

# Lecture 2 - finishing up from last time

The humans strike back,  
The humans-in-the-loop

# Course logistics

Discussion sections:

- We will discuss two papers
- We will combine both papers together across roles to save time

**Project proposals are due Jan 24 at 11:59pm**

# The humans-in-the-loop: two perspectives



## Artificial Intelligence

**Goal:** To produce high quality labels as efficiently as possible

**Artifact:** training data for models

Impacts across **short time horizon**



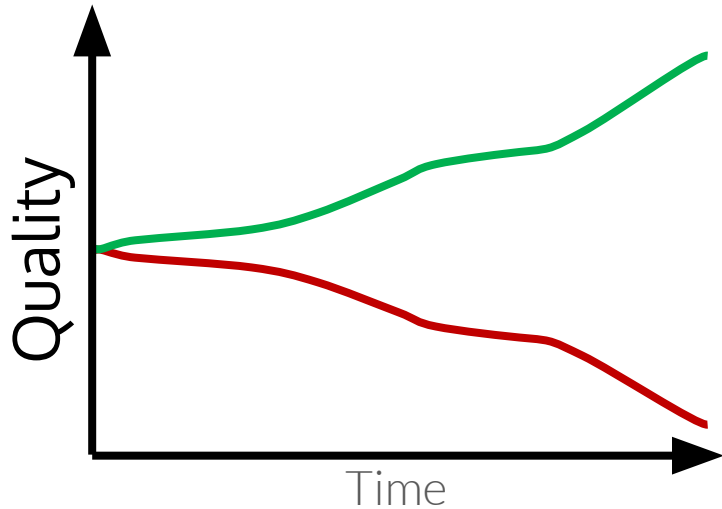
## Human-Computer Interaction

**Goal:** To support a labor force achieve their financial and career goals

**Artifact:** automations that structure work

Impacts across **long time horizon**

# Studying long term annotator quality



Quality increases over time:

Quality decreases over time:

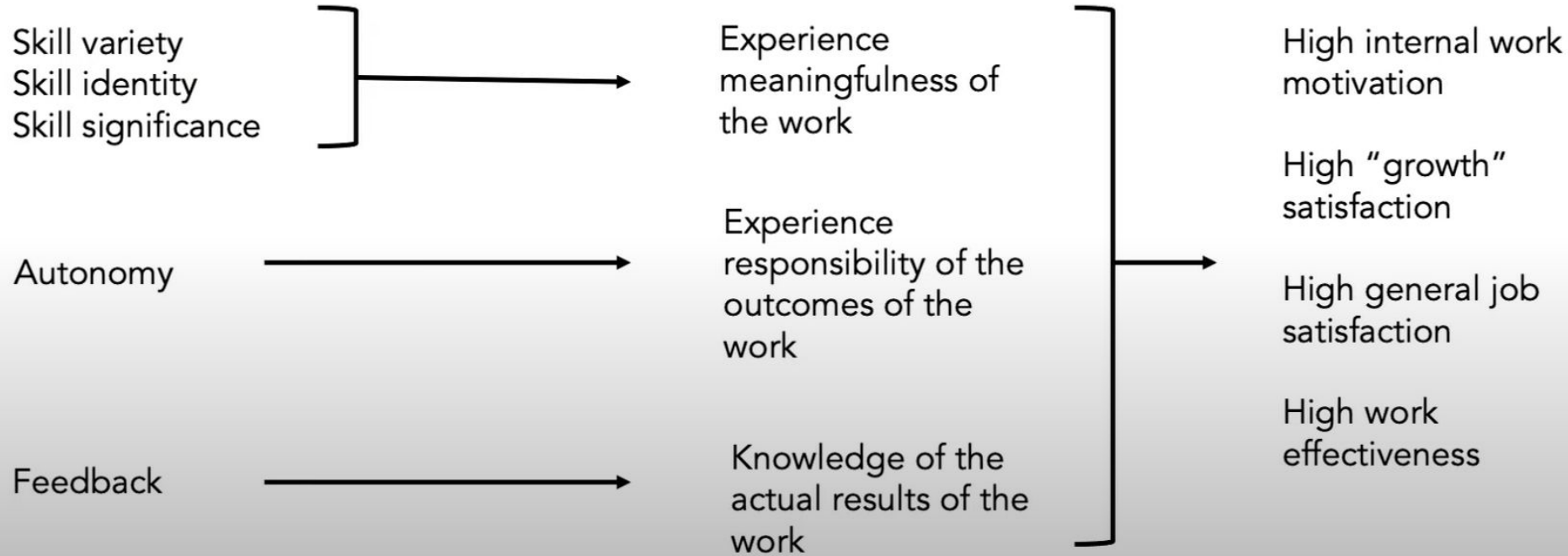


# Speeding up annotation

# Job Characteristic Model

Hackman & Oldham, 1980

Core Job Characteristics → Critical Psychological States → Outcomes



# Existing platforms do not support these job characteristics

Requester	Title	Hits	Reward	Created	Actions	
<a href="#">James Billings</a>	Market Research Survey	25,571	\$0.05	9m ago	<a href="#">Preview</a>	<a href="#">Accept &amp; Work</a>
<a href="#">Research Rewards</a>	Quick Market Research Survey	22,826	\$0.02	6m ago	<a href="#">Preview</a>	<a href="#">Accept &amp; Work</a>
<a href="#">Mayanksoniphd</a>	Generate praise, given a persona.	6,655	\$0.03	15d ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">Shopping Receipts</a>	Extract General Data & Items From Shopping Receipt	1,150	\$0.01	11s ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">Shopping Receipts</a>	Extract General Data & Items From Shopping Receipt	1,121	\$0.02	4h ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">minsVA</a>	Draw a polygon around the tailgate of the requested cars	915	\$0.10	4h ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">Shopping Receipts</a>	Extract General Data & Items From Shopping Receipt	811	\$0.03	3h ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">VacationRentalAPI CA</a>	Address Identification - 10207 - Kelowna, BC	676	\$7.50	5h ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">Shopping Receipts</a>	Extract General Data & Items From Shopping Receipt	628	\$0.05	16h ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">minsVA</a>	Draw a polygon around the front hood of the requested cars	616	\$0.10	4h ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">Shopping Receipts</a>	Extract General Data & Items From Shopping Receipt	554	\$0.04	12h ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">VacationRentalAPI</a>	Address Identification - 10227 - Minneapolis, MN	405	\$2.50	5h ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">VacationRentalAPI</a>	Address Identification - 10243 - New Listing Mix	371	\$2.00	3h ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">str11223344</a>	Tell us what this item is - General Contents - Batch ID #44814	353	\$0.08	6d ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">VacationRentalAPI</a>	Address Identification - 10242 - New Listing Mix	353	\$2.00	4h ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">Alexander Gulin</a>	Run a query in ChatGPT	326	\$0.02	11d ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">VacationRentalAPI CA</a>	Address Identification - 10200 - Brampton, ON	321	\$7.50	5h ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">Company</a>	Company Logos	297	\$0.01	17s ago	<a href="#">Preview</a>	<a href="#">Accept &amp; Work</a>
<a href="#">Shopping Receipts</a>	Extract Data From Shopping Receipt	294	\$0.01	1m ago	<a href="#">Preview</a>	<a href="#">Quality</a>
<a href="#">VacationRentalAPI CA</a>	Address Identification - 10201 - Burnaby, BC	258	\$7.50	5h ago	<a href="#">Preview</a>	<a href="#">Quality</a>

# Existing platforms do not support these job characteristics

Requester	Task	Hits	Reward	Created	Actions
J		25,571	\$0.05	9m ago	<a href="#">Preview</a> <a href="#">Accept &amp; Work</a>
R		22,826	\$0.02	6m ago	<a href="#">Preview</a> <a href="#">Accept &amp; Work</a>
M		6,655	\$0.03	15d ago	<a href="#">Preview</a> <a href="#">Quality</a>
Shopping Receipts	Extract General Data & Items From Shopping Receipt	1,150	\$0.01	11s ago	<a href="#">Preview</a> <a href="#">Quality</a>
Shopping Receipts	Extract General Data & Items From Shopping Receipt	1,121	\$0.02	4h ago	<a href="#">Preview</a> <a href="#">Quality</a>
minsVA		915	\$0.10	4h ago	<a href="#">Preview</a> <a href="#">Quality</a>
Shopping Receipts		811	\$0.03	3h ago	<a href="#">Preview</a> <a href="#">Quality</a>
VacationRentalAPI CA		676	\$7.50	5h ago	<a href="#">Preview</a> <a href="#">Quality</a>
Shopping Receipts	Extract General Data & Items From Shopping Receipt	628	\$0.05	16h ago	<a href="#">Preview</a> <a href="#">Quality</a>
minsVA	Draw a polygon around the front hood of the requested cars	616	\$0.10	4h ago	<a href="#">Preview</a> <a href="#">Quality</a>
Shopping Receipts	Extract General D	554	\$0.04	12h ago	<a href="#">Preview</a> <a href="#">Quality</a>
VacationRentalAPI	Address Identifica	405	\$2.50	5h ago	<a href="#">Preview</a> <a href="#">Quality</a>
VacationRentalAPI	Address Identifica	371	\$2.00	3h ago	<a href="#">Preview</a> <a href="#">Quality</a>
str11223344	Tell us what this it	353	\$0.08	6d ago	<a href="#">Preview</a> <a href="#">Quality</a>
VacationRentalAPI	Address Identification - 10242 - New Listing Mix	353	\$2.00	4h ago	<a href="#">Preview</a> <a href="#">Quality</a>
Alexander Gutin	Run a query in ChatGPT				<a href="#">Quality</a>
VacationRentalAPI CA	Address Identification - 10200 - Brampton, ON				<a href="#">Quality</a>
Company	Company Logos				<a href="#">Quality</a> <a href="#">Accept &amp; Work</a>
Shopping Receipts	Extract Data From Shopping Receipt				<a href="#">Quality</a>
VacationRentalAPI CA	Address Identification - 10201 - Burnaby, BC	258	\$7.50	5h ago	<a href="#">Preview</a> <a href="#">Quality</a>

Does this task design even work?

What skills does this task require or help me develop?

Why does Amazon take between 20-40% of overhead?

Bad ratings hurt my future earning potential. Can't even rate the requestors



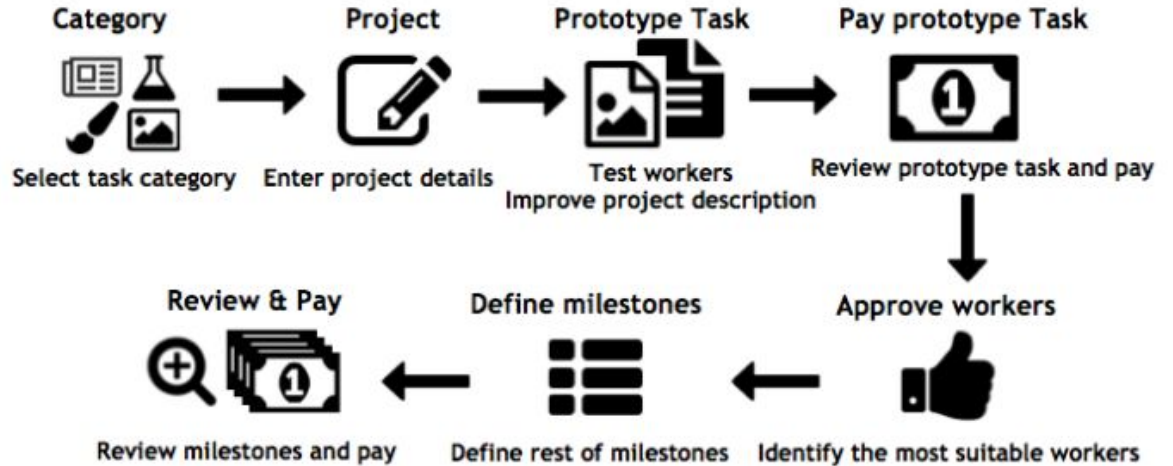
Humans-in-the-loop from an HCI perspective:  
Can we develop a platform that supports worker  
needs?

# Daemo: a Self-Governed Crowdsourcing Marketplace

V1 launched with :

## Prototype tasks

- Workers improve task design

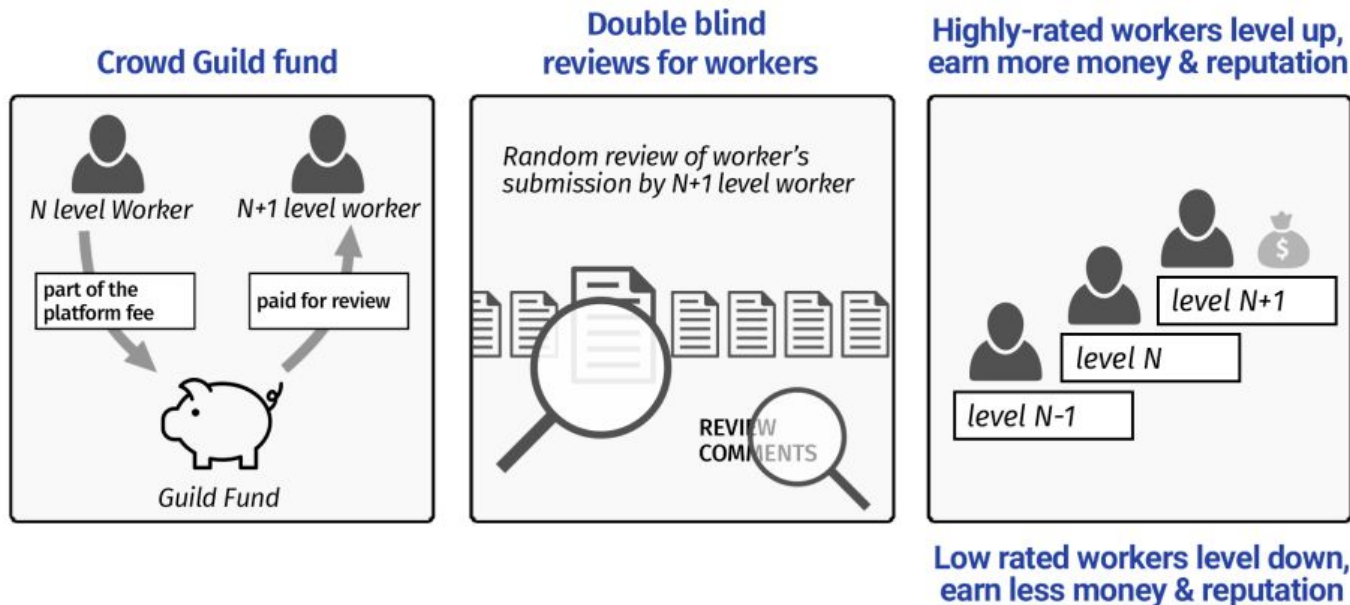


## Open governance

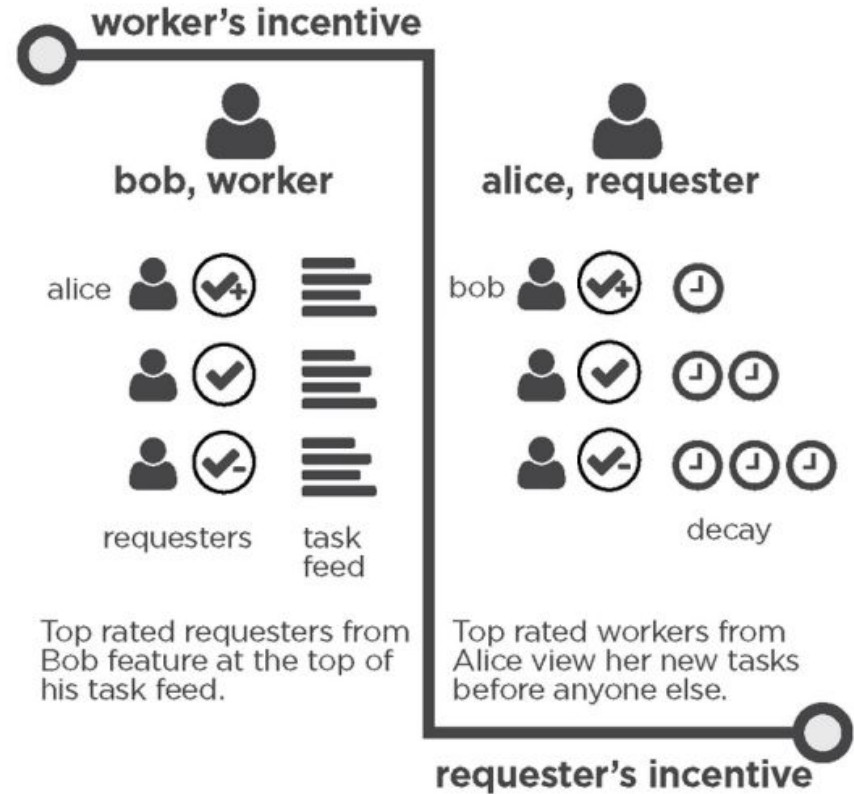
- 3 workers
- 3 requesters
- 1 researcher

Gaikwad et al. Daemo: a Self-Governed Crowdsourcing Marketplace. UIST 2017

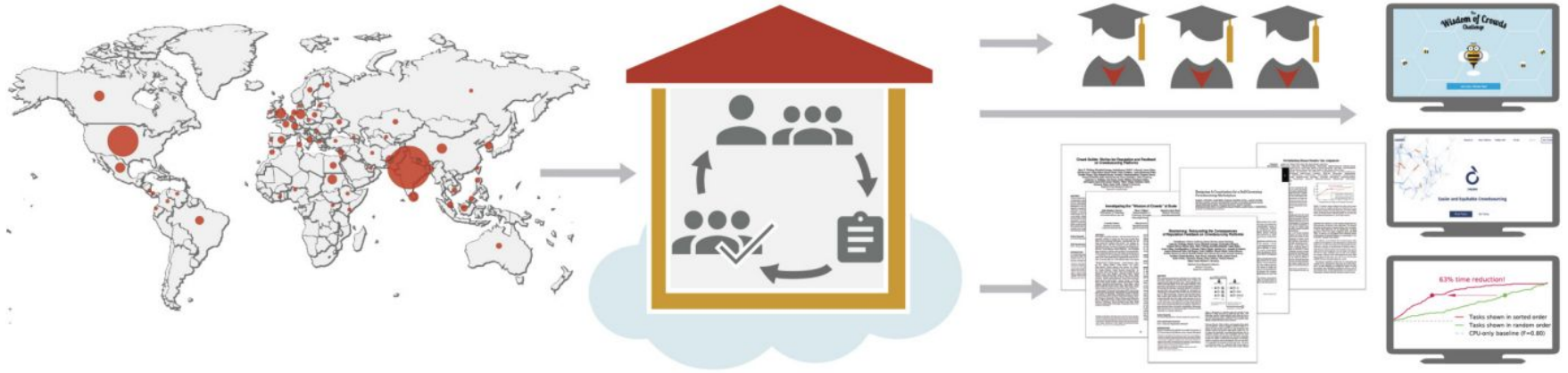
# A reputation protocol: workers received **feedback**



A rating system:  
To trade off skill  
variety of identity



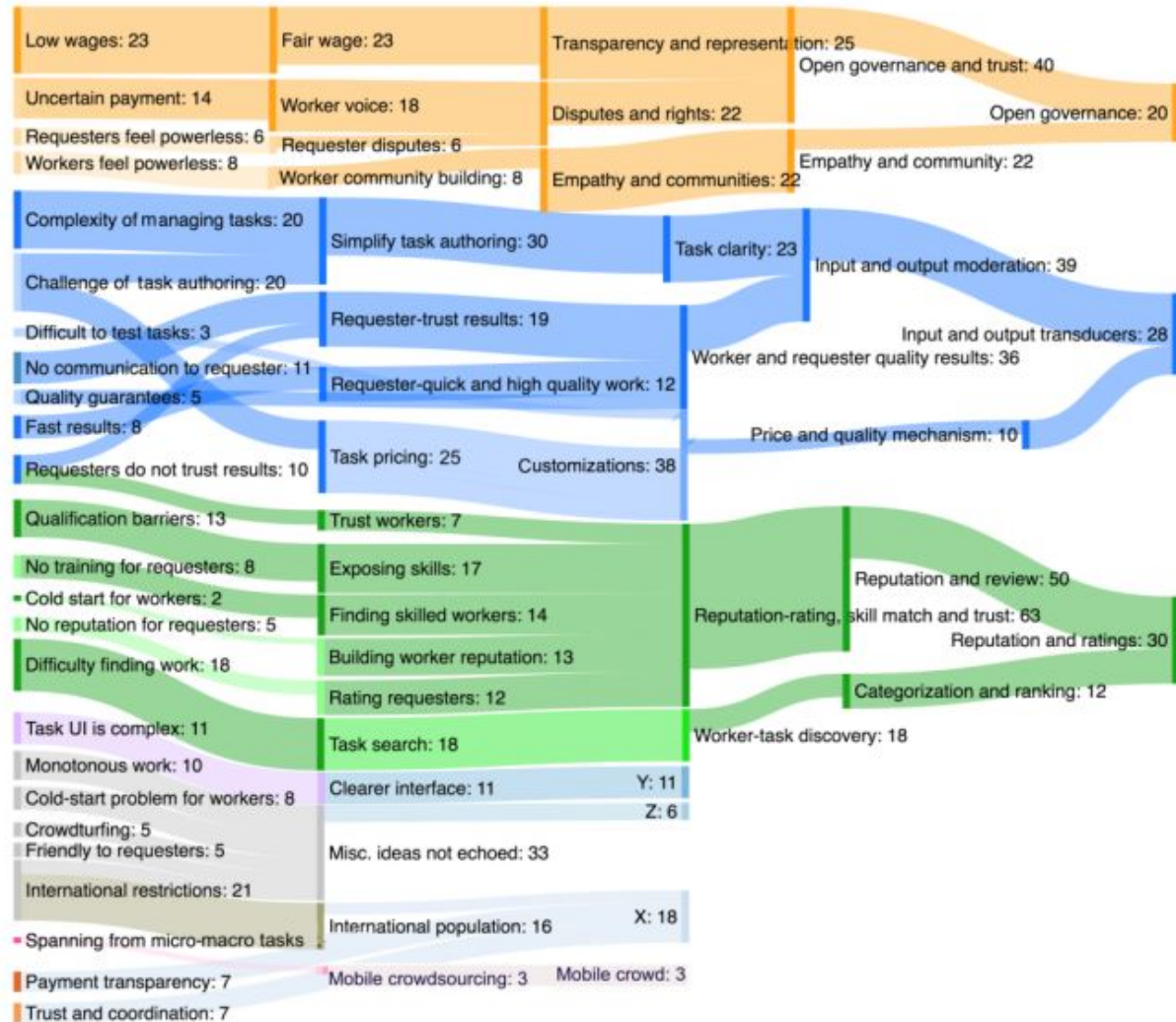
# Building a new decentralized crowdsourcing system with a crowd of researchers



Achieve upward educational mobility while creating research systems and co-authoring papers

# Ideas

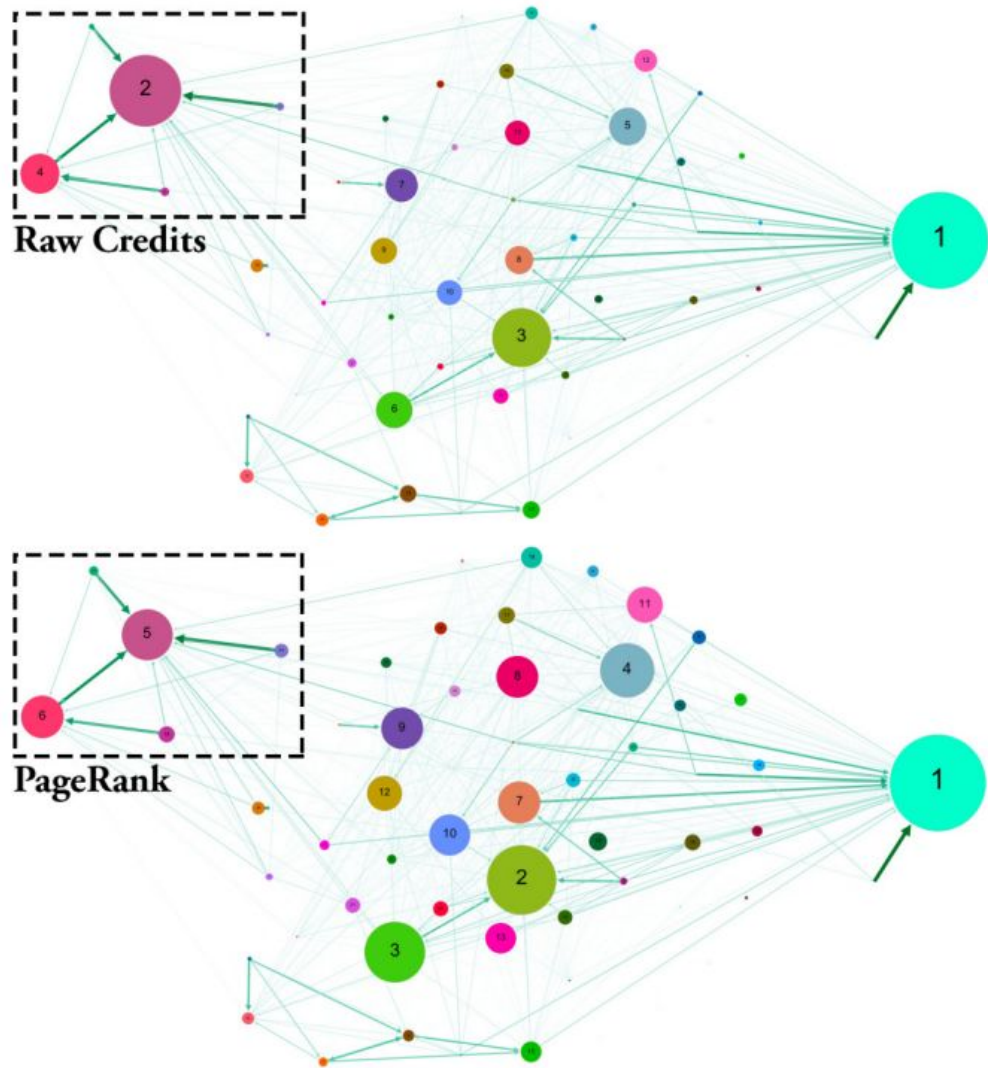
Changes to the platform were ideated on transparently and collectively prioritized



Author order determined using crowdsourced points and page rank

Potential challenges:

- Link ring
- Quid-proquo strategy



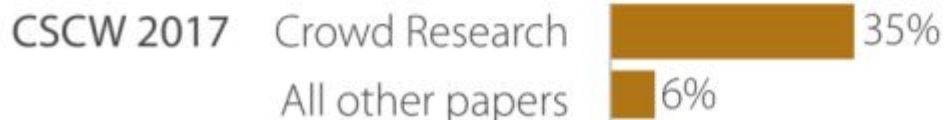
# Supporting upward mobility

Our authors were more diverse than those from other papers at the same venue

Coauthors' universities that are ranked below 500 worldwide



Coauthors whose countries are ranked below 50 worldwide in GDP per capita

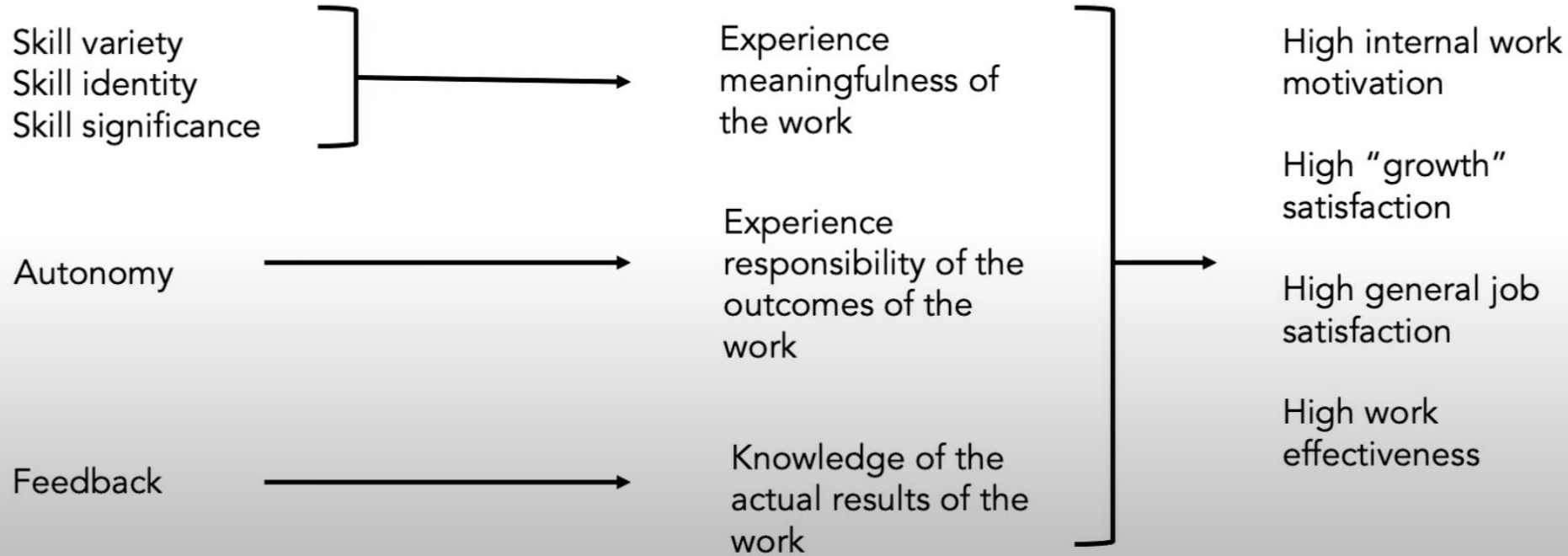




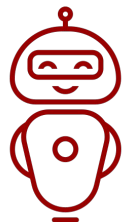
# Job Characteristic Model

Hackman & Oldham, 1980

Core Job Characteristics → Critical Psychological States → Outcomes



# The humans-in-the-loop: two perspectives



## Artificial Intelligence

**Goal:** To produce high quality labels as efficiently as possible

**Artifact:** training data for models

Impacts across **short time horizon**



## Human-Computer Interaction

**Goal:** To support a labor force achieve their financial and career goals

**Artifact:** automations that structure work

Impacts across **long time horizon**

# Lecture 3

Return of the metrics,  
The challenges with evaluating models

Today's questions: **We will take an AI perspective today**

Is a model good enough for deployment?

Which model is better?

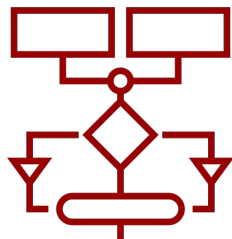
How do we design effective evaluation metrics?

How do we utilize these metrics within an appropriate evaluation protocols?

Main take away from today's lecture

Machine learning evaluation is a challenging  
unsolved problem.

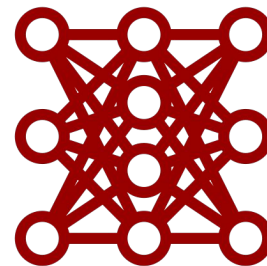
# A shift in AI: From algorithms to machine learning



Classical algorithms

**Problems:** precisely defined algebraically

**Example:** Graphcut algorithm

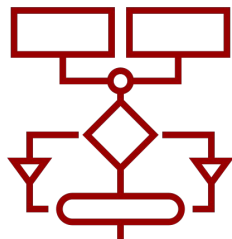


Empirical machine learning

**Problems:** loosely defined by datasets

**Example:** ResNet50 trained on ImageNet 1K

# A shift in AI: From algorithms to machine learning



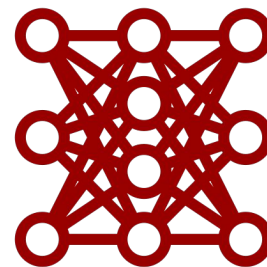
## Classical algorithms

**Problems:** precisely defined algebraically

**Example:** Graphcut algorithm

**Accuracy:** measured by correctness

**Artifact:** provably correct, transparent process



## Empirical machine learning

**Problems:** loosely defined by datasets

**Example:** ResNet50 trained on ImageNet 1K

**Accuracy:** measured using test set

**Artifact:** stochastic black box model

# Object Classification: The ImageNet task



[This image](#) by [Nikita](#) is licensed under [CC-BY 2.0](#)

(assume given a set of possible labels)  
{dog, cat, truck, plane, ...}

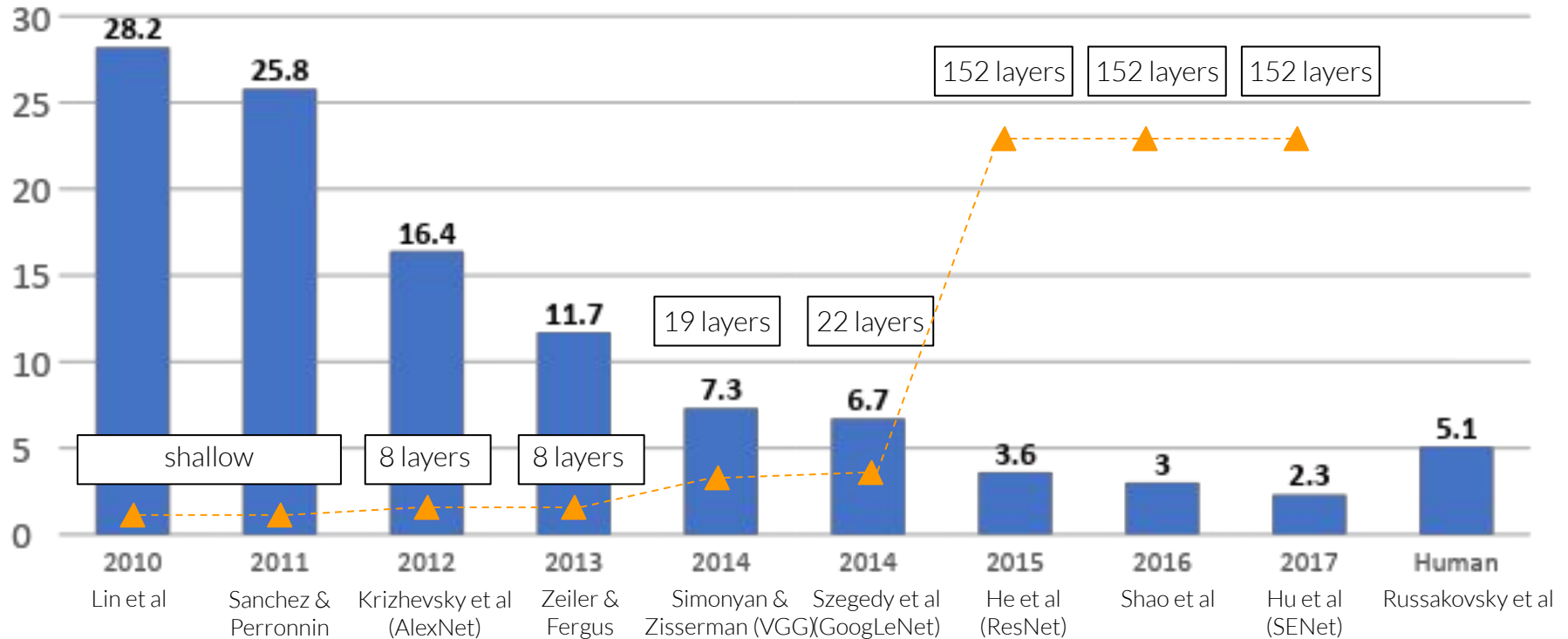


**cat**

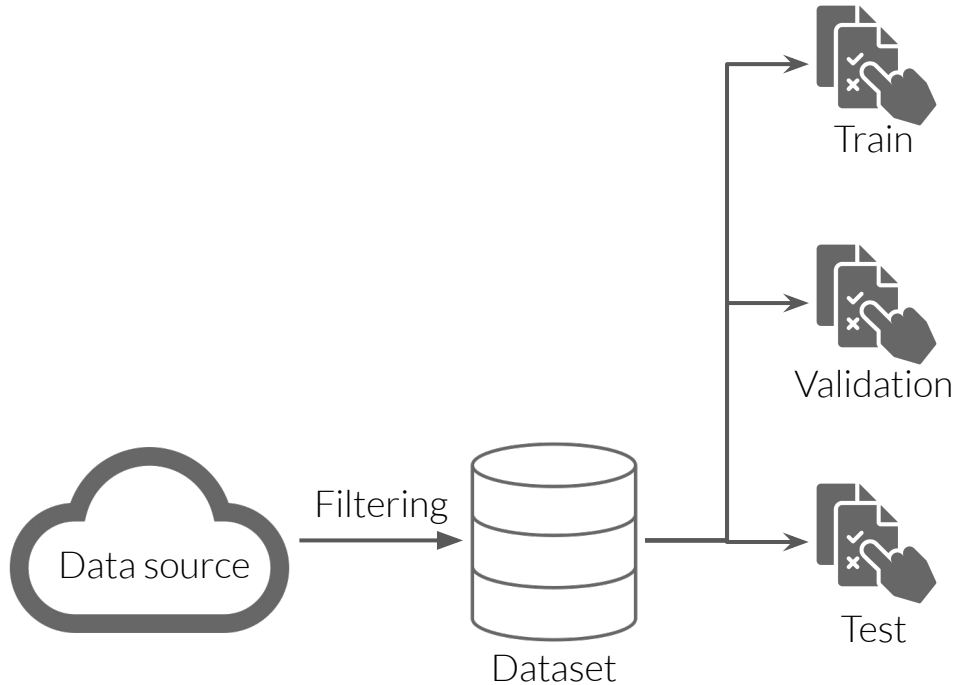
Evaluated using either top-1 or top-5 accuracy



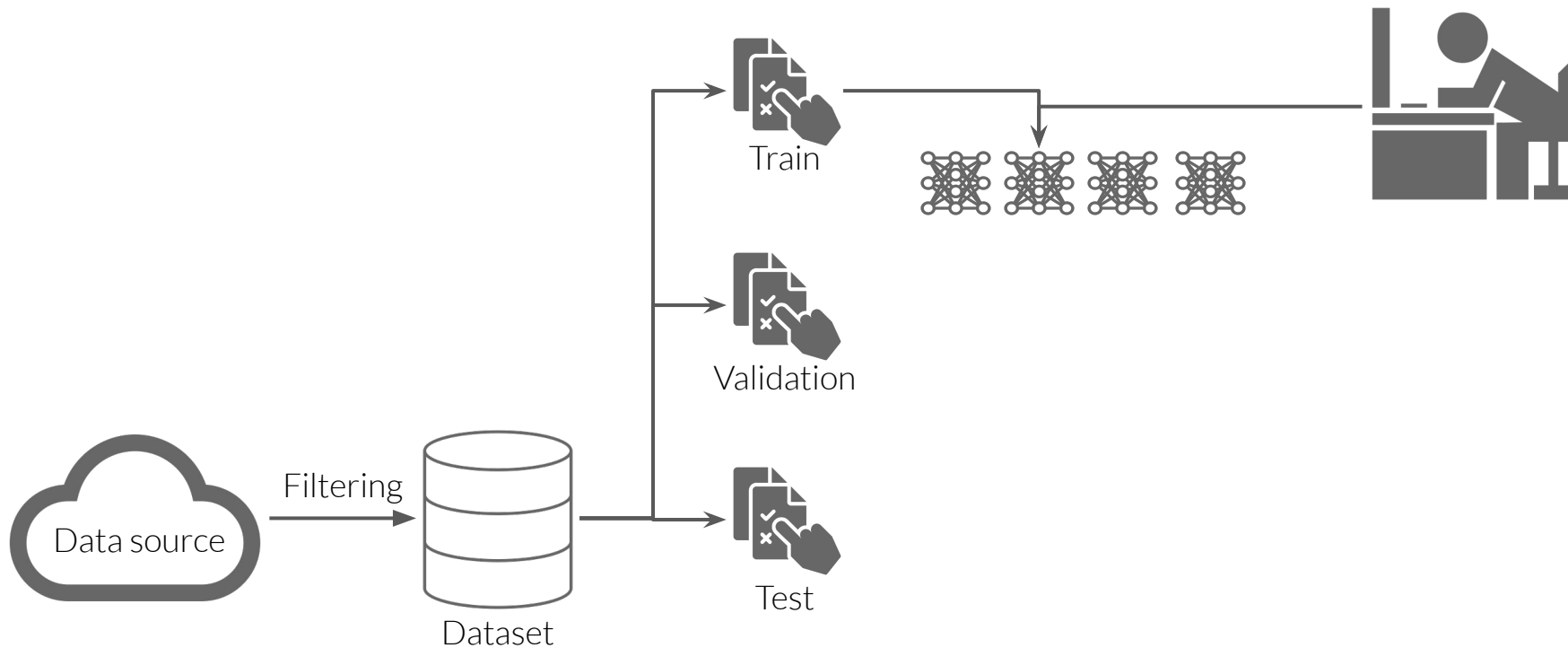
# Top-5 accuracy on ImageNet challenge over the years



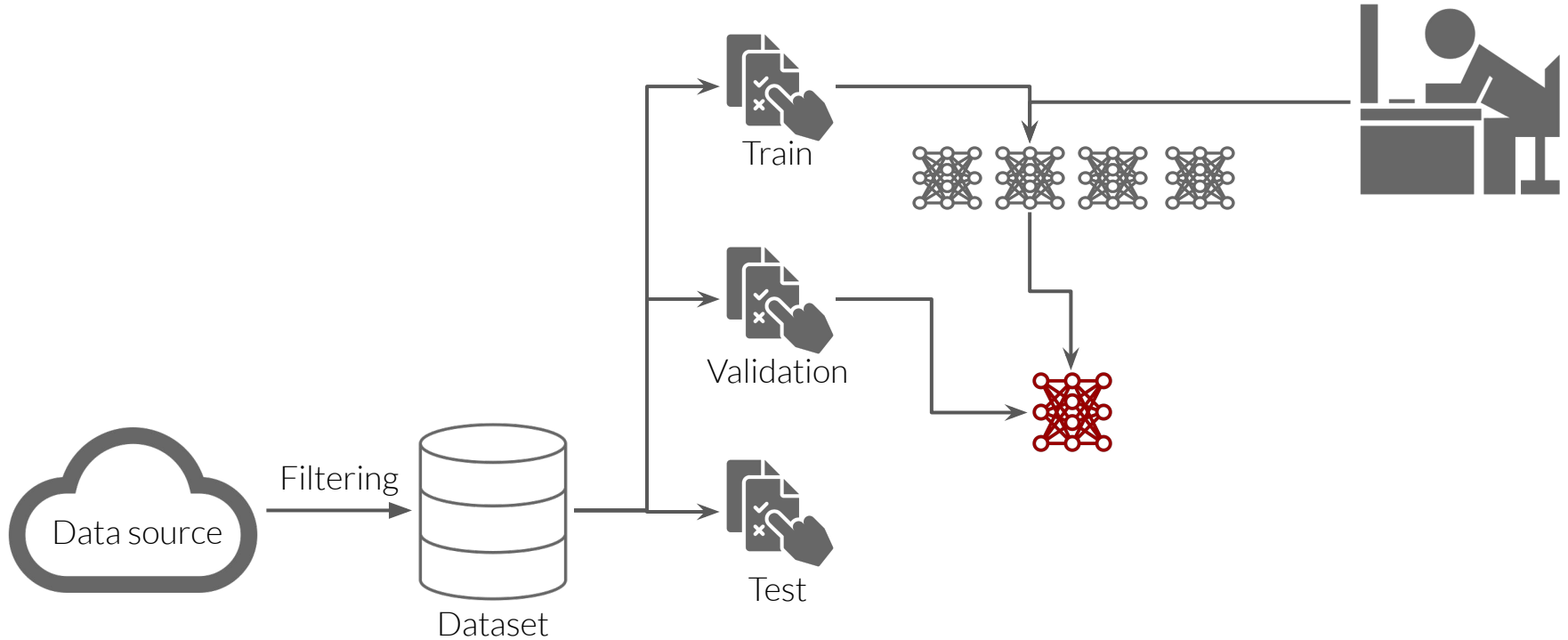
# High level evaluation protocol for empirical machine learning



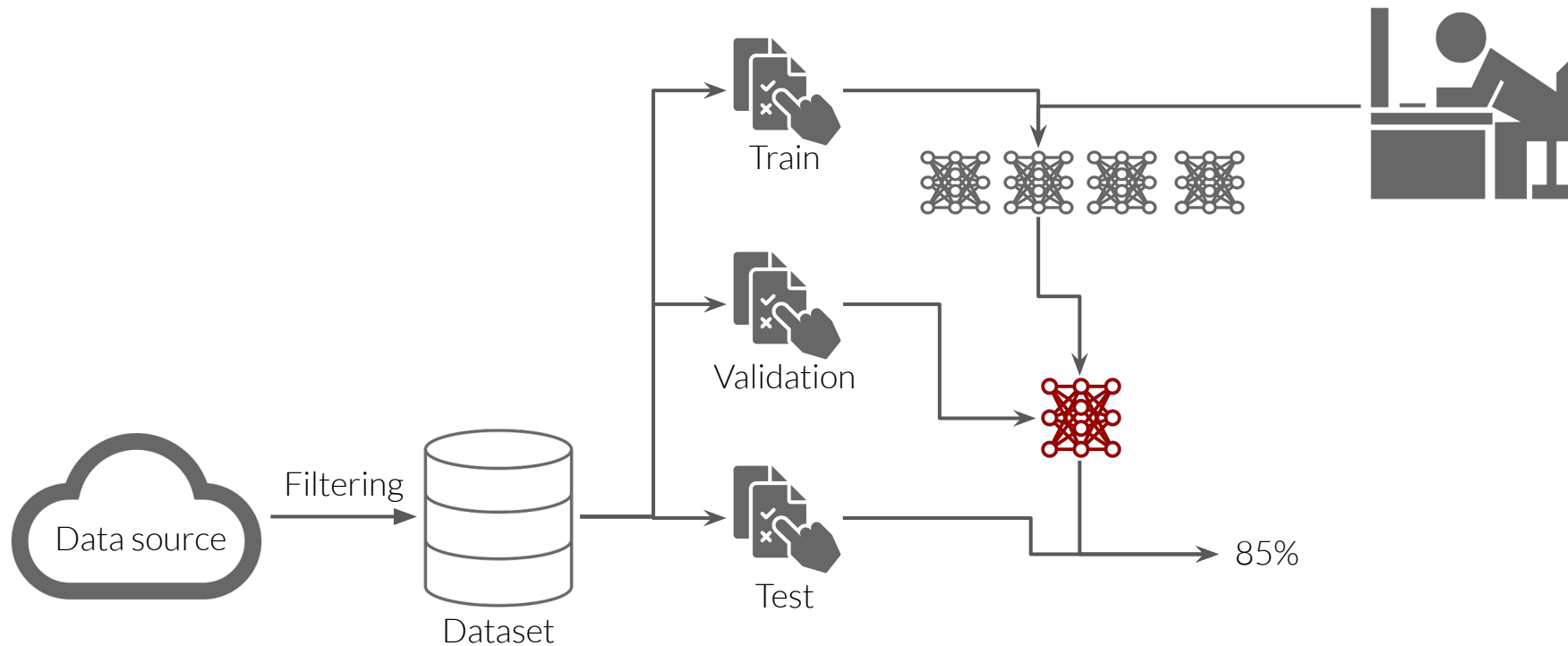
# High level evaluation protocol for empirical machine learning



# High level evaluation protocol for empirical machine learning



# High level evaluation protocol for empirical machine learning

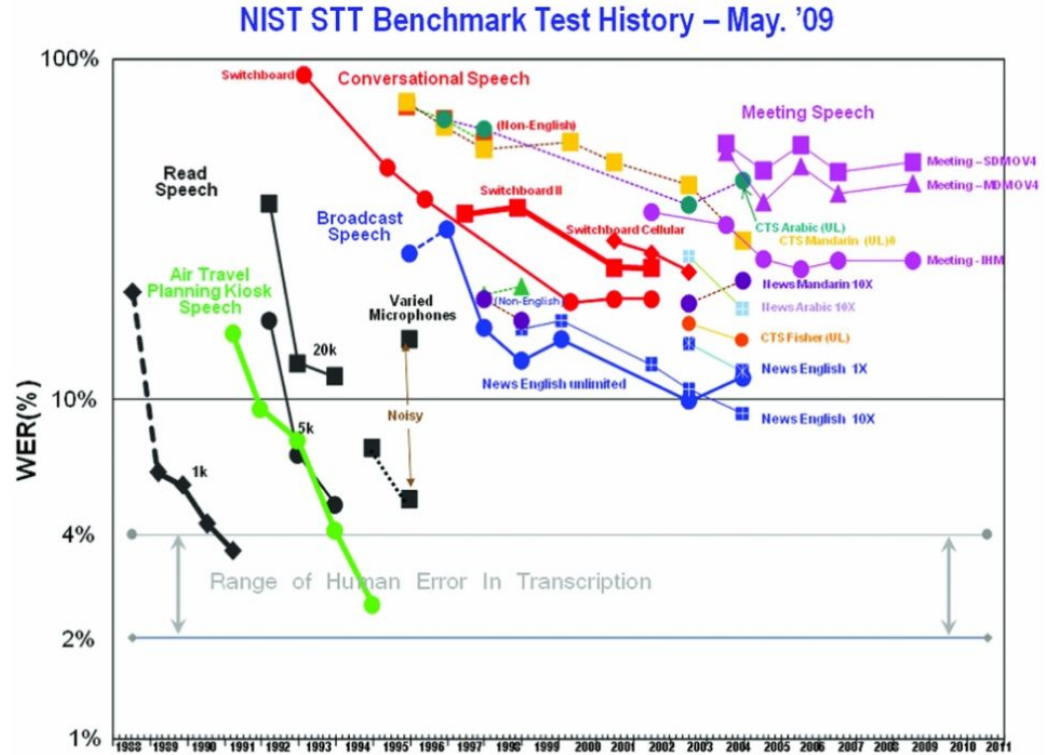


# Use of benchmark test datasets and common metrics

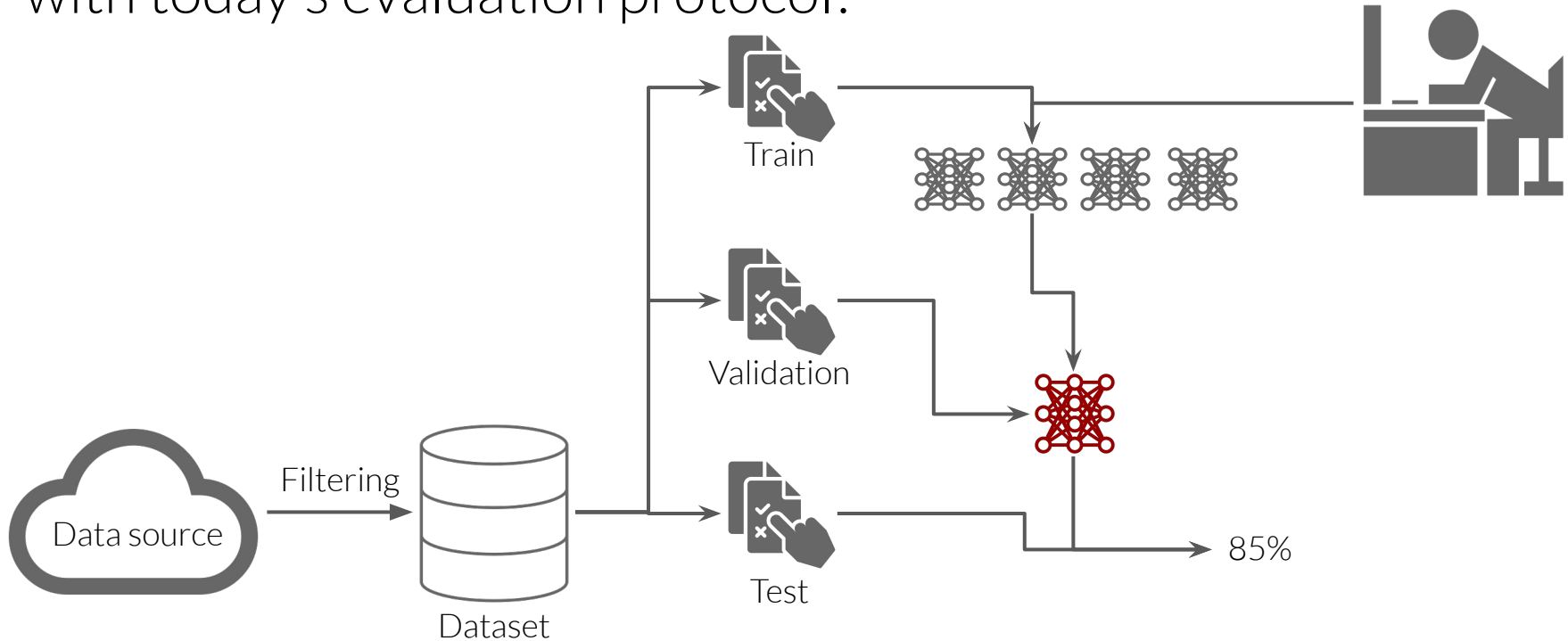
- Dates back to 1980s.
- Funded by DAPRA and led by IBM
- Goal: solve general diction problem
  
- Metric: Word error rate (WER)
- Artifact: a shared set of datasets, evaluation protocols, common metric, etc.

# UCI machine learning collection of datasets

Started in 1987 by David Aha and fellow graduate students at UC Irvine.



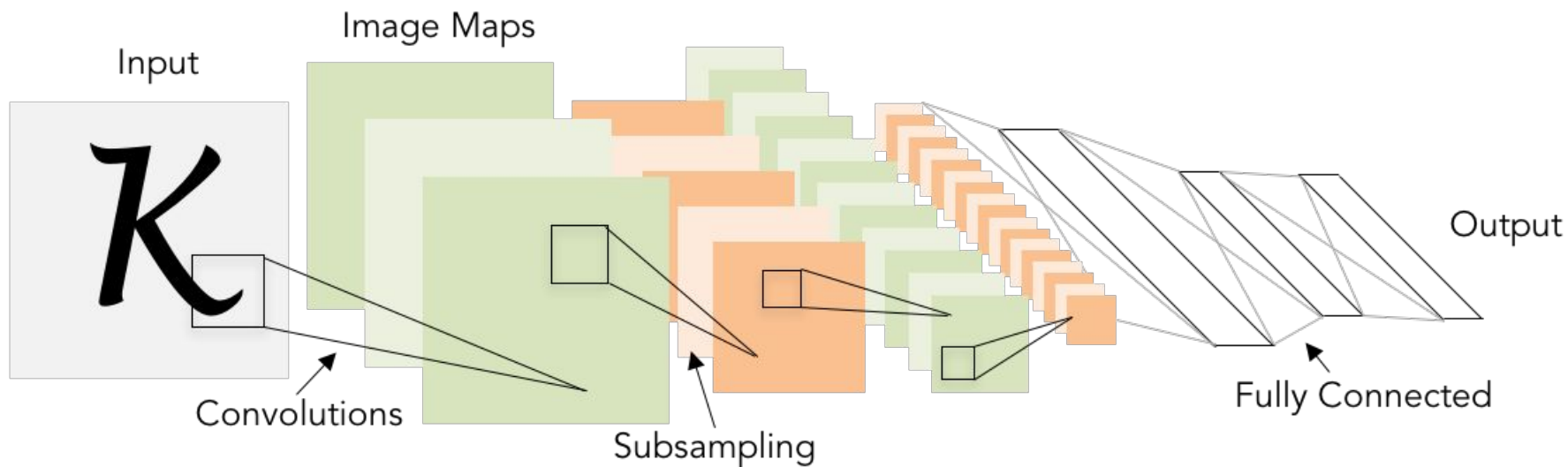
**Class activity:** So if things are working, can you think of issues with today's evaluation protocol?



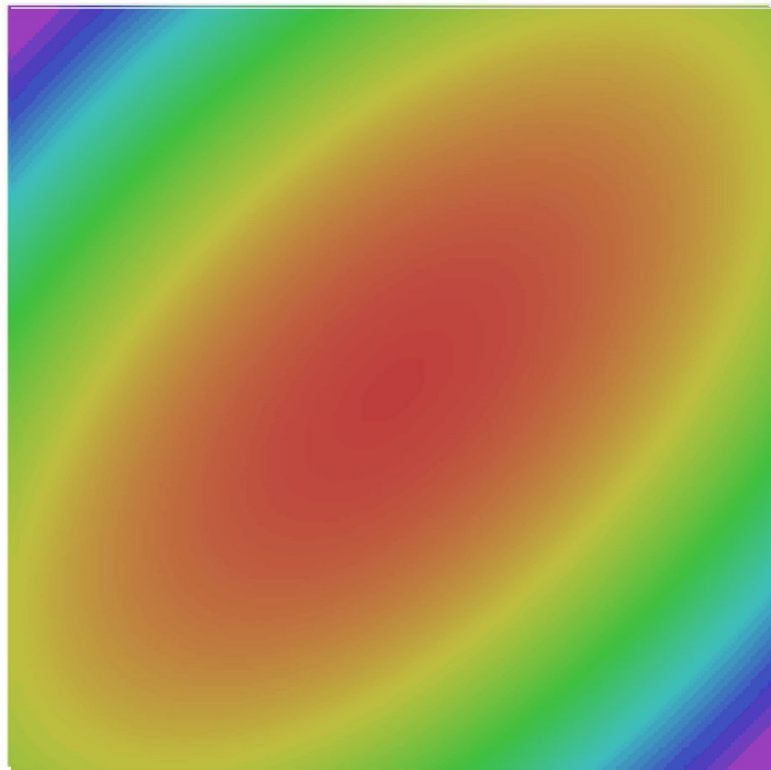


# #1: The replication crisis

Take a basic convolution neural network to solve object classification for instance



# Optimization options

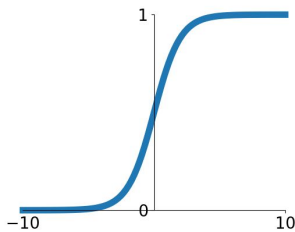


- SGD
- SGD+Momentum
- RMSProp
- Adam

# So many choices of activation Functions

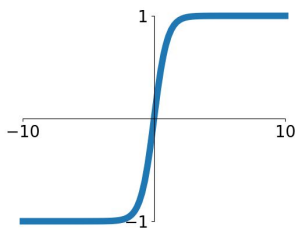
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



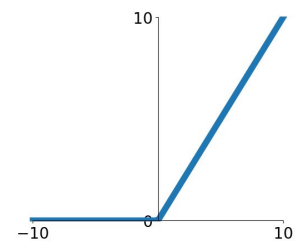
tanh

$$\tanh(x)$$



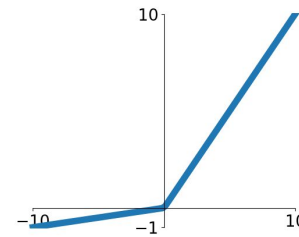
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

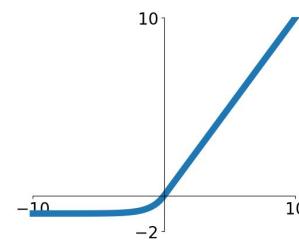


Maxout

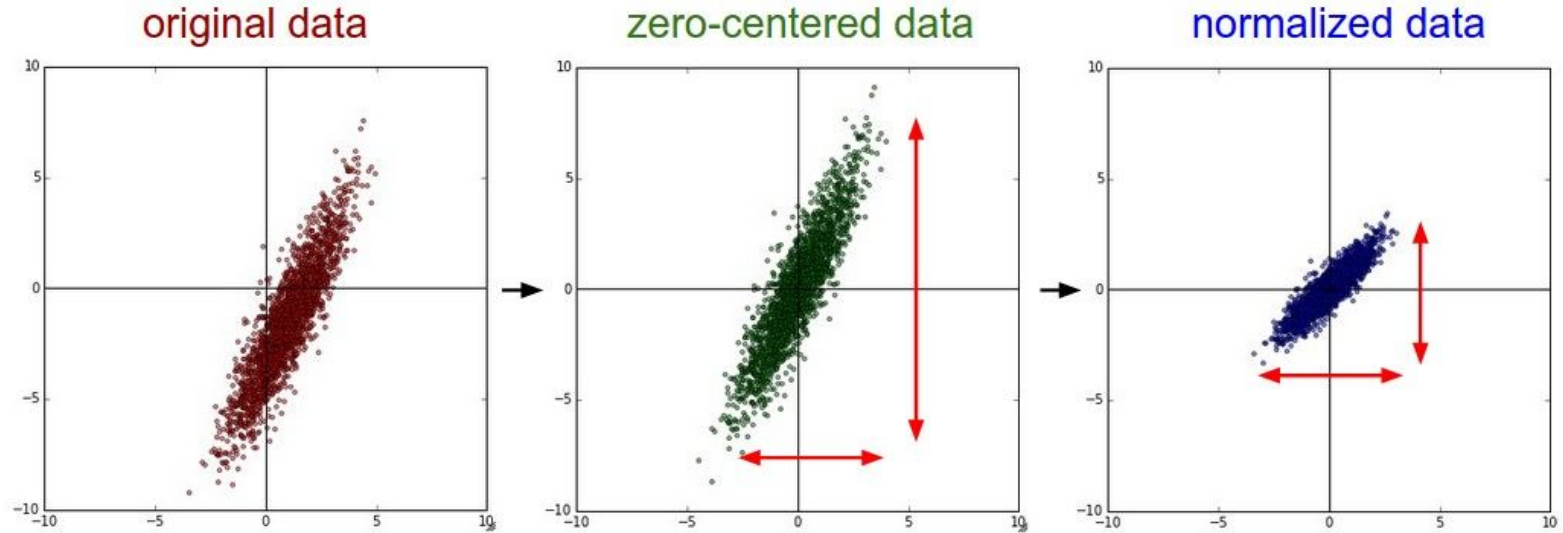
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# Data preprocessing



# Regularization options: e.g. Mixup

Training: Train on random blends of images

Testing: Use original images

## Examples:

Dropout

Batch Normalization

Data Augmentation

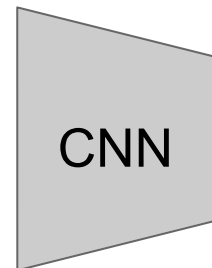
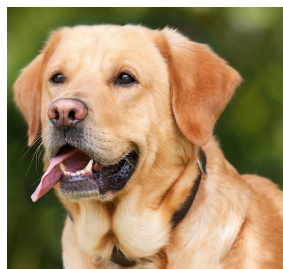
DropConnect

Fractional Max Pooling

Stochastic Depth

Cutout / Random Crop

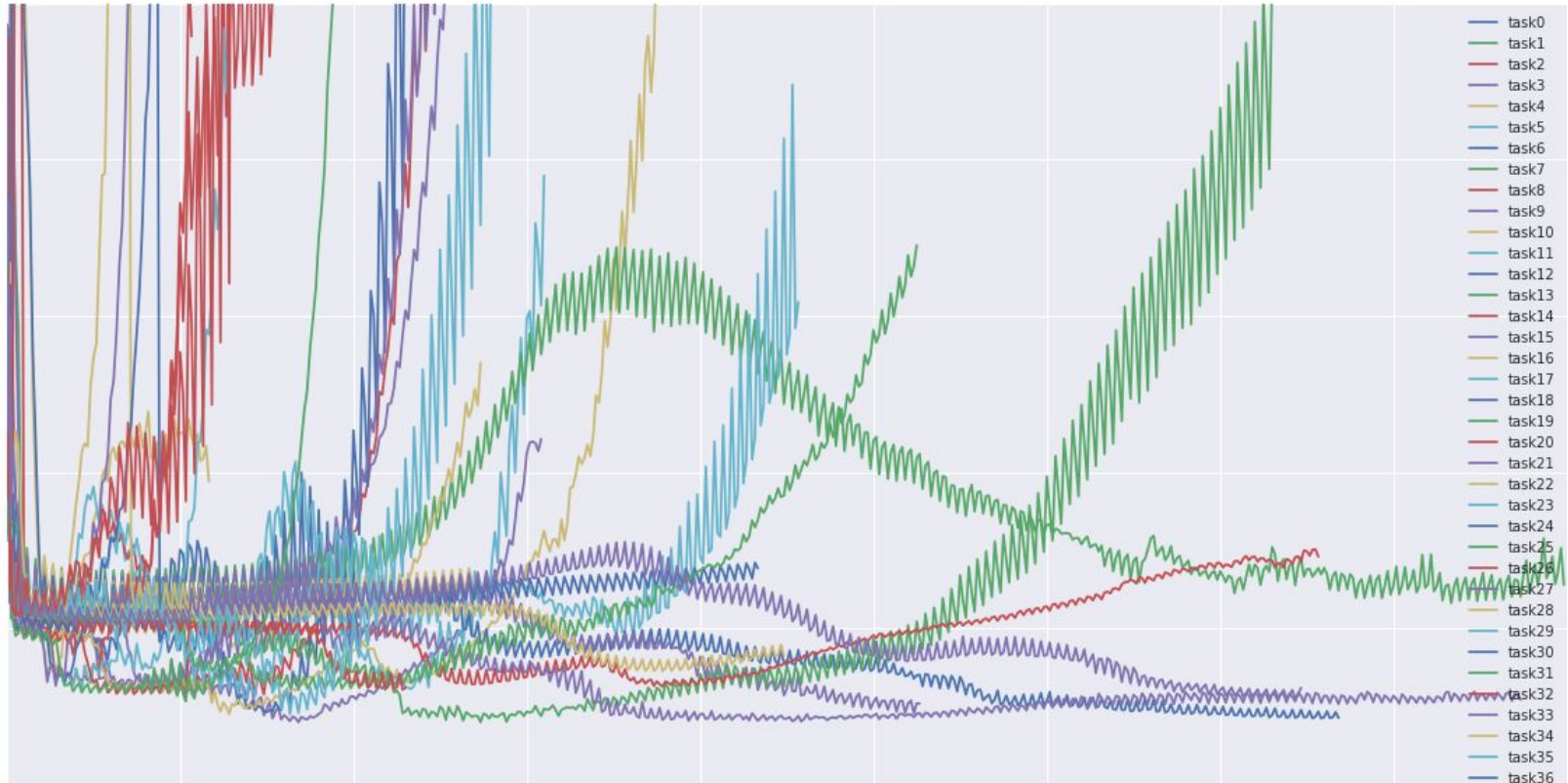
Mixup



Target label:  
cat: 0.4  
dog: 0.6

Randomly blend the pixels of pairs of training images, e.g. 40% cat, 60% dog

Loss curves are often used instead of real metrics to make decisions



# Hardware + Software options

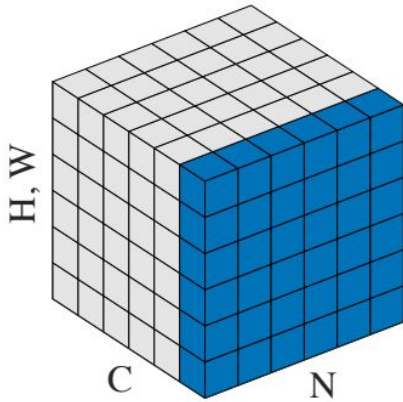


PyTorch

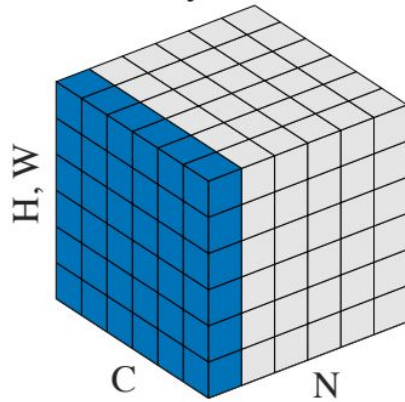
TensorFlow

# Normalization layer options

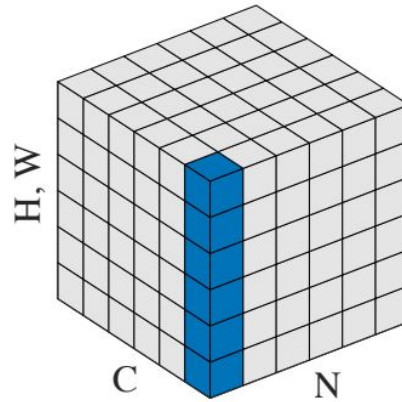
Batch Norm



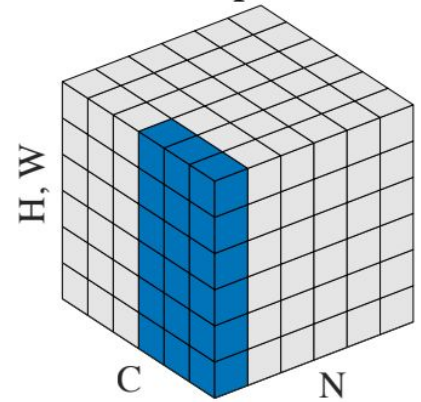
Layer Norm



Instance Norm



**Group Norm**



Wu and He, "Group Normalization", ECCV 2018

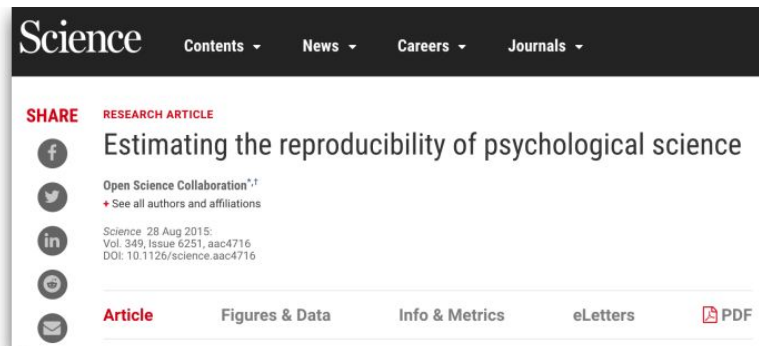
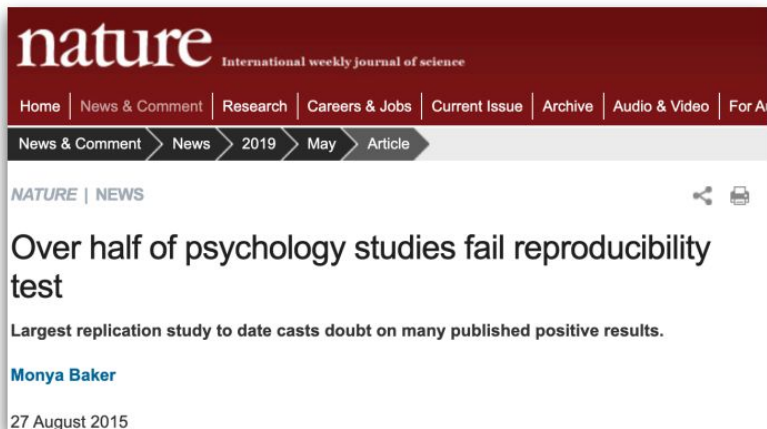


# #1: The replication crisis

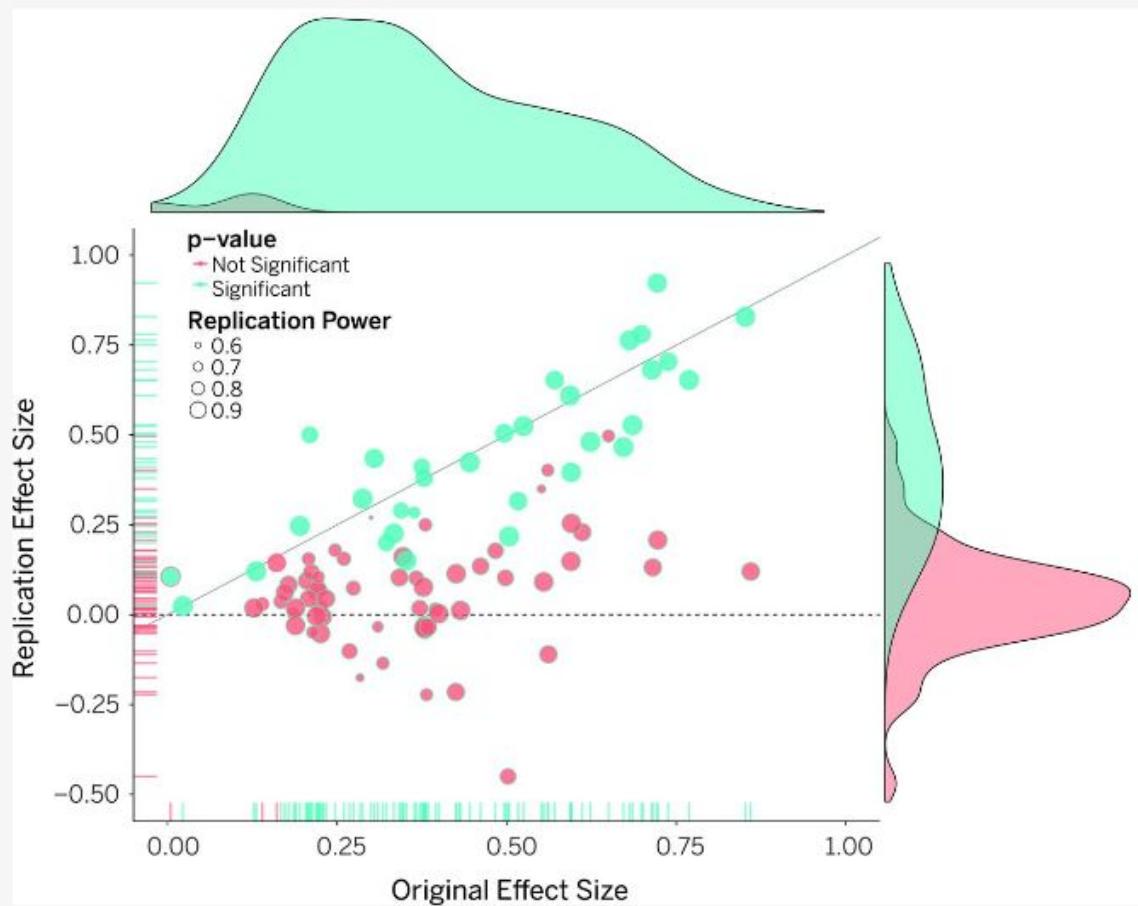
All these details are lost in appendixes or during experiments.

Anecdote: sometimes we can't reproduce our own results because of other processes interfering.

# #1: The replication crisis: not just a machine learning challenge



# Bad news



**Original study effect size versus replication effect size (correlation coefficients).**

Diagonal line represents replication effect size equal to original effect size. Dotted line represents replication effect size of 0. Points below the dotted line were effects in the opposite direction of the original. Density plots are separated by significant (blue) and nonsignificant (red) effects.

Open Science Collaboration.  
Estimating the reproducibility of  
psychological science. Science 2015

# #2: Labeling errors

MNIST

CIFAR-10

CIFAR-100

Caltech-256

ImageNet

QuickDraw

correctable



given: 8  
corrected: 9



given: cat  
corrected: frog



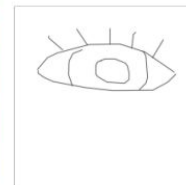
given: lobster  
corrected: crab



given: dolphin  
corrected: kayak



given: white stork  
corrected: black stork



given: tiger  
corrected: eye

multi-label

(N/A)

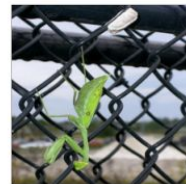
(N/A)



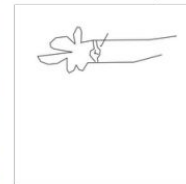
given: hamster  
also: cup



given: laptop  
also: people



given: mantis  
also: fence



given: wristwatch  
also: hand

neither



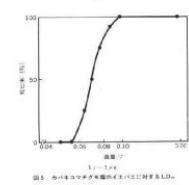
given: 6  
alt: 1



given: deer  
alt: bird



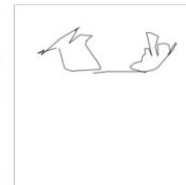
given: rose  
alt: apple



given: house-fly  
alt: ladder



given: polar bear  
alt: elephant



given: pineapple  
alt: raccoon

non-agreement



given: 4  
alt: 9



given: automobile  
alt: airplane



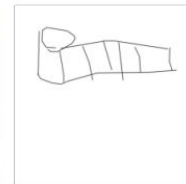
given: dolphin  
alt: ray



given: yo-yo  
alt: frisbee



given: eel  
alt: flatworm



given: bandage  
alt: roller coaster

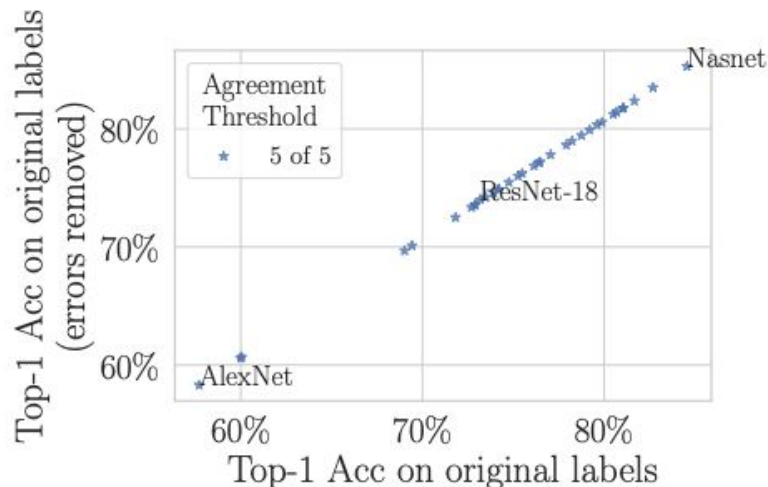
Northcutt et al. Pervasive  
Label Errors in Test Sets  
Destabilize Machine  
Learning Benchmarks.  
NeurIPS 2021

## #2: Labeling errors: % errors in test sets

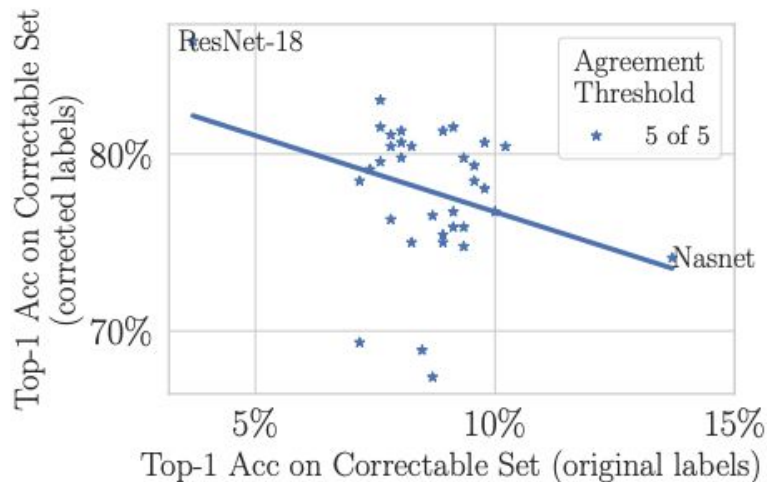
Dataset	Modality	Size	Model	Test Set Errors				% error
				CL guessed	MTurk checked	validated	estimated	
MNIST	image	10,000	2-conv CNN	100	100 (100%)	15	-	0.15
CIFAR-10	image	10,000	VGG	275	275 (100%)	54	-	0.54
CIFAR-100	image	10,000	VGG	2,235	2,235 (100%)	585	-	5.85
Caltech-256 <sup>†</sup>	image	29,780	Wide ResNet-50-2	2,360	2,360 (100%)	458	-	1.54
ImageNet*	image	50,000	ResNet-50	5,440	5,440 (100%)	2,916	-	5.83
QuickDraw <sup>†</sup>	image	50,426,266	VGG	6,825,383	2,500 (0.04%)	1870	5,105,386	10.12
20news	text	7,532	TFIDF + SGD	93	93 (100%)	82	-	1.09
IMDB	text	25,000	FastText	1,310	1,310 (100%)	725	-	2.90
Amazon Reviews <sup>†</sup>	text	9,996,437	FastText	533,249	1,000 (0.2%)	732	390,338	3.90
AudioSet	audio	20,371	VGG	307	307 (100%)	275	-	1.35

\*Because the ImageNet test set labels are not publicly available, the ILSVRC 2012 validation set is used.

## #2: Labeling errors: Errors make larger models overfit

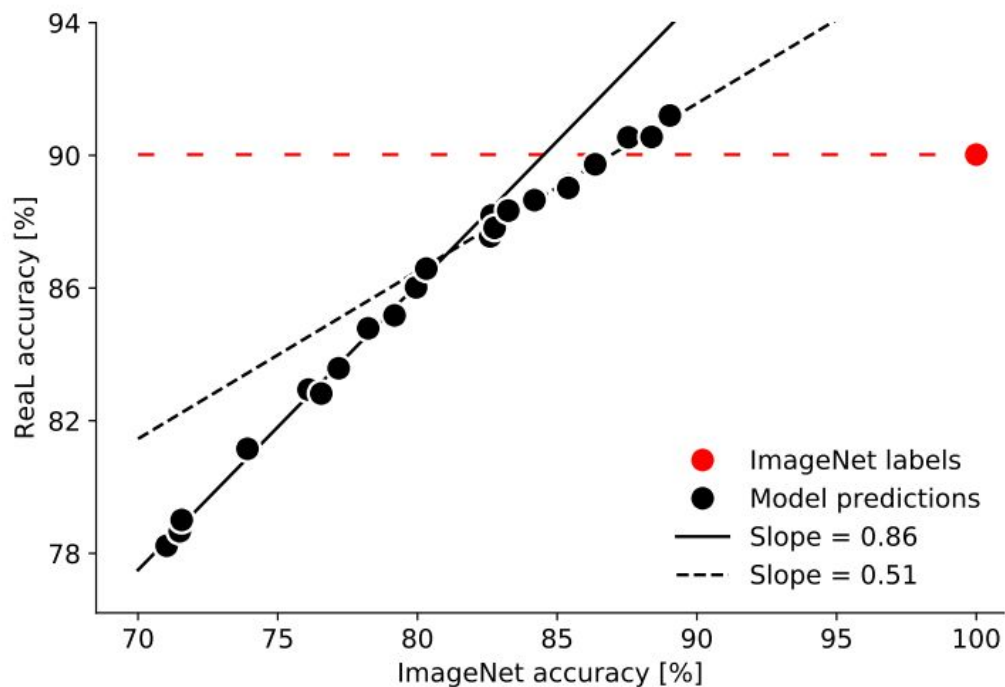


(a) ImageNet val set acc.



(b) ImageNet correctable set acc.

## #2: Labeling errors: Relabeled ImageNet test set

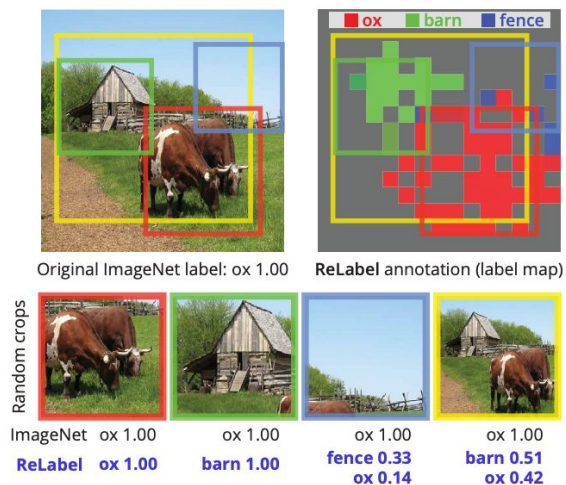


Gains reported using fixed labels is smaller than those with original ImageNet labels is

Beyer et al. Are we done with ImageNet?  
2020

## #2: Labeling errors: Of course the training set also has errors

Gains reported using fixed labels is smaller than those with original ImageNet labels is



Variants	ImageNet top-1 (%)
ReLabel (localized mutli-labels)	78.9
Localized single labels	78.4 (-0.5)
Global multi-labels	78.5 (-0.4)
Global single labels	77.5 (-1.4)
Original ImageNet labels	77.5 (-1.4)

Yun et al. Re-labeling ImageNet: from Single to Multi-Labels, from Global to Localized Labels. CVPR 2021



# #3: Generalization errors: Test sets represent a small slice of the real world.

Models may have seen:

- people,
- phones,
- bottles,
- people holding bottles



Can they generalize to:

- People holding phones?

Example compositional spatio-temporal questions:

Q: What did the person **hold** after **putting a phone somewhere**?

A: **bottle**

Q: Were they **taking a picture** or **holding a bottle** for longer?

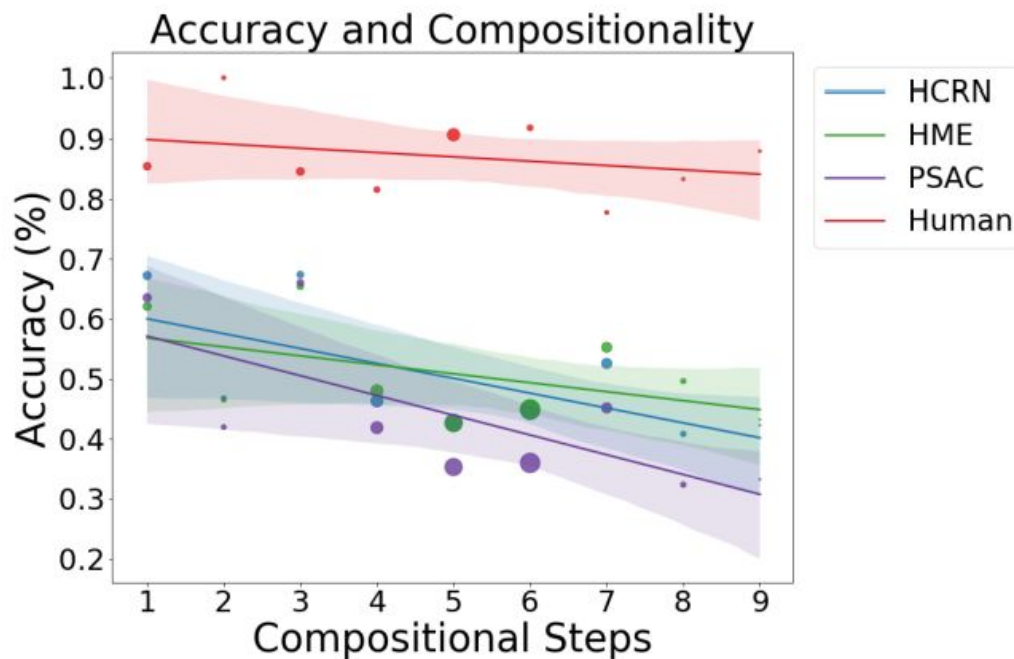
A: **holding a bottle**

Q: Did they **take a picture** before or after they did **the longest action**?

A: **before**

### #3: Generalization errors: Systematic generalization in video understanding decreases as composition steps increase

- Human performance: 86%
- Best model performance: 48%




# #3: Generalization errors: Maybe videos are too hard... what about images?


CREPE: a benchmark to test for compositional generalization of CLIP and other image-text models

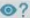
Can models at the very least generalize to new compositions of seen concepts?

Ma et al. CREPE: Can Vision-Language Foundation Models Reason Compositionally? ArXiv 2023


 **CREPE-Systematicity**

Unseen atoms  
Seen compounds  
Unseen compounds




✓ Crepe on a skillet. 

- ✗ Boats on a skillet.
- ✗ Crepe under a skillet.
- ✗ Crepe on a dog.

 **CREPE-Productivity**

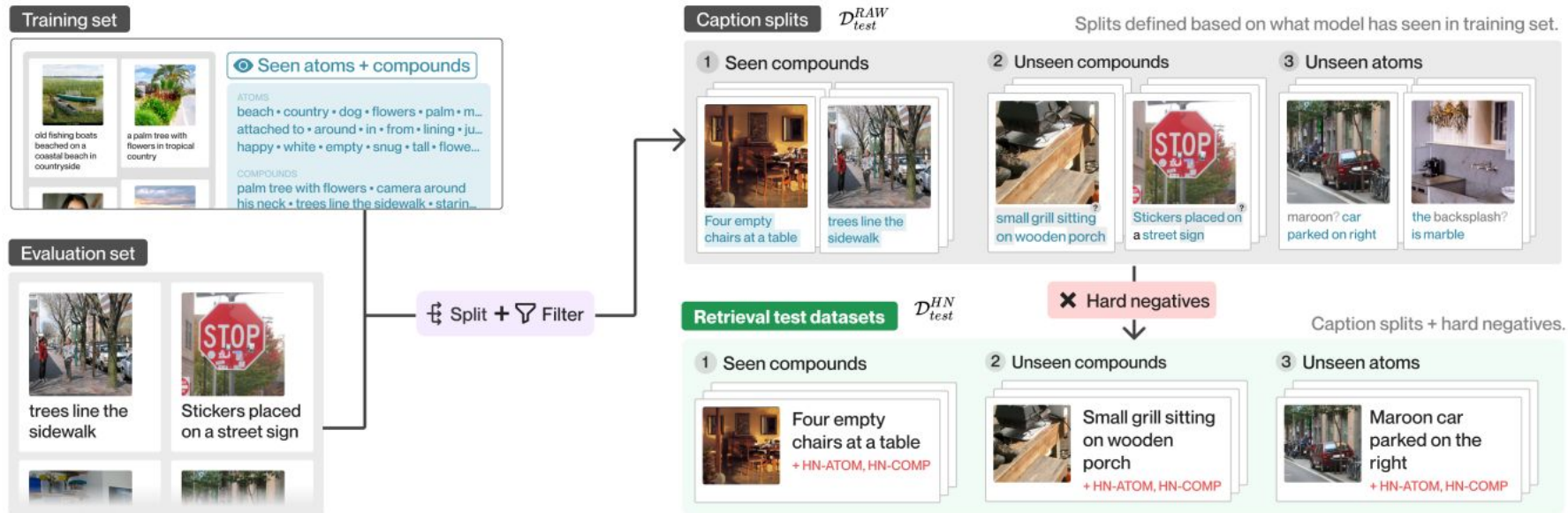
other  $n +$  negative types  
 $n = 9$  • swap negatives  
 $n = 8$  • atomic negatives



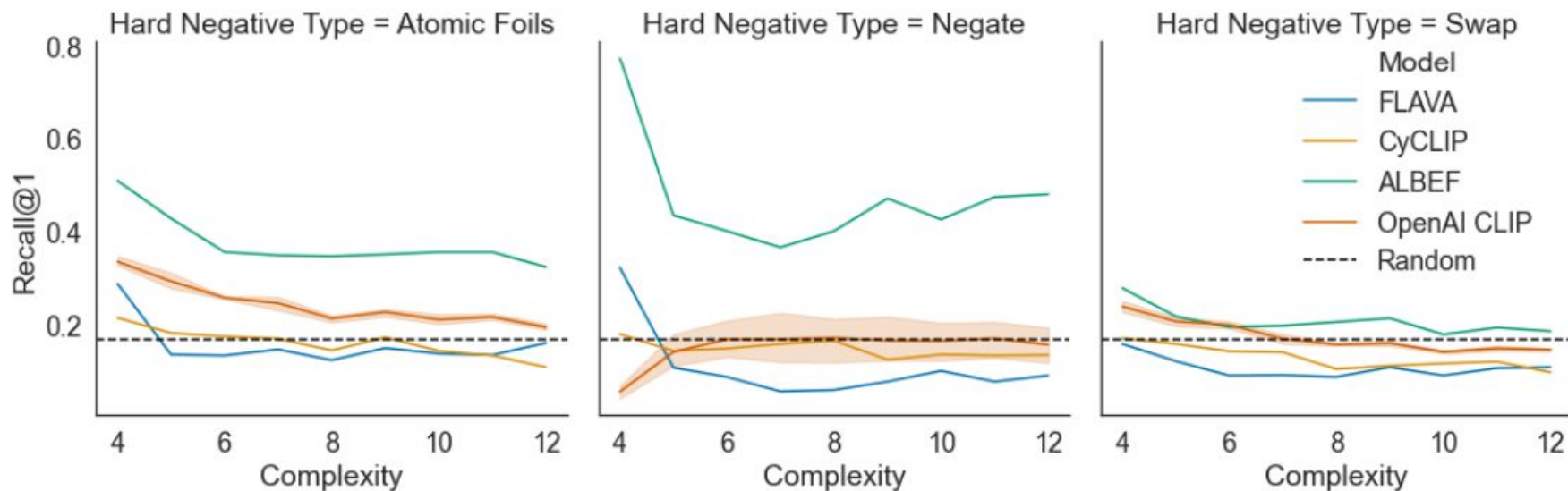
✓ Browned crepe next to leafy salad and in front of metal fork.

- ✗ Blue crepe next to leafy salad and in front of metal fork.
- ✗ Browned chair next to leafy salad and in front of metal fork.

# #3: Generalization errors: compositional generalization



### #3: Generalization errors: today's models can't represent composition in language or vision



### #3: Generalization errors: Increasing model size or increasing dataset size doesn't improve compositional generalization

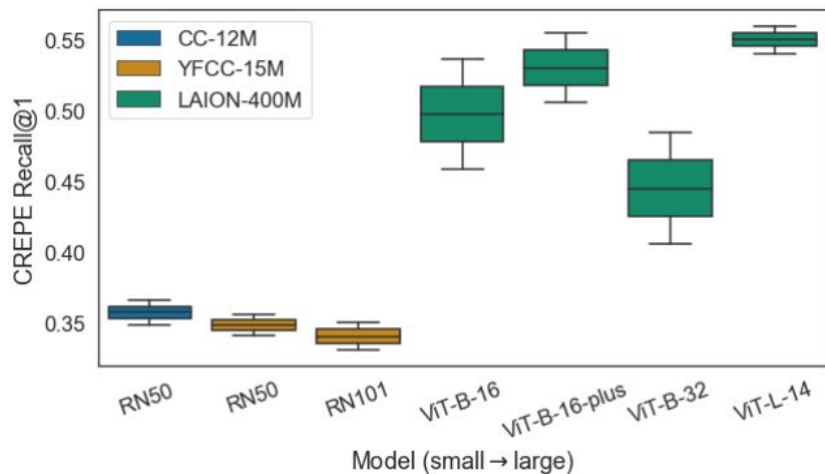


Figure 6. *Systematicity Analysis*. We plot the retrieval Recall@1 of all models pre-trained on the three datasets and observe no particular correlation with model size within datasets.

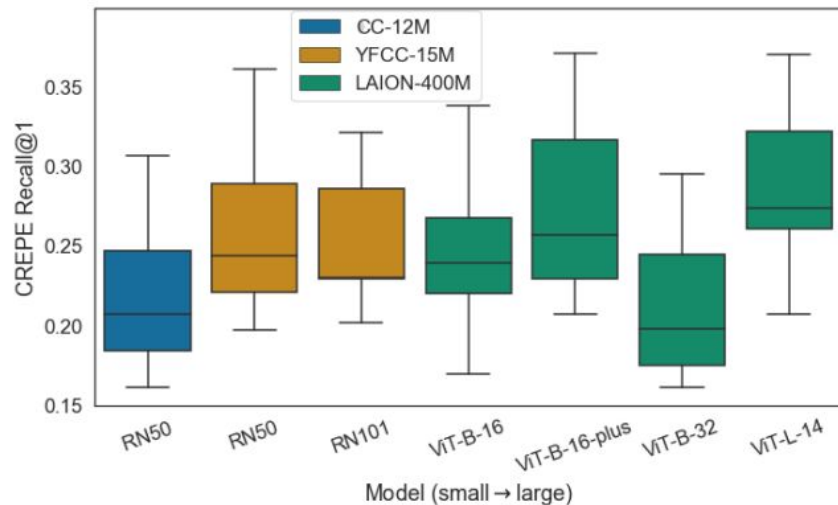
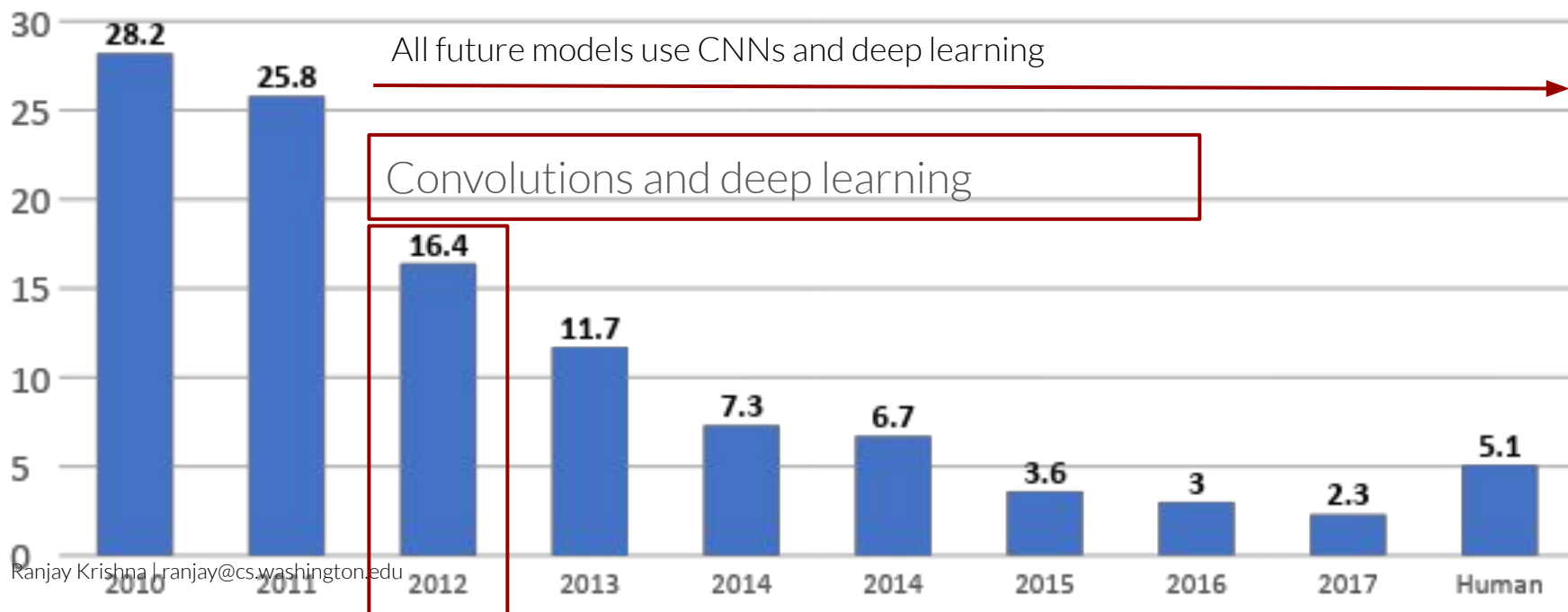
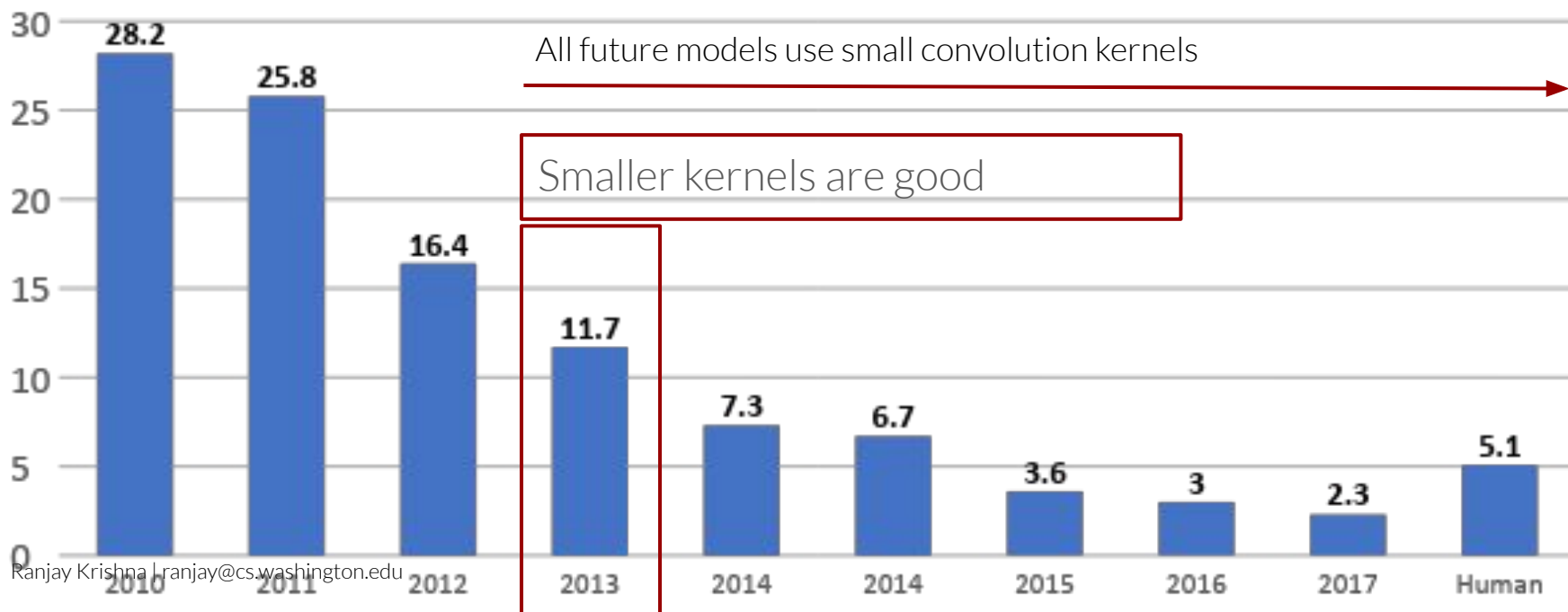


Figure 7. *Productivity Analysis*. We plot the retrieval Recall@1 of all models trained on all three datasets. We observe that there is no consistent correlation with model size within datasets.

## #4: Static test sets: Reusing test sets every year?

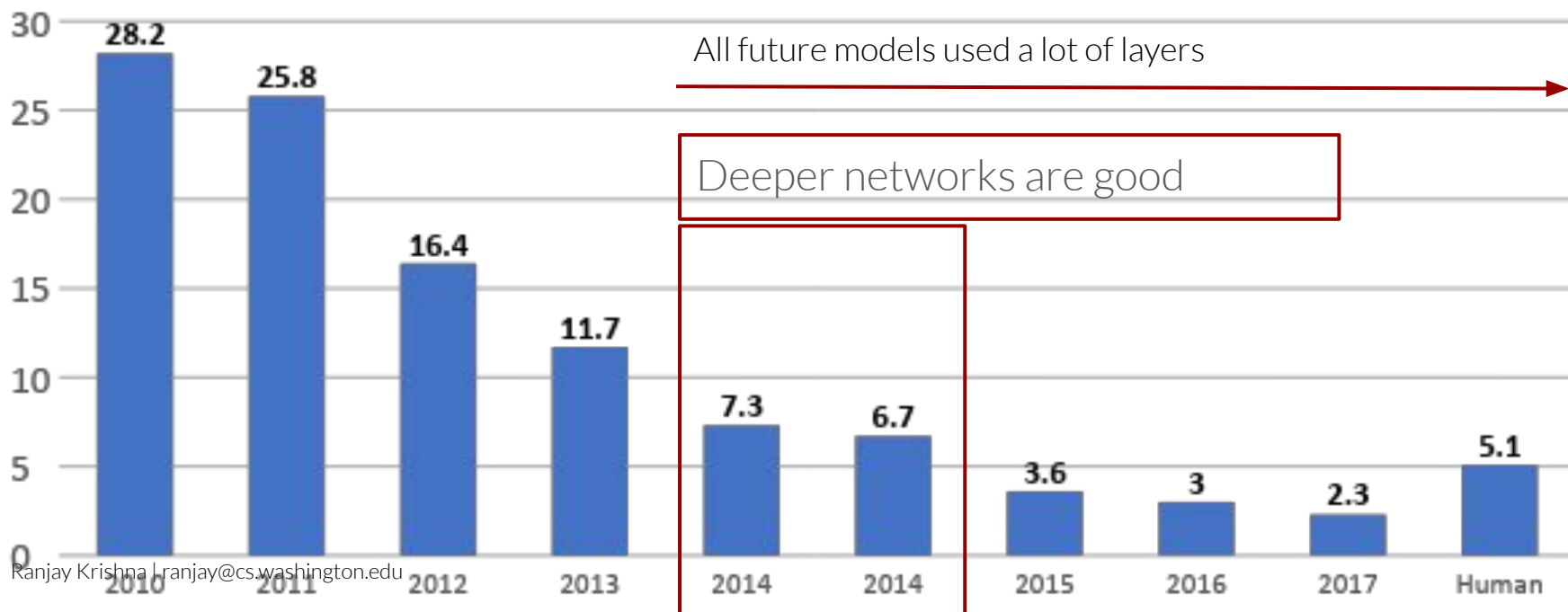


## #4: Static test sets: Reusing test sets every year?

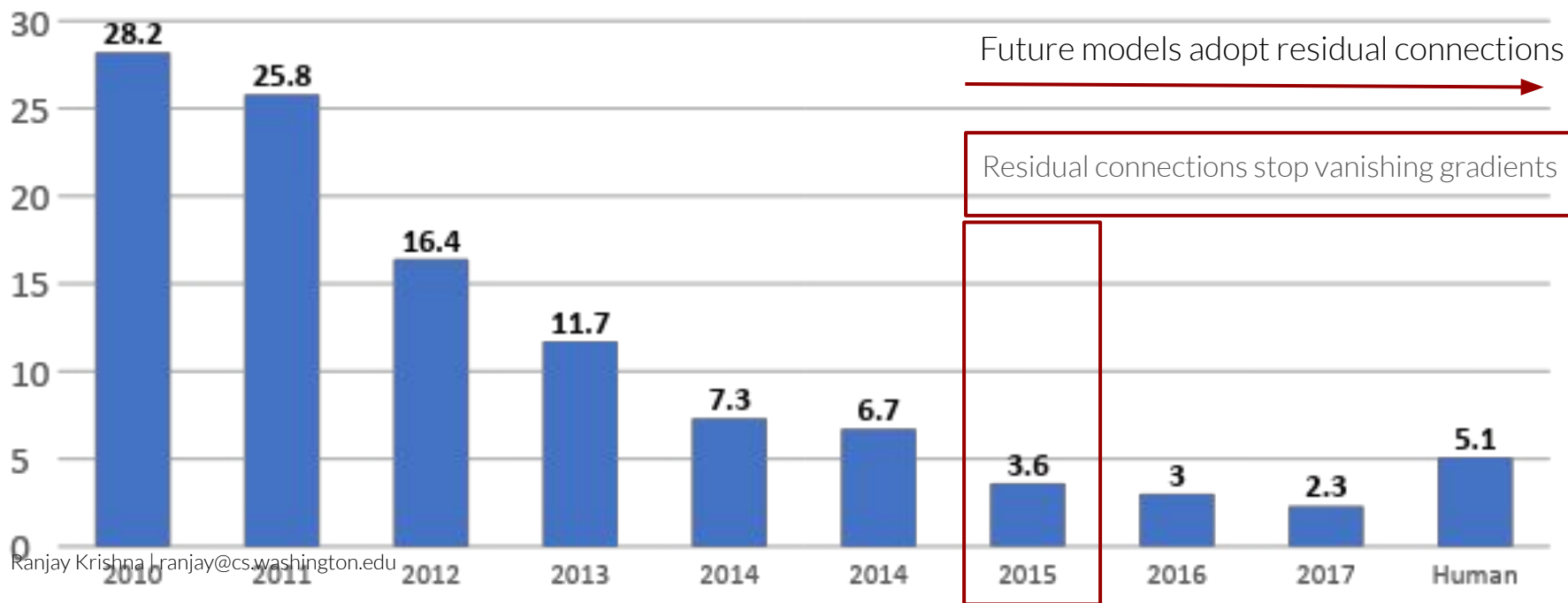




## #4: Static test sets: Reusing test sets every year?

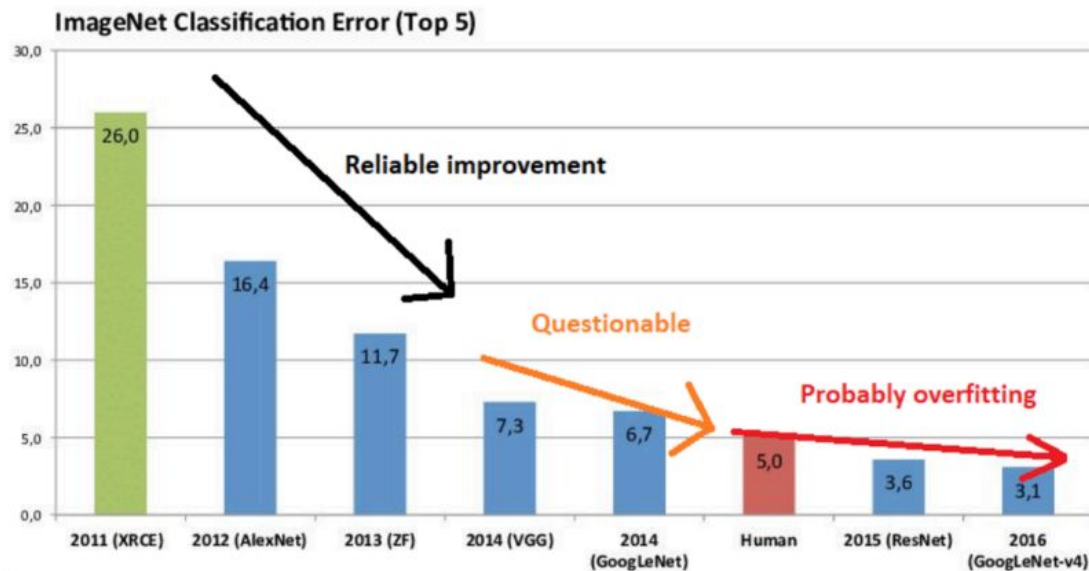


## #4: Static test sets: Reusing test sets every year?



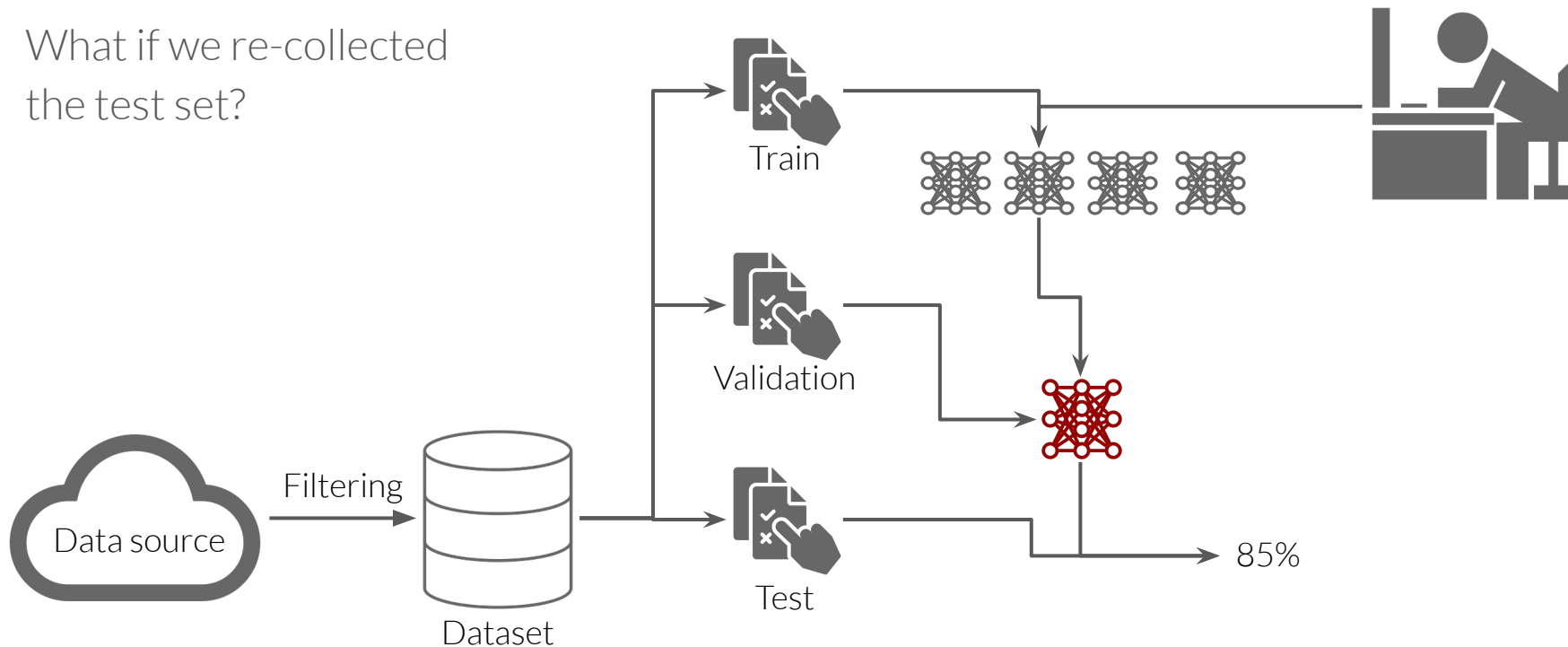
## #4: A static dataset: Are models overfitting to the test set?

AI competitions don't produce useful models



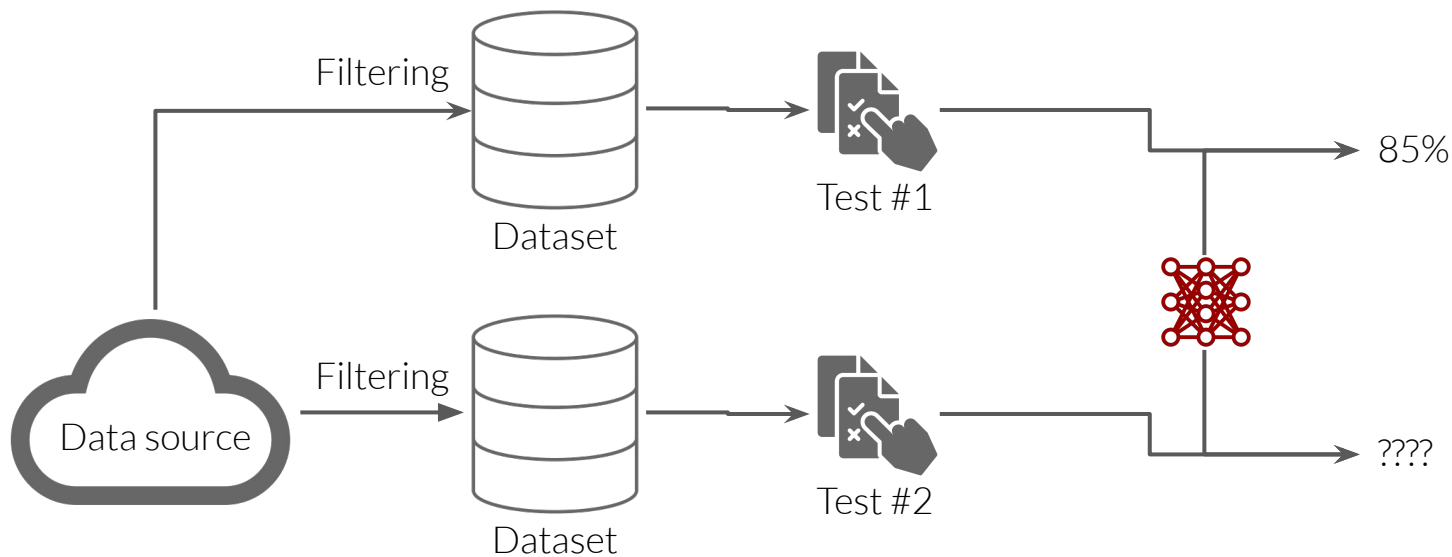
## #4: A static dataset: Let's collect a new test set

What if we re-collected the test set?

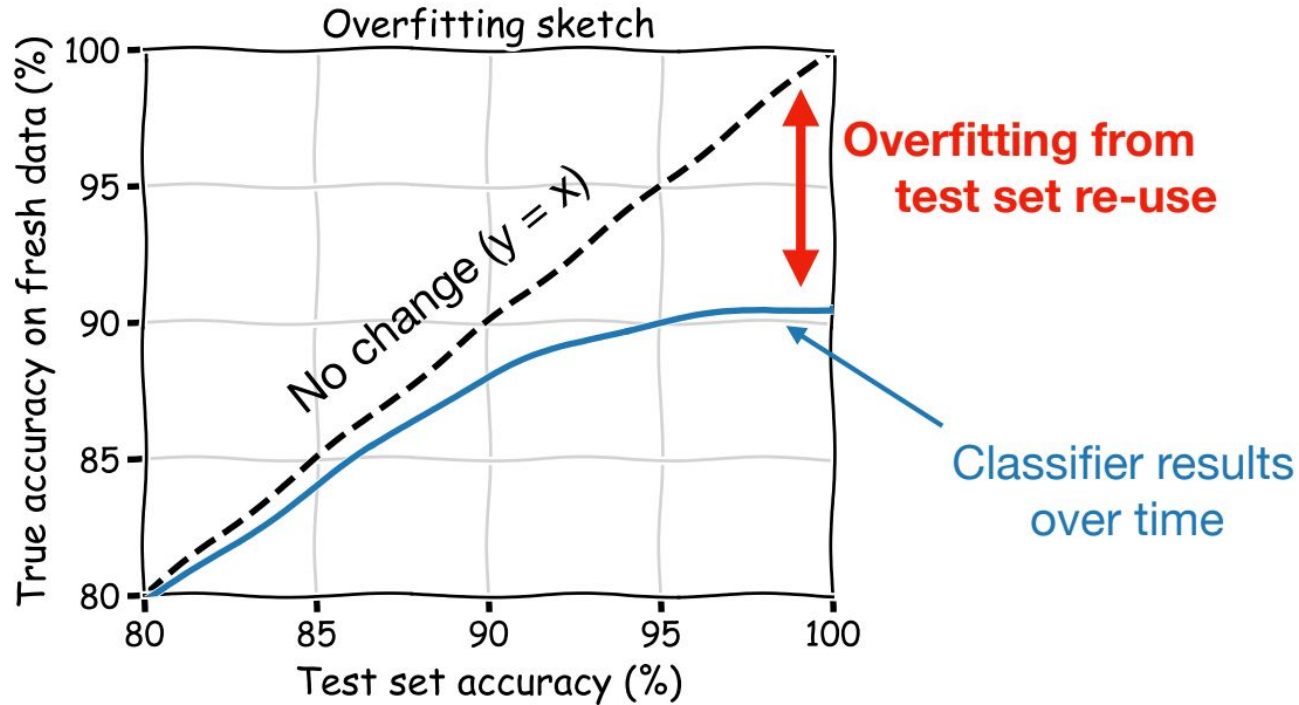


## #4: A static dataset: Checking for overfitting

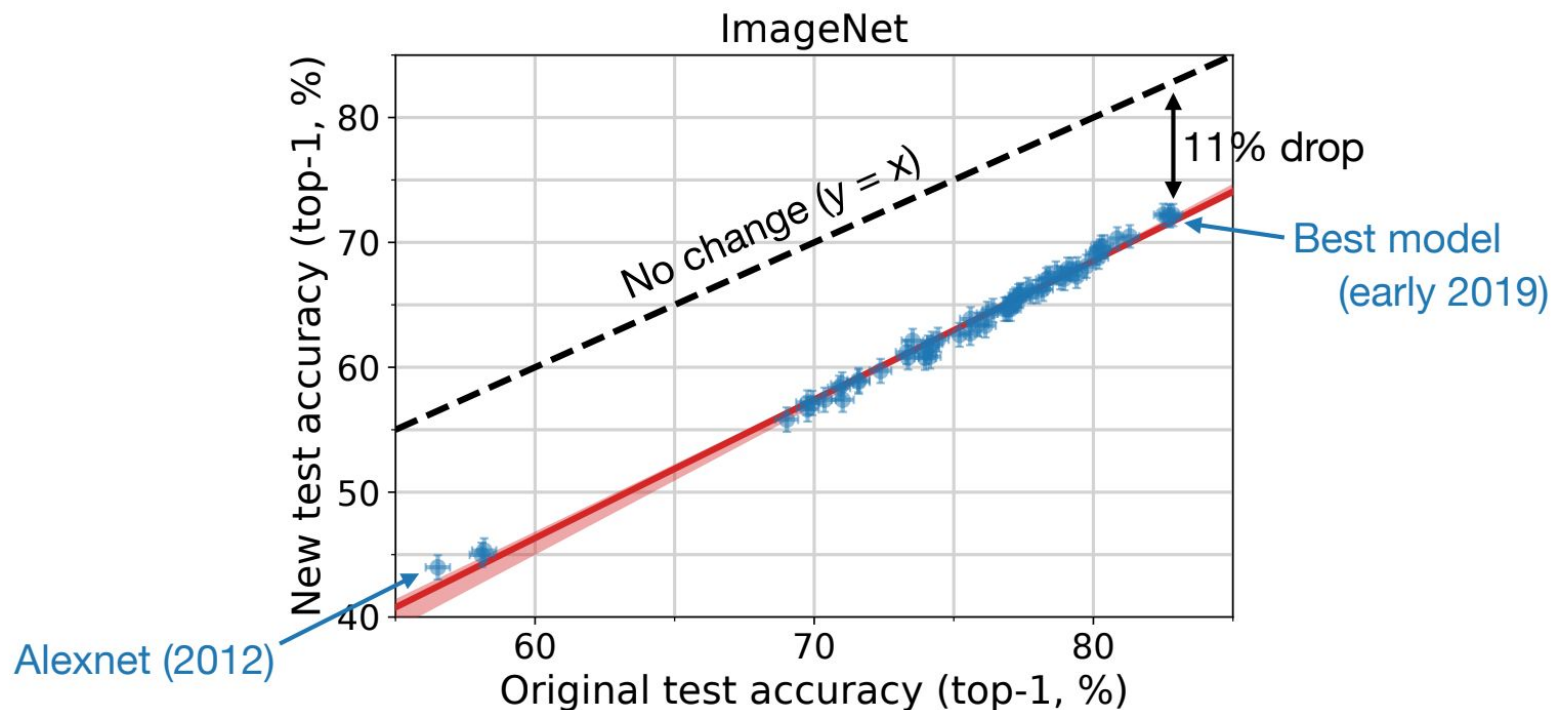
What if we re-collected the test set?



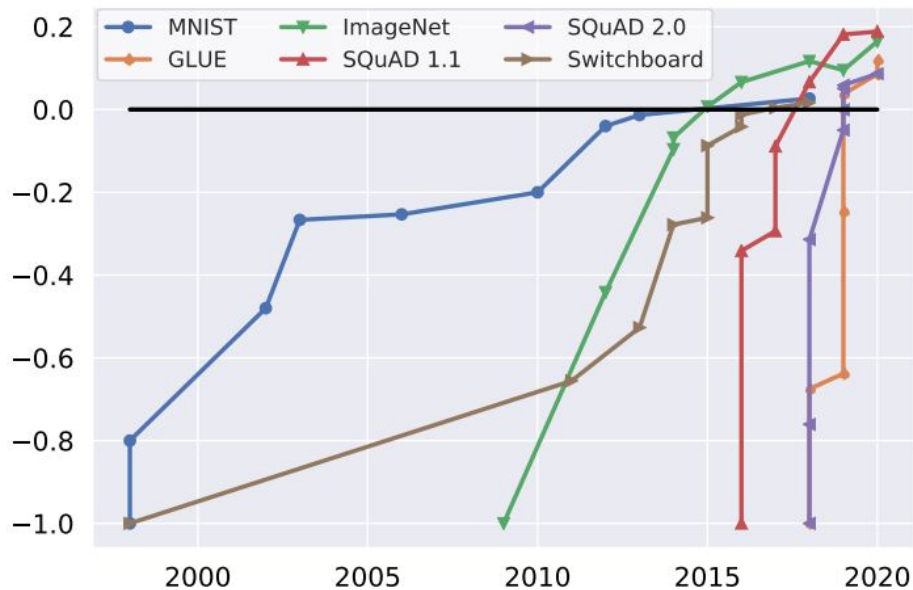
## #4: A static dataset: If models are overfitting to test set



## #4: A static dataset: Surprisingly no overfitting



## #4: A static dataset: creating dynamic benchmarks



Benchmark saturation over time for popular benchmarks, normalized with initial performance at minus one and human performance at zero.



## Find examples that fool the model

Your goal: enter a **negative** statement that fools the model into predicting positive.

Please pretend you are reviewing a place, product, book or movie.

This year's NAACL was very different because of Covid

Model prediction: **positive**

**Well done!** You fooled the model.

Optionally, provide an explanation for your example: [Draft](#). [Click out of input box to save.](#)

Covid is clearly not a good thing

The model probably doesn't know what Covid is

Model Inspector

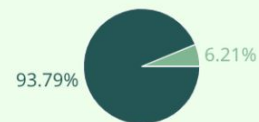
#s This year 's NA AC L was very different because of Cov id #/s

The model inspector shows the [layer integrated gradients](#) for the input token layer of the model.

Retract

Flag

Inspect



This year's NAACL was very different because of Covid

Live Mode

Switch to next context

Submit

#4: A static dataset: adversarial training only helps improve performance on adversarial test sets

Adversarially collected training data did not improve model performance

So far, dynamic adversarial testing hasn't resulted in new insights

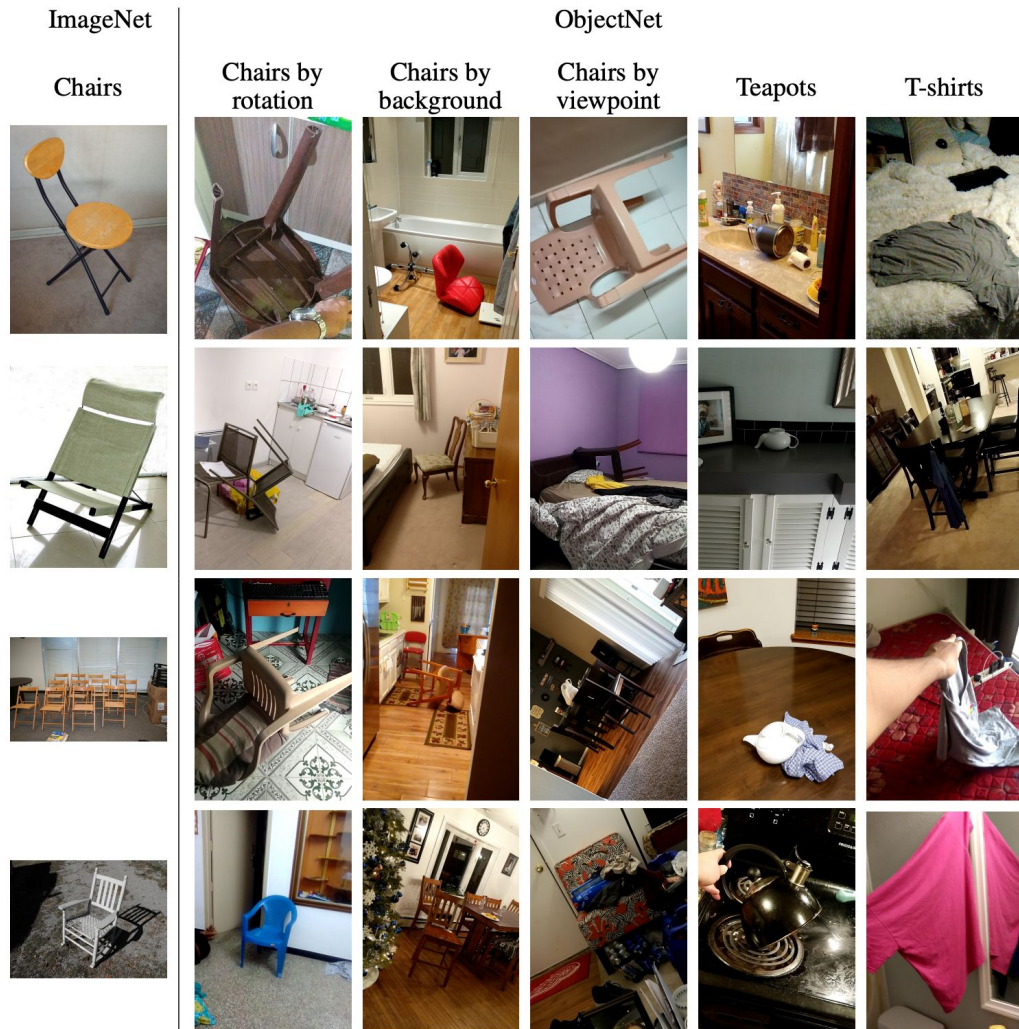
## #4: A static dataset: New guidelines for developing test sets

1. Good performance on the benchmark should imply robust in-domain performance on the task.  
↳ *We need more work on dataset design and data collection methods.*
2. Benchmark examples should be accurately and unambiguously annotated.  
↳ *Test examples should be validated thoroughly enough to remove erroneous examples and to properly handle ambiguous ones.*
3. Benchmarks should offer adequate statistical power.  
↳ *Benchmark datasets need to be much harder and/or much larger.*
4. Benchmarks should reveal plausibly harmful social biases in systems, and should not incentivize the creation of biased systems.  
↳ *We need to better encourage the development and use of auxiliary bias evaluation metrics.*

# #5: Distribution shifts

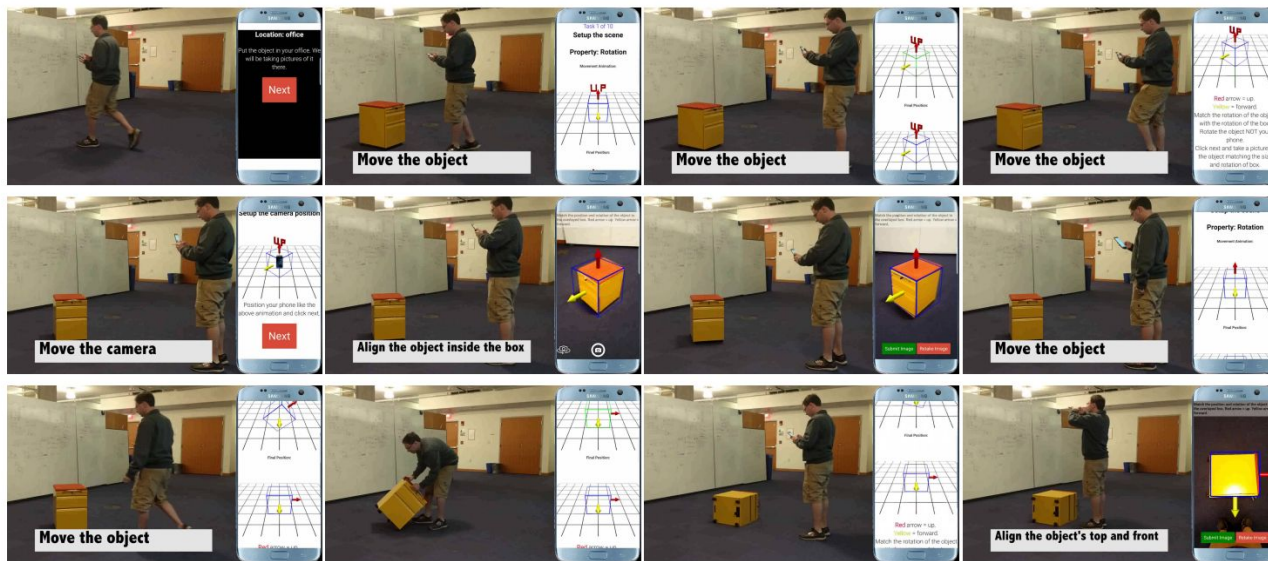
Differences between images in dataset versus images in the real world

Barbu et al. ObjectNet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. NeurIPS 2019



# #5: Distribution shifts:

Differences between images in dataset versus images in the real world

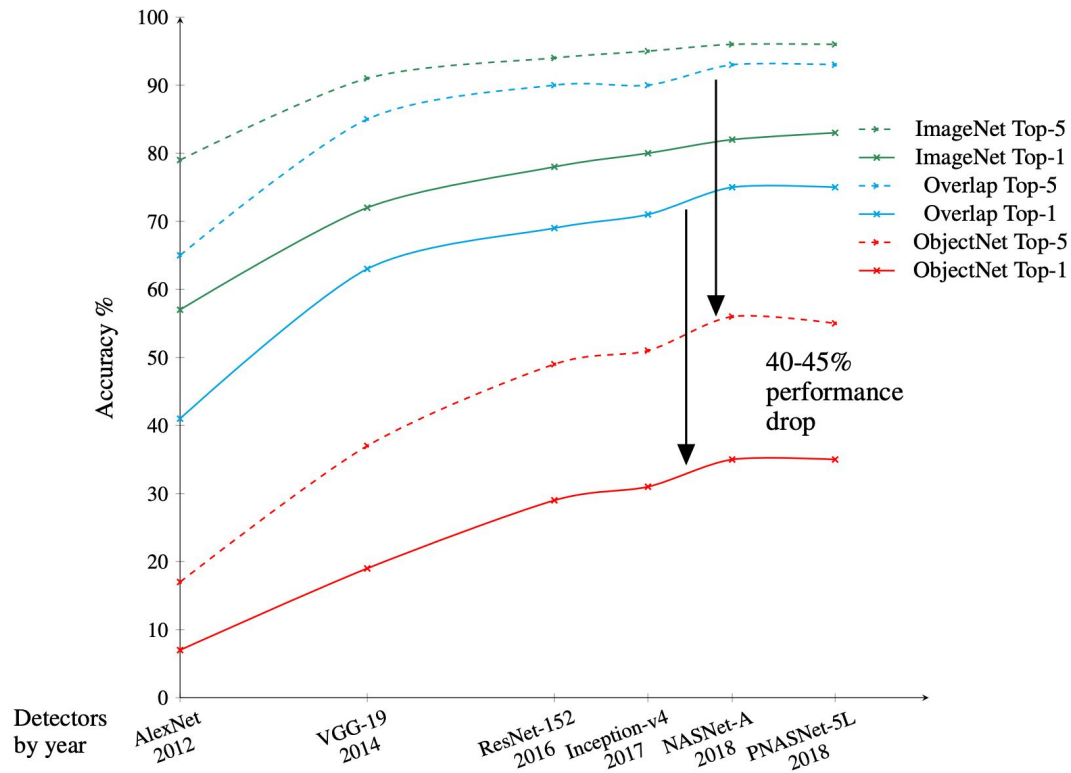


Barbu et al. ObjectNet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. NeurIPS 2019

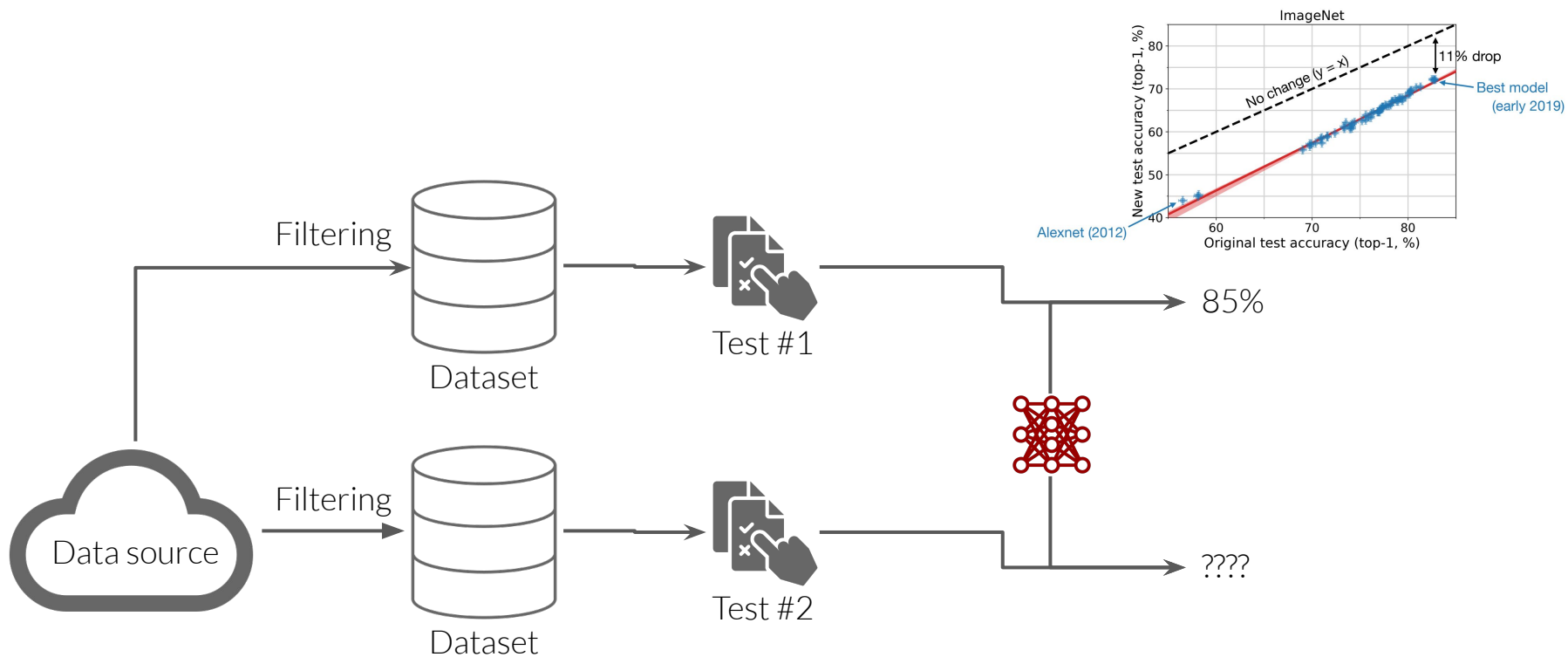
## #5: Distribution shifts

Differences between images in dataset versus images in the real world

Barbu et al. ObjectNet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. NeurIPS 2019



# #5: Distribution shifts: in data collection can explain this

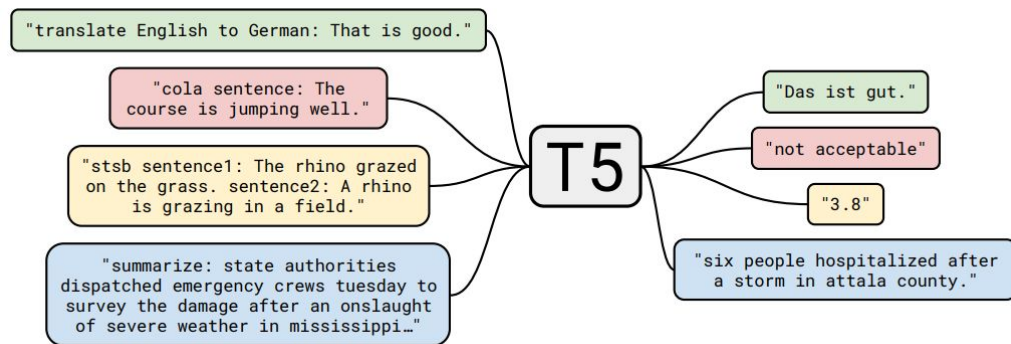
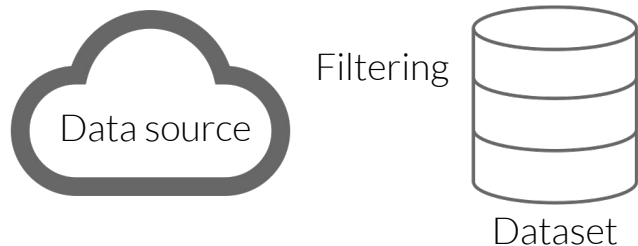


## #6: Marginalization: Filtering

T5 trained on Colossal Clean Crawled Corpus

400 words from the [List of filtered words](#)

- E.g. **swastika**, **white power** - implications?
- E.g. **twink** - implications?



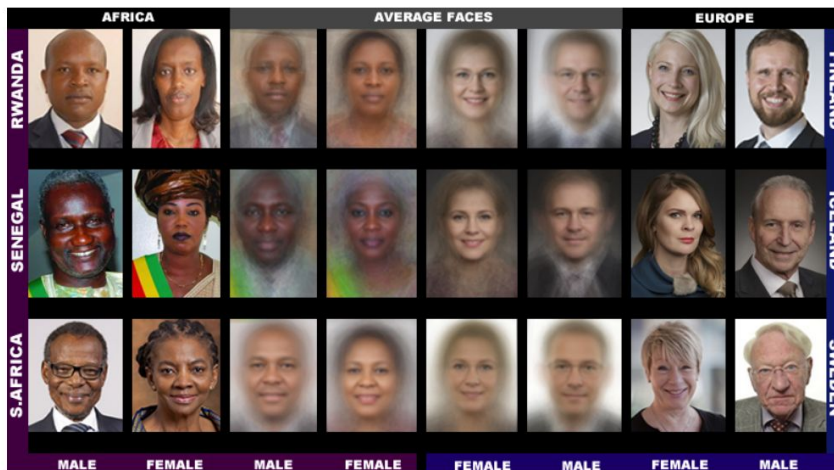
Raffel et al. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. IJML 2020

Dodge et al. Documenting Large Webtext Corpora: A Case Study on the Colossal Clean Crawled Corpus. ArXiv 2021



## #7: Bias in data source

- Then: What was not curated caused bias
- Today: More media coverage = more training data instances



Buolamwini et al. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. FAccT 2018

Bender et al. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? FAccT 2021

## #8: Environmental and financial costs

Energy for a flight from NY to SF:

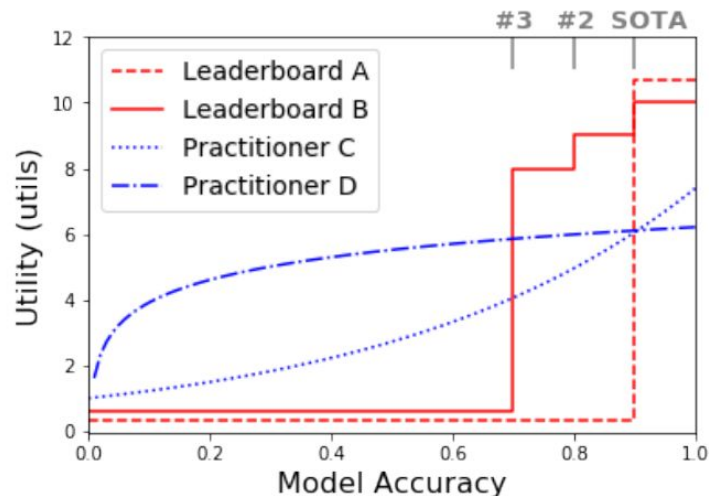


Model	Hardware	Power (W)	Hours	kWh·PUE	CO <sub>2</sub> e	Cloud compute cost
Transformer <sub>base</sub>	P100x8	1415.78	12	27	26	\$41–\$140
Transformer <sub>big</sub>	P100x8	1515.43	84	201	192	\$289–\$981
ELMo	P100x3	517.66	336	275	262	\$433–\$1472
BERT <sub>base</sub>	V100x64	12,041.51	79	1507	1438	\$3751–\$12,571
BERT <sub>base</sub>	TPUv2x16	—	96	—	—	\$2074–\$6912
NAS	P100x8	1515.43	274,120	656,347	626,155	\$942,973–\$3,201,722
NAS	TPUv2x1	—	32,623	—	—	\$44,055–\$146,848
GPT-2	TPUv3x32	—	168	—	—	\$12,902–\$43,008

## #9: Leaderboard with one metric is not enough

Utility of a new AI model:

- is NOT smooth w.r.t. Accuracy for a **leaderboard**
- Any improvement along any dimension is good for a **practitioner**



# #10: Open ended tasks: Generative models are very hard to evaluate

Research question:

How do you evaluate the output of an image generation model?





















# It used to be easy to measure progress



2014



2015



2016



2017



2018

# It's much harder now



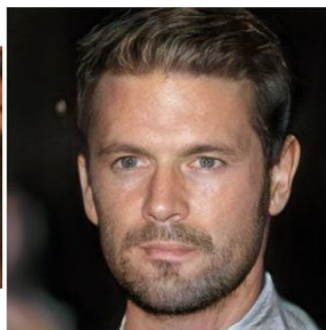
2014



2015



2016



2017



2018



Ian Goodfellow @goodfellow\_ian

# We don't even have corresponding pairs



2014



2015



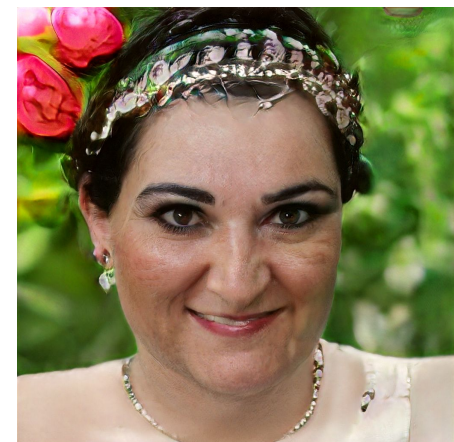
2016



2017



2018

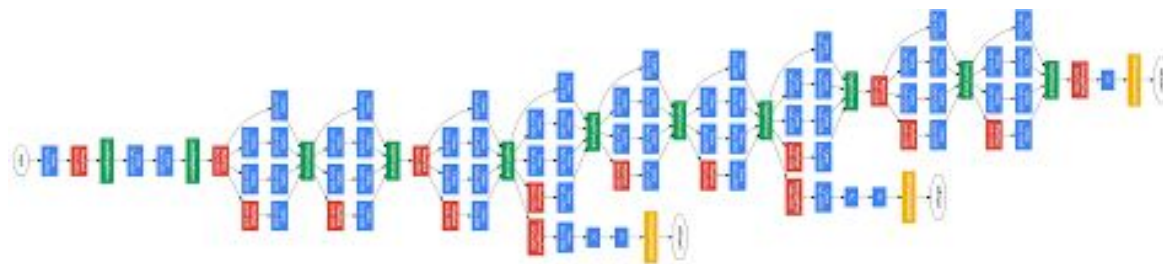


Ian Goodfellow @goodfellow\_ian



# How are models evaluated today?

Inception score, FID.



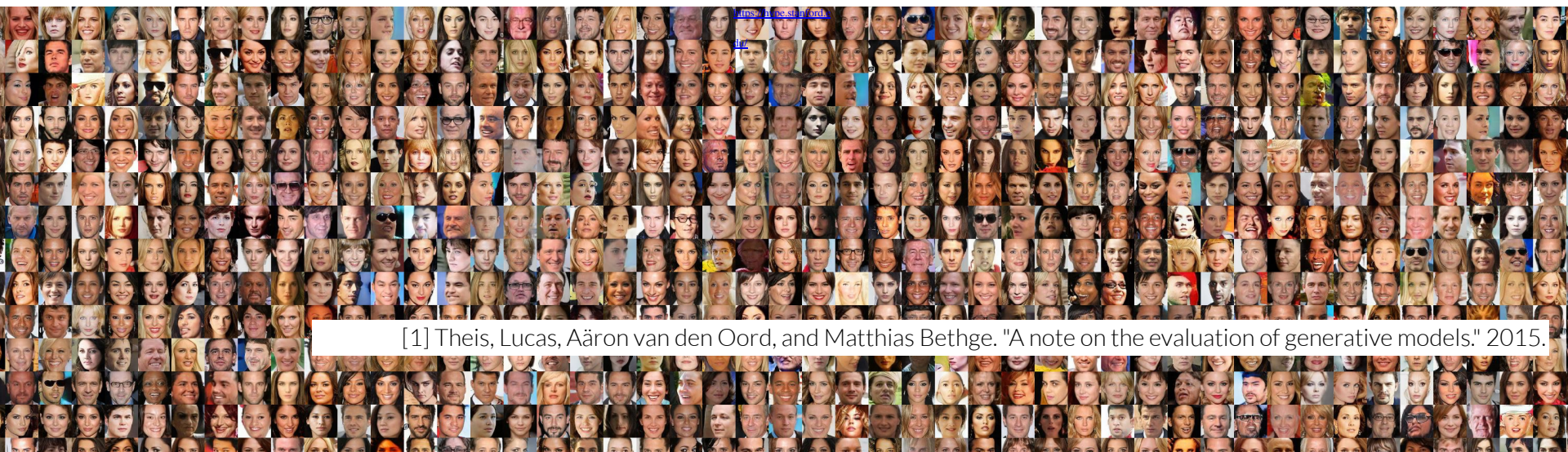
- Trained on imagenet
- **Inception score** is maximized when entropy of predicted output is low
  - Meaning if Inception says with high certainty that it's a "person", the score will be higher
- **FID** calculates distributions from activations of an Inception-v3 layer
- **What is the problem with this approach?**

# Why not use automated metrics?



# Why not use automated metrics?

Density estimation has even been shown to be misleading [1].



[1] Theis, Lucas, Aäron van den Oord, and Matthias Bethge. "A note on the evaluation of generative models." 2015.

# Why not use automated metrics?

Density estimation has even been shown to be misleading [1].

Automated evaluation metrics on sampled outputs (Inception Score [2], FID [3], Precision [4], etc.) rely on ImageNet embeddings.



[1] Theis, Lucas, Aäron van den Oord, and Matthias Bethge. "A note on the evaluation of generative models." 2015.

[2] Salimans, Tim, et al. "Improved techniques for training GANs." 2016.

[3] Heusel, Martin, et al. "GANs trained by a two time-scale update rule converge to a local nash equilibrium." 2017.

[4] Sajjadi, Mehdi SM, et al. "Assessing generative models via precision and recall." 2018.

# Why not use automated metrics? Or human metrics?

Density estimation has even been shown to be misleading [1].

Automated evaluation metrics on sampled outputs (Inception Score [2], FID [3], Precision [4], etc.) rely on ImageNet embeddings.



[1] Theis, Lucas, Aäron van den Oord, and Matthias Bethge. "A note on the evaluation of generative models." 2015.

[2] Salimans, Tim, et al. "Improved techniques for training GANs." 2016.

[3] Heusel, Martin, et al. "GANs trained by a two time-scale update rule converge to a local nash equilibrium." 2017.

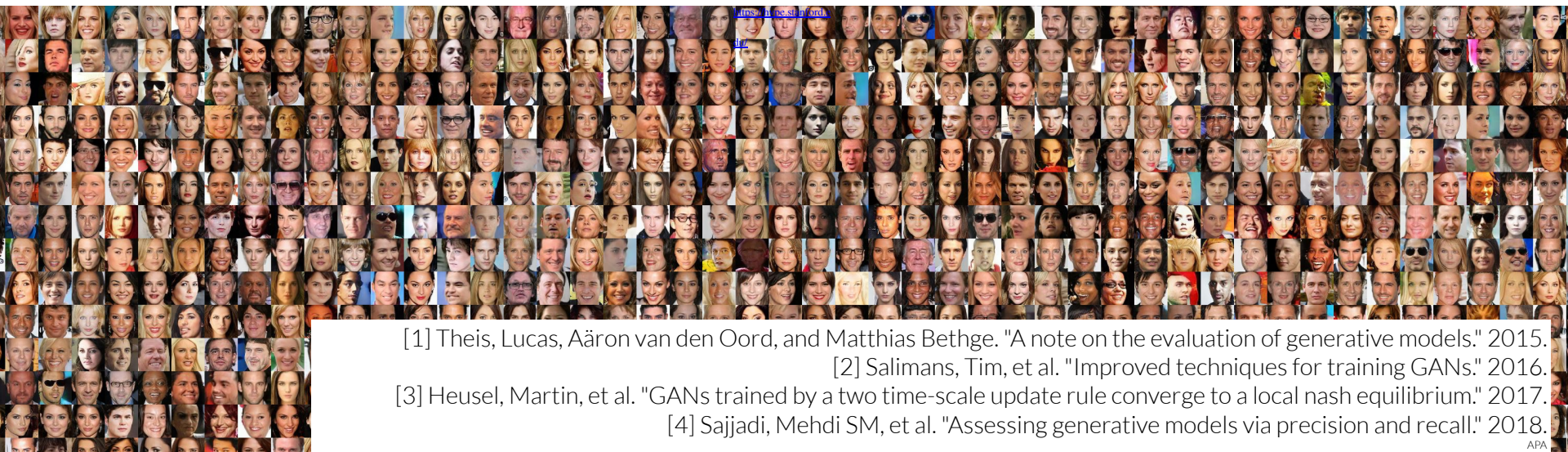
[4] Sajjadi, Mehdi SM, et al. "Assessing generative models via precision and recall." 2018.

# Why not use automated metrics? Or human metrics?

Density estimation has even been shown to be misleading [1].

Automated evaluation metrics on sampled outputs (Inception Score [2], FID [3], Precision [4], etc.) rely on ImageNet embeddings.

Human evaluation metrics are ad-hoc — unreliable and costly.



[1] Theis, Lucas, Aäron van den Oord, and Matthias Bethge. "A note on the evaluation of generative models." 2015.

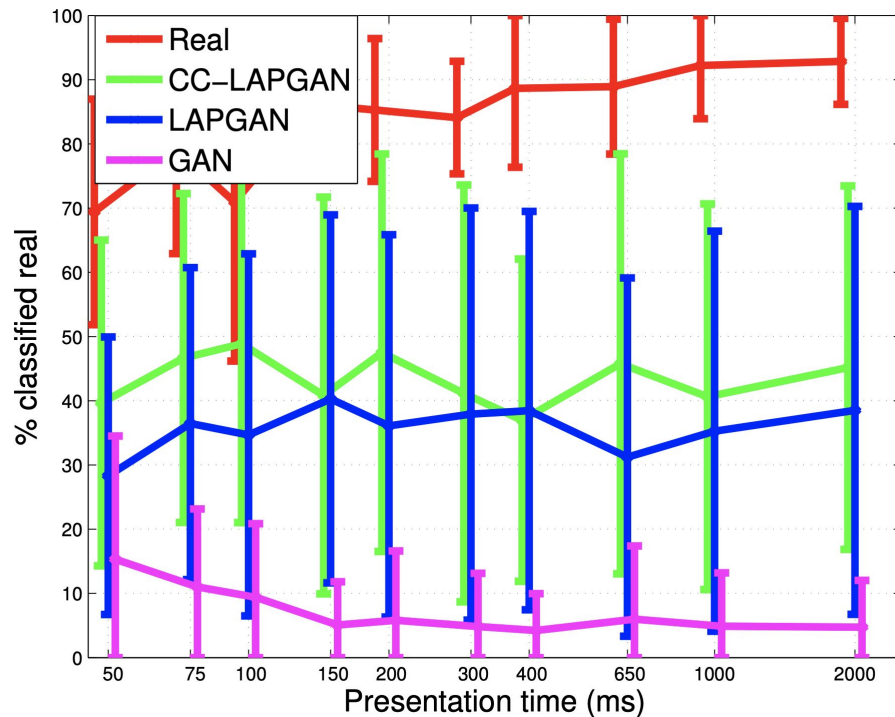
[2] Salimans, Tim, et al. "Improved techniques for training GANs." 2016.

[3] Heusel, Martin, et al. "GANs trained by a two time-scale update rule converge to a local nash equilibrium." 2017.

[4] Sajjadi, Mehdi SM, et al. "Assessing generative models via precision and recall." 2018.

# Why not use human evaluation?

1. *Ad-hoc*, each executed in idiosyncrasy without proof of reliability or grounding to theory.
2. High *variance* in their estimates.
3. Lack clear *separability* between models.
4. Expensive and *time-consuming*

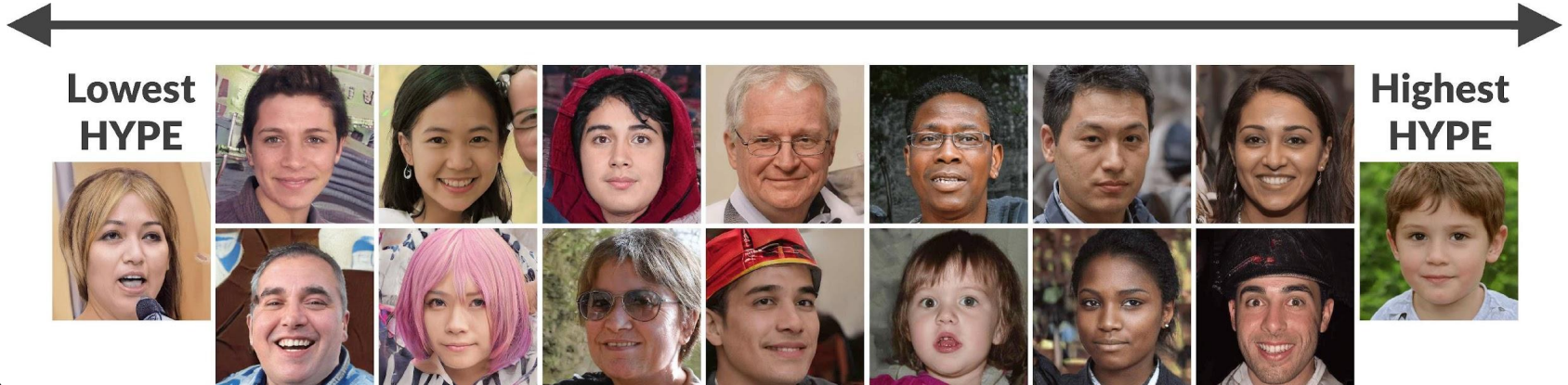


HYPE measures this progress using human evaluation that is consistent, efficient, and grounded in theory



# HYPE is designed to address these problems:

1. **Grounded** method inspired by psychophysics methods in perceptual psychology.
2. **Reliable** and consistent estimator.
3. Statistically **separable** to enable a comparative ranking.
4. Cost and time **efficient**.



# Psychophysics method: adaptive staircase procedure

- Staircase methods can determine human perceptual thresholds efficiently and reliably (Cornsweet, 1962).

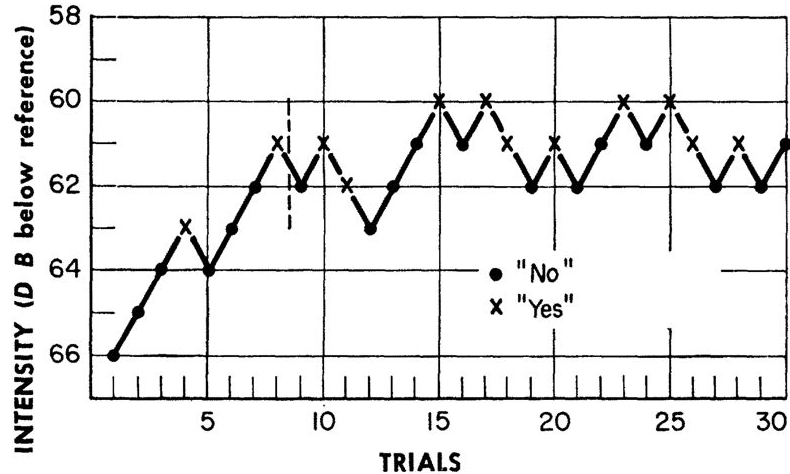
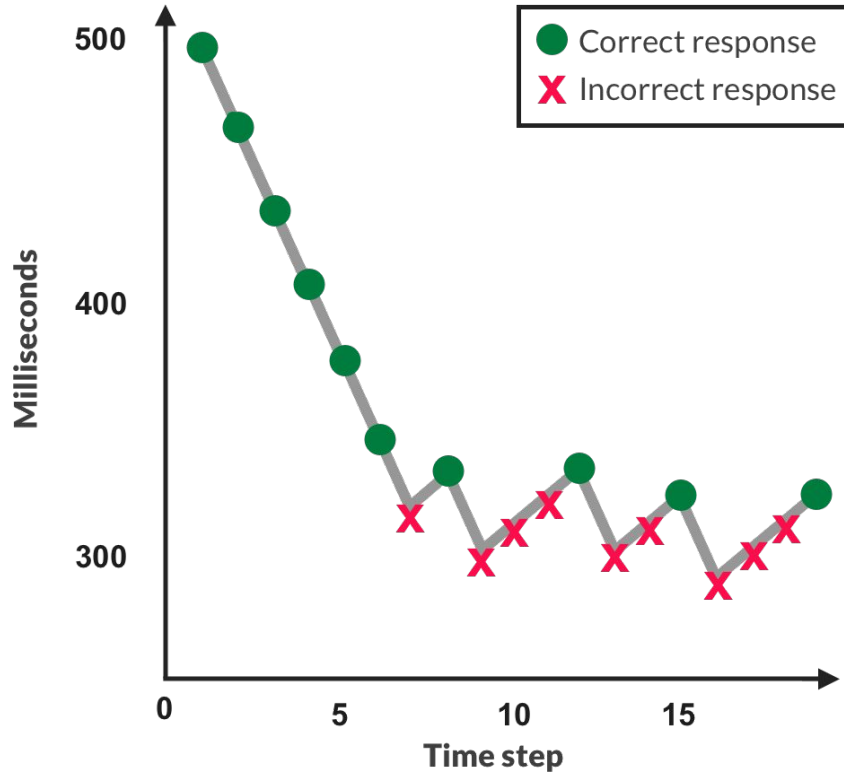


FIG. 1. DATA FROM THE DETERMINATION OF A TYPICAL AUDITORY THRESHOLD BY THE STAIRCASE-METHOD

# HYPE: adaptive staircase procedure

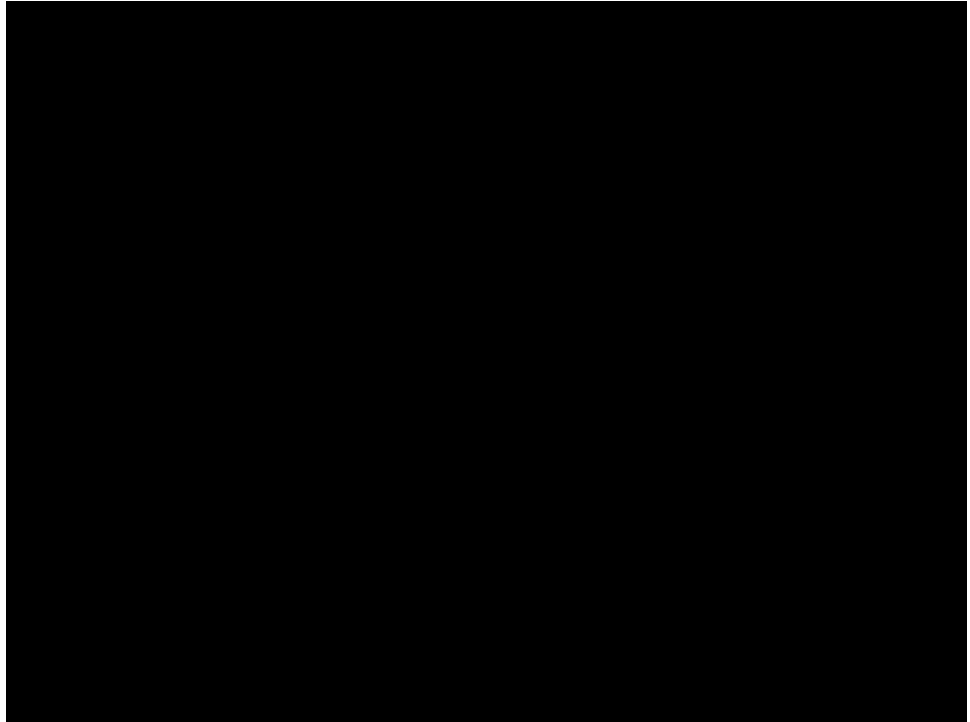


Time: 375ms

real

or

fake

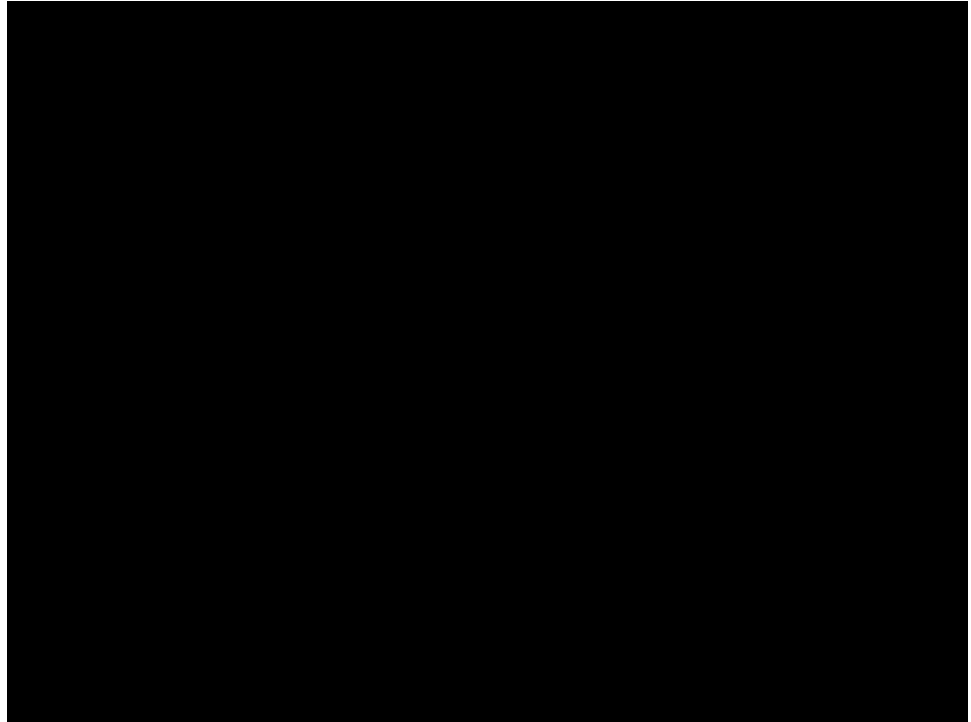


Time: 500ms

real

or

fake

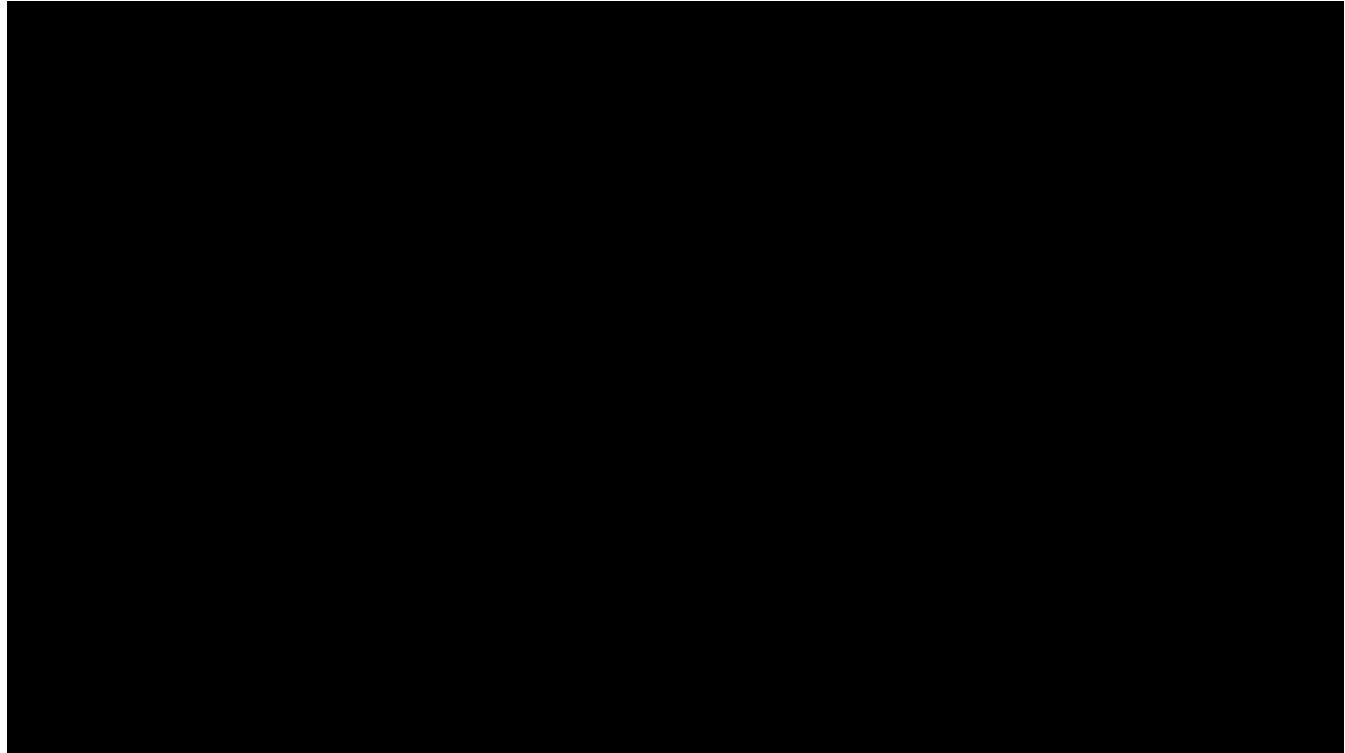


Time: 250ms

real

or

fake

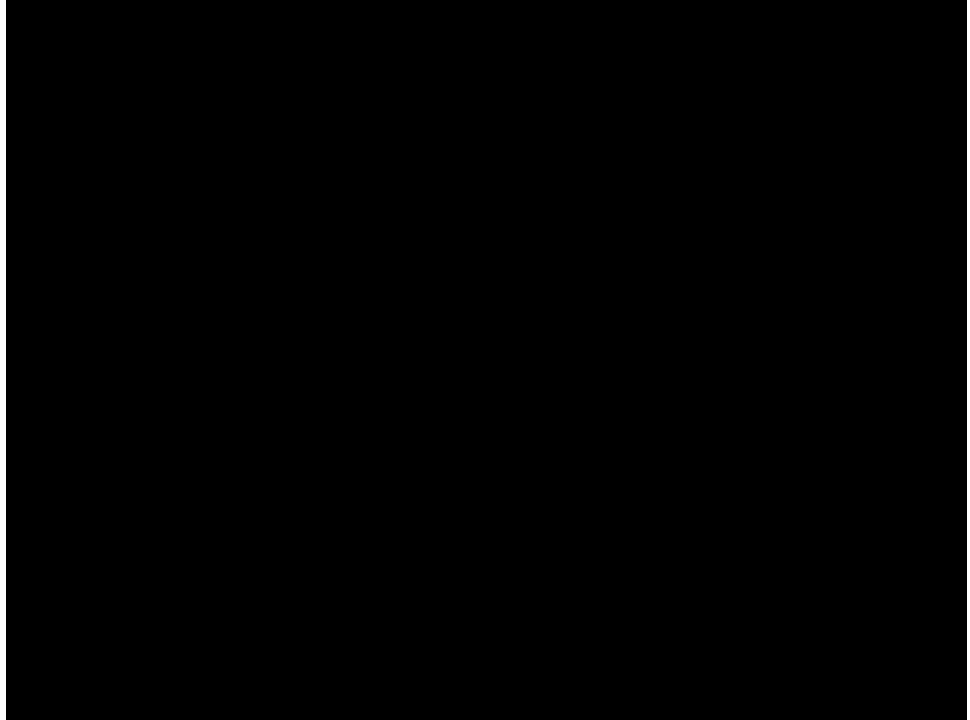


Time: 125ms

real

or

fake



# Creating a reliable score

To ensure reliability, we need to:

1. Hire and train/filter a sufficient number of evaluators.
2. Sample sufficient outputs.
3. Aggregate.



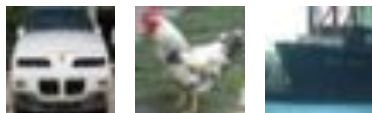
# Experiments

# Datasets

.CelebA



.CIFAR-10



.ImageNet-5

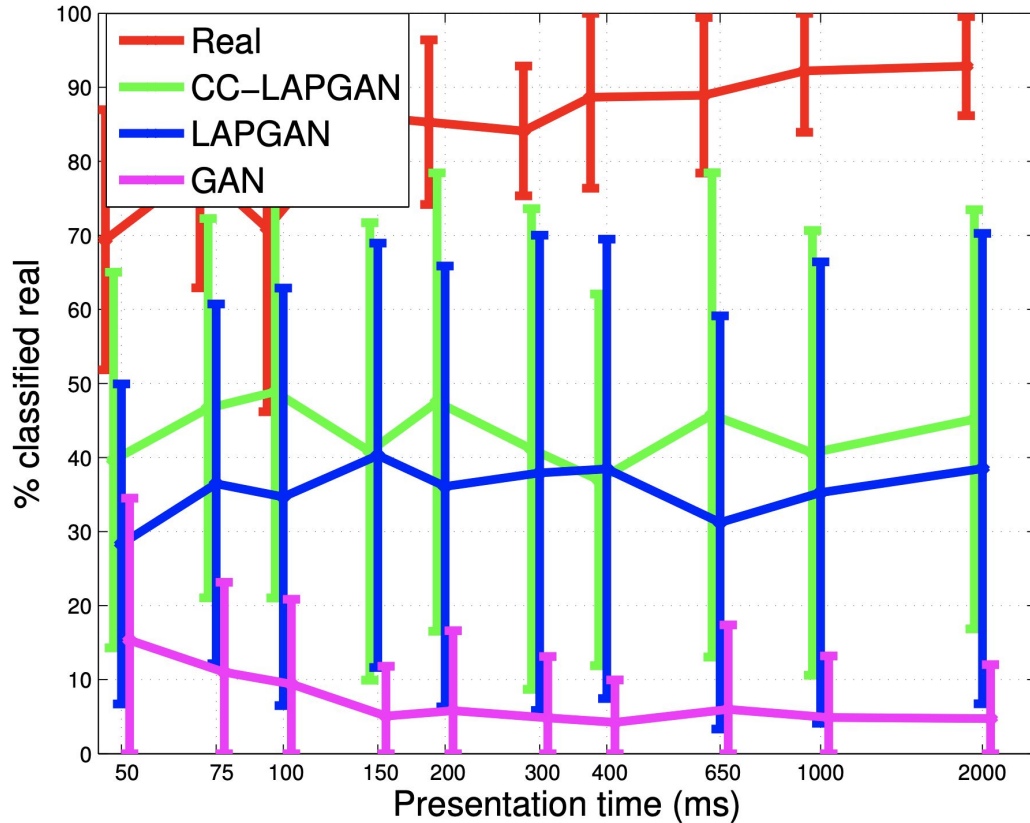


.FFHQ

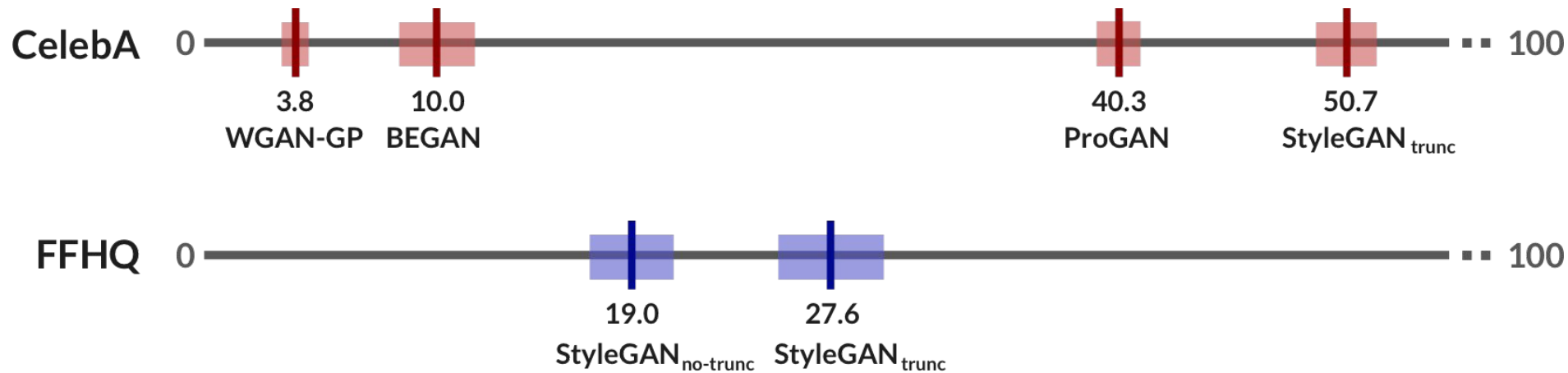


# Results

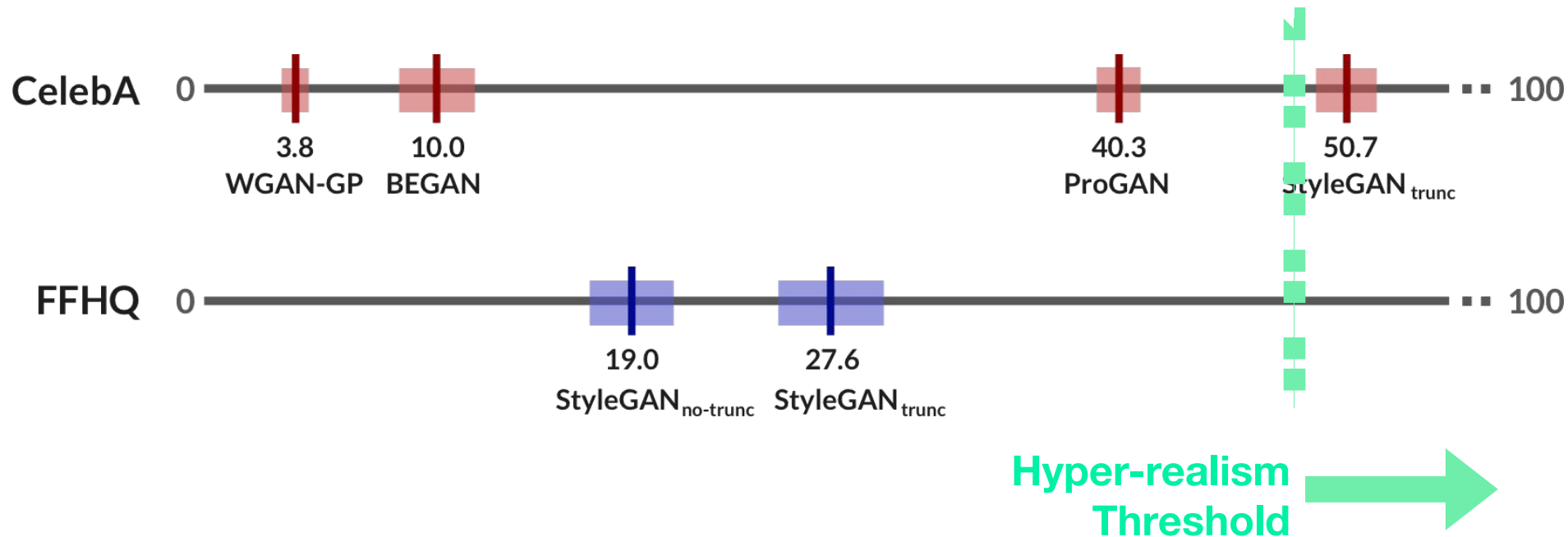
# Are HYPE's results statistically separable?



# Are HYPE's results statistically separable?

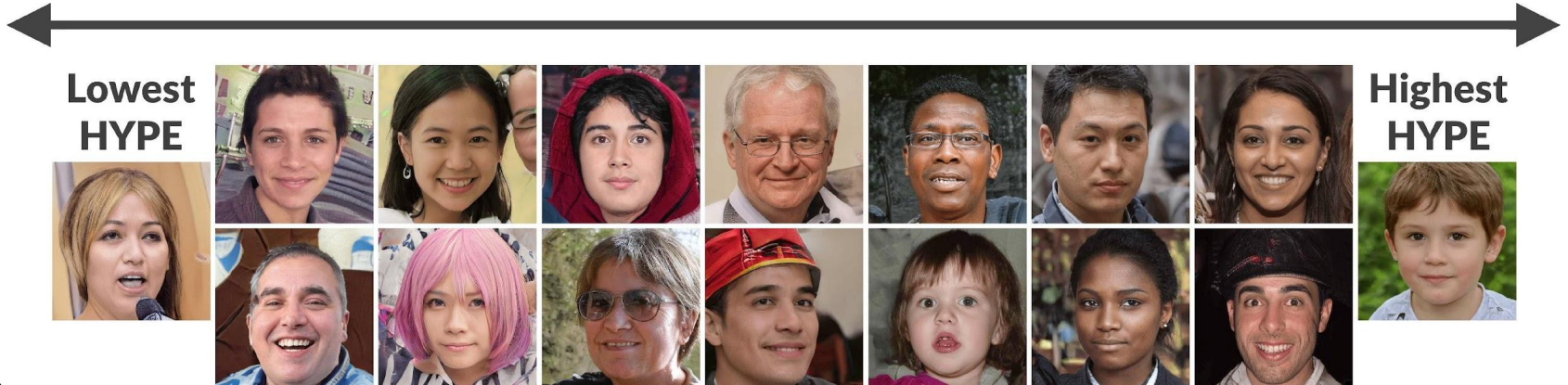


# Are HYPE's results statistically separable?



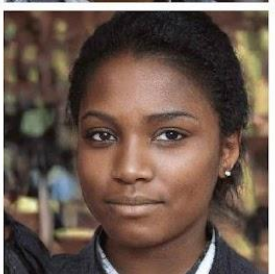
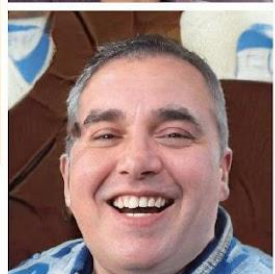
# HYPE achieves:

1. **Grounded** method inspired by psychophysics methods in perceptual psychology.
2. **Reliable** and consistent estimator.
3. Statistically **separable** to enable a comparative ranking.
4. Cost and time **efficient**.









Next time:  
evaluations with real users from an  
AI+HCI perspective