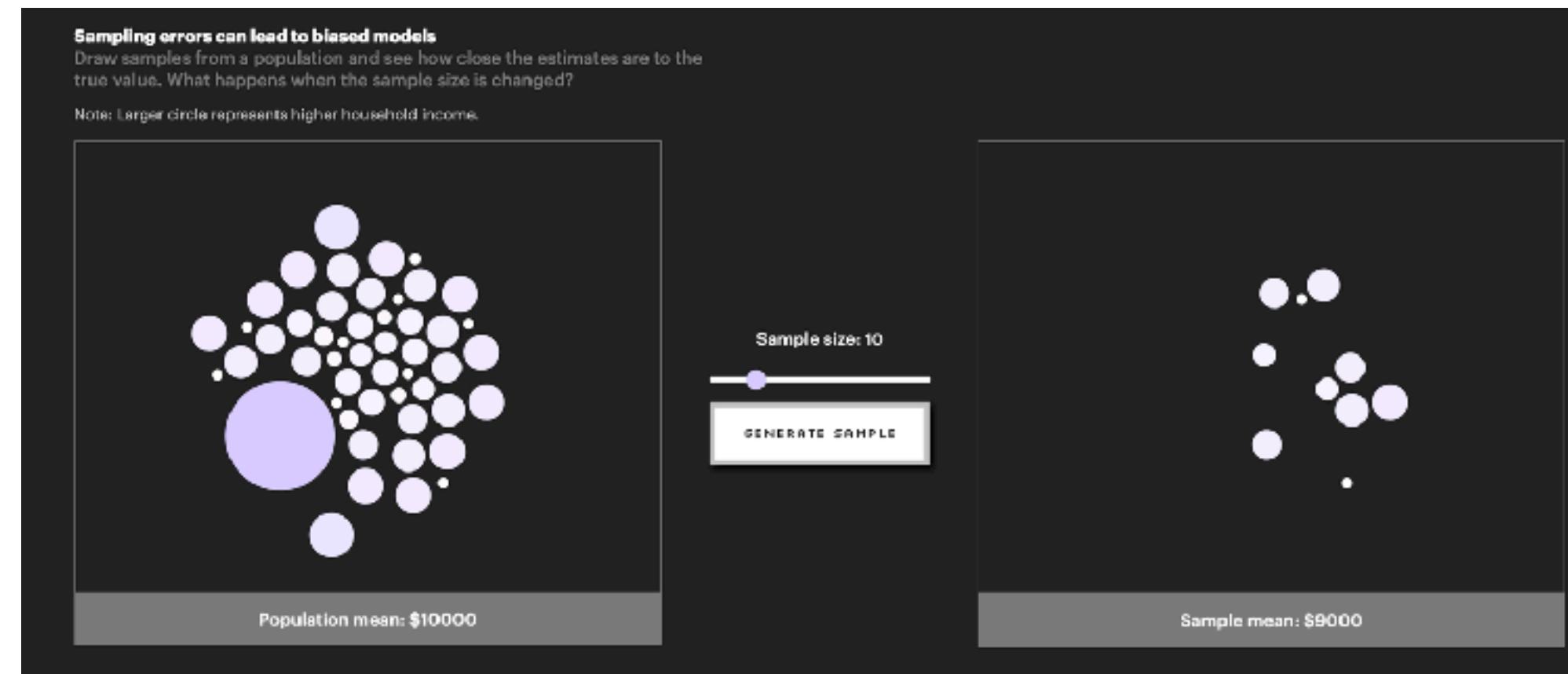
Rest of world

US

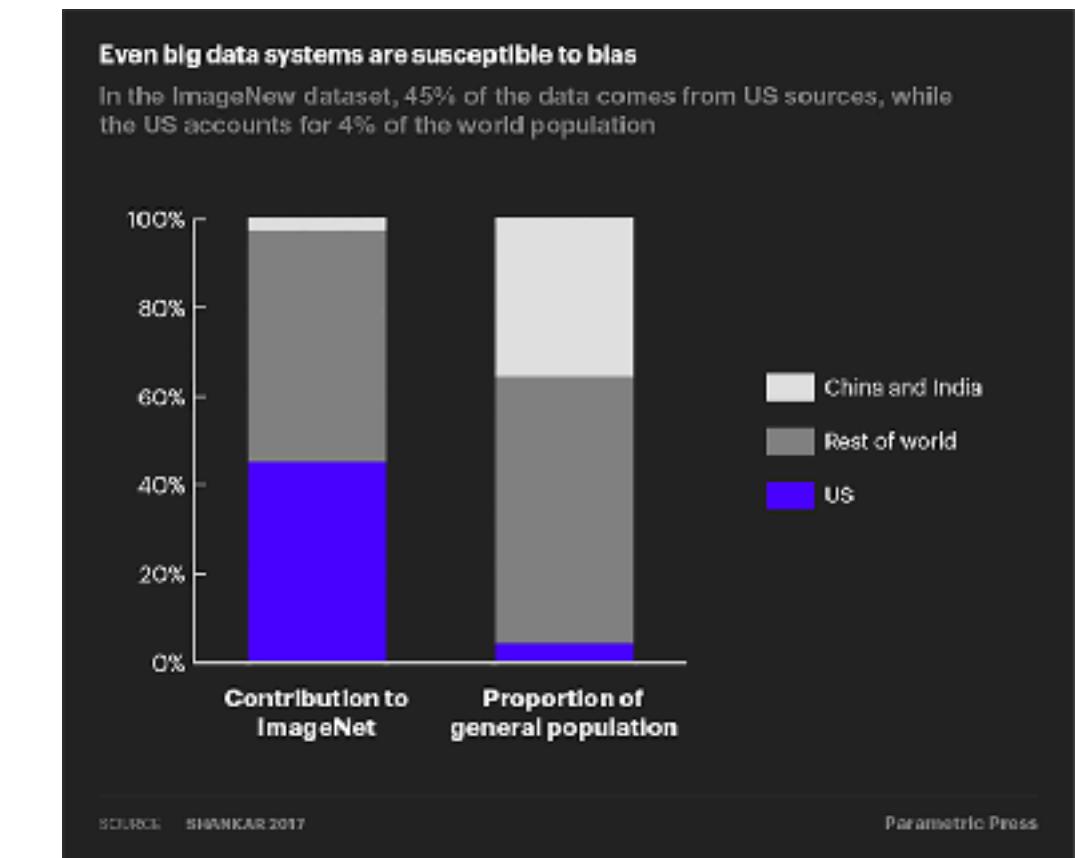
SOURCE: SHANKAR 2017

Parametric Press

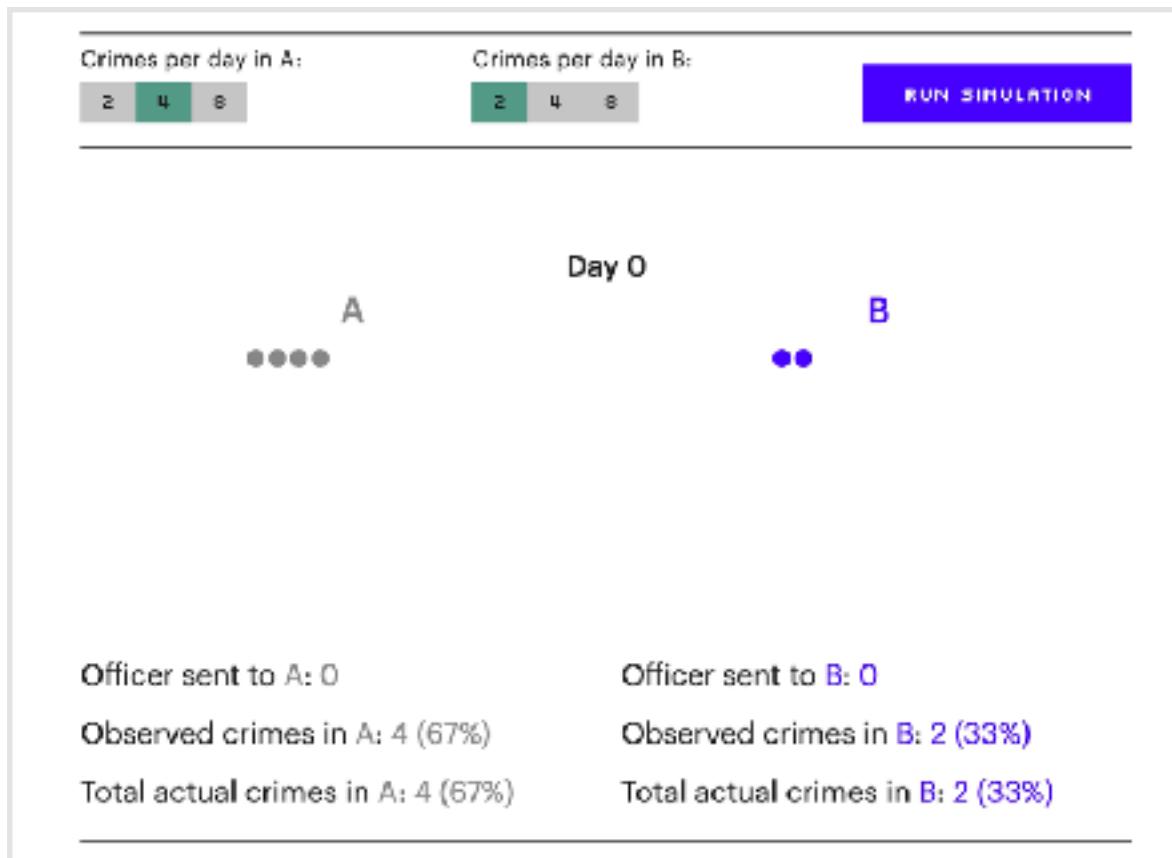
# Interactive Graphics



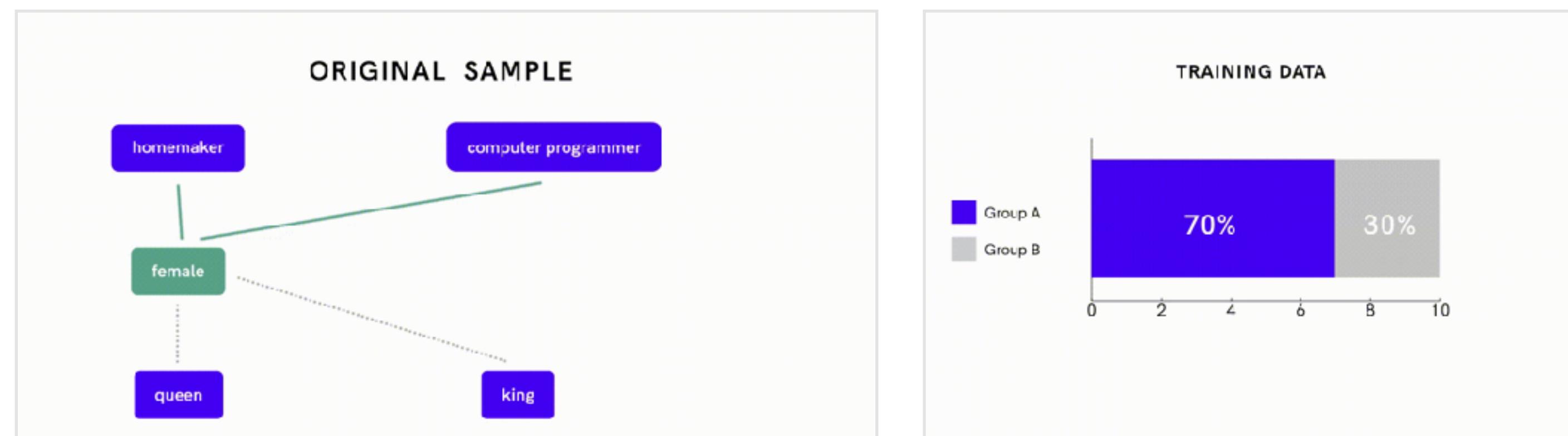
# Visualization



# Simulation

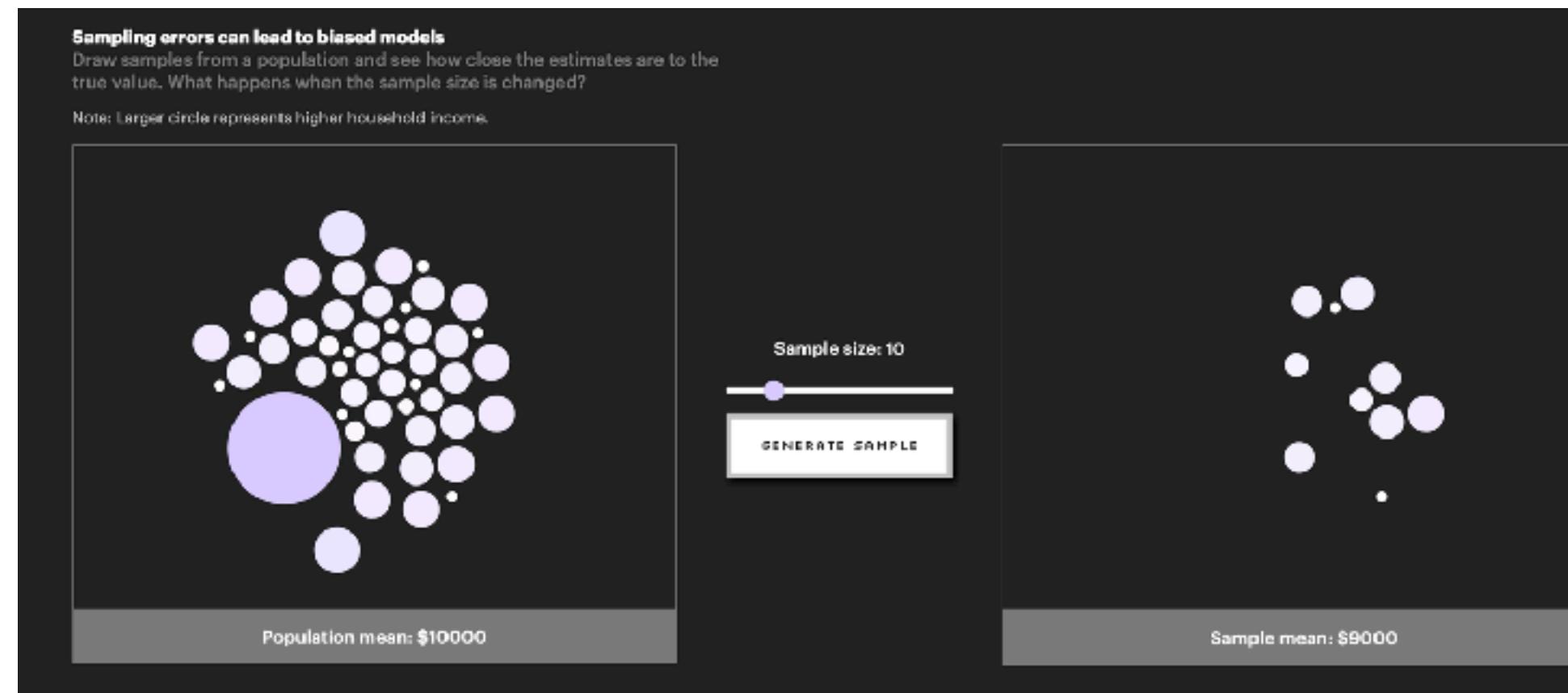


# Animation

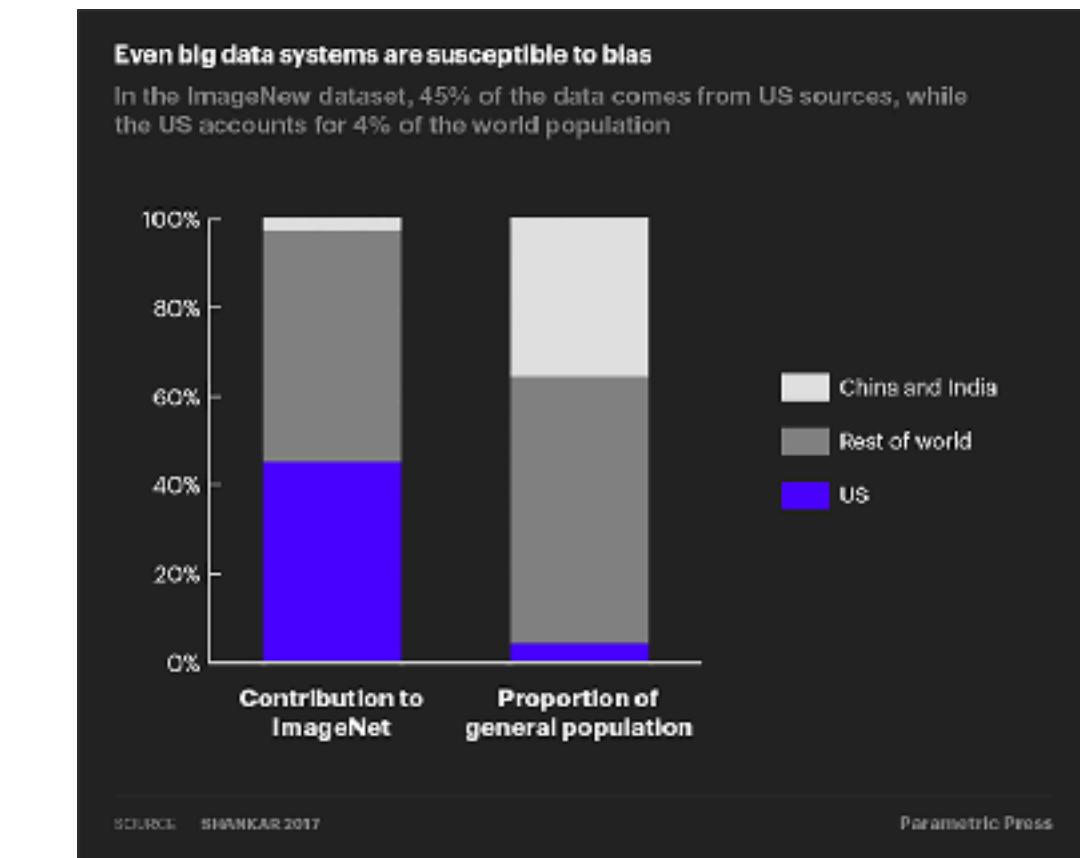




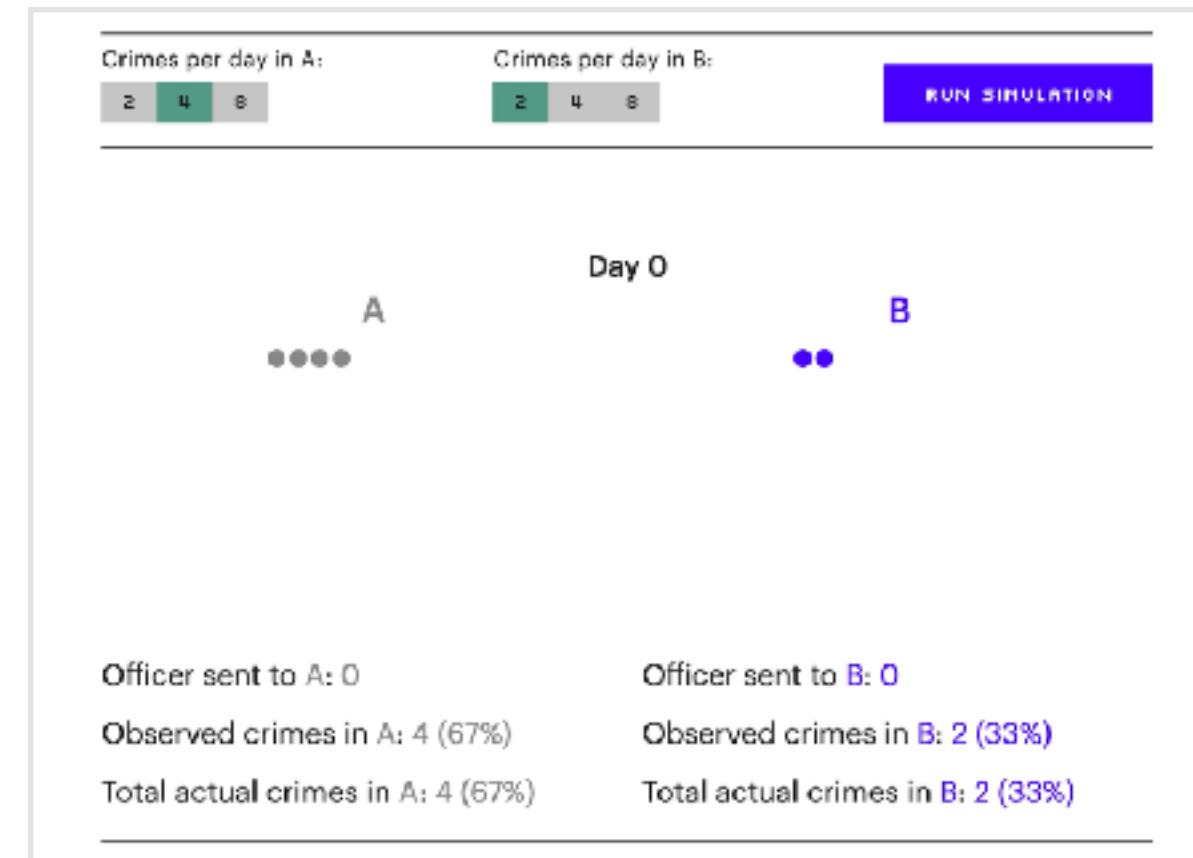
# Interactive Graphics



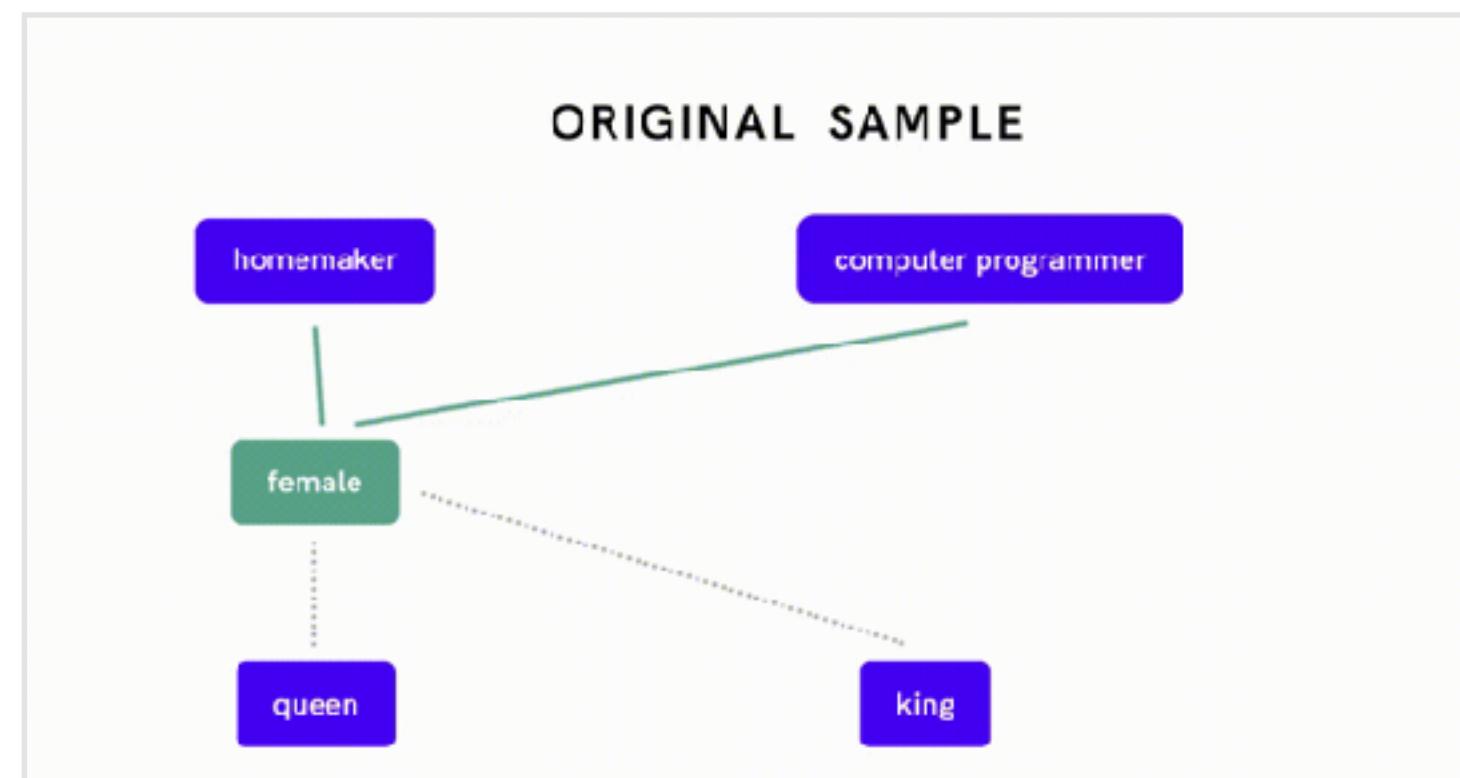
# Visualization



# Simulation



# Animation



# Modeling

# Co-launched Parametric Press

*Interactive publishing platform & zine to break down complex topics.*



**Matthew Conlen**

University of Washington



**Fred Hohman**

Georgia Tech

parametric.press

VISCOMM 2019

Parametric Press  
Issue 01  
Science + Society  
Spring 2019

Unraveling the JPEG  
The Myth of the Impartial Machine  
Data Science for Fair Housing  
Flatland Follies: An Adjunct Simulator  
Parametric Press

FAST COMPANY

CO.DESIGN TECH WORK LIFE CREATIVITY IMPACT AUDIO VIDEO NEWS

05.29.19

## The secret life of a JPEG

What's inside all those cat photos? Let this interactive explainer be your guide.

stackoverflow

Latest Newsletter Podcast stackoverflow.com

Explore the interactive demo [here](#). [Image: Parametric Press]

BY JOHN PAVLUS 3 MINUTE READ

FLOWINGDATA

Membership Tutorials Courses Projects

June 11, 2019

## Myth of the impartial machine

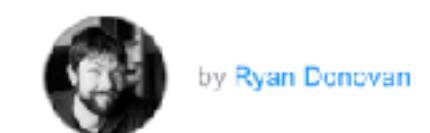
Sampling errors can lead to biased models. Under what circumstances does this happen? What happens when the sample size is changed?

Topic Statistics / bias, machine

Sample size: 15

Sample mean: 83733

The Interactive News Platform for Everyone



by Ryan Donovan



on November 18, 2019

Parametric Press describes [how bias can](#) [work with data](#):

ceptible to non-sampling errors. A study by found that the United States (which accounts pulation) contributed over 45% of the data for of more than 14 million labelled images.

People loved it!

# Helped Teach People About ML's Impact



jonathan.ai  
@jonathandinu

First issue of @ParametricPress is 🔥🔥🔥

"The Myth of the Impartial Machine" is one of the pithiest and accessible treatments of bias in machine learning... and the interactivity puts it in a league of its own



Jon Fisher  
@sciencejon

This is a fascinating and really important read about how AI and machine learning can amplify bias if not corrected:



Patrick Cleary  
@PatJ.Cleary

An excellent overview of how many biases can factor into "impartial" #ai systems, even despite the best of intentions.



Brett2point0  
@Brett2point0

Sharing my Top 5 favorite #AI reads of 2019, crucial education for 2020 and beyond: "The Myth of the Impartial Machine"



Zak Rogoff  
@Zakkai

Interactive article with sliders and infographics that teaches how #AI bias works.

The Myth of the Impartial Machine

The Myth of the Impartial Machine

Wide-ranging applications of data science bring utopian proposals of a world free from bias, but in reality, machine learning models reproduce the inequalities that ...

parametric.press



Frank Ostermann  
@f\_ostermann

Great interactive introduction to AI/ML and what can possibly go wrong (and how to fix it):



Dheeraj Tommandru  
@realtlsdheeraj

This is great article explaining biases from data gathering to model productionization, with a lot of interesting visual and interactive graphs.



Brent Roraback  
@doublepower

Really good piece - with great interactive visuals - from @ParametricPress on biases in machine learning models.



Callum Flack  
@callumflack

Yum, @ParametricPress. A pertinent article, visualised with D3. And then they have all their articles on Github: open publishing, file an issue! This is smart and it looks great.



Chelsea Waite  
@chelseawaite

It's easy to say AI isn't impartial; it's harder to explain exactly why. Loved this (interactive!) article from @ParametricPress with nuanced & accessible answers

People loved it!

# Helped Teach People About ML's Impact



**jonathan.ai**

@jonathandinu

First issue of [@ParametricPress](#) is 🔥🔥🔥

"The Myth of the Impartial Machine" is one of the pithiest and accessible treatments of bias in machine learning... and the interactivity puts it in a league of its own



**Patrick Cleary**

@PatJCleary

An excellent overview of how many biases can factor into "impartial" #ai systems, even despite the best of intentions.



**Brett2point0**

@Brett2point0

Sharing my Top 5 favorite #AI reads of 2019, crucial education for 2020 and beyond: "The Myth of the Impartial Machine"



The Myth of the Impartial Machine  
Wide-ranging applications of data sci  
from bias, but in reality, machine learn  
parametric.press



**Chelsea Waite**  
@chelseawaite

It's easy to say AI isn't impartial; it's harder to explain exactly why. Loved this (interactive!) article from [@ParametricPress](#) with nuanced & accessible answers



**Frank Ostermann**  
@f\_ostermann

Great interactive introdu  
possibly go wrong (and



**Dheeraj Tommandru**

@realtlsdheeraj

This is great article explaining biases from data gathering to model productionization, with a lot of interesting visual and interactive graphs.



**Brent Roraback**

@doublepower

Really good piece - with great interactive visuals - from [@ParametricPress](#) on biases in machine learning models.



**Callum Flack**

@callumflack

Yum, [@ParametricPress](#). A pertinent article, visualised with D3. And then they have all their articles on GitHub!

People loved it!

# Glimpse at an Interactive Reading Future



Aw this is so cool, you can interactively corrupt a JPEG. People, this is how textbooks should have worked since I was a kid!

This block contains a screenshot of an interactive reading interface. It features a map of a city with various data points highlighted in purple. A sidebar on the left contains text about JPEG corruption and a diagram of a particle's path through a magnetic field. Another sidebar on the right contains text about accelerating particles.

This block contains a screenshot of an interactive reading interface. It features a map of Southeast Asia with various data points highlighted in purple. A sidebar on the left contains text about JPEG corruption and a diagram of a particle's path through a magnetic field. Another sidebar on the right contains text about accelerating particles.

This block contains a screenshot of an interactive reading interface. It features a map of Laos with various data points highlighted in purple. A sidebar on the left contains text about JPEG corruption and a diagram of a particle's path through a magnetic field. Another sidebar on the right contains text about accelerating particles.

This block contains a screenshot of an interactive reading interface. It features a map of the United States with various data points highlighted in purple. A sidebar on the left contains text about JPEG corruption and a diagram of a particle's path through a magnetic field. Another sidebar on the right contains text about accelerating particles.

This block contains a screenshot of an interactive reading interface. It features a map of the United States with various data points highlighted in purple. A sidebar on the left contains text about JPEG corruption and a diagram of a particle's path through a magnetic field. Another sidebar on the right contains text about accelerating particles.

This block contains a screenshot of an interactive reading interface. It features a map of the United States with various data points highlighted in purple. A sidebar on the left contains text about JPEG corruption and a diagram of a particle's path through a magnetic field. Another sidebar on the right contains text about accelerating particles.

*Interact with systems—no setup required. 100% open-source. Article archival.*



Paul Ford  
@ftrain

This publication features great nerd content (see the JPEG editor built into this page) and each file can be downloaded as a WARC file to encourage decentralized archiving of interactive content. It's...just really good and thoughtful.

Interactive Articles: Contribution 1  
Help broad audiences learn about machine learning's  
impact on their lives.

# Understanding (Interactive) Reading Behavior

## *The Beginner's Guide to Dimensionality Reduction*



Matthew Conlen  
University of Washington



Fred Hohman  
Georgia Tech

**The Beginner's Guide to Dimensionality Reduction**

Explore the methods that data scientists use to visualize high-dimensional data.

By: [Matthew Conlen](#) and [Fred Hohman](#)  
July 16, 2018

Dimensionality reduction is a powerful technique used by data scientists to look for hidden structure in data. The method is useful in a number of domains, for example document categorization, protein disorder prediction, and machine learning model debugging<sup>[2]</sup>.

The results of a dimensionality reduction algorithm can be visualized to reveal patterns and clusters of similar or dissimilar data. Even though the data is displayed in only two or three dimensions, structures present in higher dimensions are maintained, at least roughly<sup>[7]</sup>.

The technique is available in many applications, for example Google's [Embedding Projector](#)<sup>[10]</sup> let's you view high-dimensional datasets embedded in two or three dimensions under a variety of different projections.

This guide will teach you how to think about these embeddings, and provide a comparison of some of the most popular dimensionality reduction algorithms used today.

VISXAI 2018

Recorded reader interactions

7,000+ views over 1<sup>st</sup> month

#1 of r/datascience subreddit

# Understanding (Interactive) Reading Behavior

## *The Beginner's Guide to Dimensionality Reduction*



Matthew Conlen  
University of Washington



Fred Hohman  
Georgia Tech



**The Beginner's Guide to Dimensionality Reduction**

Explore the methods that data scientists use to visualize high-dimensional data.

By: [Matthew Conlen](#) and [Fred Hohman](#)  
July 16, 2018

Dimensionality reduction is a powerful technique used by data scientists to look for hidden structure in data. The method is useful in a number of domains, for example document categorization, protein disorder prediction, and machine learning model debugging<sup>[2]</sup>.

The results of a dimensionality reduction algorithm can be visualized to reveal patterns and clusters of similar or dissimilar data. Even though the data is displayed in only two or three dimensions, structures present in higher dimensions are maintained, at least roughly<sup>[7]</sup>.

The technique is available in many applications, for example Google's [Embedding Projector](#)<sup>[10]</sup> let's you view high-dimensional datasets embedded in two or three dimensions under a variety of different projections.

This guide will teach you how to think about these embeddings, and provide a comparison of some of the most popular dimensionality reduction algorithms used today.

VISXAI 2018 

Recorded reader interactions

7,000+ views over 1<sup>st</sup> month

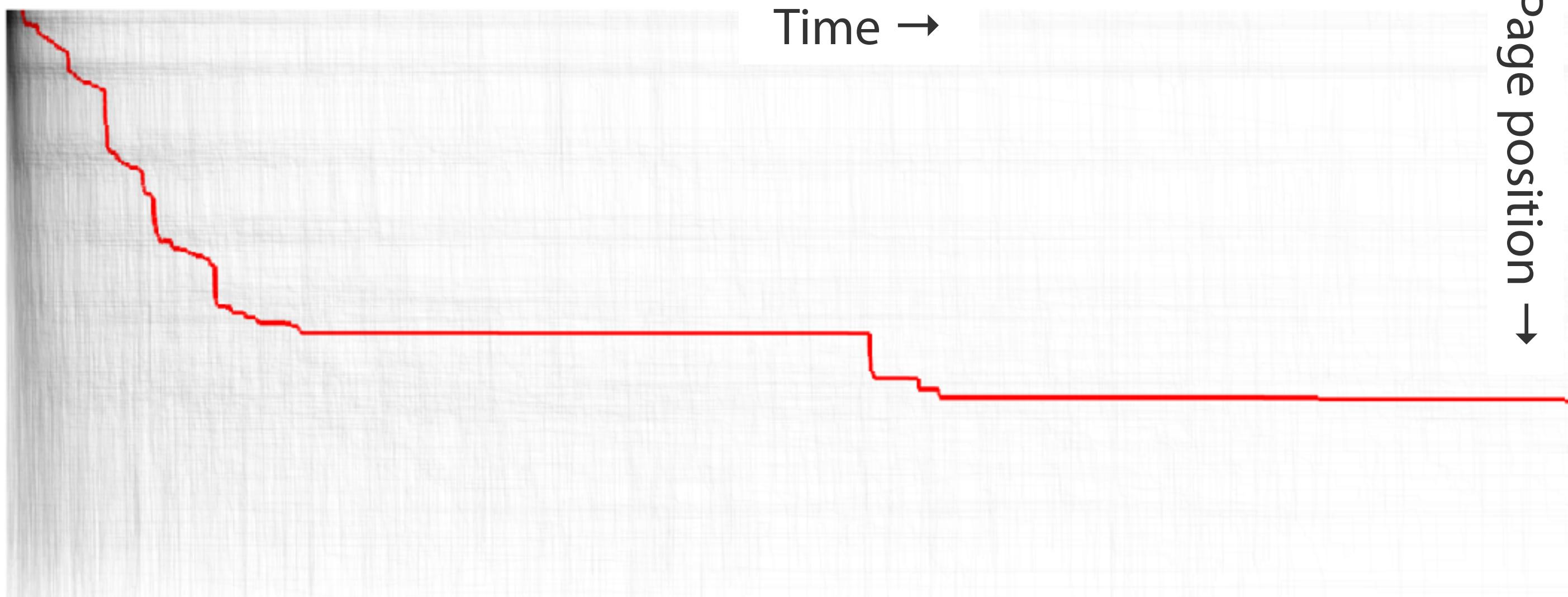
#1 of r/datascience subreddit

Discovery #1

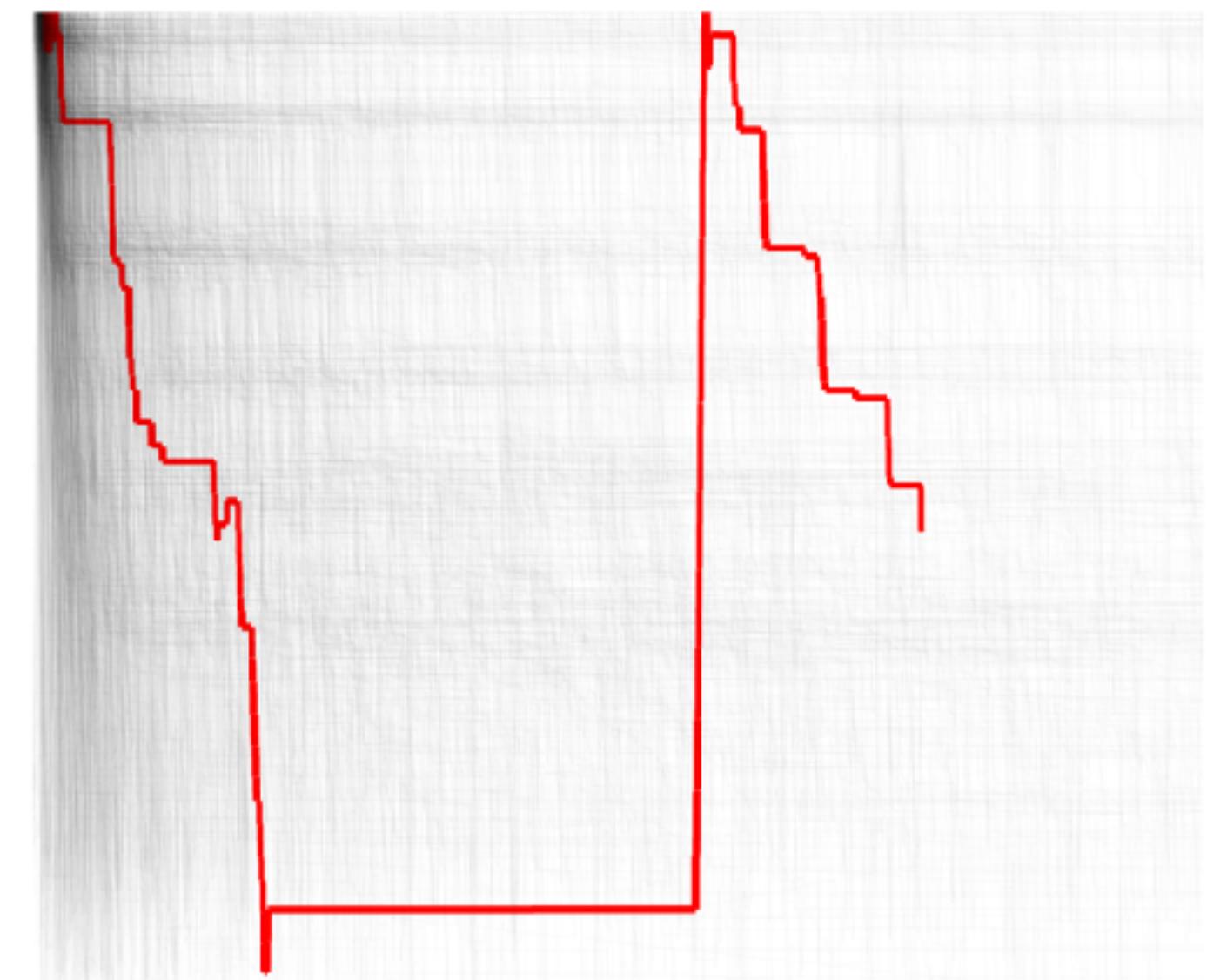
# Revealing Different Reader Patterns

*Critical, yet understudied: How do people read interactive content?*

A. “Scroll & stop”



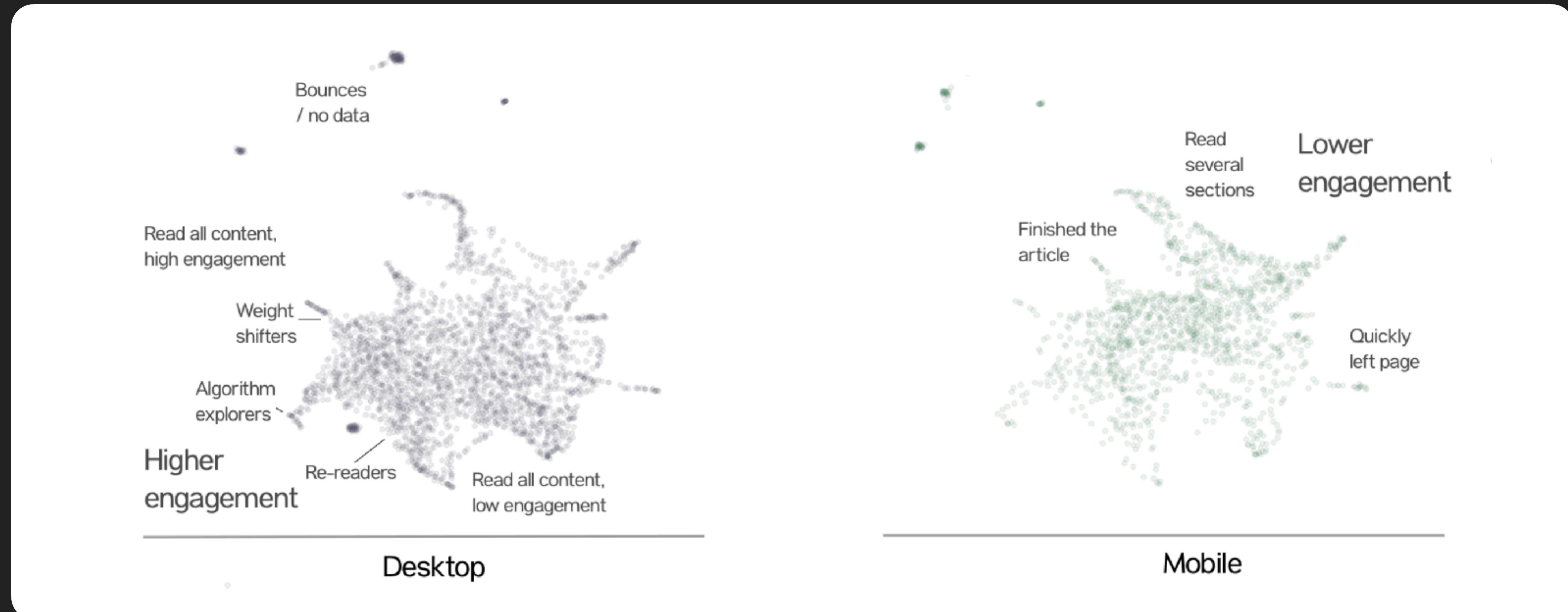
B. “Preview & read”



Discovery #2

# Enable Macro-analysis of Readership

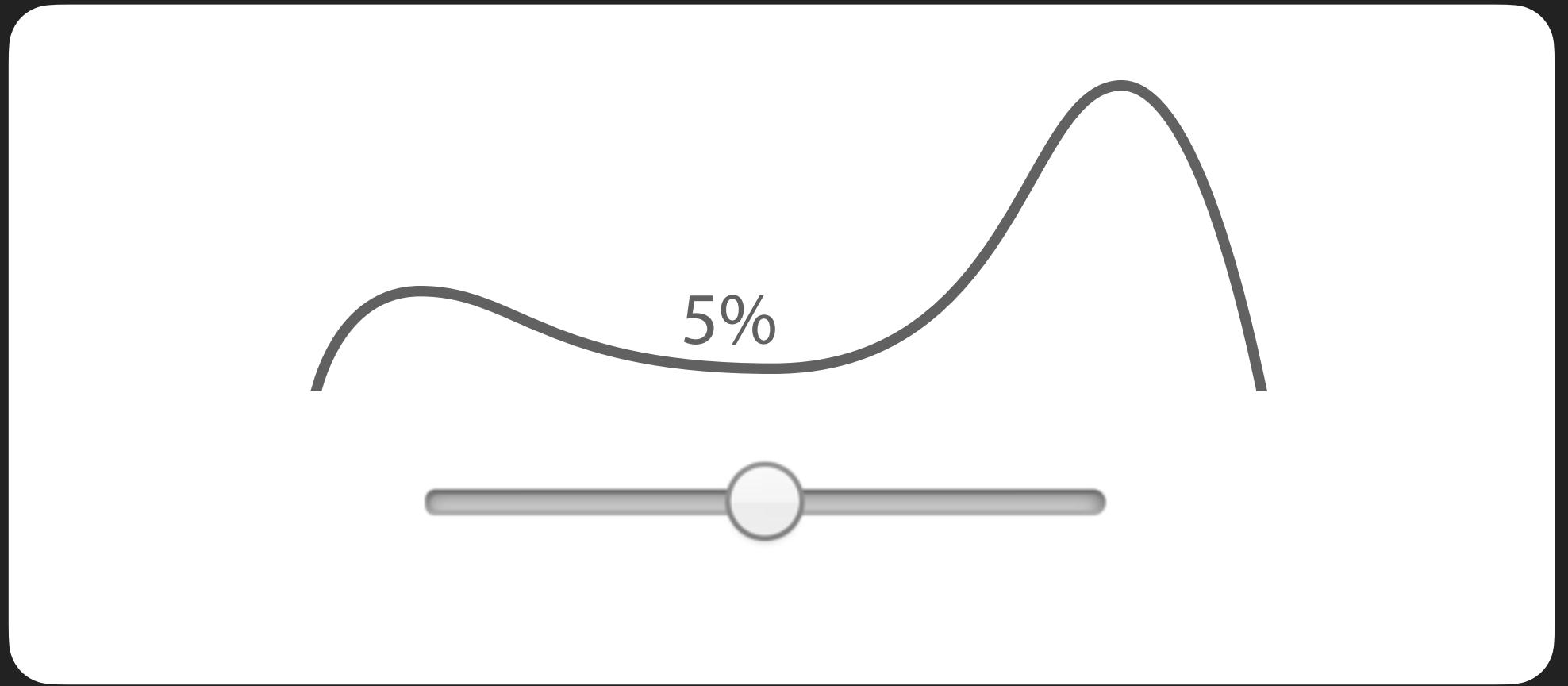
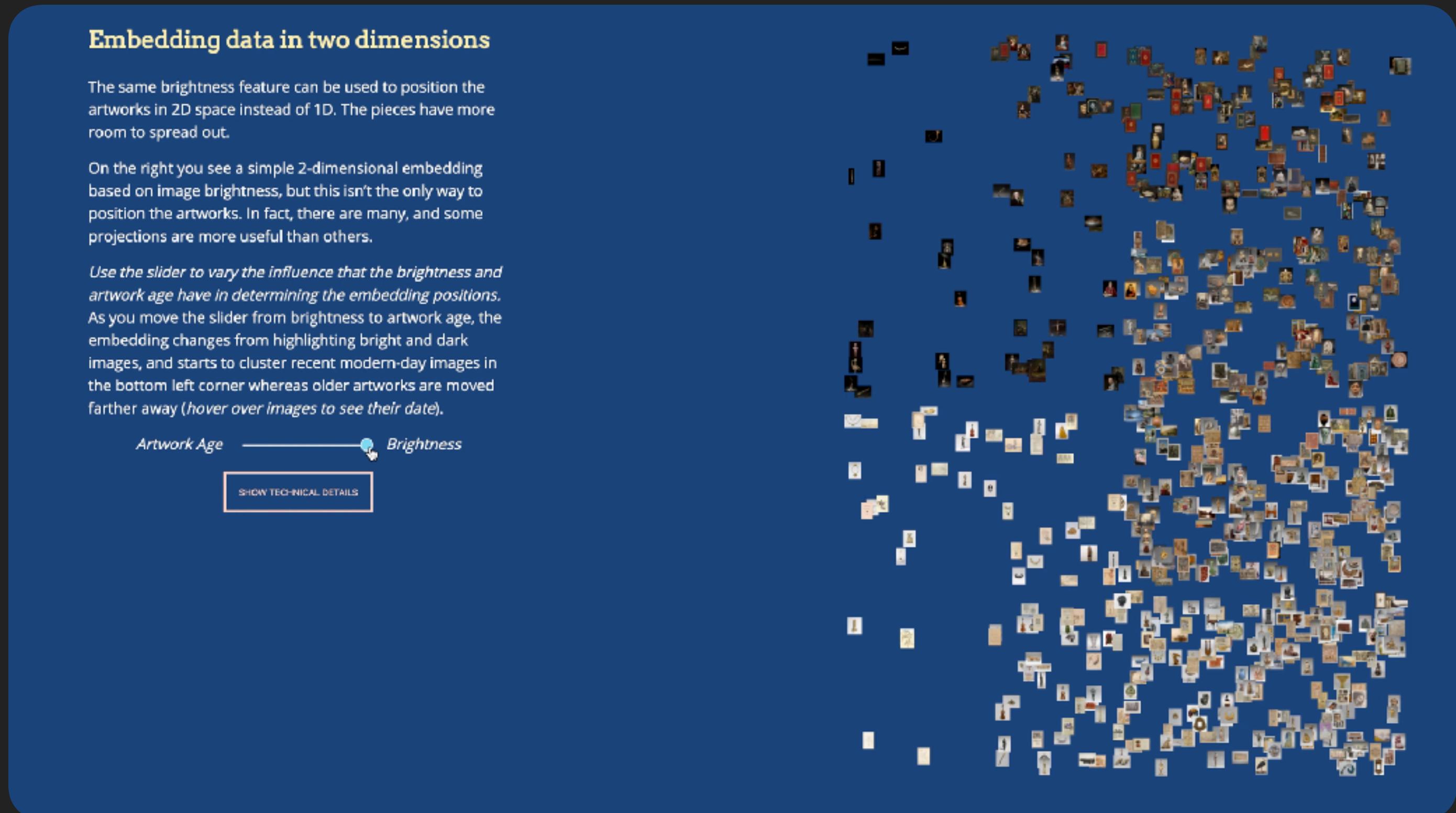
*Reveals broad clusters of reading behavior across device formats.*



Discovery #3

# Content-specific Design Evaluation

*People can read in ways that don't match original design intention.*



**Minimum value distribution**

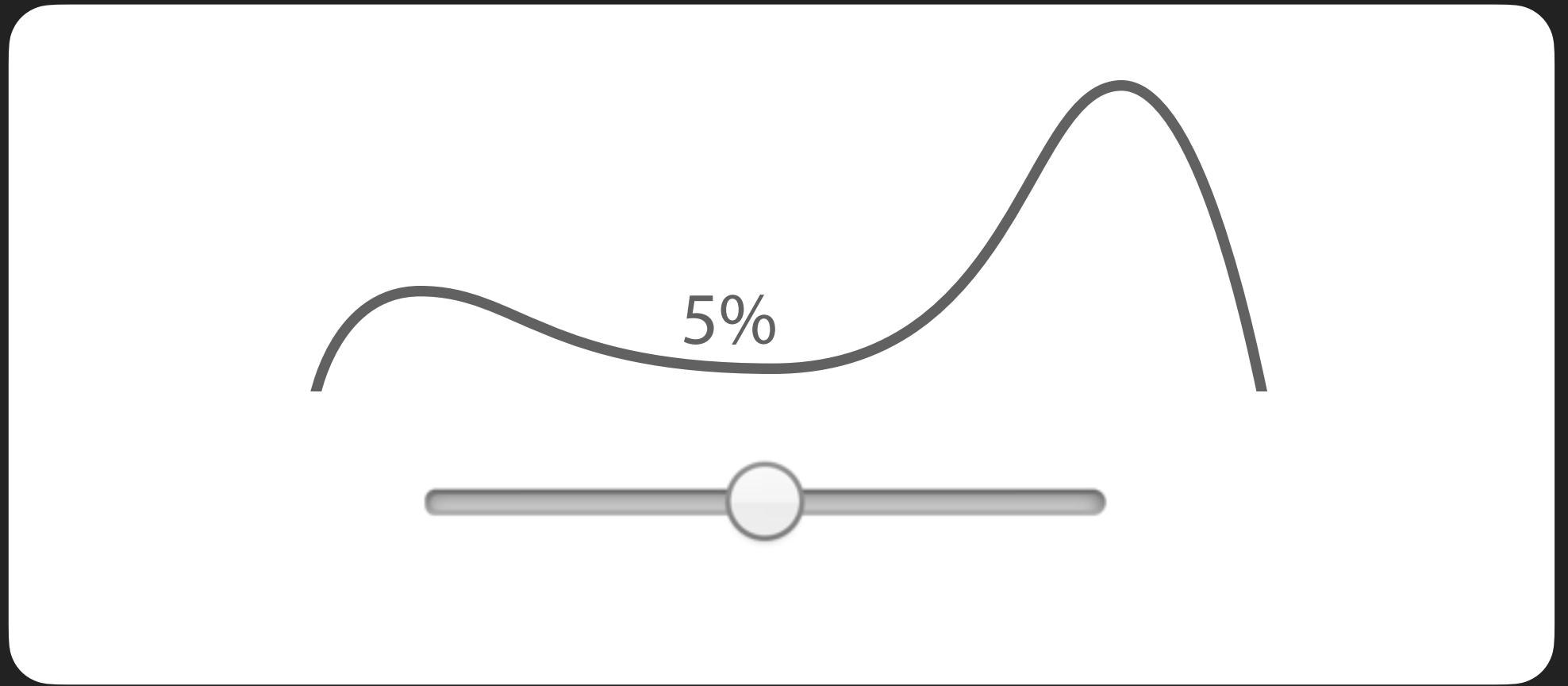
→ *This could've been a toggle?*

→ *Nudge readers towards the middle?*

Discovery #3

# Content-specific Design Evaluation

*People can read in ways that don't match original design intention.*



**Minimum value distribution**  
→ *This could've been a toggle?*  
→ *Nudge readers towards the middle?*

Interactive Articles: Contribution 1

Help broad audiences learn about machine learning's impact on their lives.

Interactive Articles: Contribution 2

Encourage active reading and engage people with learning material.

# Theory + Practice for Interactive Articles



**Fred Hohman**  
Georgia Tech



**Matthew Conlen**  
Univ. of Washington



**Jeffrey Heer**  
Univ. of Washington



**Polo Chau**  
Georgia Tech

## Why care about interactives?

- Better learning & engagement
- Highlight important topics
- Attract broad readership & acclaim

### Communicating with Interactive Articles

Examining the design of interactive articles by synthesizing theory from disciplines such as education, journalism, and visualization.

AUTHORS	AFFILIATIONS	PUBLISHED	DOI
Teddi Iribarren Matthew Conlen Jeffrey Heer Dvir Homan (Polo Chau)	Georgia Tech University of Washington University of Washington Georgia Tech	Not published yet	No DOI yet

#### Contents

Interactive Articles: Theory and Practice  
Engagement and Cognition  
Interactivity with Stories  
Imagining Smart  
Personalization  
Content on Devices  
Critical Perspectives  
Challenges of Interactive Writing  
Leading Figures

Computing has changed how people communicate. The transmission of news, messages, and ideas is instant. Anyone's voice can be heard. In fact, access to digital communication technologies such as the internet is so fundamental to daily life that their disruption by government is condemned by the United Nations Human Rights Council [1]. But while the technology to distribute our ideas has grown in scope and bounds, the interfaces have remained largely the same.

Parallel to the development of the internet, researchers like Alan Kay and Douglas Engelbart worked to build technology that would empower individuals and enhance cognition. Kay imagined the Dynabook [2] in the hands of children across the world. Engelbart, while best remembered for his "mother of all demos" was more interested in the ability of computation to augment human intellect [3]. Neal Stephenson wrote speculative fiction that imagined interactive paper that could display videos and interfaces, and books that could teach and respond to their readers [4].

More recent (yet still historical by personal computing standards) designs point to a future where computers are connected and assist people in decision-making and communicating using rich graphics and interactive user interfaces [5]. While some technologies have seen mainstream adoption, such as Hypertext [6], unfortunately, many others have not. The most popular publishing platforms (for example, WordPress and Medium) choose to prioritize social features and ease-of-use while forfeiting the ability for authors to communicate using the dynamic features of the web.

In the spirit of various computer-assisted cognition technologies, a new form of computational communication medium has emerged that leverages active reading techniques to make ideas more accessible to a broad range of people. These interactive articles build on a long history from Padoa [7] to PhD [8] to reproducible publications [9]. They have been used to better engage, can help improve recall and learning, and attract broad readership and acclaim<sup>1</sup>, yet we did not know that much about them.

Because diverse communities create interactive content, this medium goes by many different names and has not yet settled on a standardized format (or definition)<sup>2</sup>. In research, data journalists, developers, and designers work together to make complex news and investigative reporting clear and engaging using interactive stories [10]. Researchers have proposed artifacts such as explorable multiverse analyses [11], explanations [12], and explanations [13] to more effectively disseminate their work, communicate their results to the public, and remove research debt [14]. Educators use interactive textbooks as an alternative learning format to give students hands-on experience with learning material [15].

Besides these groups, others such as academics, game developers, web developers, and designers blend editorial, design, and programming skills to create and publish aspirational explorations [16], interactive fiction [17], interactive non-fiction [18], active essays [19], and interactive games [20]. While these all slightly differ in their technical approach and target audience, they all largely leverage the interactivity of the modern web.

#### Research Dissemination

#### Research Dissemination

Conducting novel research requires deep understanding and expertise in a specific area. Once achieved, researchers continue contributing new knowledge for other researchers to use and build upon. Over time, this consistent addition of new knowledge can build up, contributing to what some have called research debt. Not everyone is an expert in their field, and it can be easy to lose perspective and forget the bigger picture. Yet research can be understood by many. Interactive articles can be used to distill the latest progress in various research fields and make their methods and results accessible and understandable to a broader audience.

#### Opportunities

- Engage and excite broader audience with latest research programs
- Remove research debt, connect new researchers
- Make faster and clearer research progress

#### Challenges

- No clear incentive structure for researchers
- Limited funding for helping research dissemination and communication
- Not seen as a legitimate research contribution (e.g., to the field, or one's career)

Attacking Discrimination with Smarter Machine Learning [21]

Contests: Context-aware Programming Languages [22]

What Is Complexity Science? [23]

<sup>1</sup>A search query in common academic citation databases, created by leading experts, practitioners, and students in the field, with accompanying interactive visualizations to share data, control, and visualize different complex systems.

<sup>2</sup>A search query in common academic citation databases, created by leading experts, practitioners, and students in the field, with accompanying interactive visualizations to share data, control, and visualize different complex systems.

121

Interactive Articles: Contribution 1

Help broad audiences learn about machine learning's impact on their lives.

Interactive Articles: Contribution 2

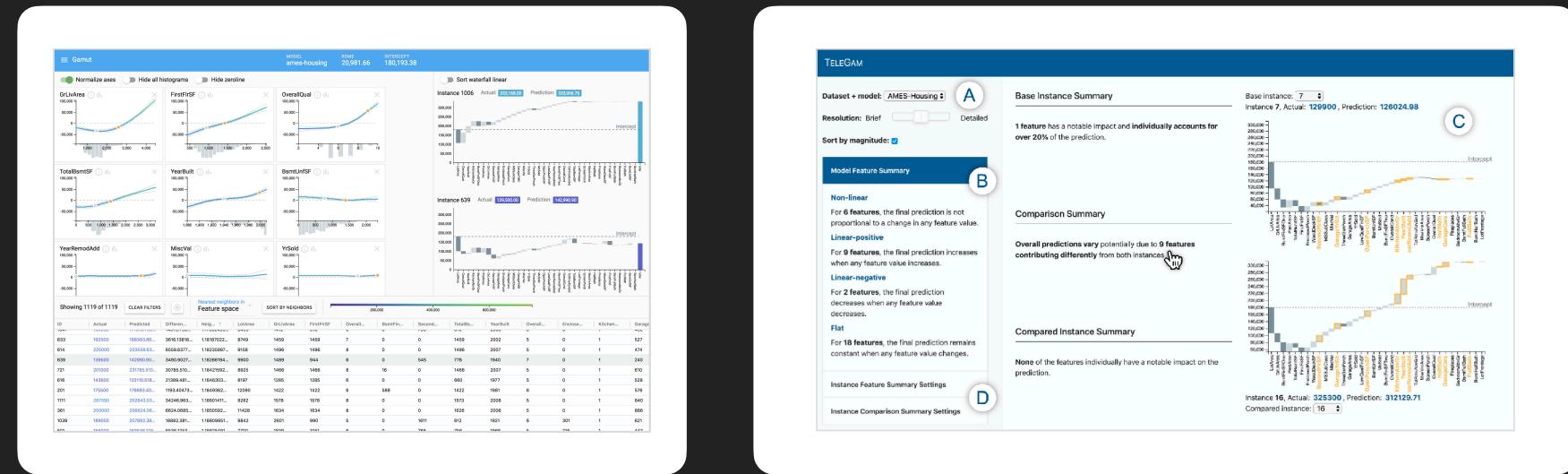
Encourage active reading and engage people with learning material.

Interactive Articles: Contribution 3

Synthesize unique capabilities and formalize research challenges in diverse domains.

# Impact

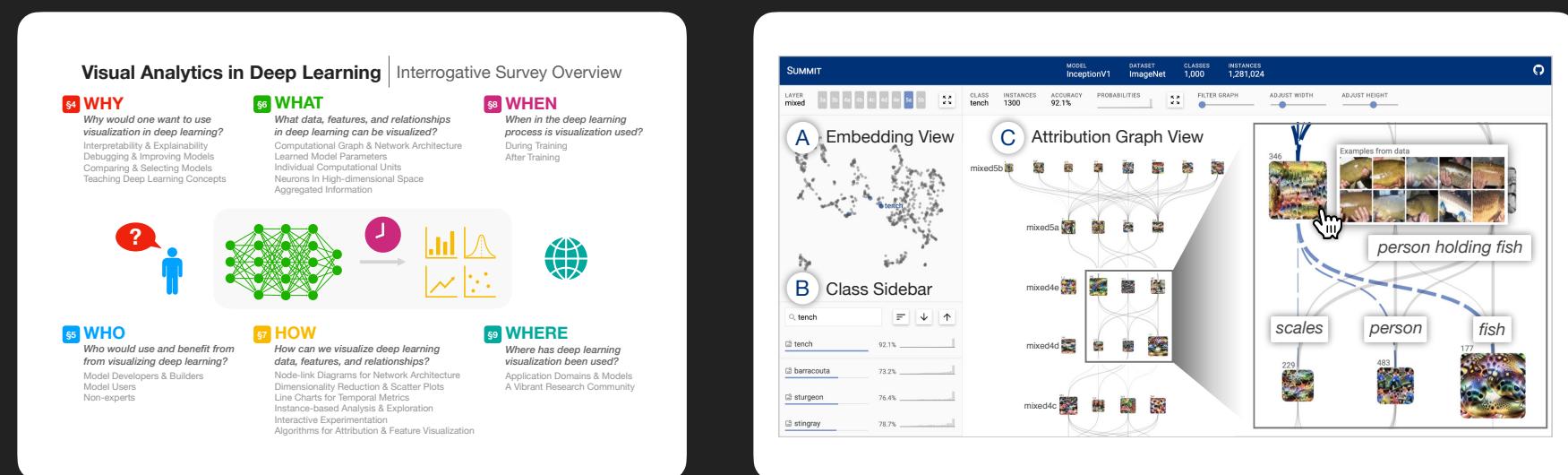
## Enable interpretability



## GAMUT + TELEGAM

Deployed and demoed at Microsoft Research  
Integrated in popular toolkits: InterpretML, SandDance

## Scale interpretability



## Interrogative Survey + SUMMIT

Recognized by NASA PhD Fellowship  
First-of-its-kind survey for deep learning visualization

## Communicate interpretability

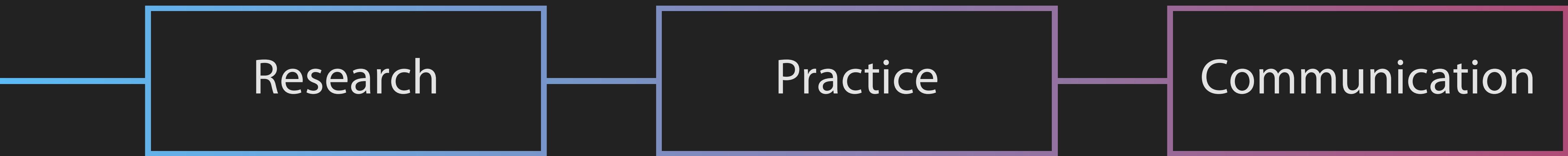


## PARAMETRIC PRESS + Interactive Articles

Helped 250,000+ people learn about machine learning  
Synthesis & formalization of interactive articles

# Future Research Vision

## HCI + AI



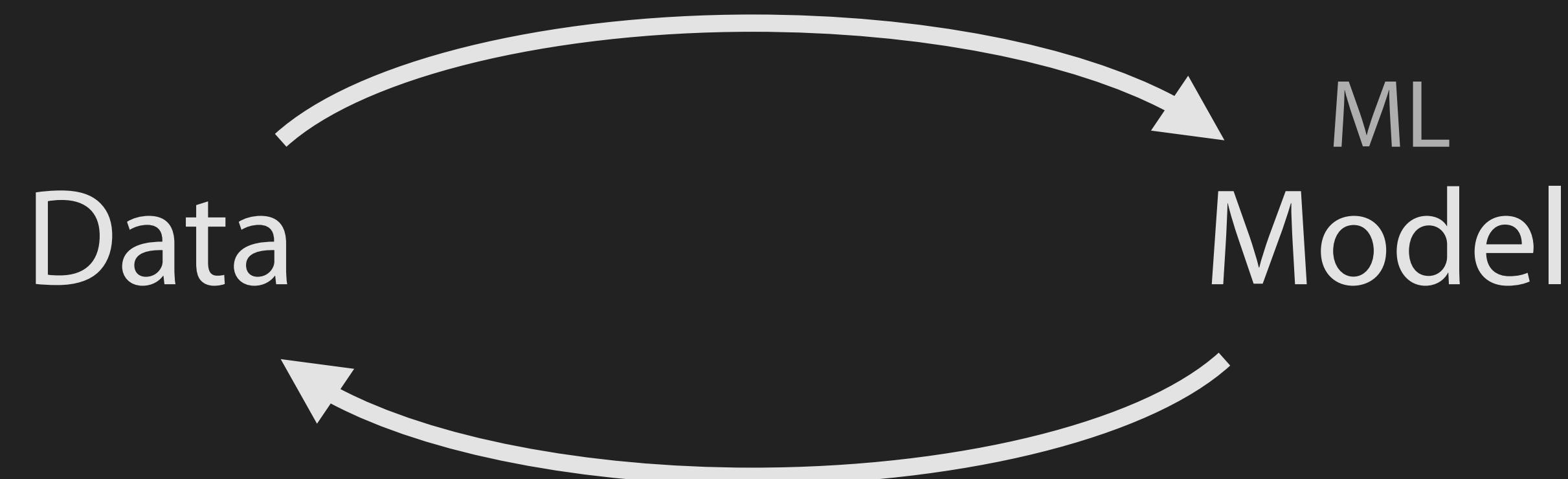
Research Thrust #1

# Mixed-initiative Model Development



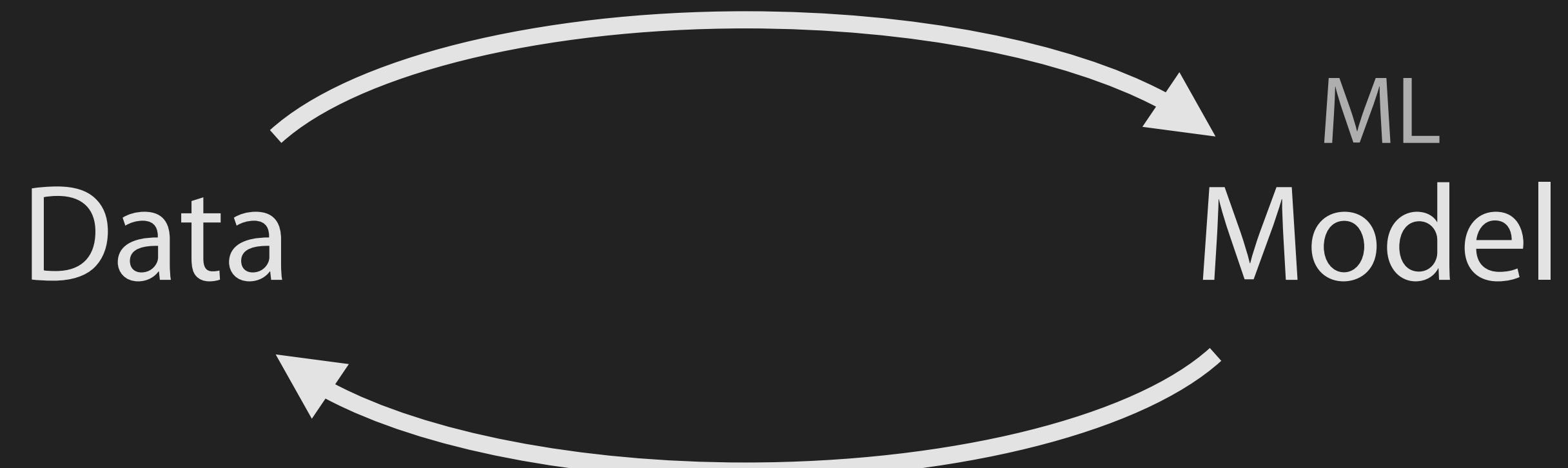
Research Thrust #1

# Mixed-initiative Model Development

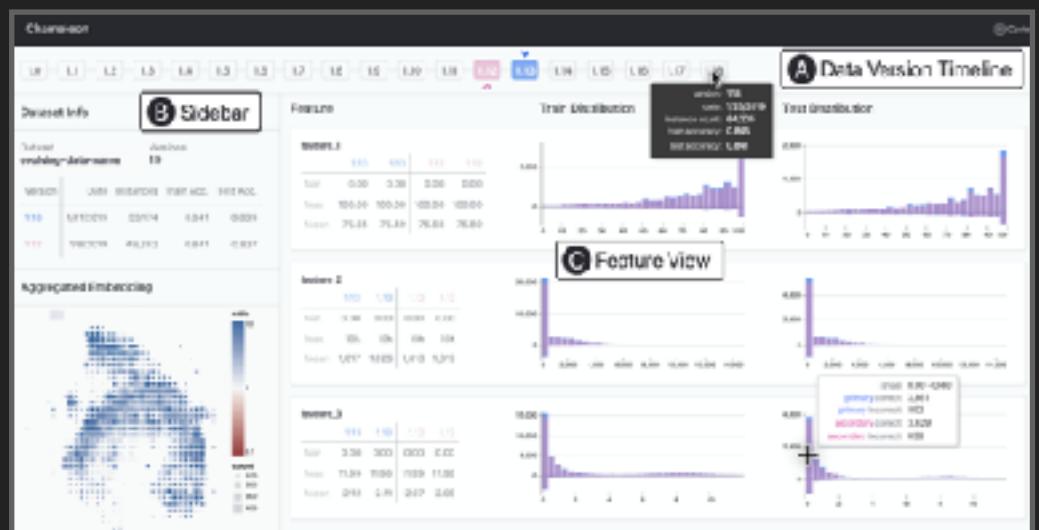


Research Thrust #1

# Mixed-initiative Model Development



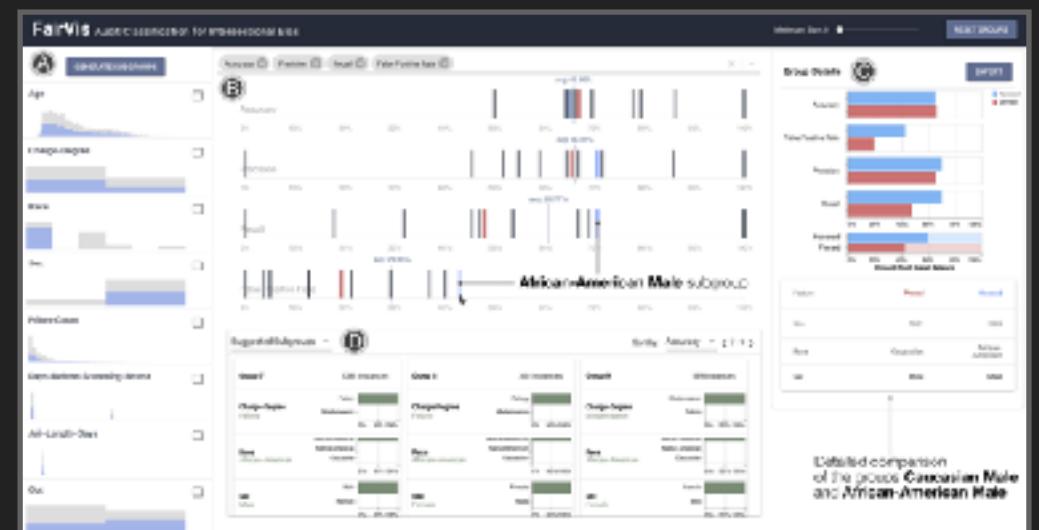
CHAMELEON *CHI 2020*



**Understanding and Visualizing  
Data Iteration in Machine Learning**

Fred Hohman, Kanit Wongsuphasawat, Mary Beth Kery,  
Kayur Patel. *CHI*, 2020.

FAIRVIS *VAST 2019*



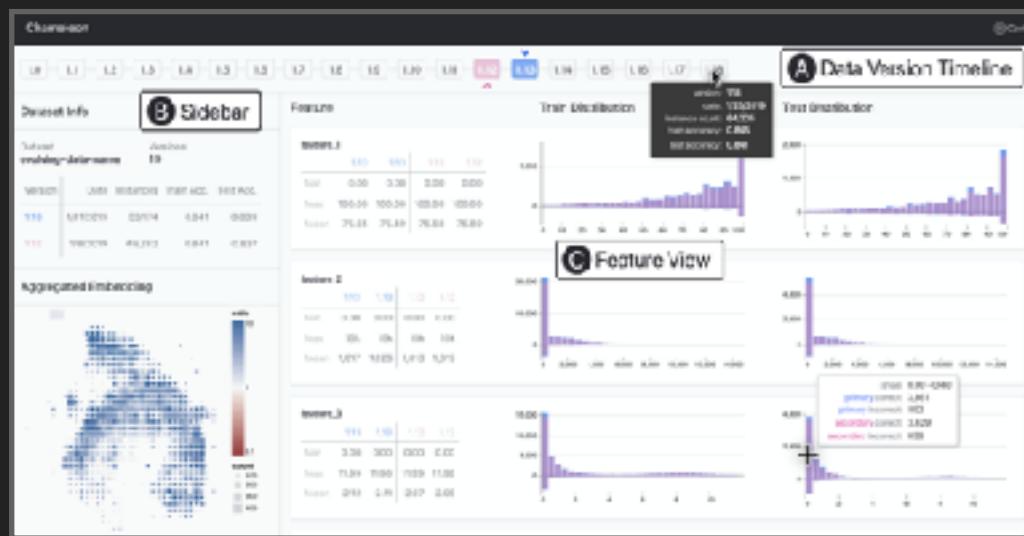
**Visual Analytics for Discovering  
Intersectional Bias in Machine Learning**

Angel Cabrera, Will Epperson, Fred Hohman, Minsuk  
Kahng, Jamie Morgenstern, Duen Horng (Polo) Chau. *VAST*, 2019.

Research Thrust #1

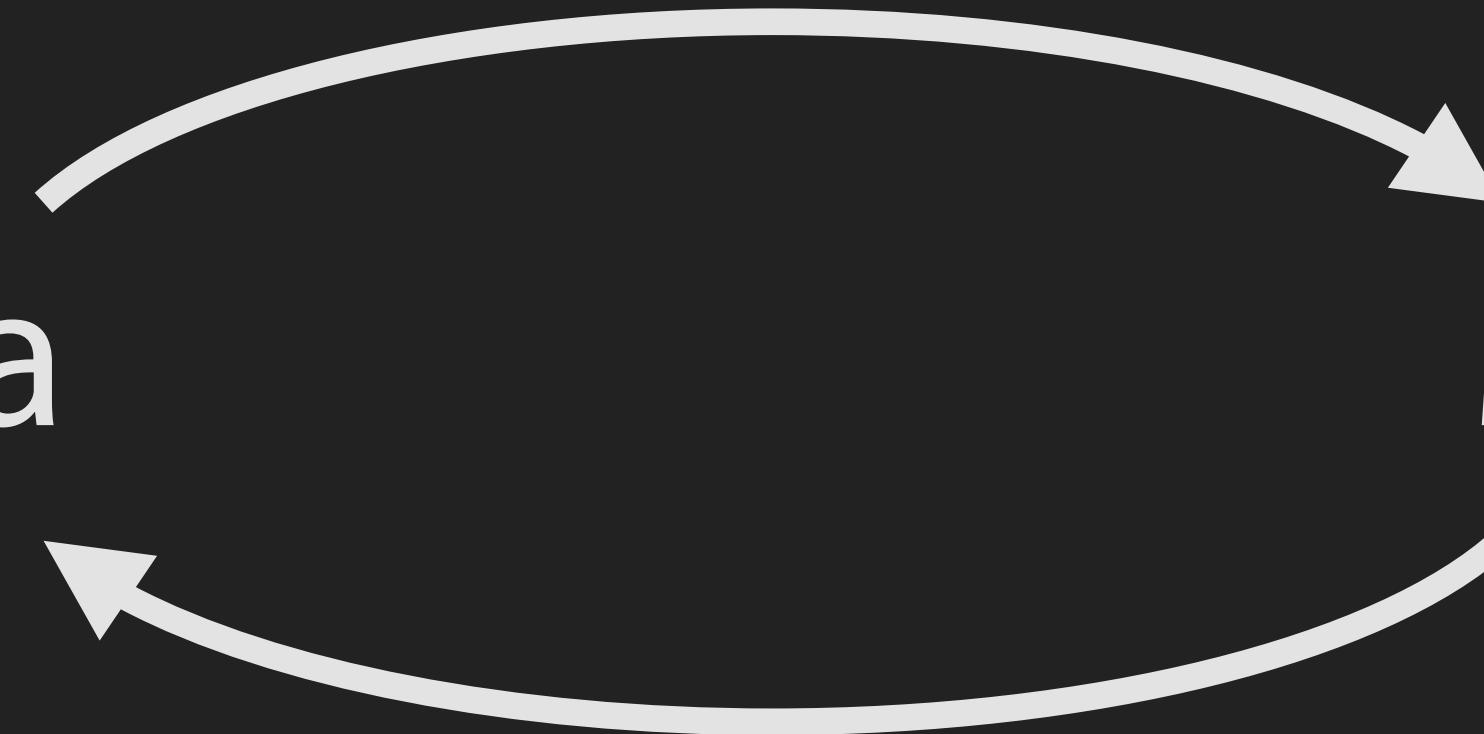
# Mixed-initiative Model Development

CHAMELEON CHI 2020



Data

ML  
Model



FAIRVIS VAST 2019



## RESEARCH DIRECTIONS

**Interactive Data Programming:** Visualize and track updating labels.

**Augmented Models:** New models for underperforming data subgroups.

**Unit Tests for Interpretability:** Specification for learned rules.

Research Thrust #2

# Making Interpretability Common Practice

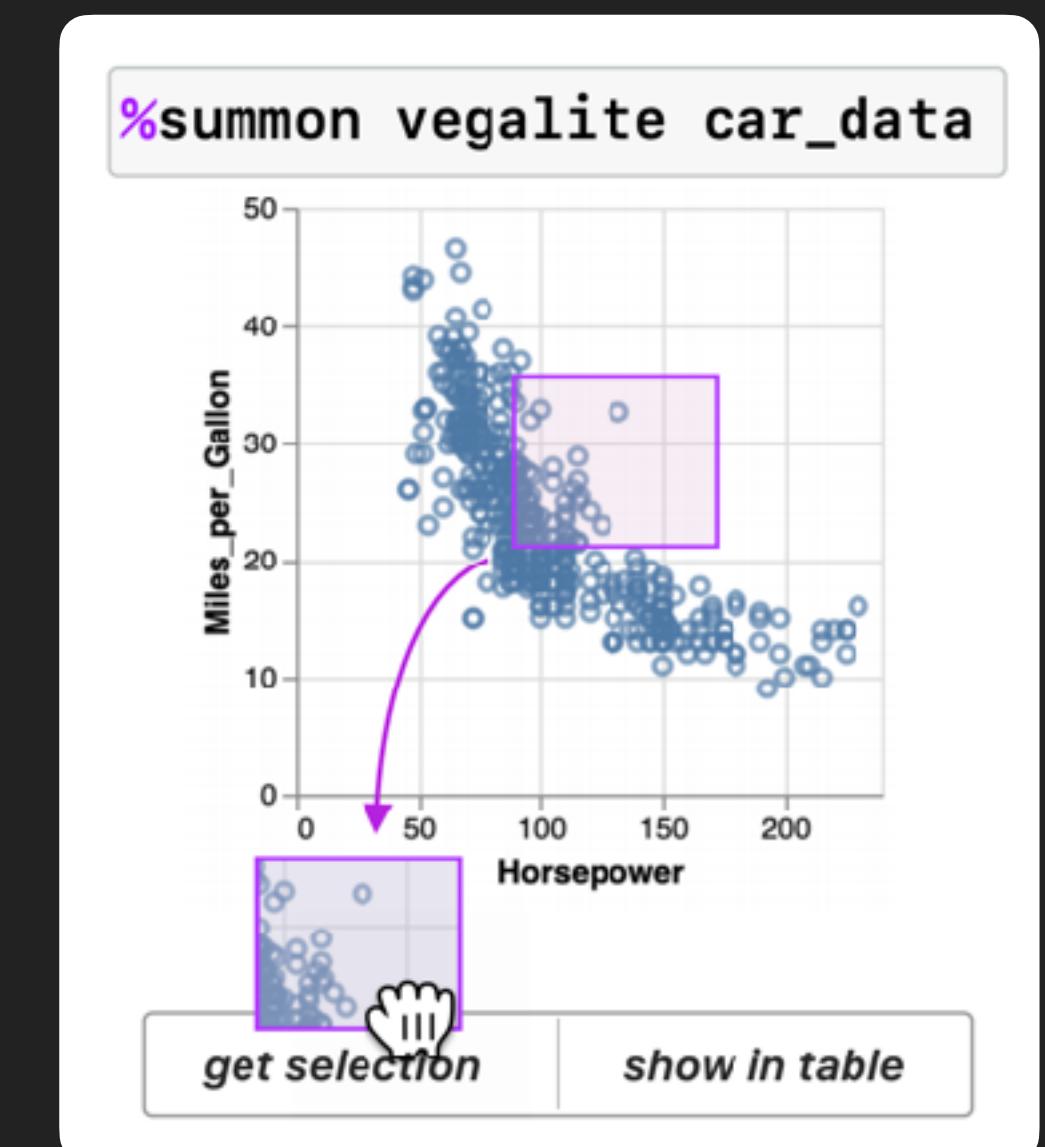
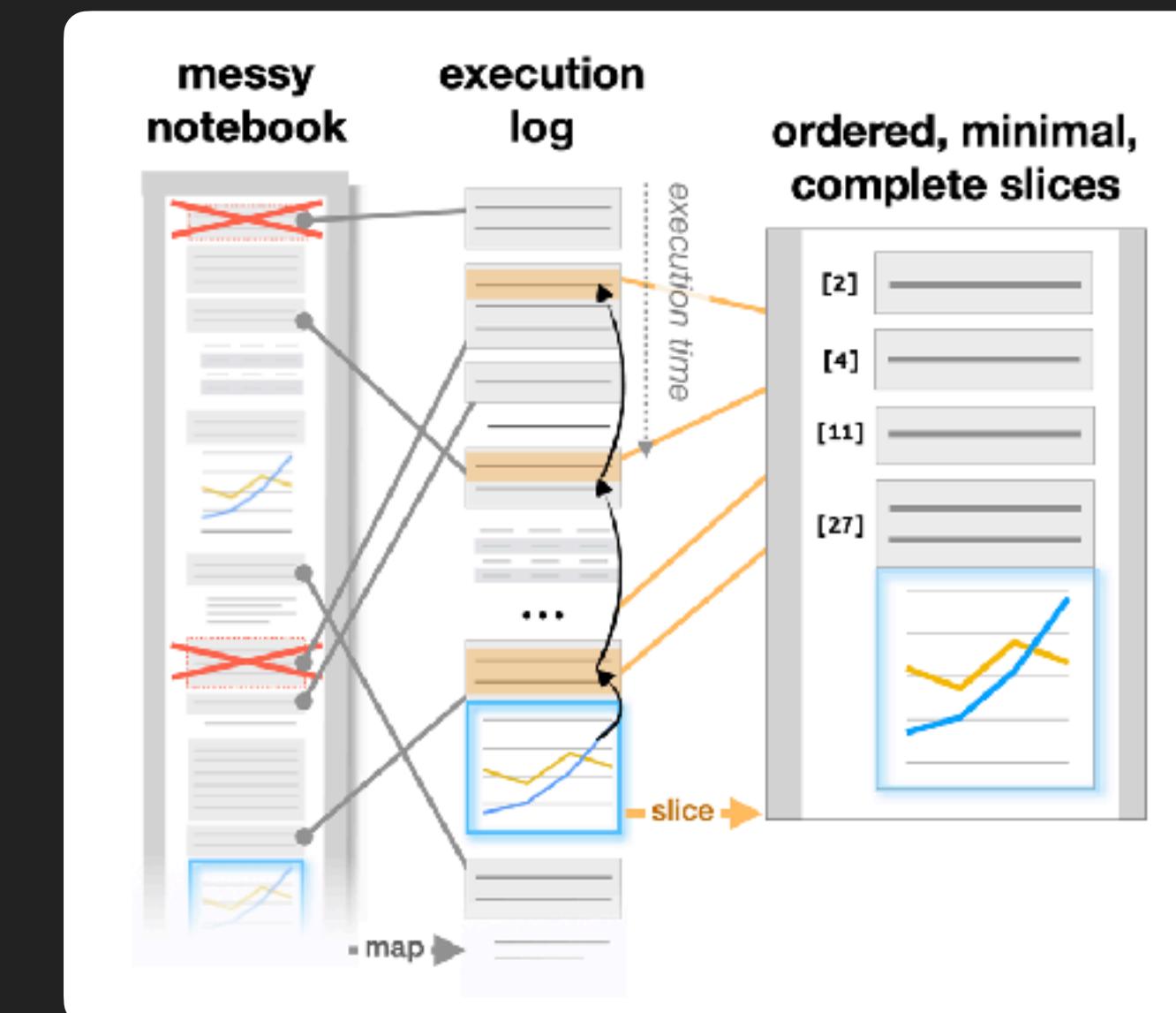
## RESEARCH DIRECTIONS

Integrate interpretability into existing code environments.

Generalize to create best practices for interpretability.

Preserve healthy analysis skepticism with explainable AI.

## PRELIMINARY WORK



## Managing Messes in Computational Notebooks

Andrew Head, Fred Hohman, Titus Barik, Steven Drucker, Robert DeLine. *CHI*, 2019.

## The Future of Notebook Programming Is Fluid

Mary Beth Kery, Donghao Ren, Kanit Wongsuphasawat, Fred Hohman, Kayur Patel. *CHI*, 2020.

# Research Dissemination

	Paper	Video	Code	Slides	Demo + Article
<b>GAMUT</b> CHI 2019	✓	✓		✓	✓
<b>TELEGAM</b> VAST 2019	✓	✓	✓	✓	✓
<b>Interrogative Survey</b> TVCG 2018	✓	✓	✓	✓	✓
<b>SUMMIT</b> TVCG 2020	✓	✓	✓	✓	✓
<b>Parametric Press</b> 2019	✓	✓	✓	✓	✓
<b>Dimensionality Reduction</b> VISxAI 2018	✓		✓	✓	✓
<b>Interactive Articles</b> Distill 2020	✓		✓		✓
<b>CHAMELEON</b> CHI 2020	✓	✓		○	○
<b>FAIRVIS</b> VAST 2019	✓	✓	✓	✓	✓
<b>Code Gathering</b> CHI 2019	✓	✓	✓	✓	✓
<b>Notebook Handoff</b> CHI 2020 poster	✓	✓	✓	○	○

Research Thrust #3

# Accessible Research Distillation

## Challenges

Technical: diverse devices, accessible archival

Design: authoring requires multiple skillsets

Legitimization: few incentives, time consuming

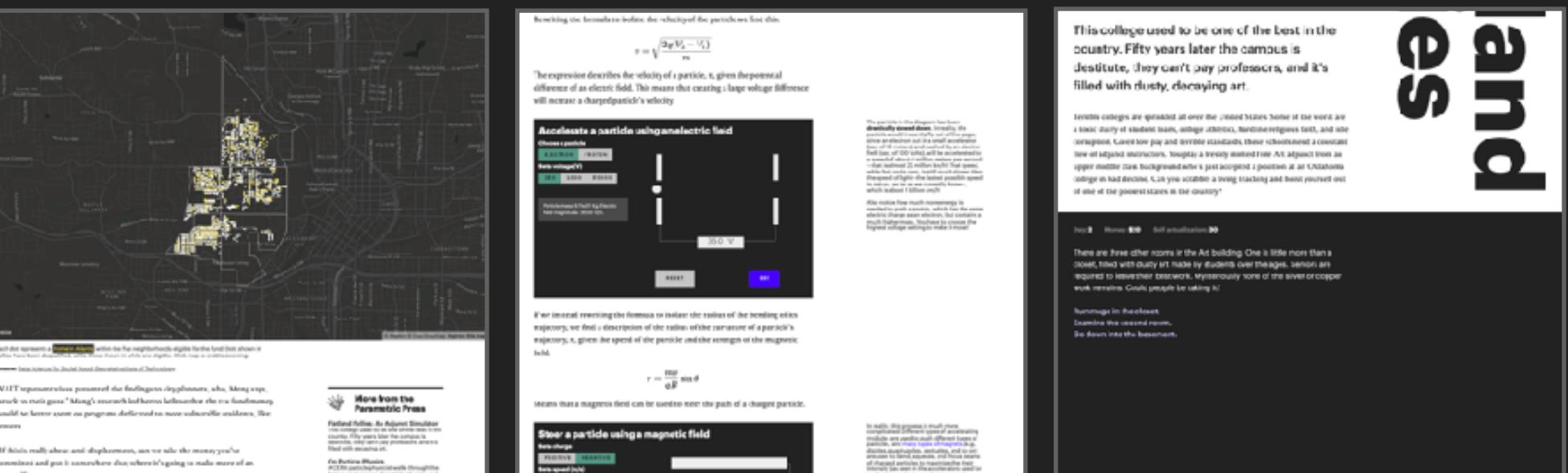
## RESEARCH DIRECTIONS

Lower cost of creation with  
authoring tools and guidelines.

Transform ML education at scale.

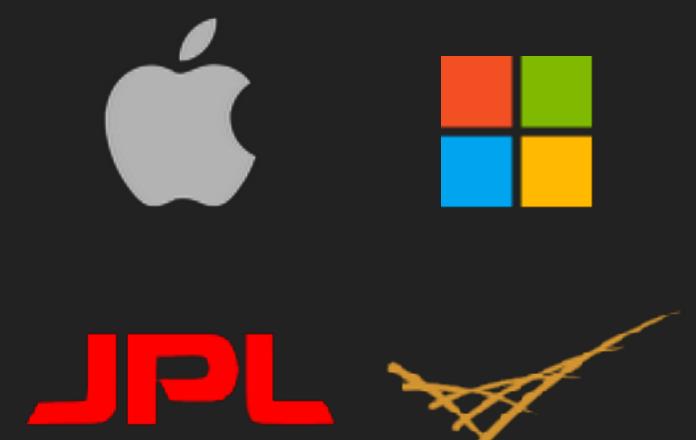
Cultivate research communication culture.

BUILD ON MY EXPERIENCE



and

***Thank you*** to collaborators from academia, industry, and the government!



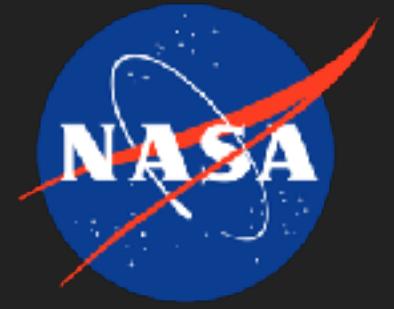
NASA PhD Fellowship

Visualization Lab

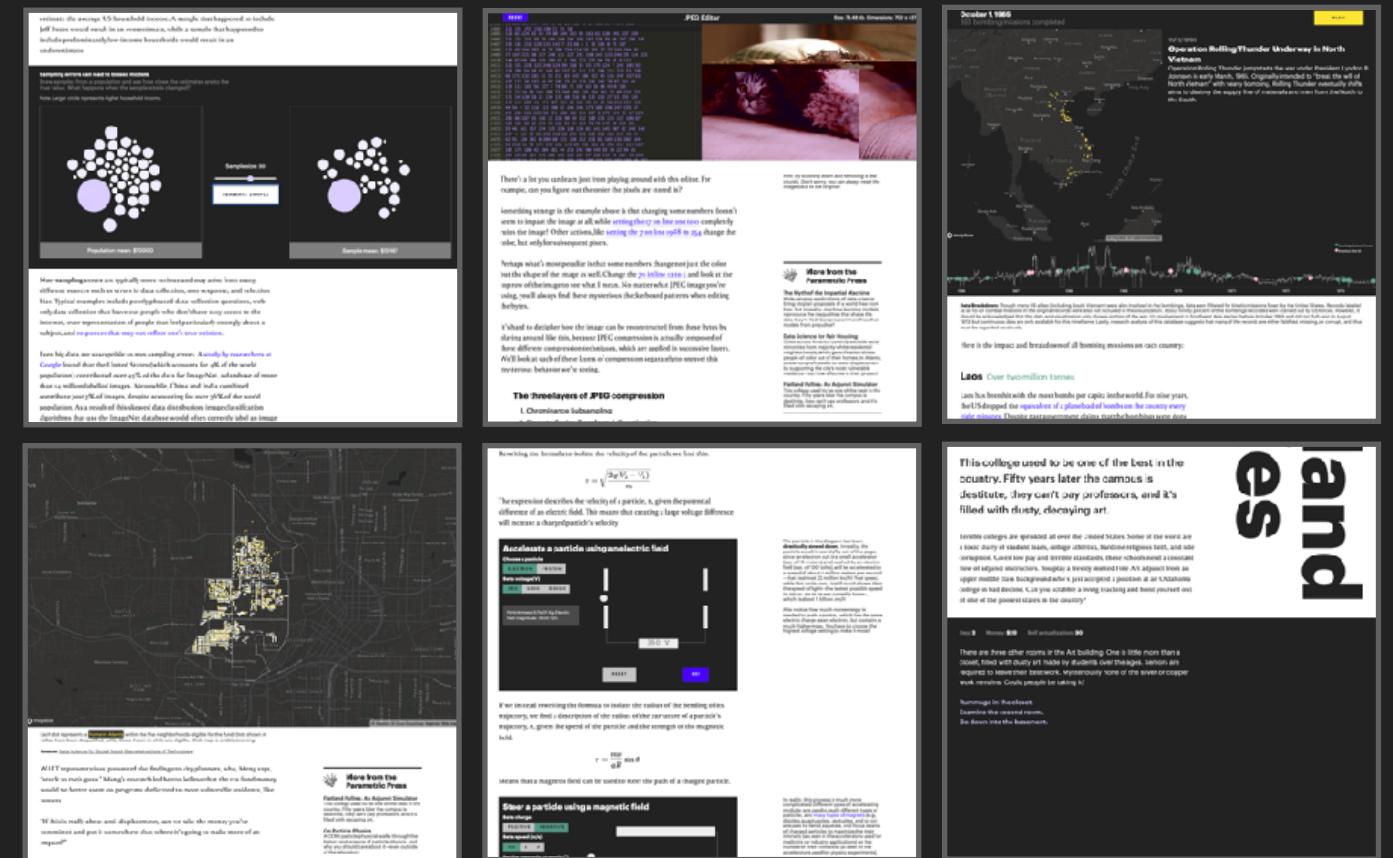
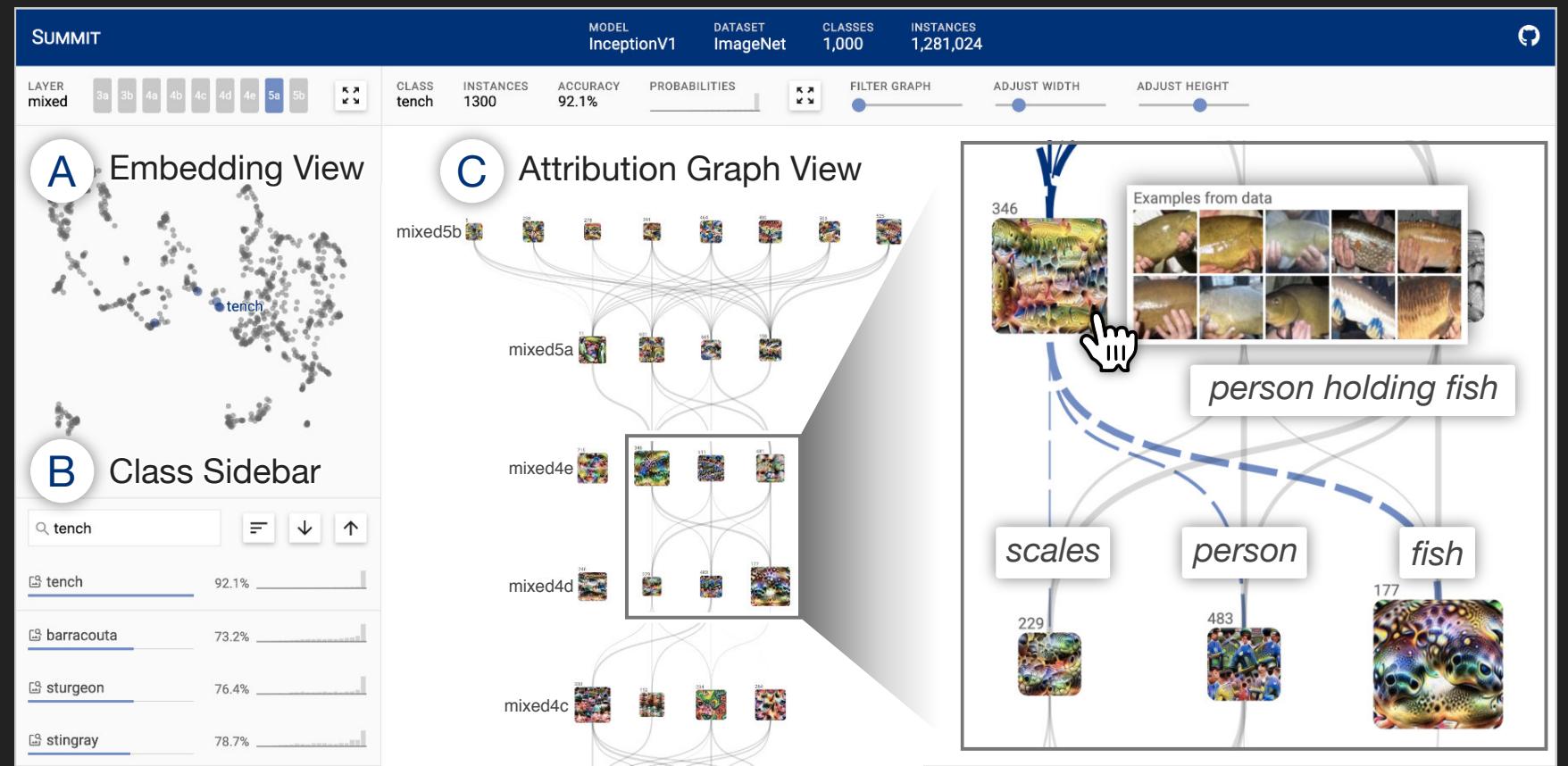
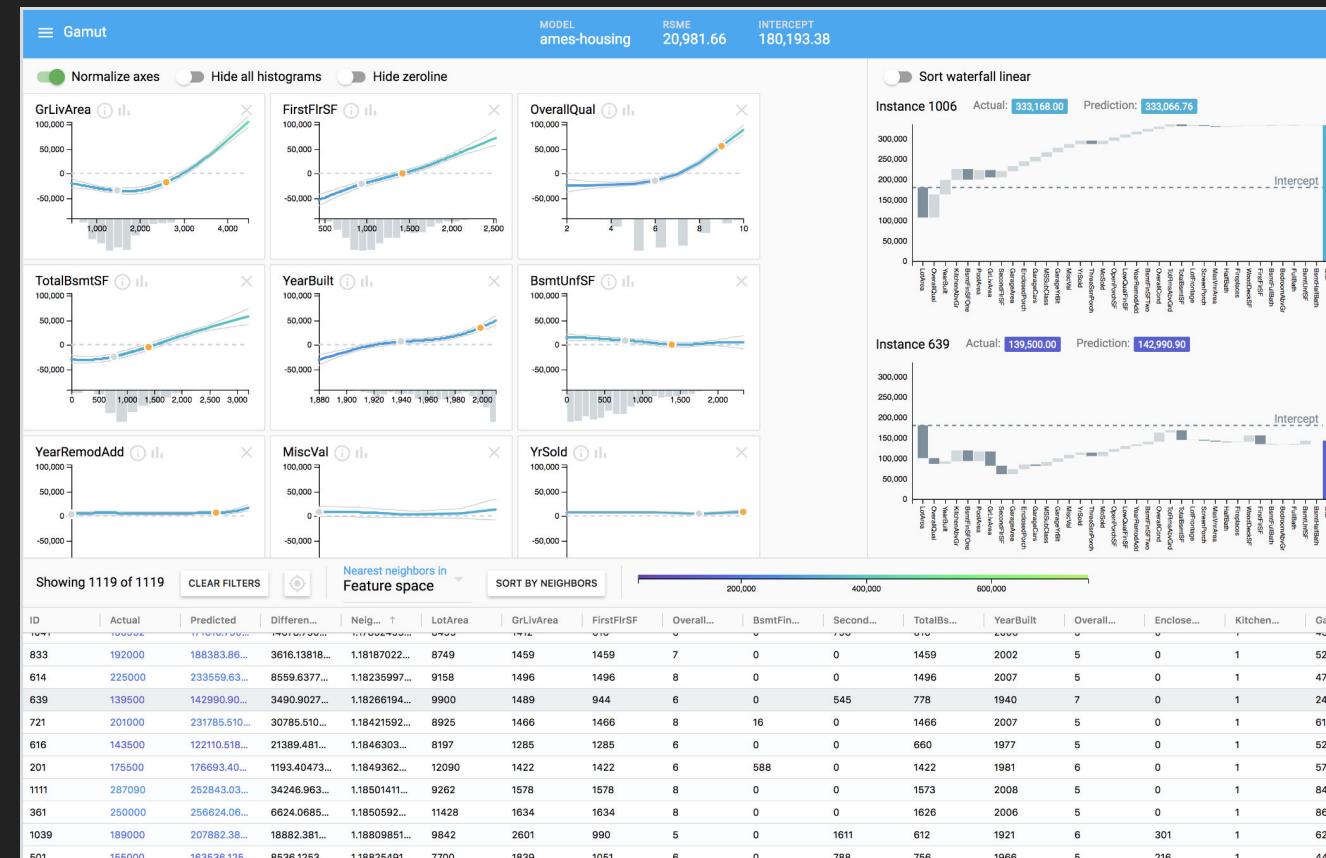
# Interactive Scalable Interfaces for Machine Learning Interpretability



**Fred Hohman** @fredhohman  
fredhohman.com



*Enabling interpretability at scale and for everyone.*



**GAMUT**  
& **TELEGAM**

**Interrogative Survey**  
& **SUMMIT**

**PARAMETRIC PRESS**  
& **Interactive Articles**