

Sprucing up a Dataset: Adversarially Filtering Dataset Artifacts

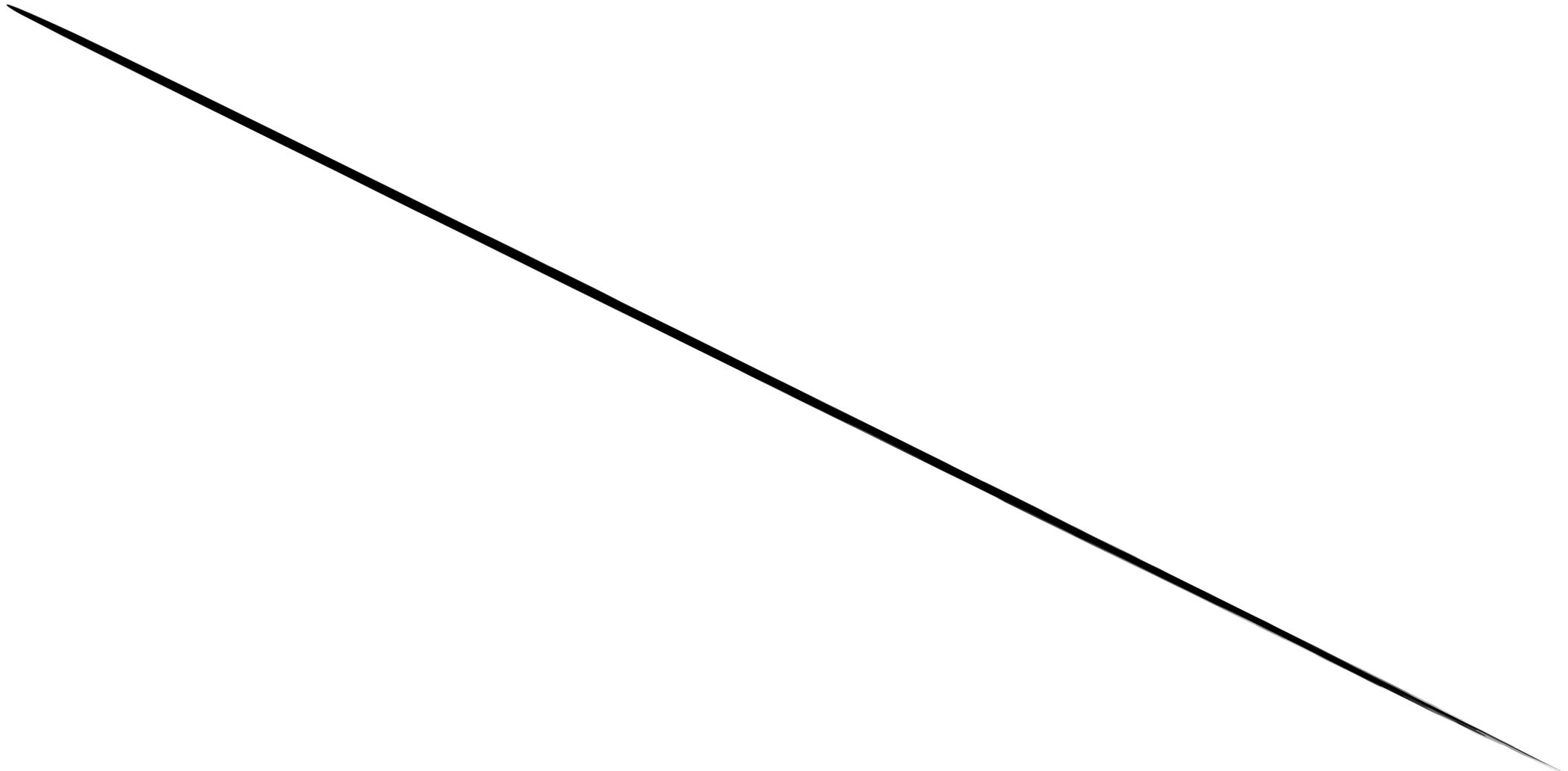


Swabha Swayamdipta
Oct 18th, 2019



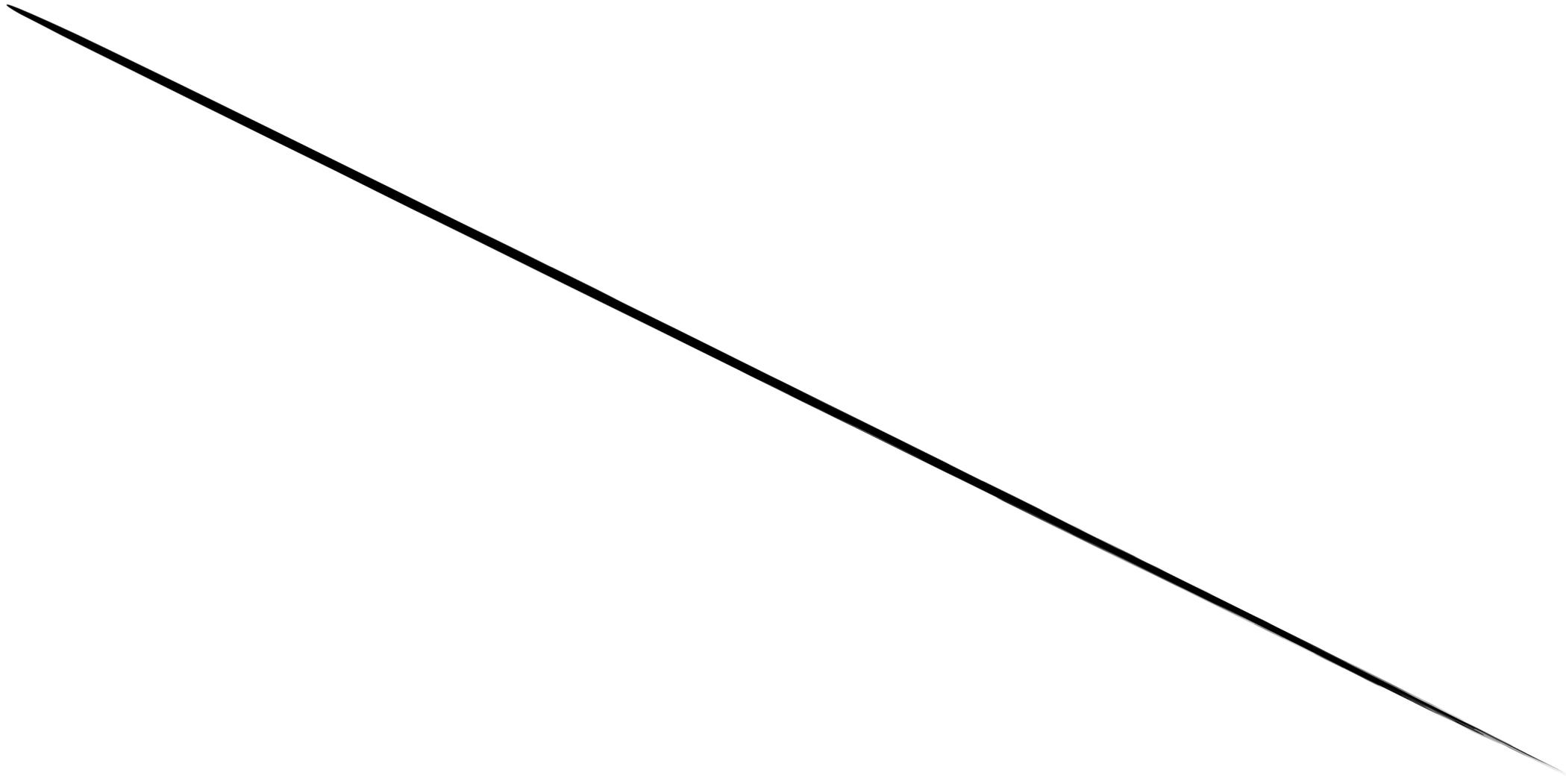
**UW Linguistics
Colloquium**

Linguistics



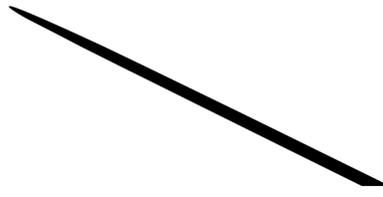
Linguistics

**Machine
Learning**

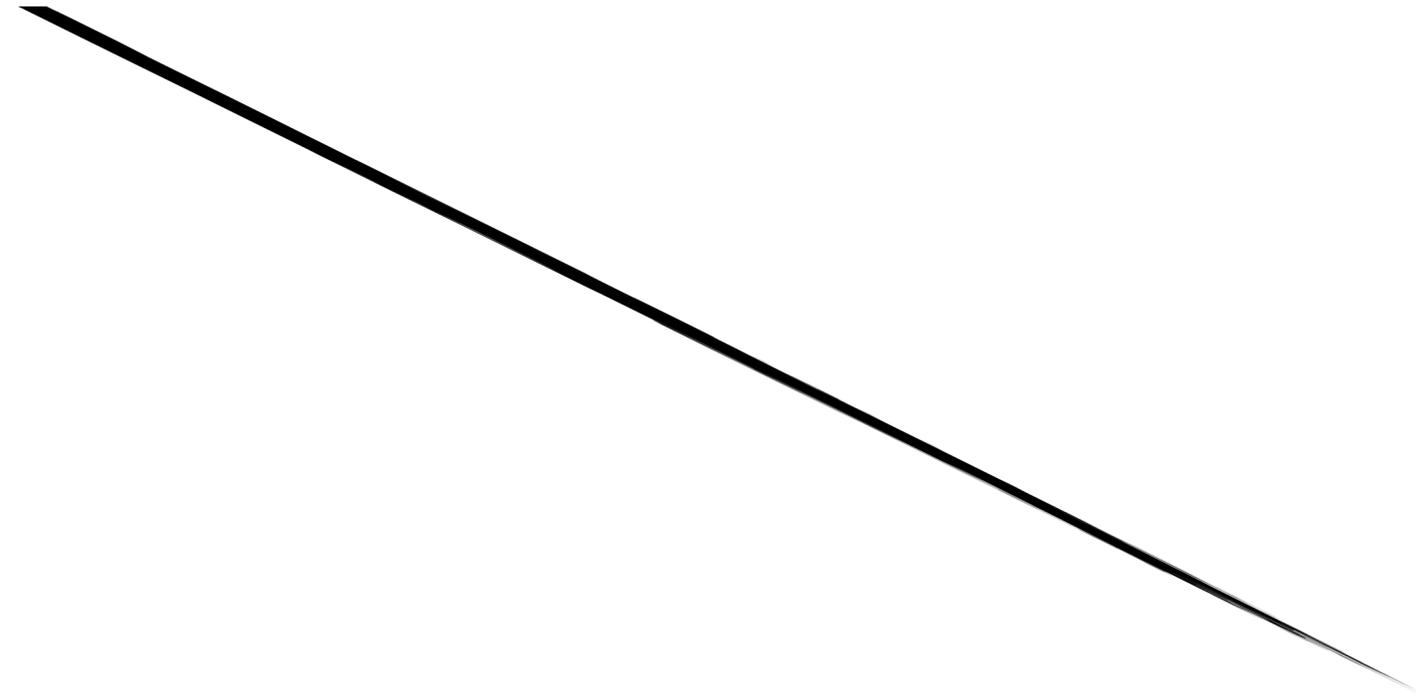


Linguistics

**Machine
Learning**



NLP



Linguistics

Machine Learning



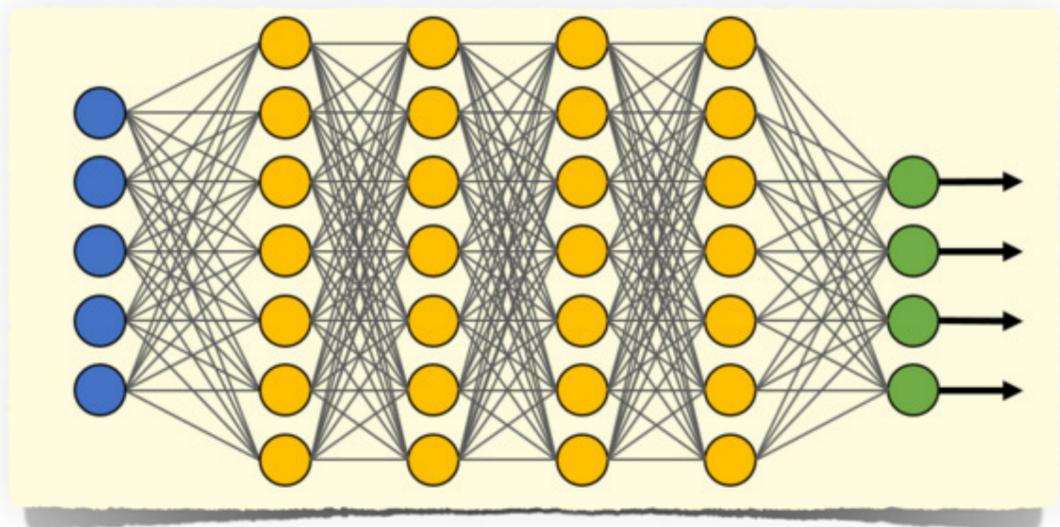
NLP

Linguistics

Machine Learning



NLP

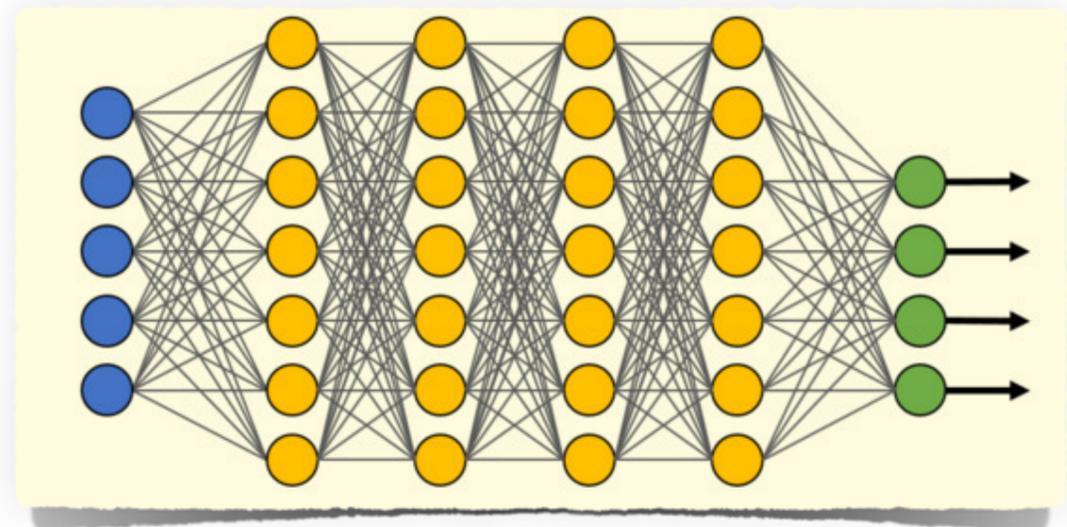
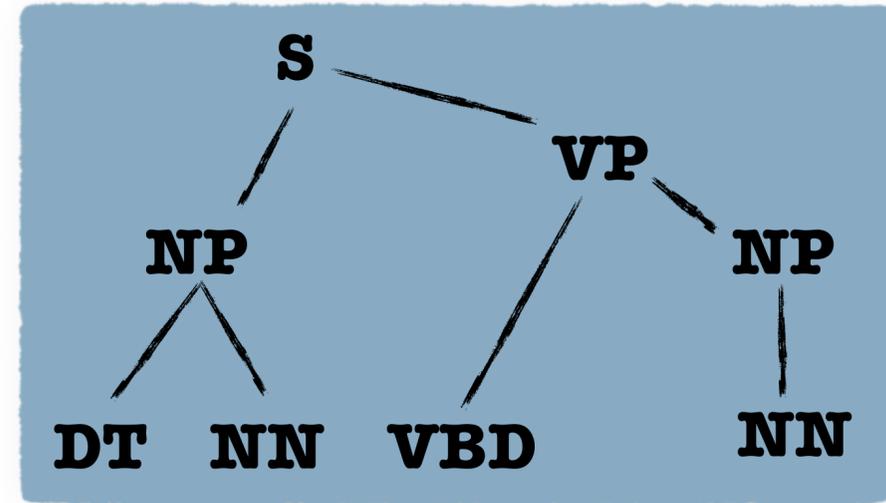


Linguistics

Machine Learning

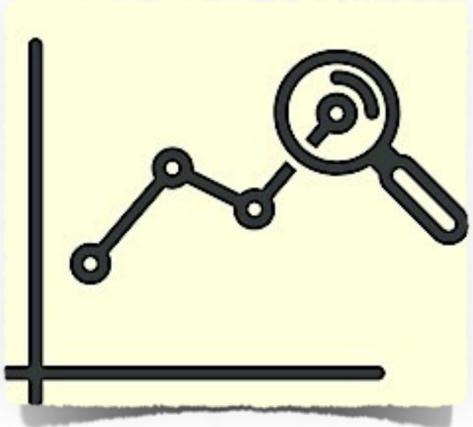


NLP

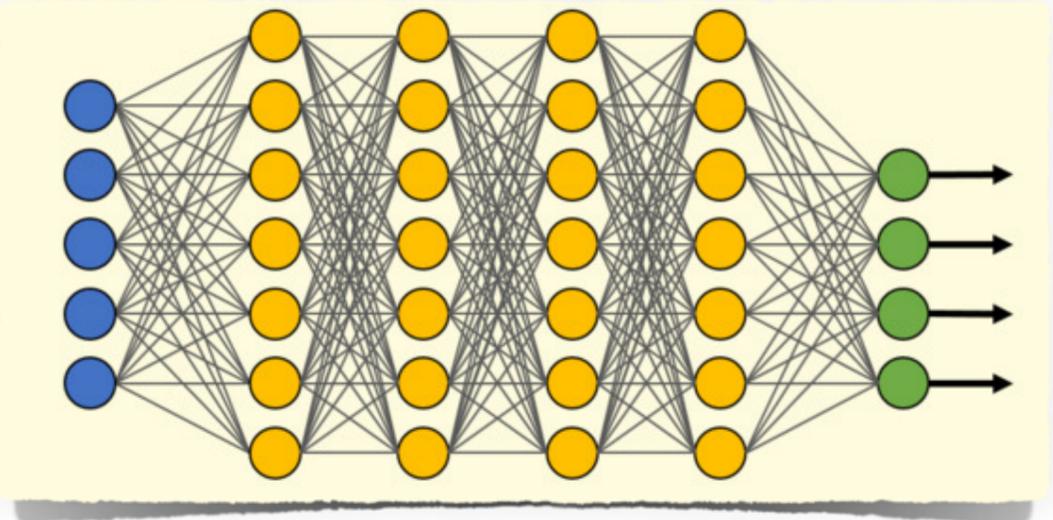
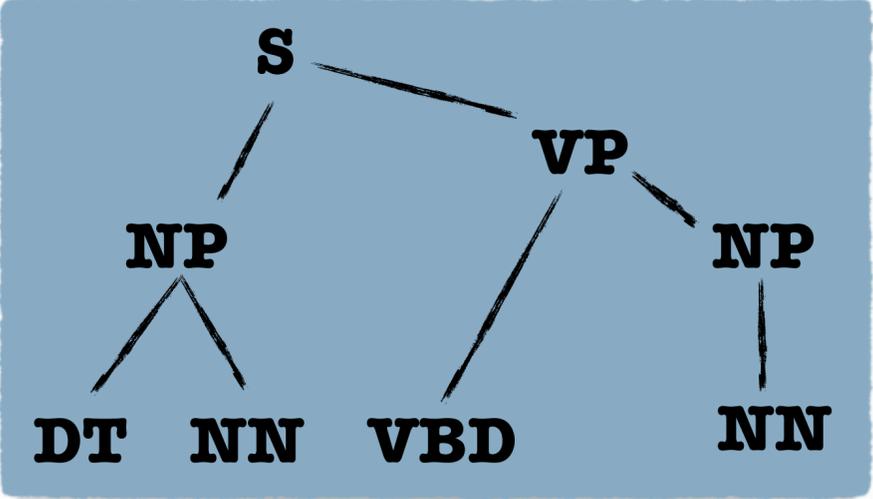


Linguistics

Machine Learning

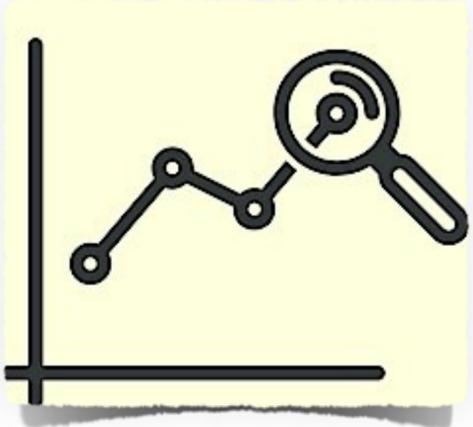


NLP

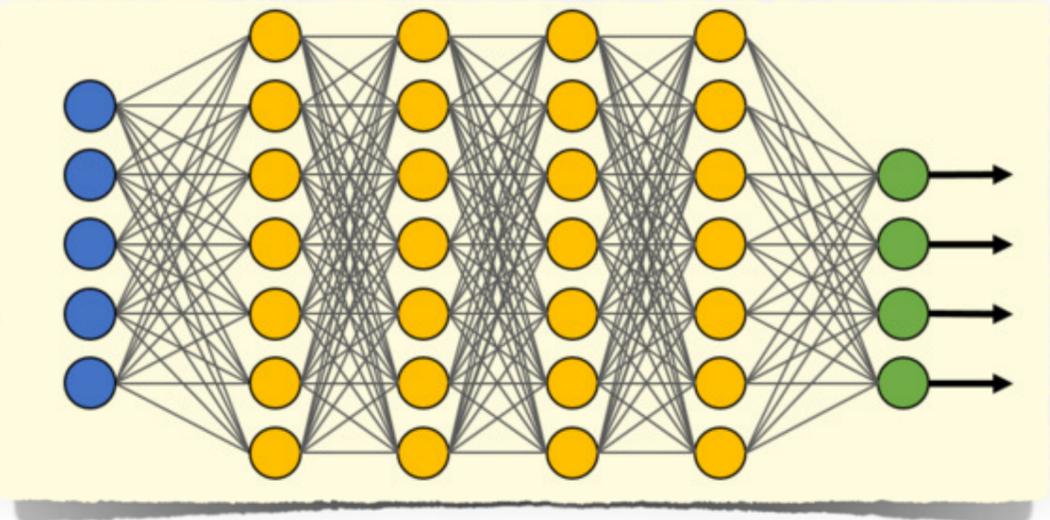
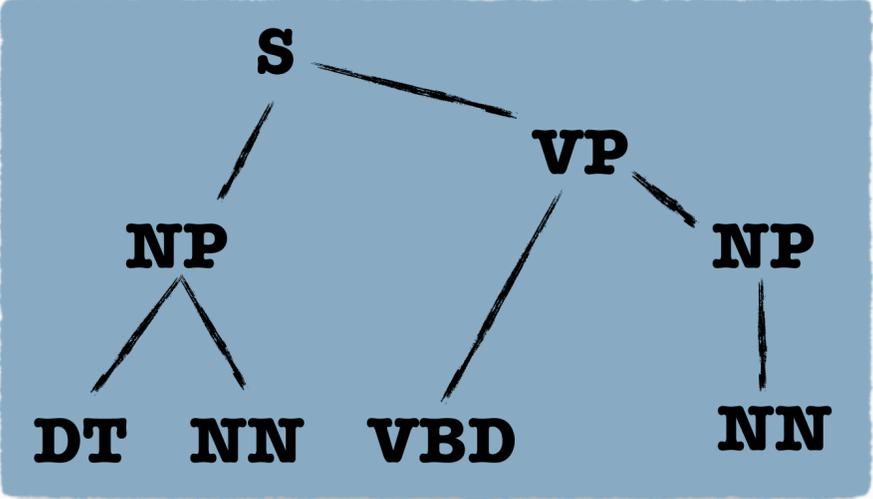


Linguistics

Machine Learning



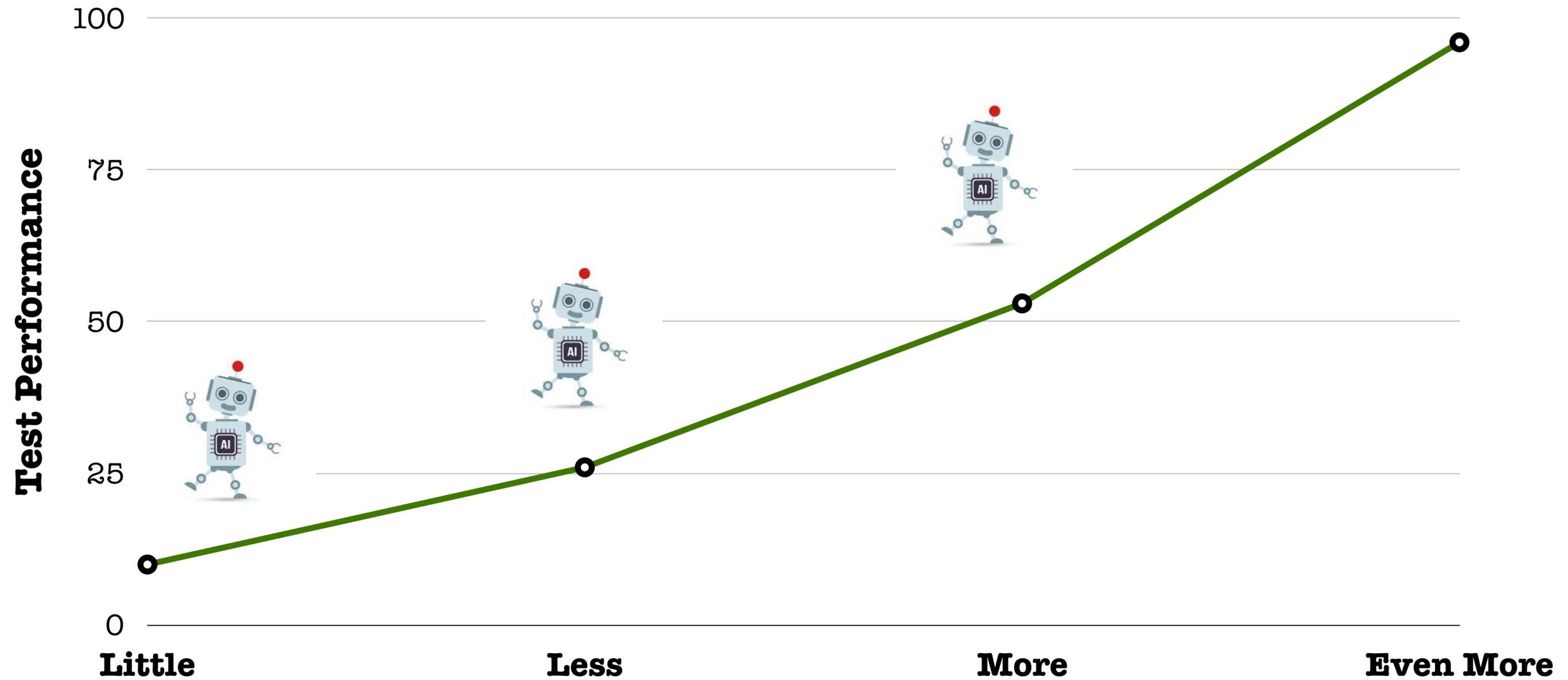
NLP



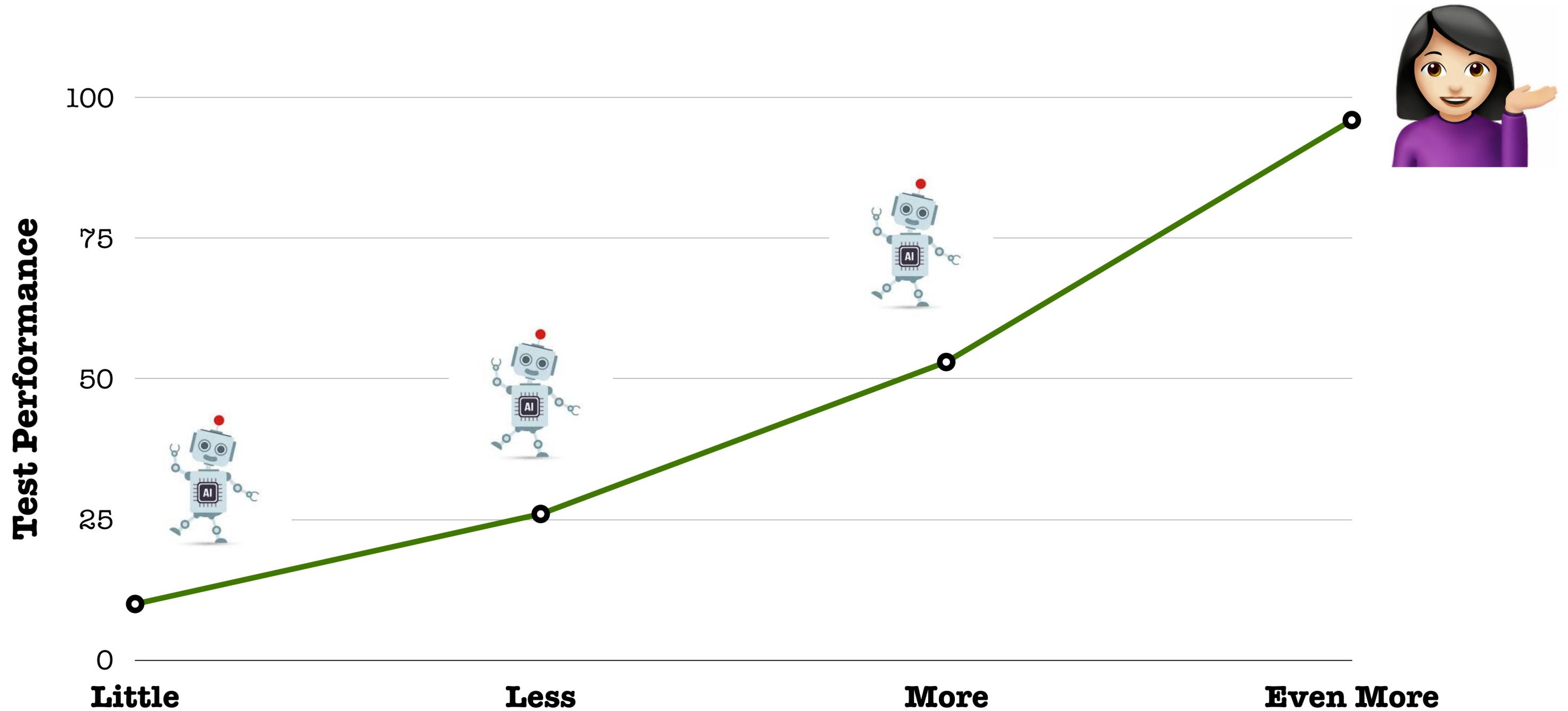
Linguistics

More data, better prediction?

More data, better prediction?



More data, better prediction?



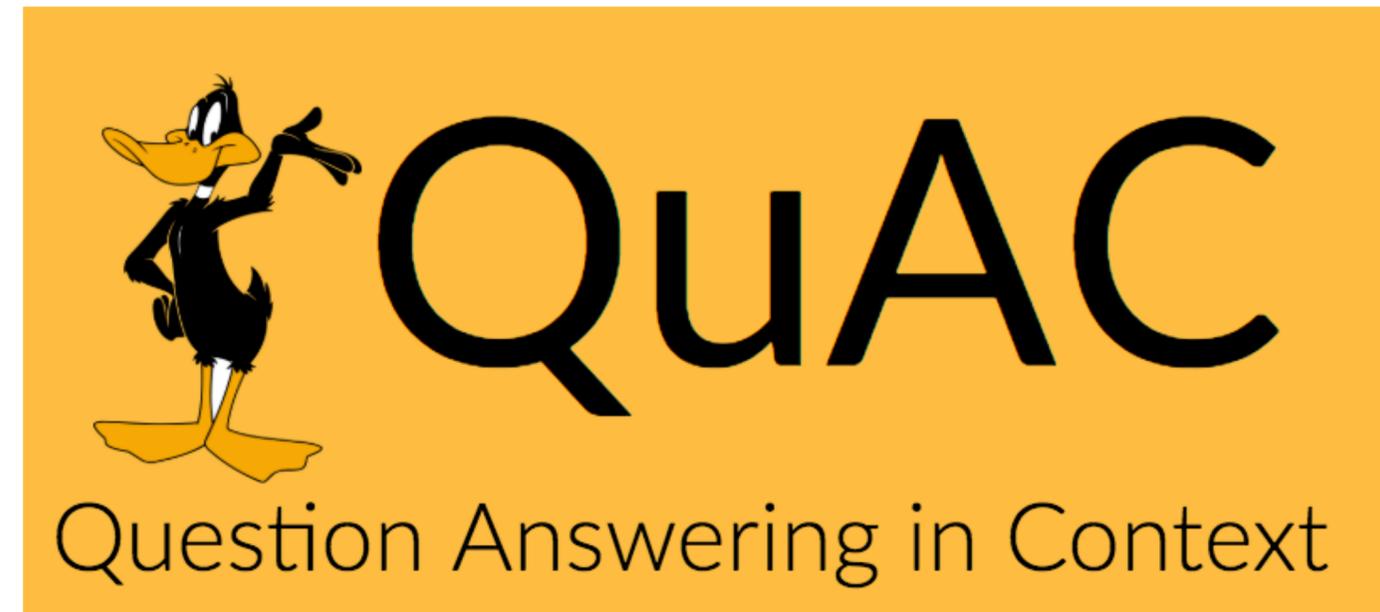
More data, better prediction?



Datasets abound!

Datasets abound!

SWAG: A Large-Scale Adversarial Dataset



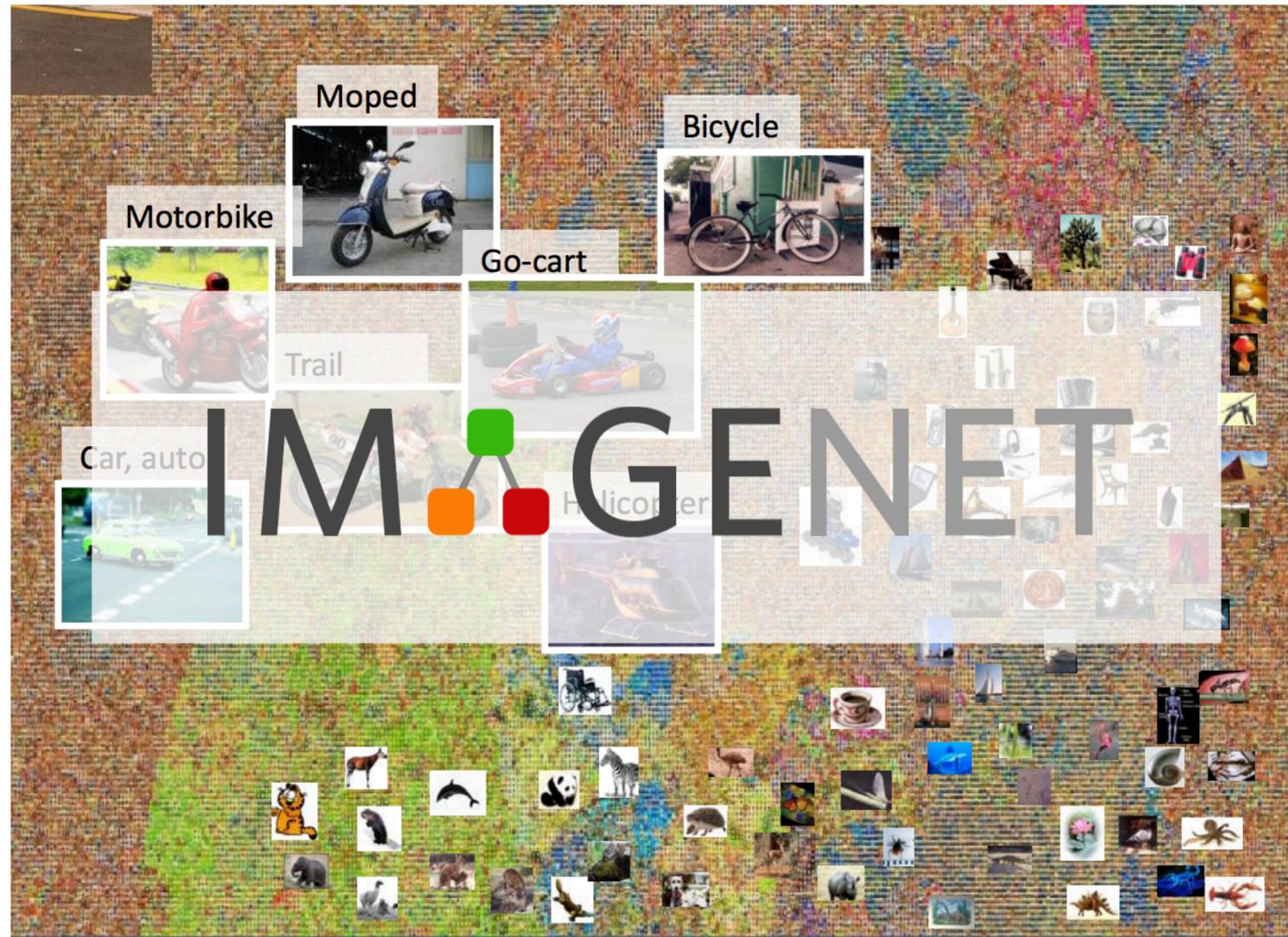
Sentiment Analysis

The Stanford Natural Language Inference (SNLI) Corpus

MultiNLI

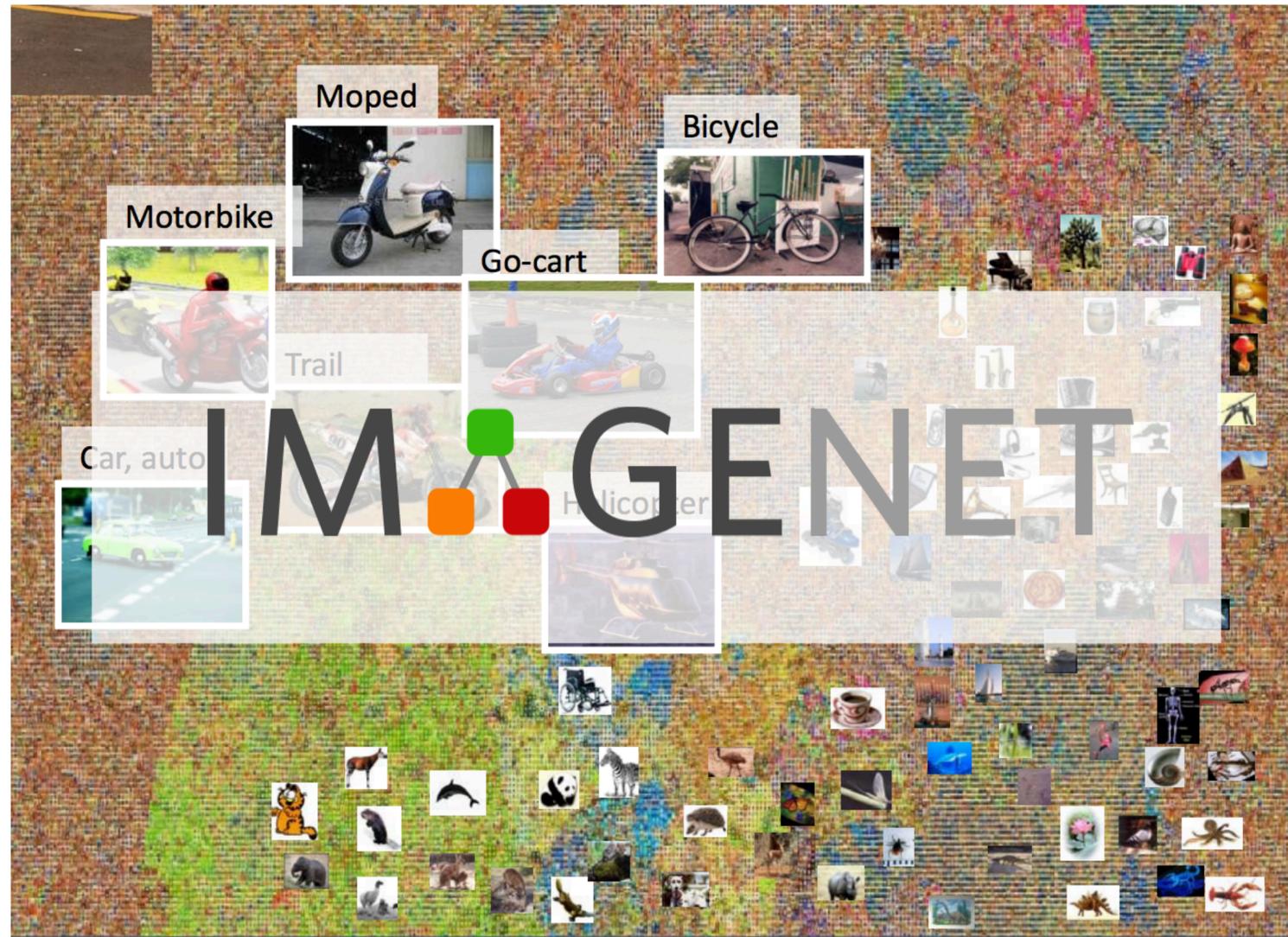
Story Cloze Test and ROCStories Corpora

Classifying images

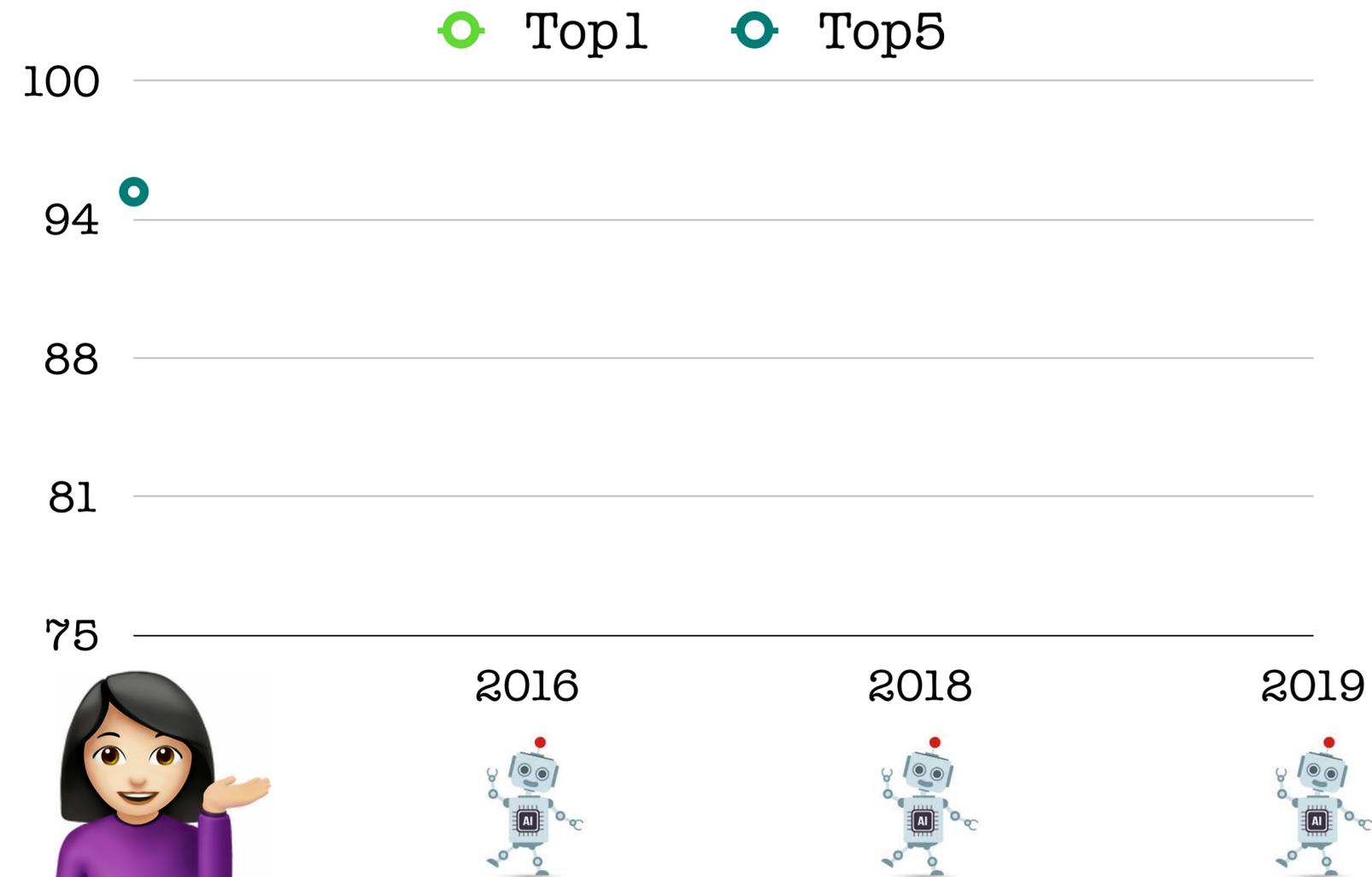


IMAGENET [Deng et al., CVPR 2009]

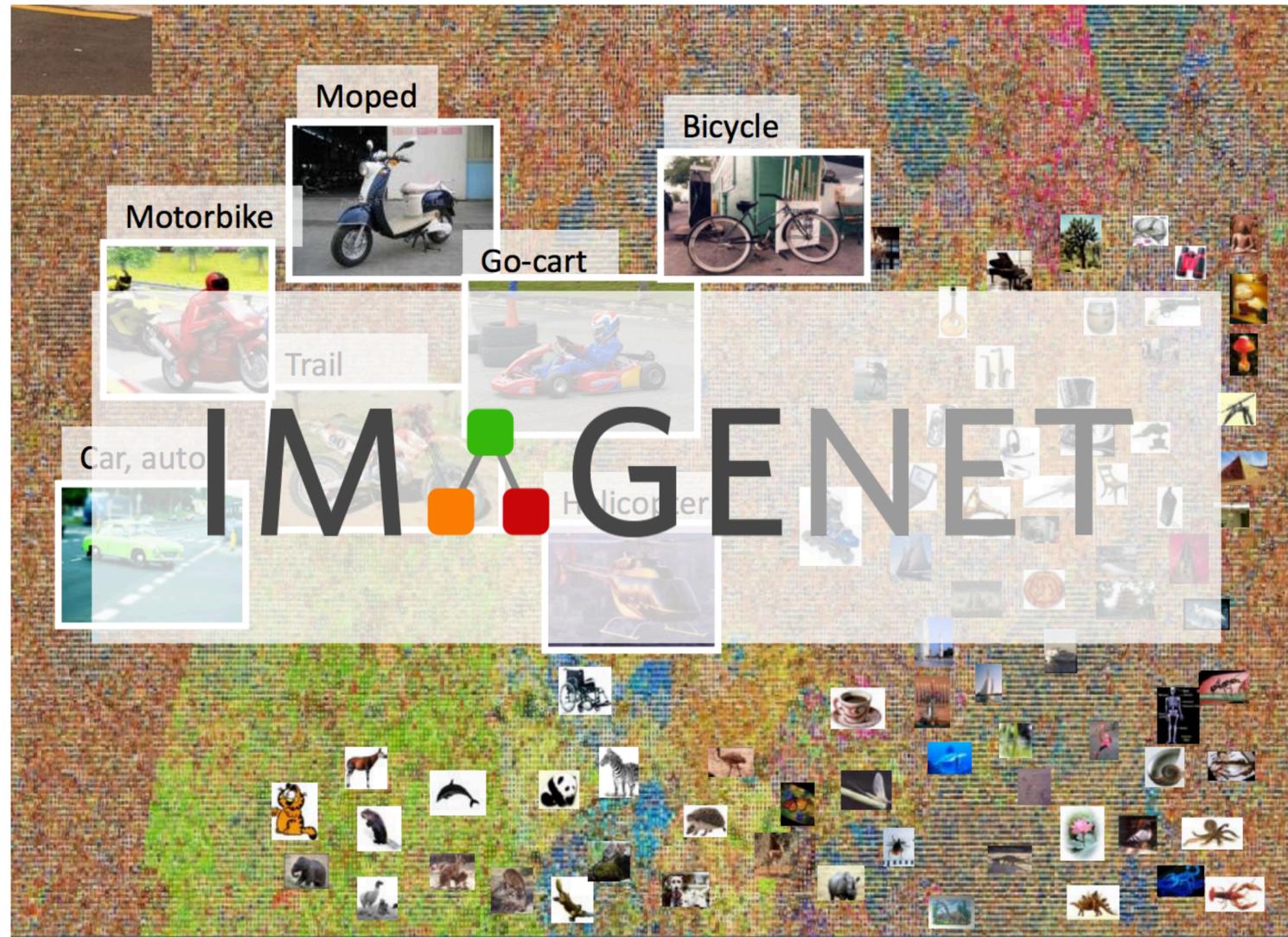
Classifying images



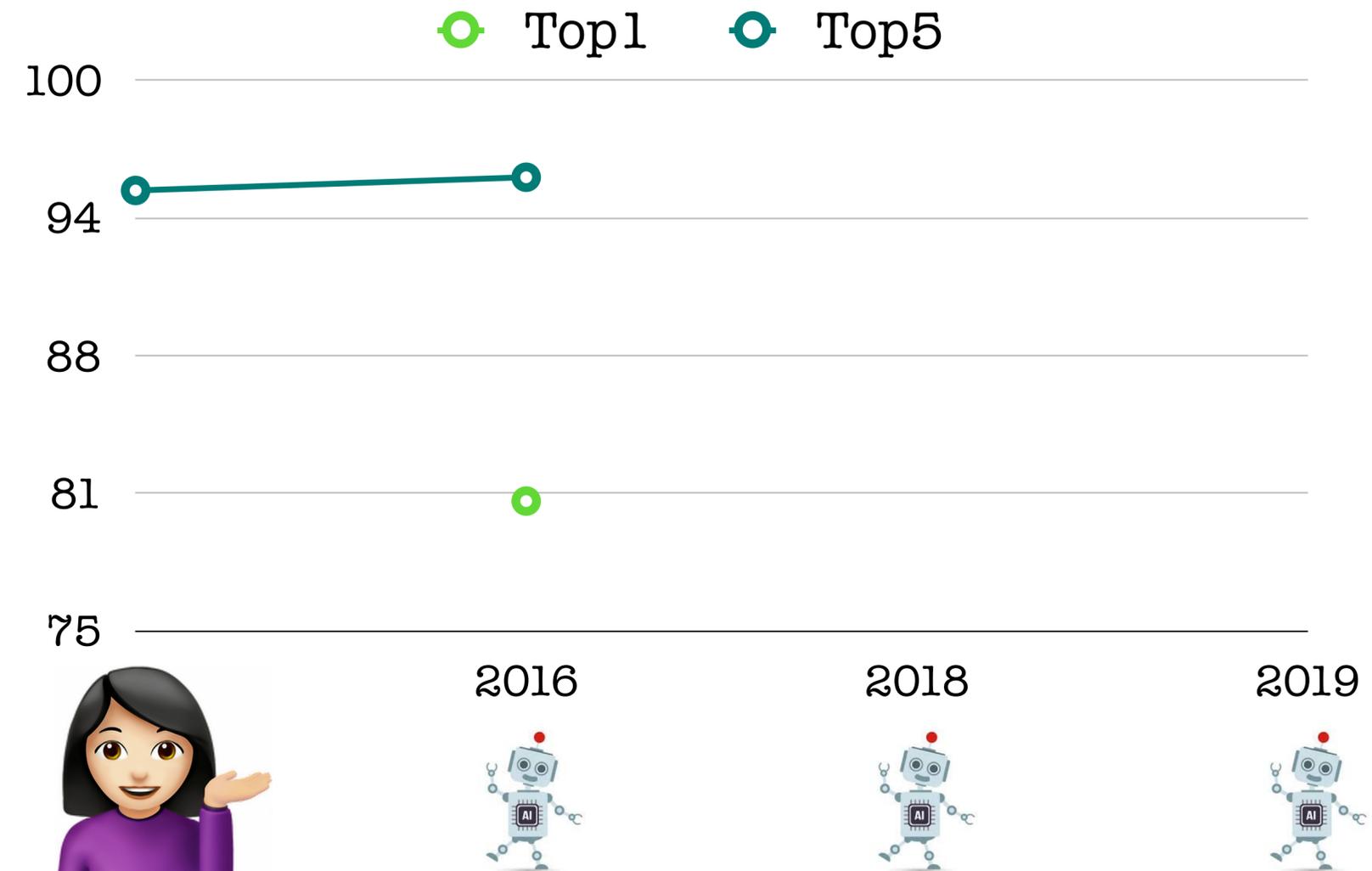
IMAGENET [Deng et al., CVPR 2009]



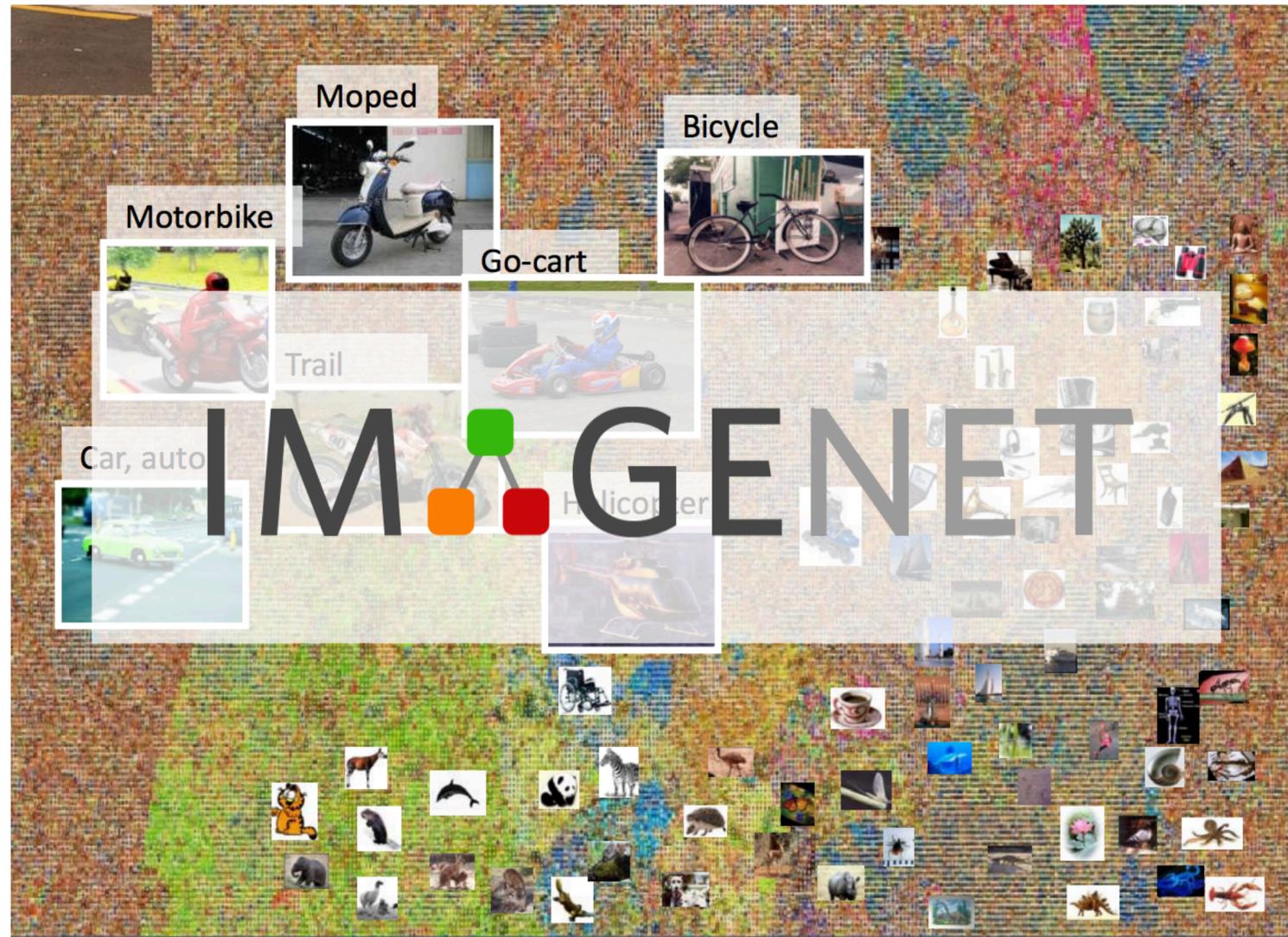
Classifying images



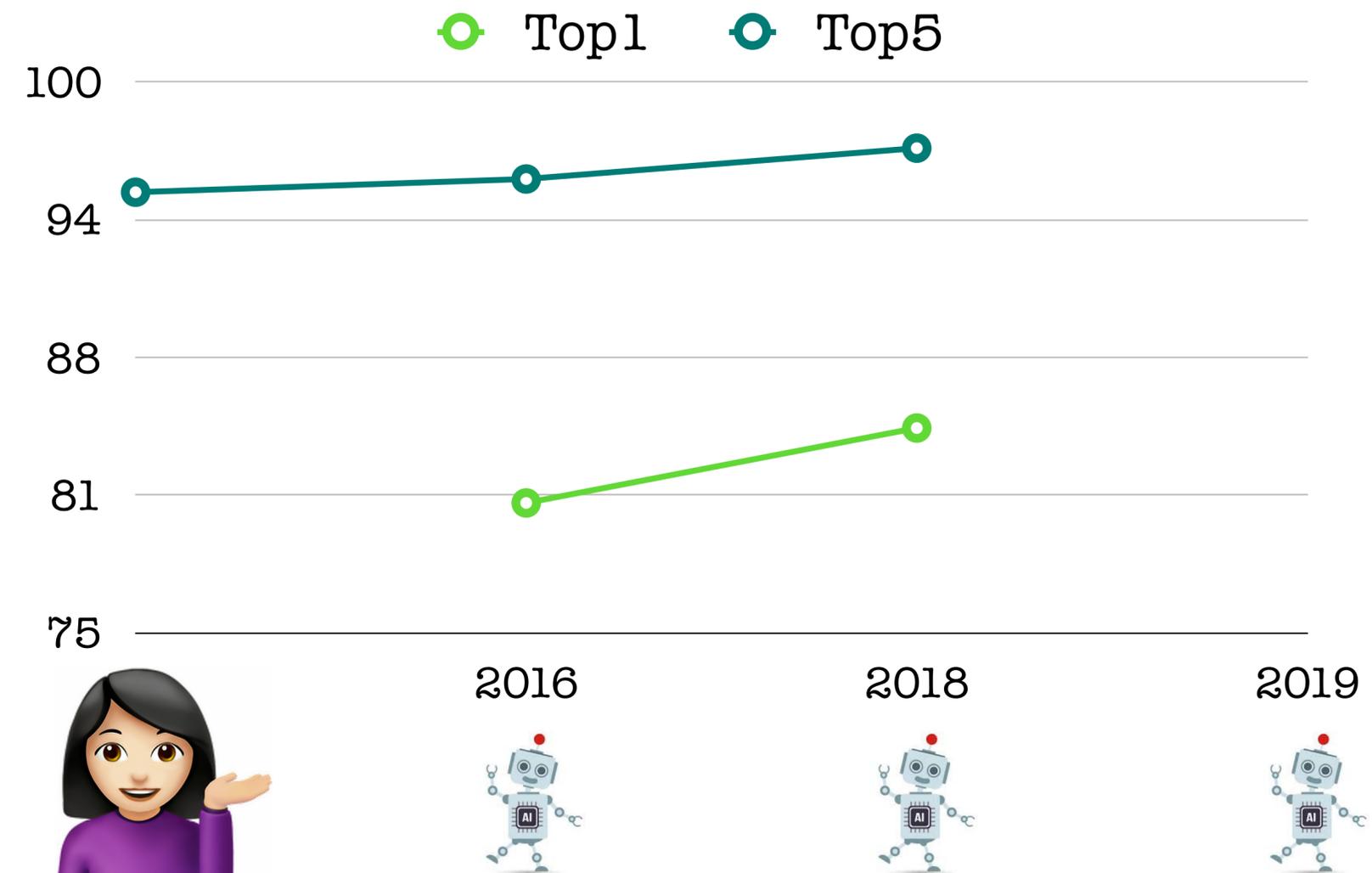
IMAGENET [Deng et al., CVPR 2009]



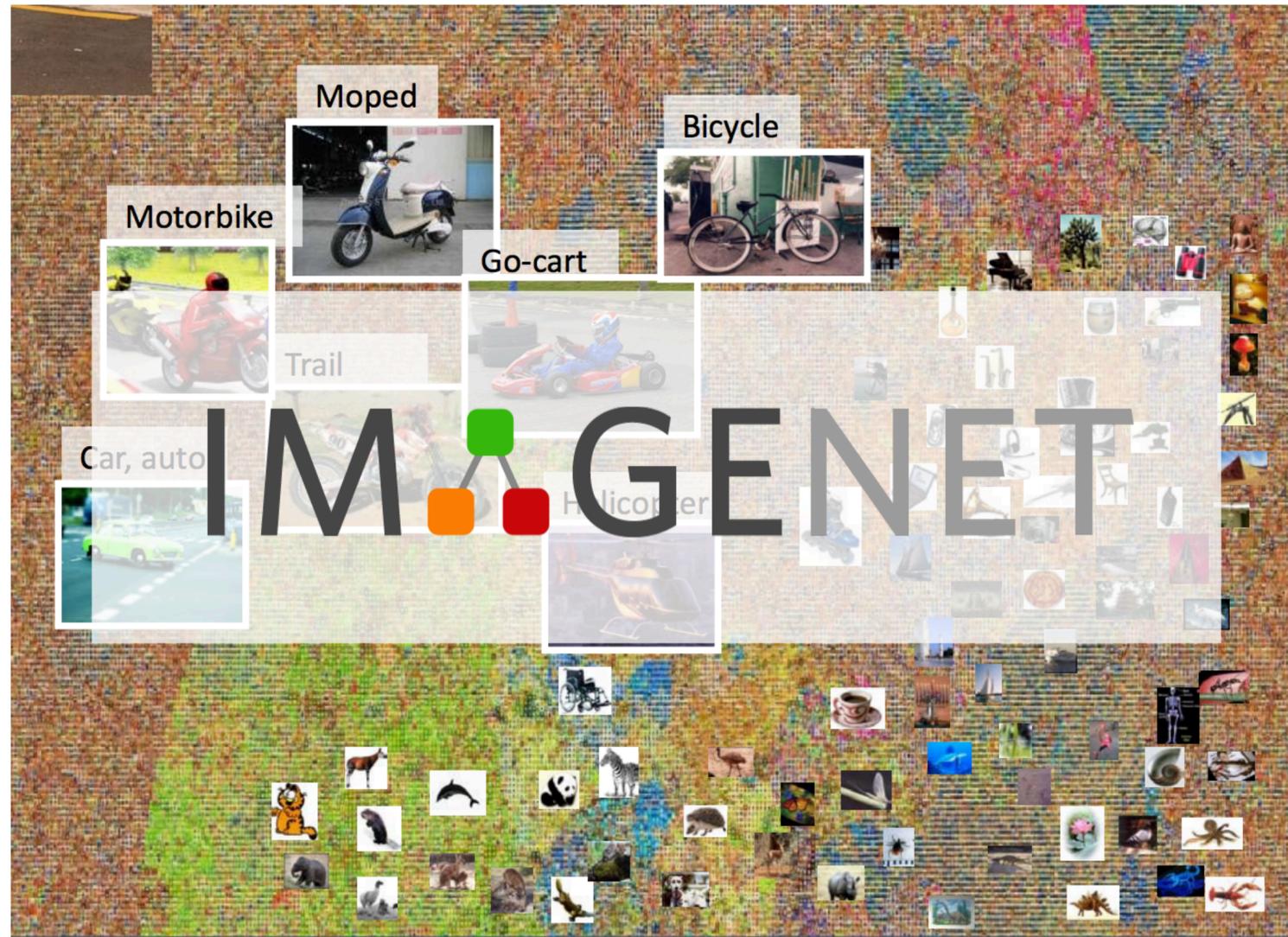
Classifying images



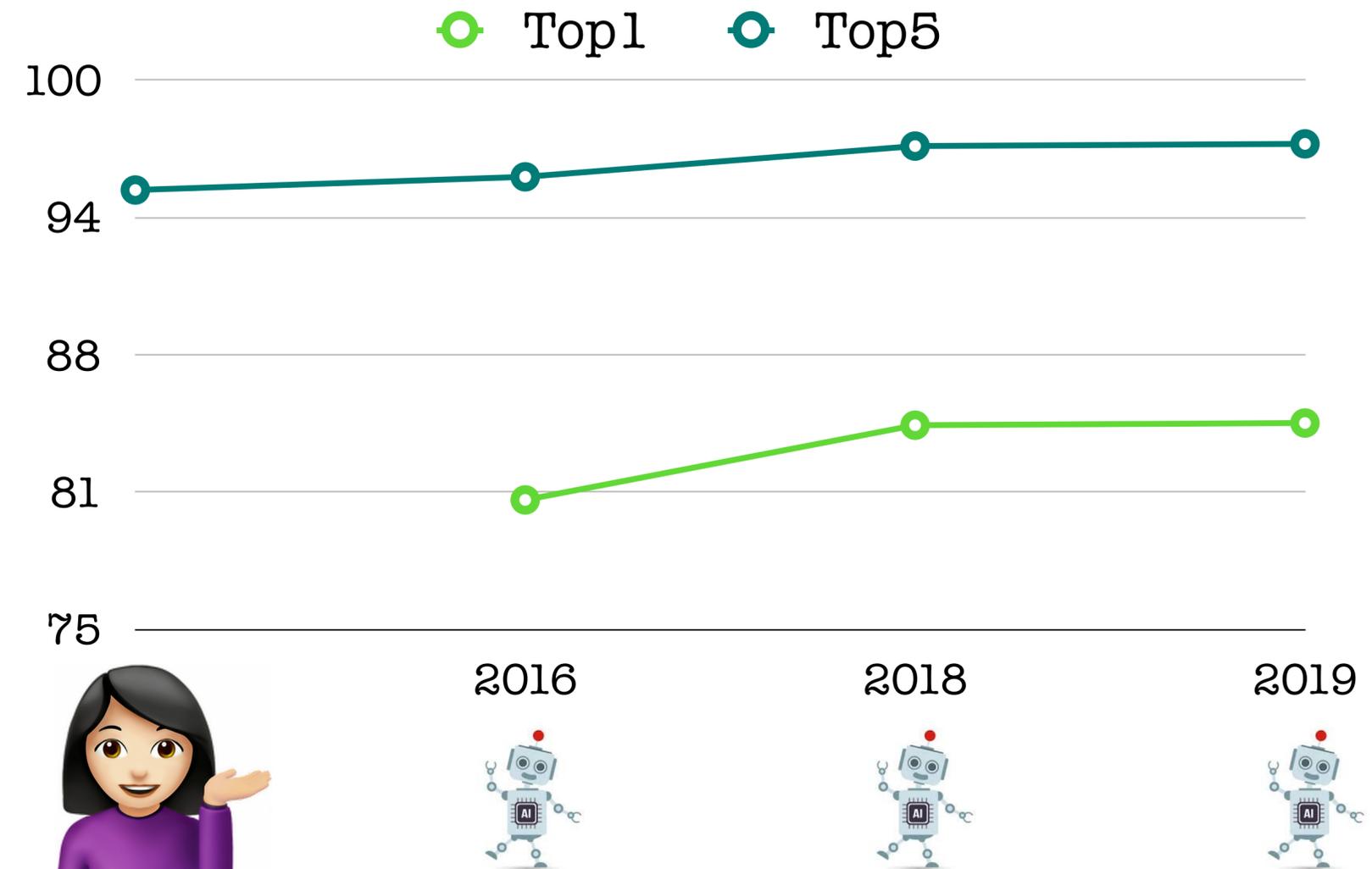
IMAGENET [Deng et al., CVPR 2009]



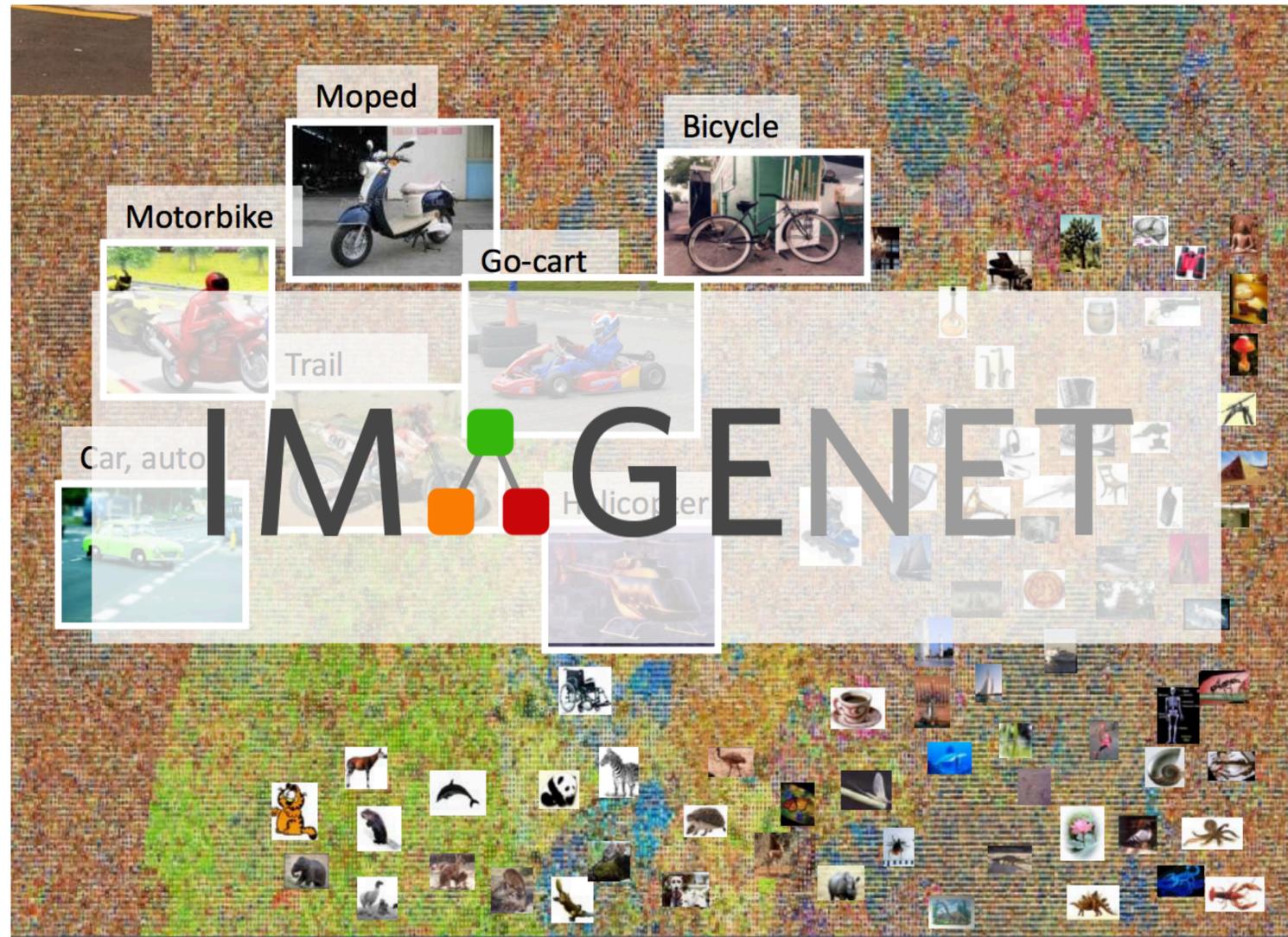
Classifying images



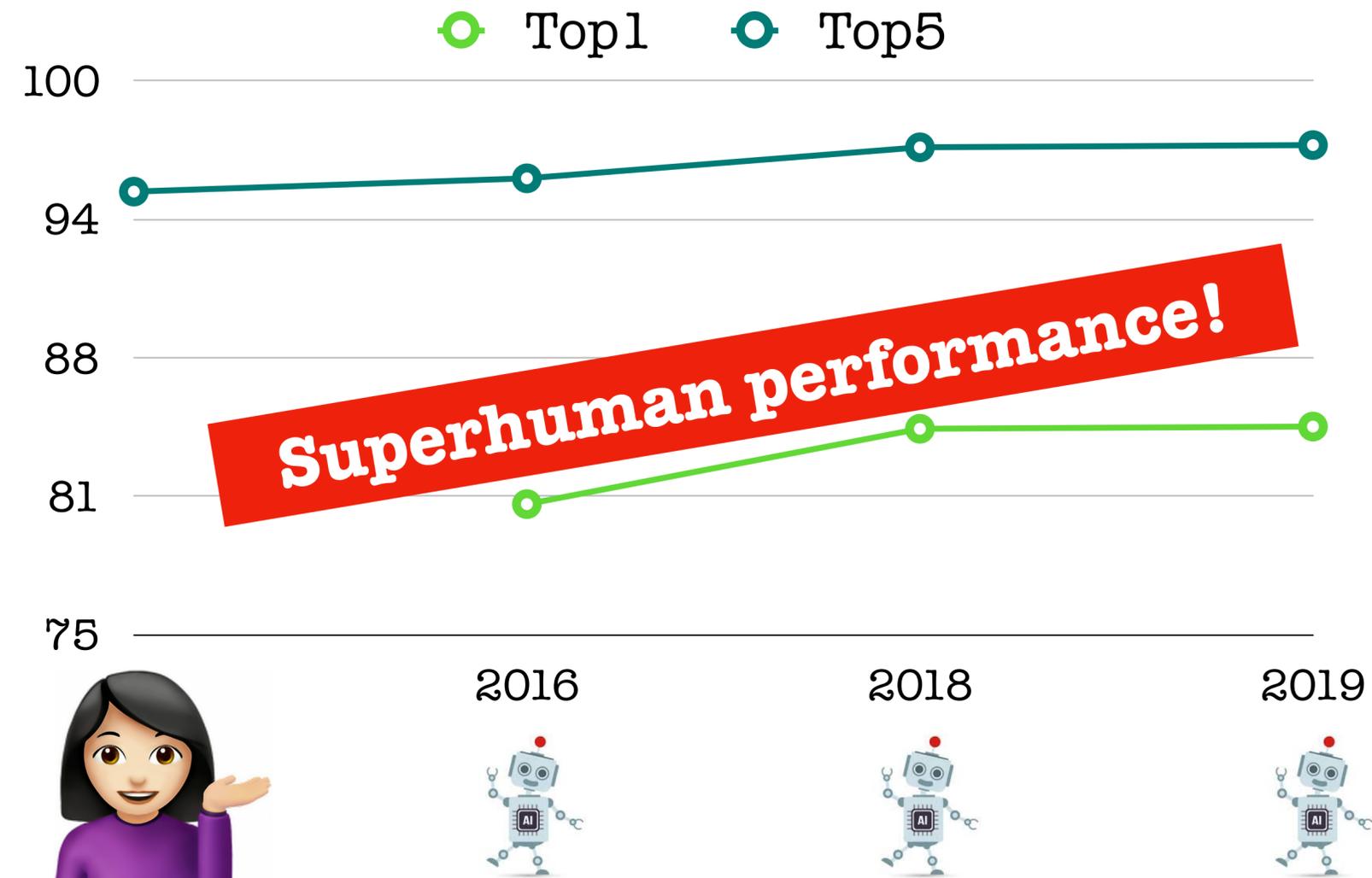
IMAGENET [Deng et al., CVPR 2009]



Classifying images

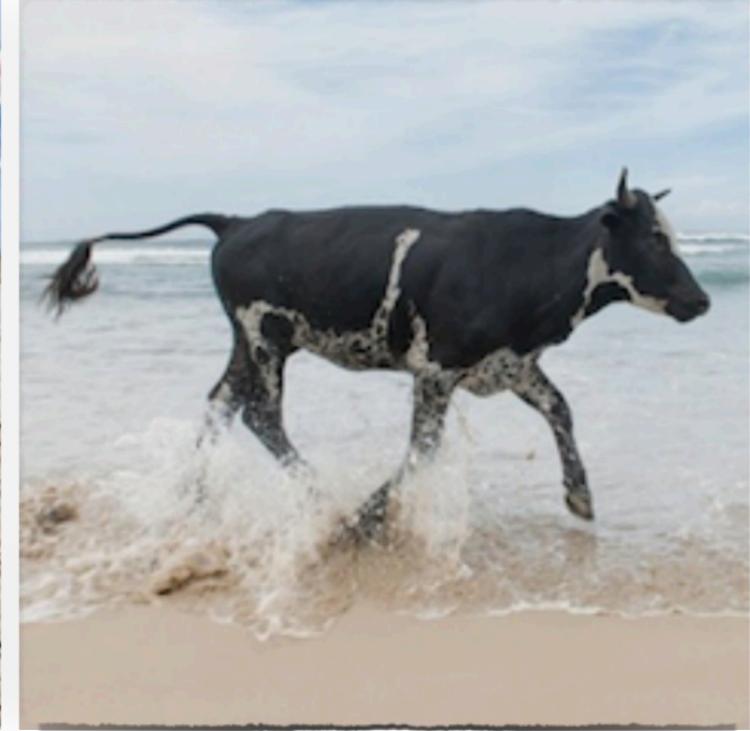


IMAGENET [Deng et al., CVPR 2009]



A simple quiz

A simple quiz



A simple quiz



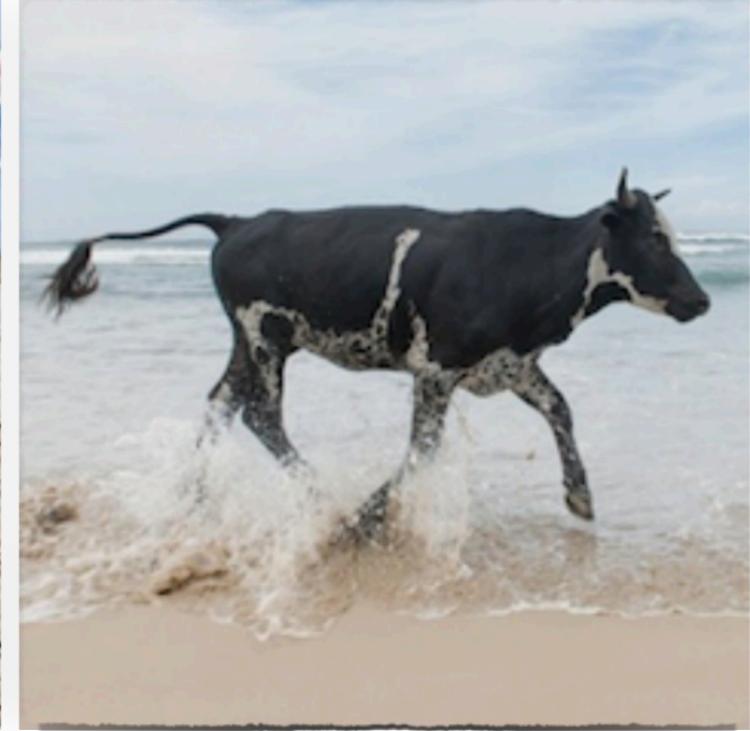
Cow



Cow



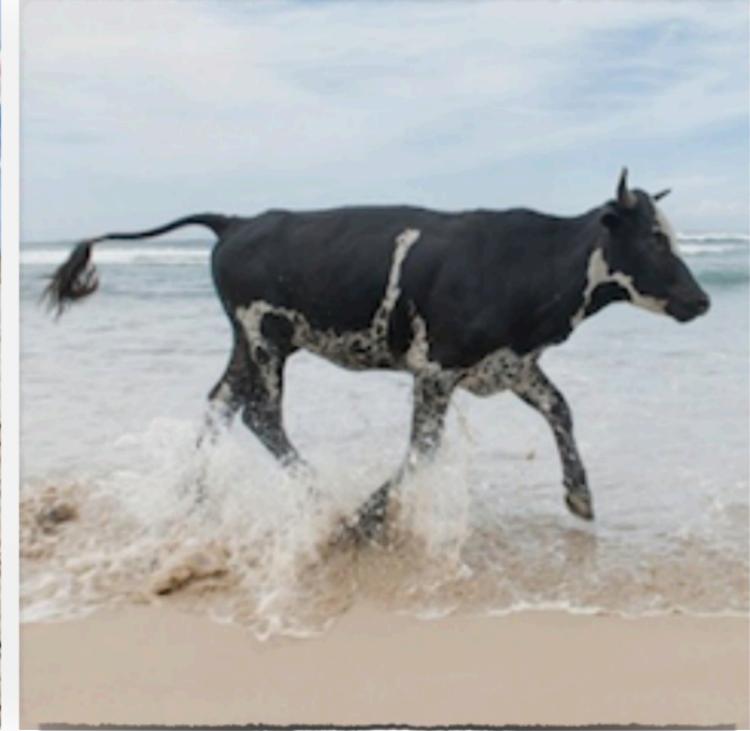
Cow



Cow



A simple quiz



Cow

Cow

Cow

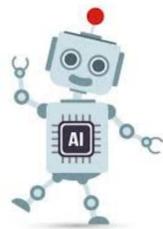
Cow

Cow

Cow

**No
person**

**No
person**



A simple quiz



Cow



Cow



Cow



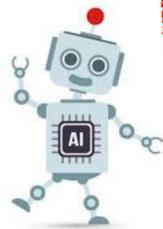
Cow

**No
person**

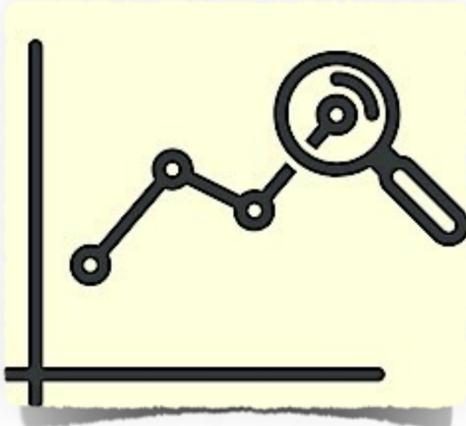
**No
person**

Cow

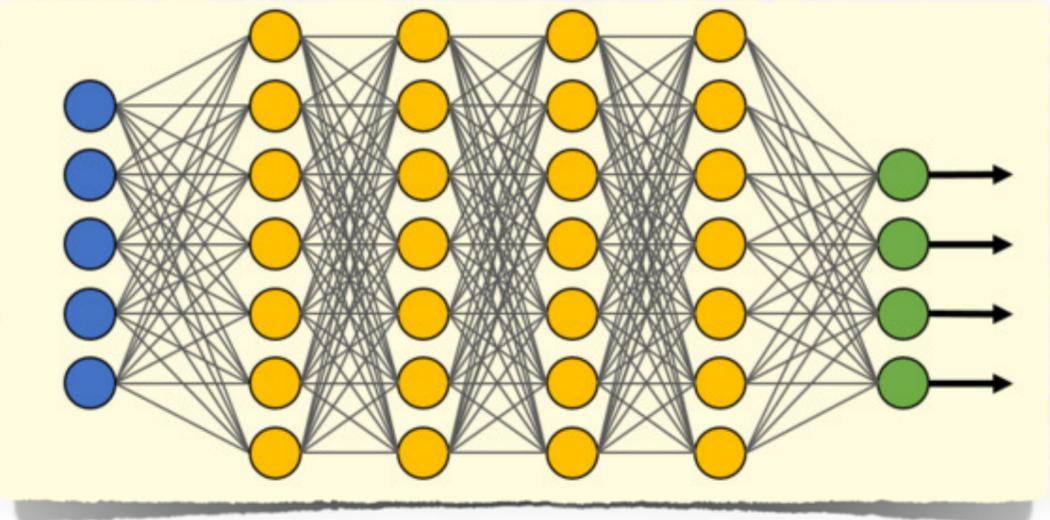
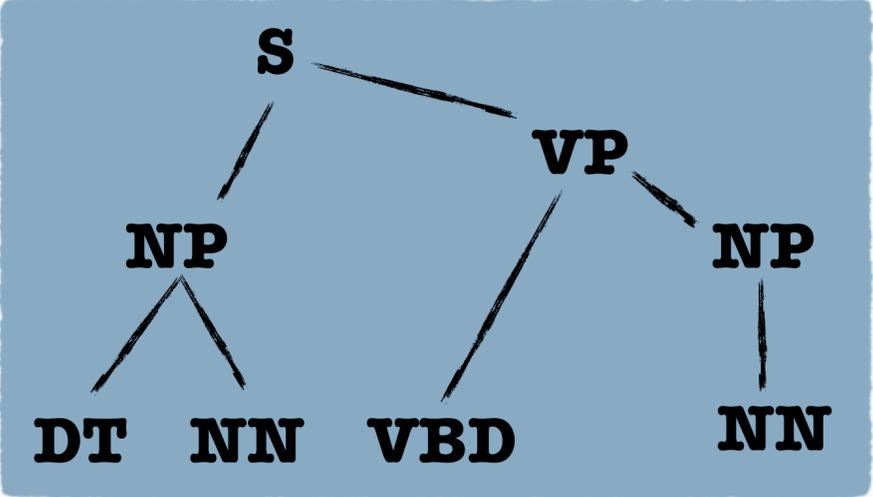
Cow



Machine Learning

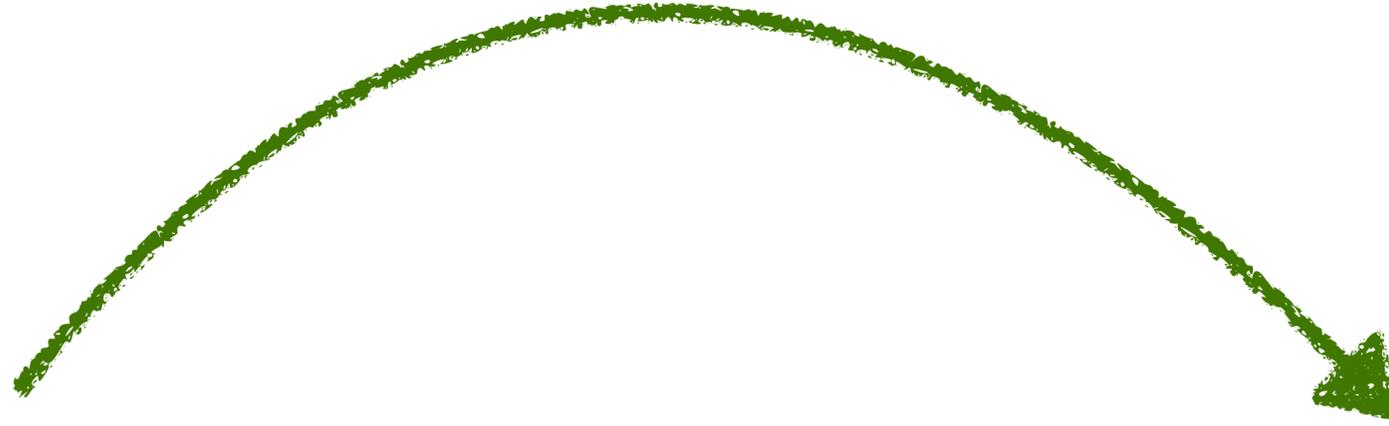
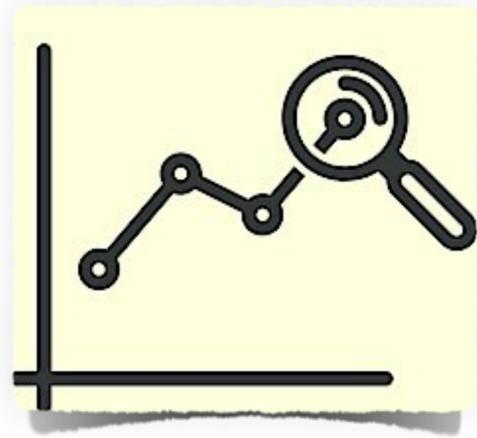


NLP



Linguistics





Question #1



**What do
predictions tell us
about the data?**

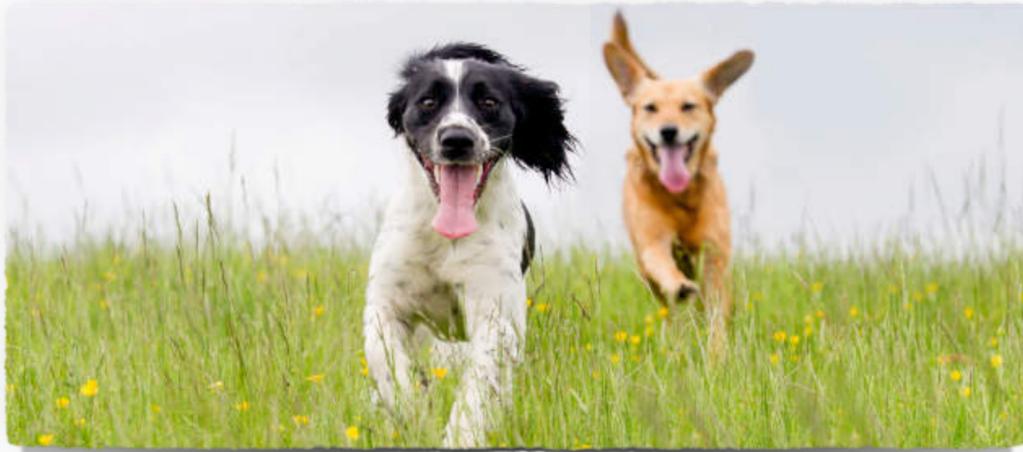


Natural Language Inference (NLI)

Natural Language Inference (NLI)

Given a premise, is a hypothesis true, false or neither?

Natural Language Inference (NLI)



Premise

Two dogs are running through a field.

Given a premise, is a hypothesis true, false or neither?



Hypothesis

The pets are sitting on a couch.

Natural Language Inference (NLI)



Premise

Two dogs are running through a field.



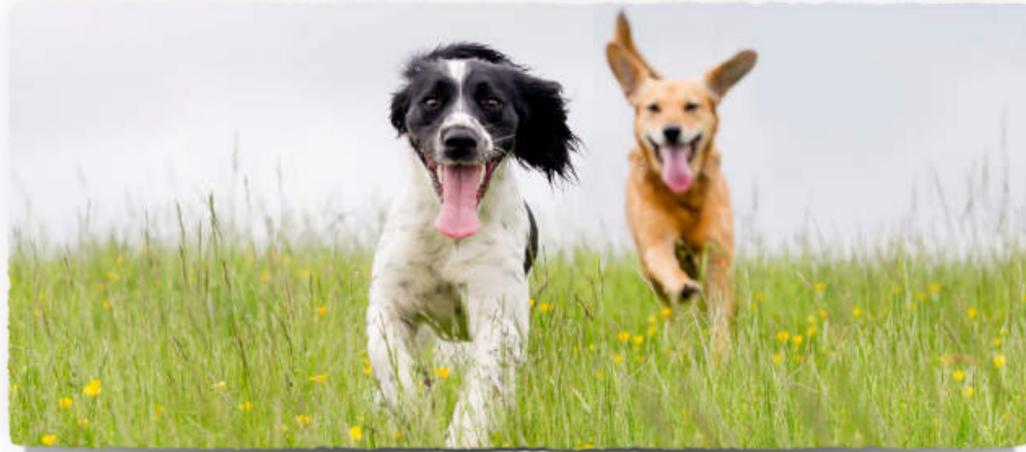
Hypothesis

The pets are sitting on a couch.

Given a premise, is a hypothesis true, false or neither?

- True → **Entailment**
- False → **Contradiction**
- Cannot Say → **Neutral**

Natural Language Inference (NLI)



Premise

Two dogs are running through a field.



Hypothesis

The pets are sitting on a couch.

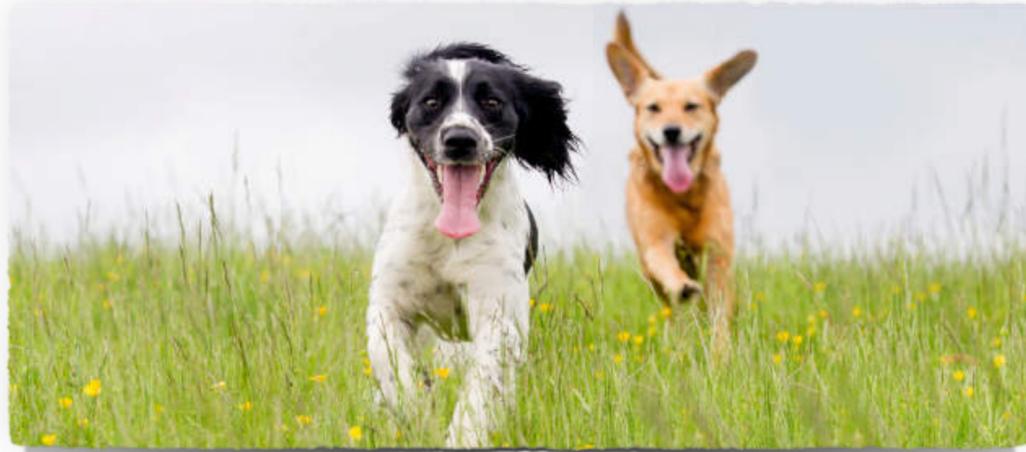
Given a premise, is a hypothesis true, false or neither?

True → **Entailment**

False → **Contradiction**

Cannot Say → **Neutral**

Natural Language Inference (NLI)



Premise

Two dogs are running through a field.



Hypothesis

The pets are sitting on a couch.

Given a premise, is a hypothesis true, false or neither?

True → **Entailment**

False → **Contradiction**

Cannot Say → **Neutral**

[Katz, 1972; van Benthem, 2008; Dagan et al., 2006]

NLI Dataset Creation Process

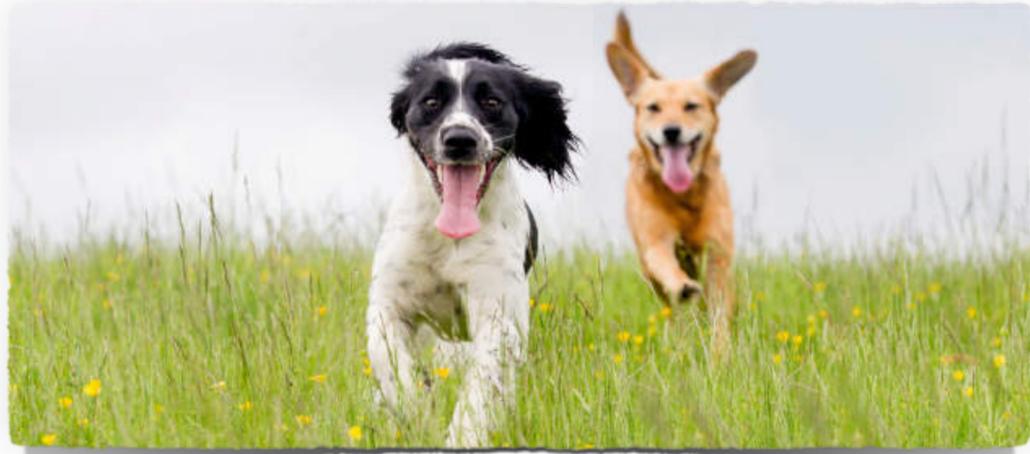
NLI Dataset Creation Process



Two dogs are
running through
a field.

Premise

NLI Dataset Creation Process

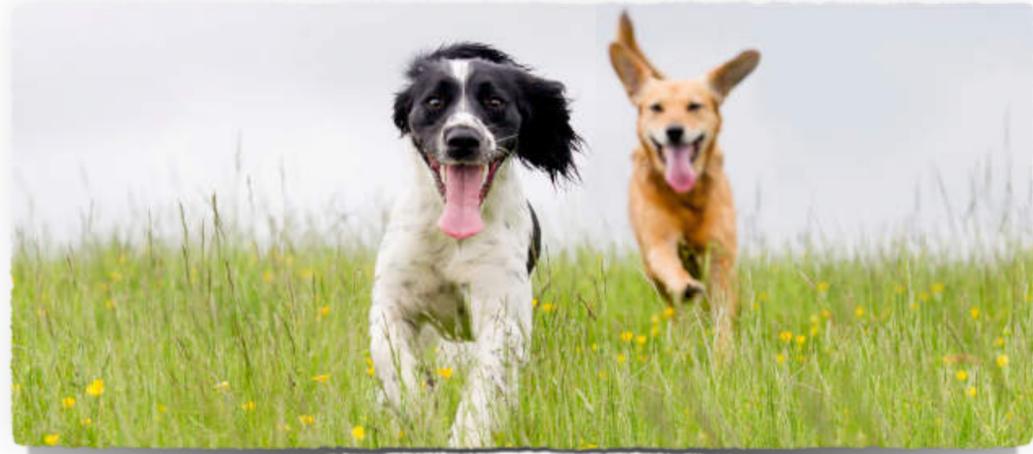


Two dogs are running through a field.

Premise

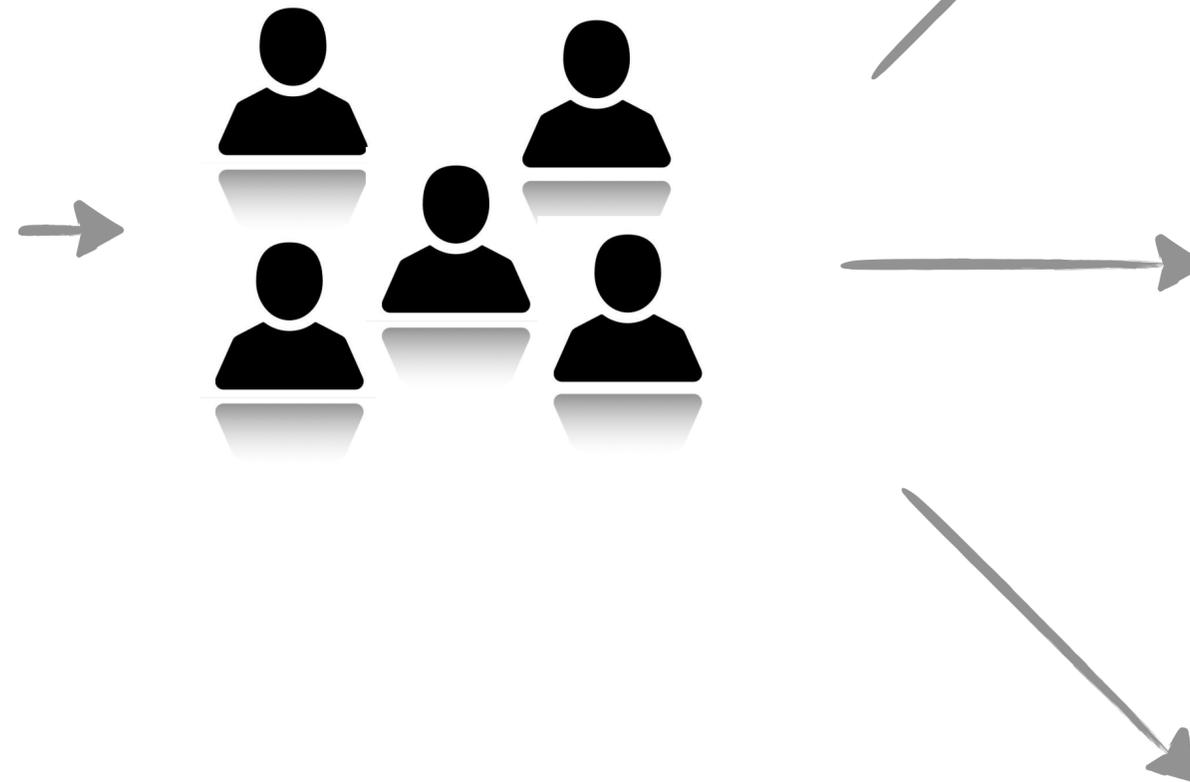


NLI Dataset Creation Process

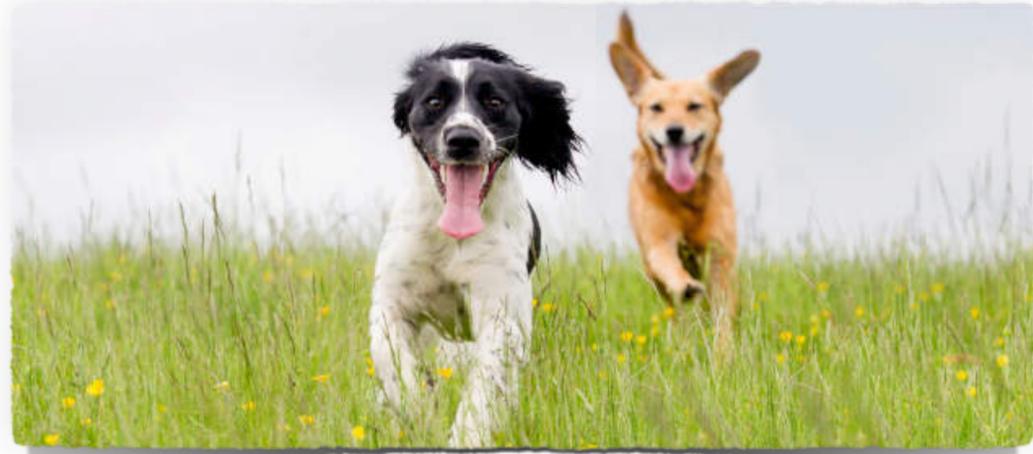


Two dogs are running through a field.

Premise



NLI Dataset Creation Process



Two dogs are running through a field.

Premise

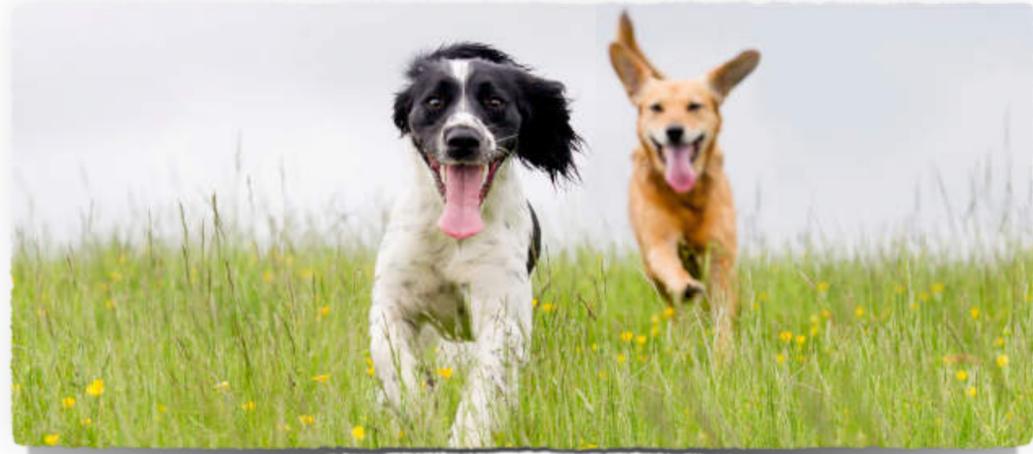


Entailment



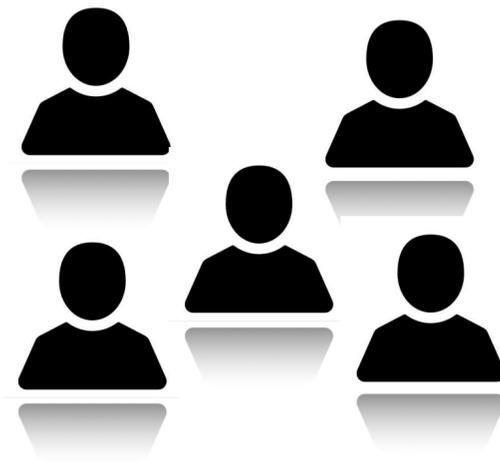
There are animals outdoors.

NLI Dataset Creation Process



Two dogs are running through a field.

Premise

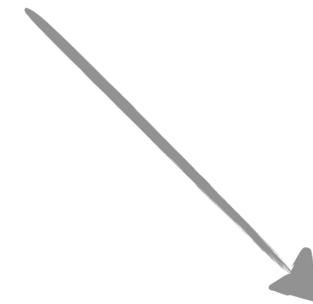


Entailment

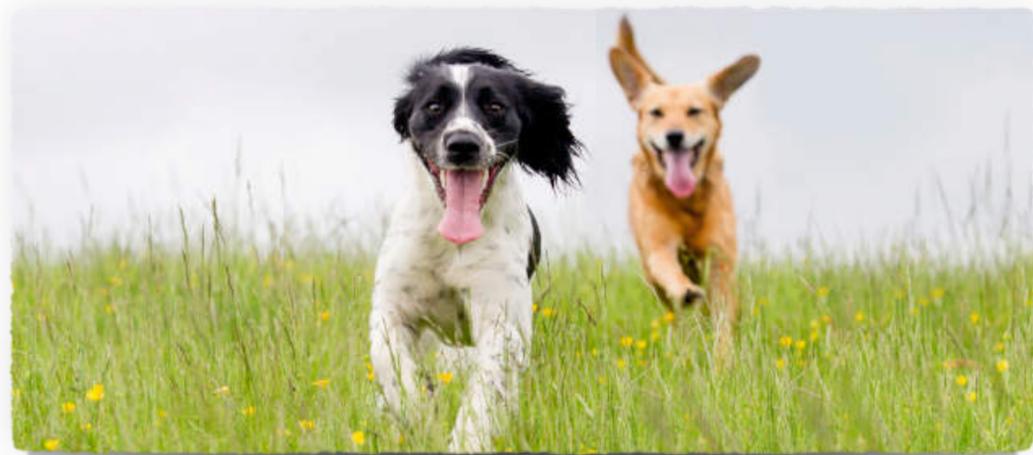
Neutral

There are animals outdoors.

Some puppies are running to catch a stick.

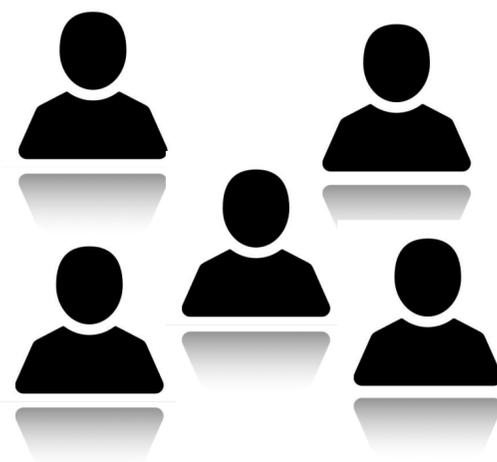


NLI Dataset Creation Process



Two dogs are running through a field.

Premise



Entailment

Neutral

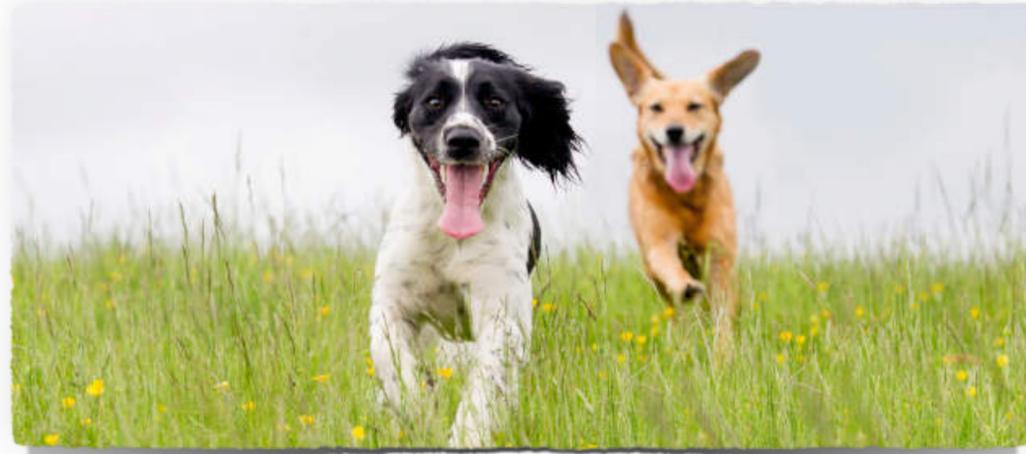
Contradiction

There are animals outdoors.

Some puppies are running to catch a stick.

The pets are sitting on a couch.

NLI Dataset Creation Process



Two dogs are running through a field.

Premise



Entailment

Neutral

Contradiction

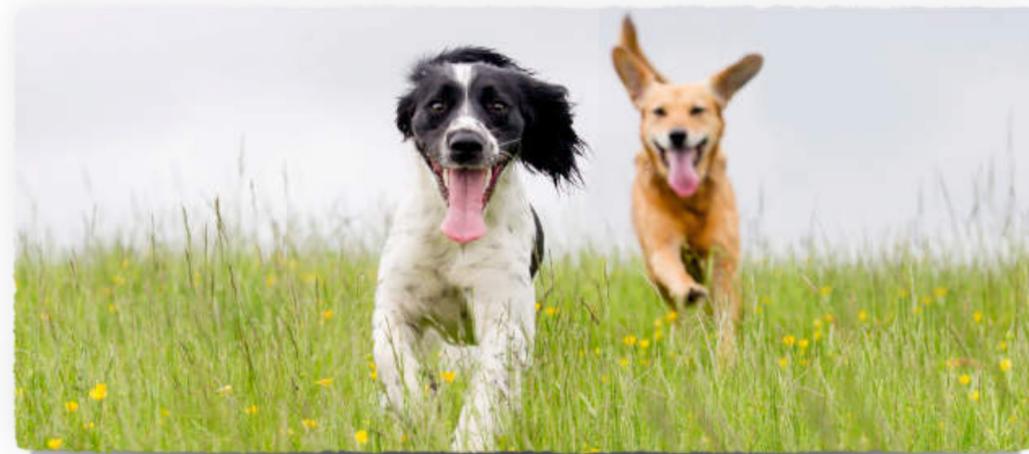
There are animals outdoors.

Some puppies are running to catch a stick.

The pets are sitting on a couch.

- **Stanford NLI** [Bowman et. al, 2015] 570 K

NLI Dataset Creation Process



Two dogs are running through a field.

Premise



Entailment

There are animals outdoors.

Neutral

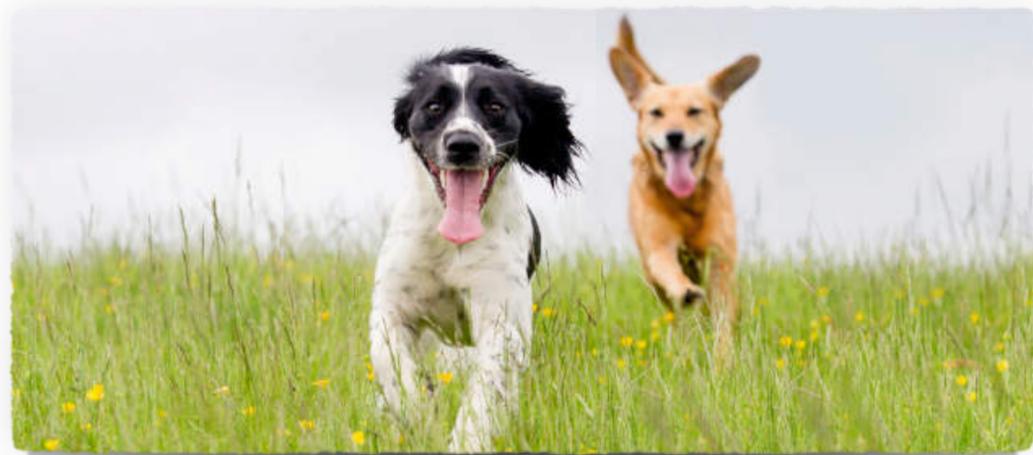
Some puppies are running to catch a stick.

Contradiction

The pets are sitting on a couch.

- **Stanford NLI** [Bowman et. al, 2015] 570 K
- **Multi-genre NLI** [Williams et. al., 2017] 433 K

NLI Dataset Creation Process



Two dogs are running through a field.

Premise



Entailment

There are animals outdoors.

Neutral

Some puppies are running to catch a stick.

Contradiction

The pets are sitting on a couch.

- **Stanford NLI** [Bowman et. al, 2015] 570 K
- **Multi-genre NLI** [Williams et. al., 2017] 433 K
- Matched and Mismatched Test Sets

Leaderboard progress

#	Team Name	Notebook	Team Members	Score ?	Entries	Last
1	bgm			0.90557	4	14d
2	Haoming Jiang			0.87923	10	1mo
3	Xiaodong Liu			0.86443	4	10mo
4	Anonymous			0.86351	2	1y
5	anonymous11111			0.85177	18	1mo
6	Ariel			0.85065	41	5mo
7	sherry77			0.85034	17	5mo
		⋮				
	 Bidirectional LSTM			0.67507		
104	gabrielalmeida			0.67313	5	8mo
105	Zippy			0.67160	2	1y
106	kudkudak			0.66435	2	1y
107	Shawn Tan			0.65271	1	6d
	 CBOW			0.65200		

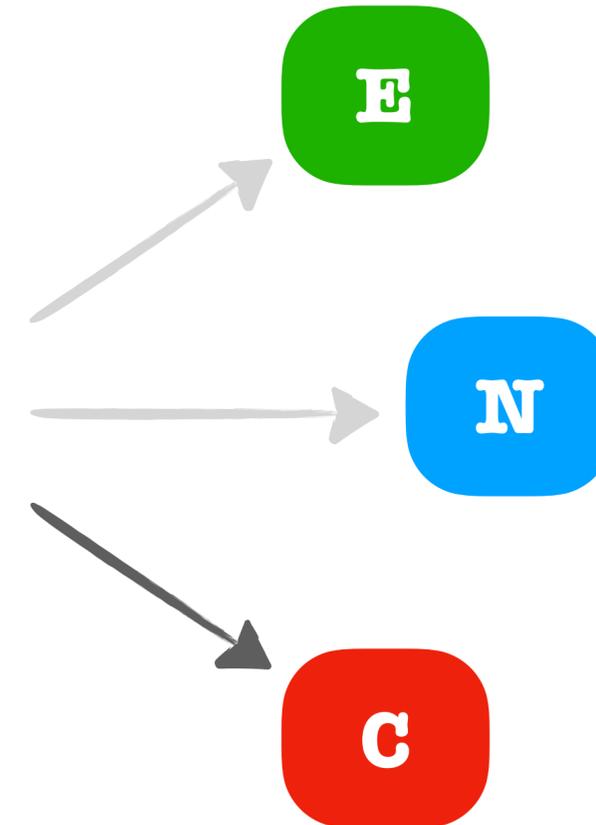
NLI as Text Classification

Two dogs are
running through
a field.

Premise

The pets are
sitting on a
couch.

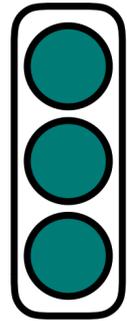
Hypothesis



NLI as Text Classification

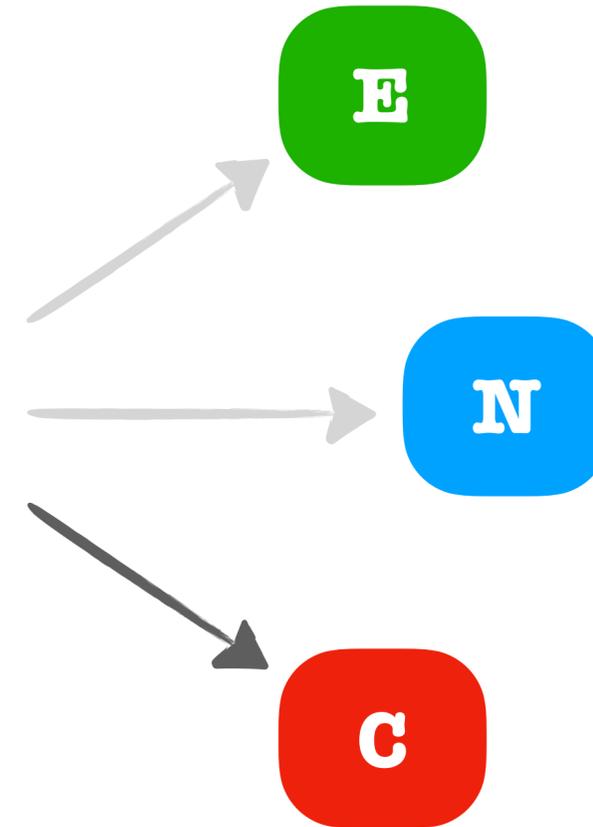
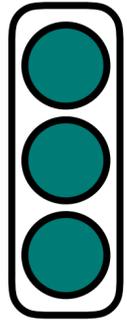
Two dogs are running through a field.

Premise

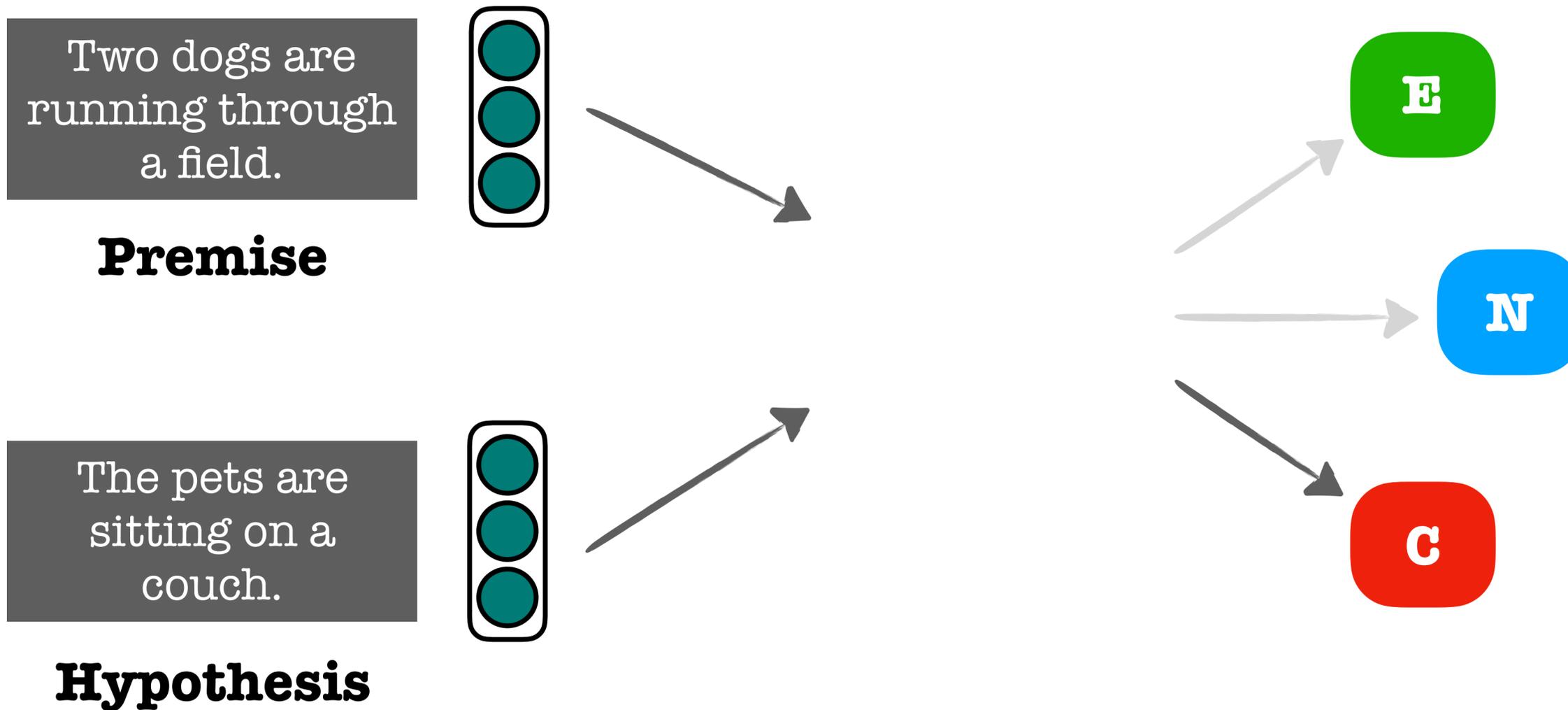


The pets are sitting on a couch.

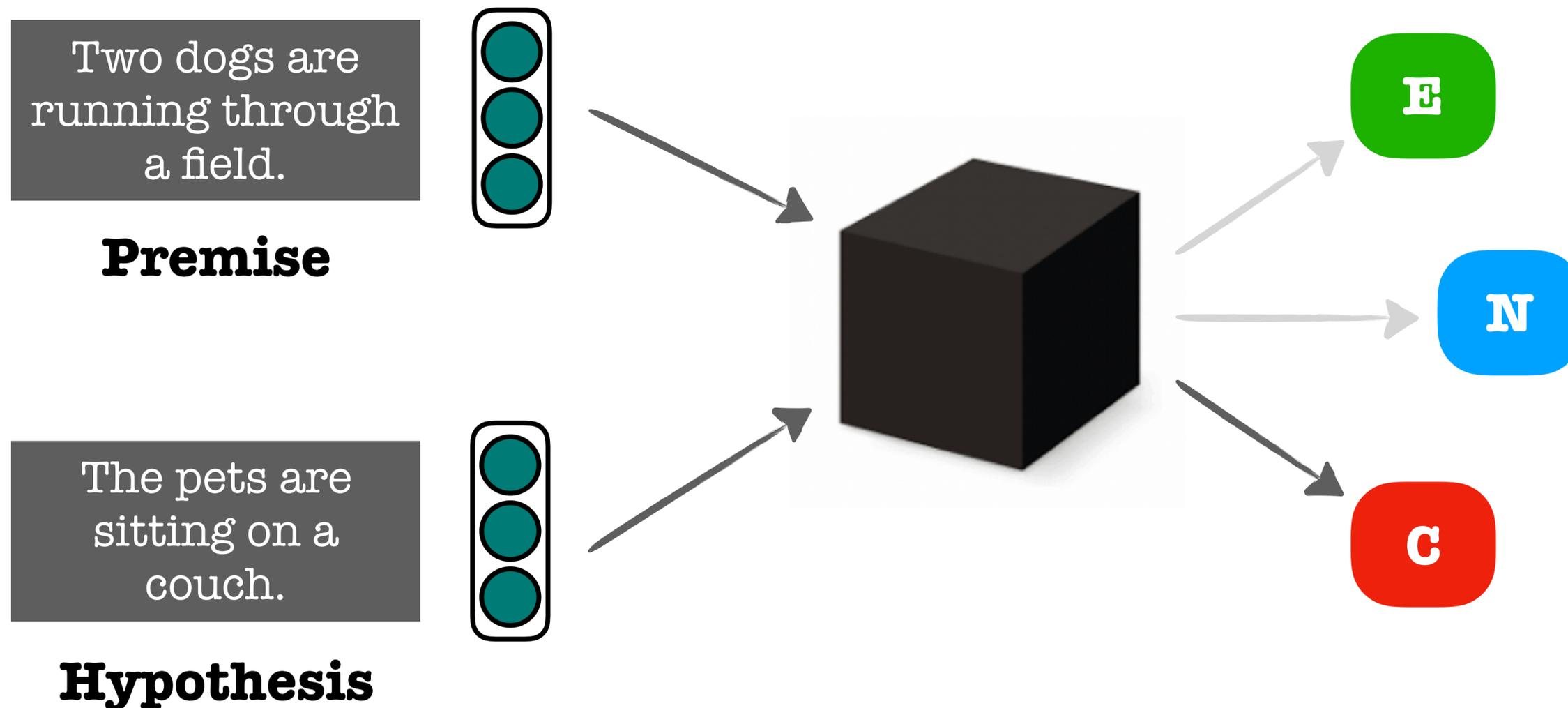
Hypothesis



NLI as Text Classification

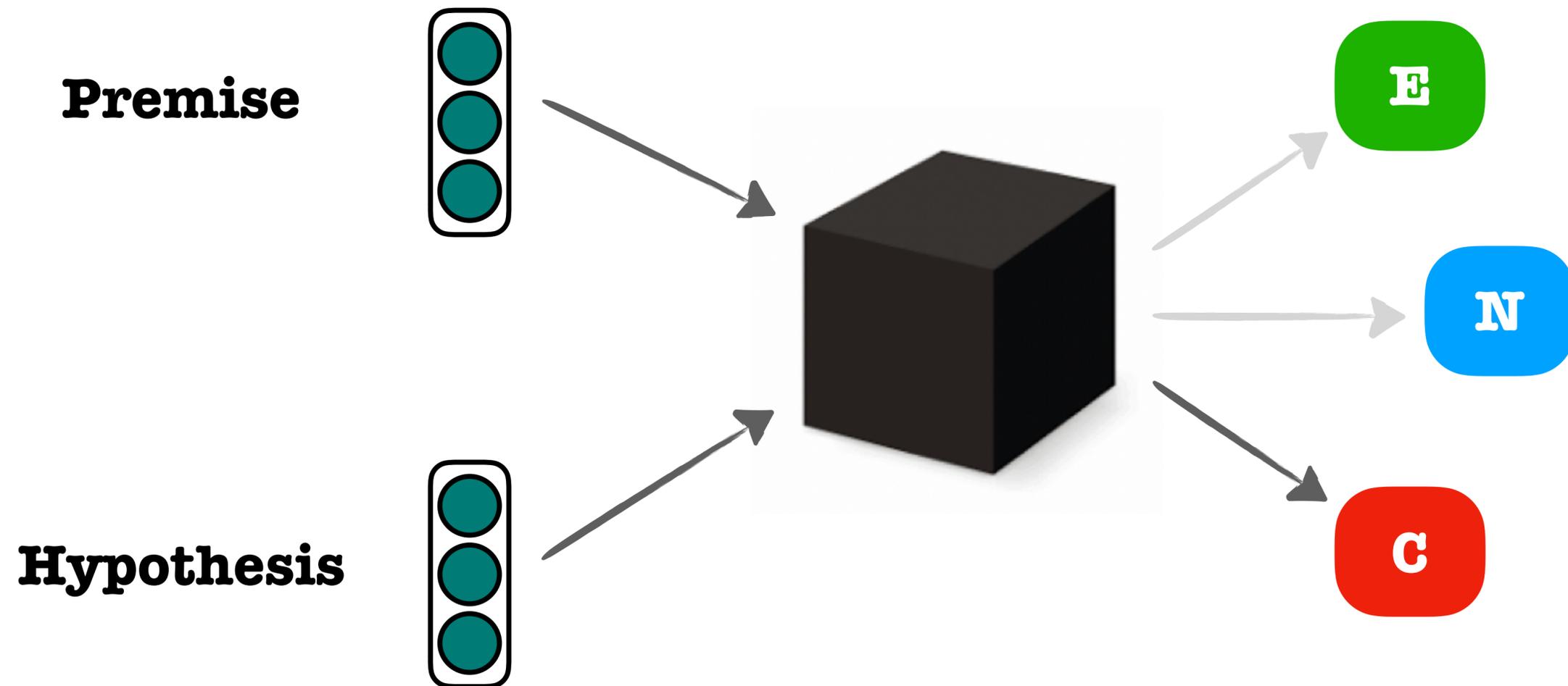


NLI as Text Classification

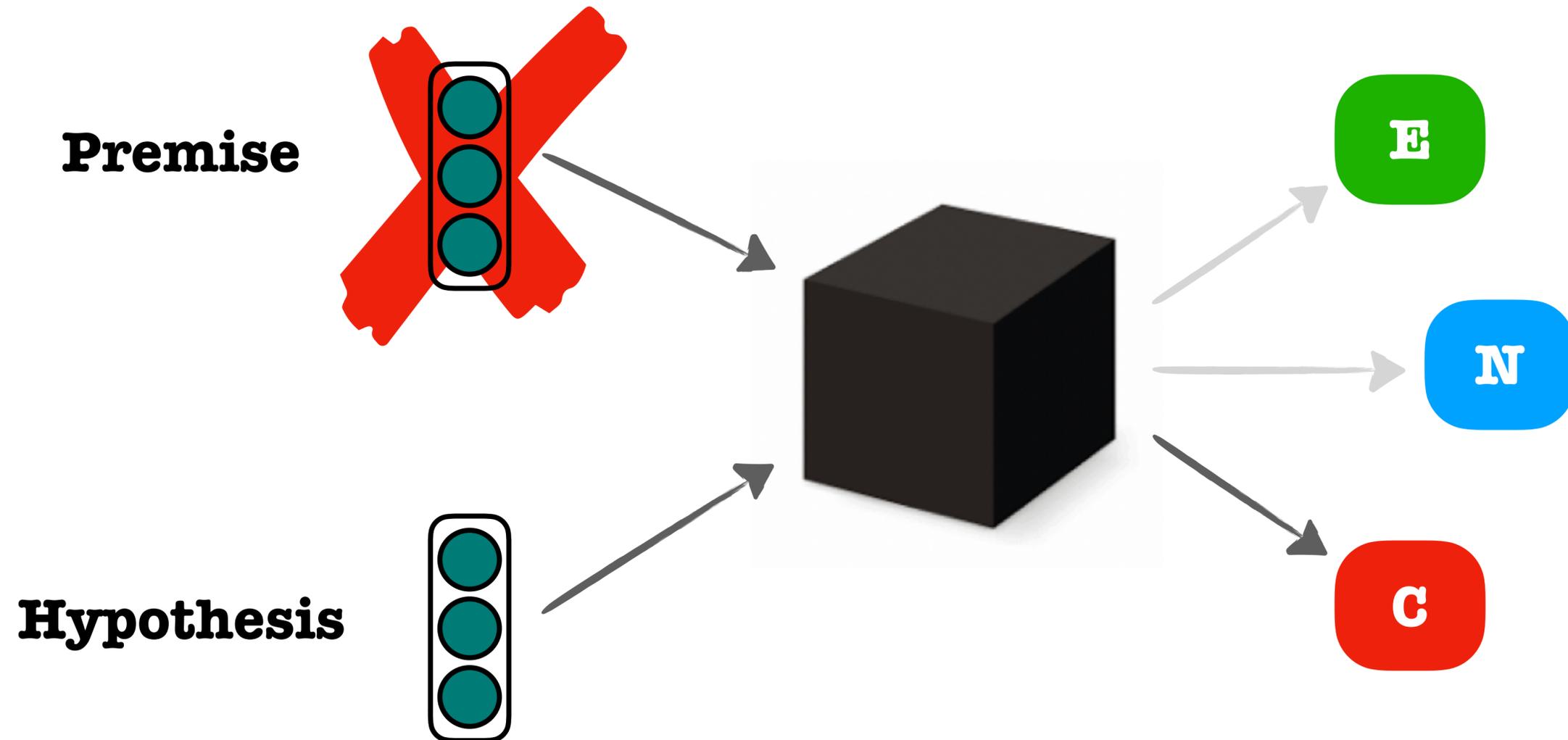


DAM - Decomposable Attention Model (Parikh et. al. 2016)
ESIM - Enhanced Sequential Inference Model (Chen et. al., 2017)
DIIN - Densely Interactive Inference Network (Gong et. al. 2018)

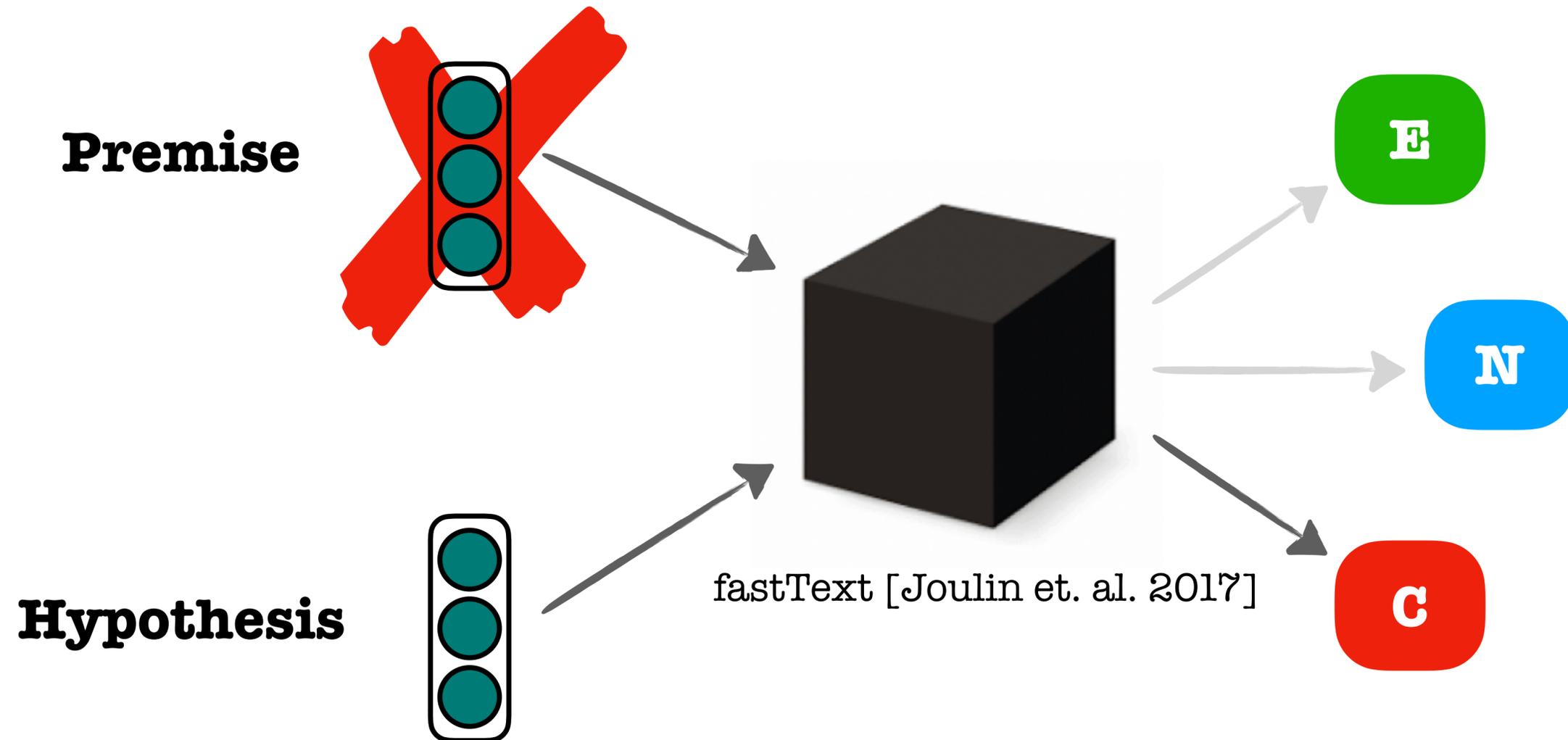
A simple experiment



A simple experiment



A simple experiment



Intuitively...

Given **no** premise, is a hypothesis true, false or neither?

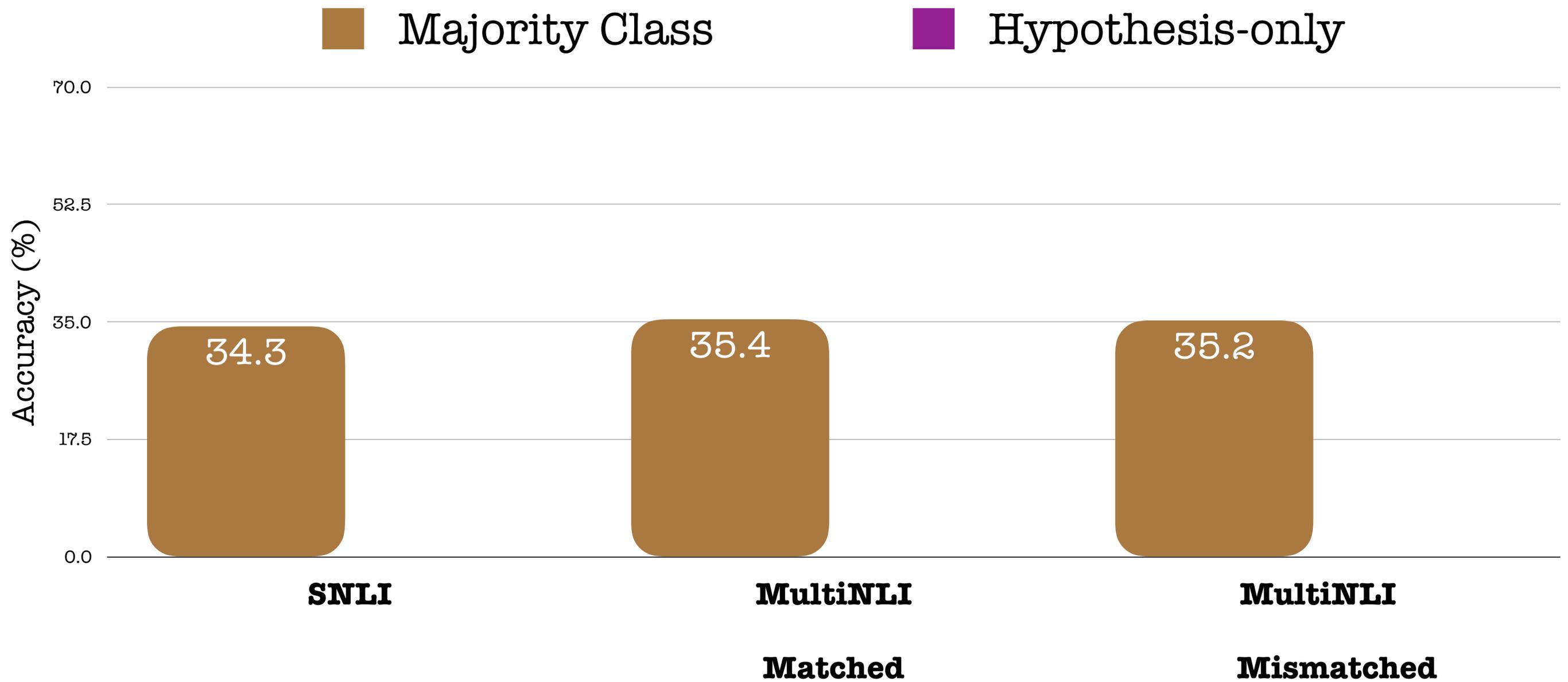
Intuitively...

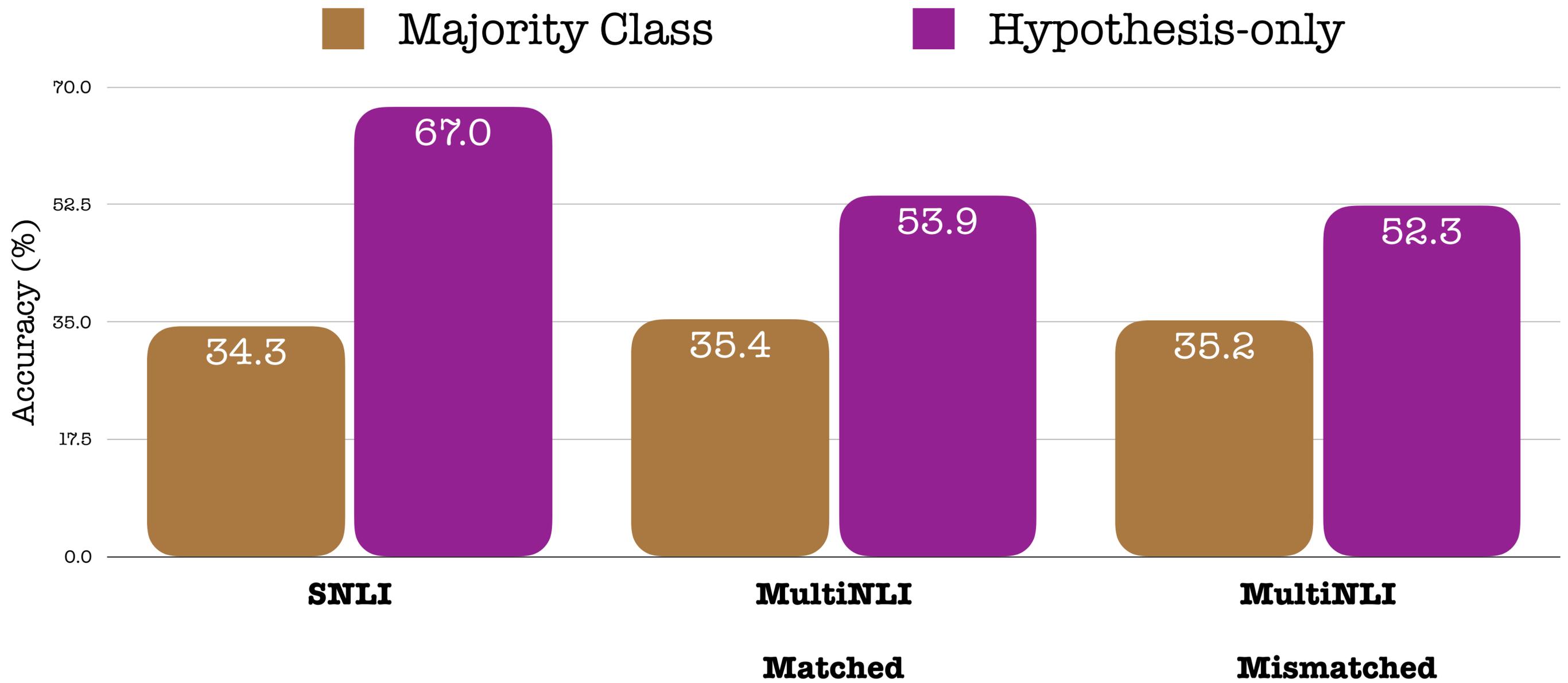
Hypothesis

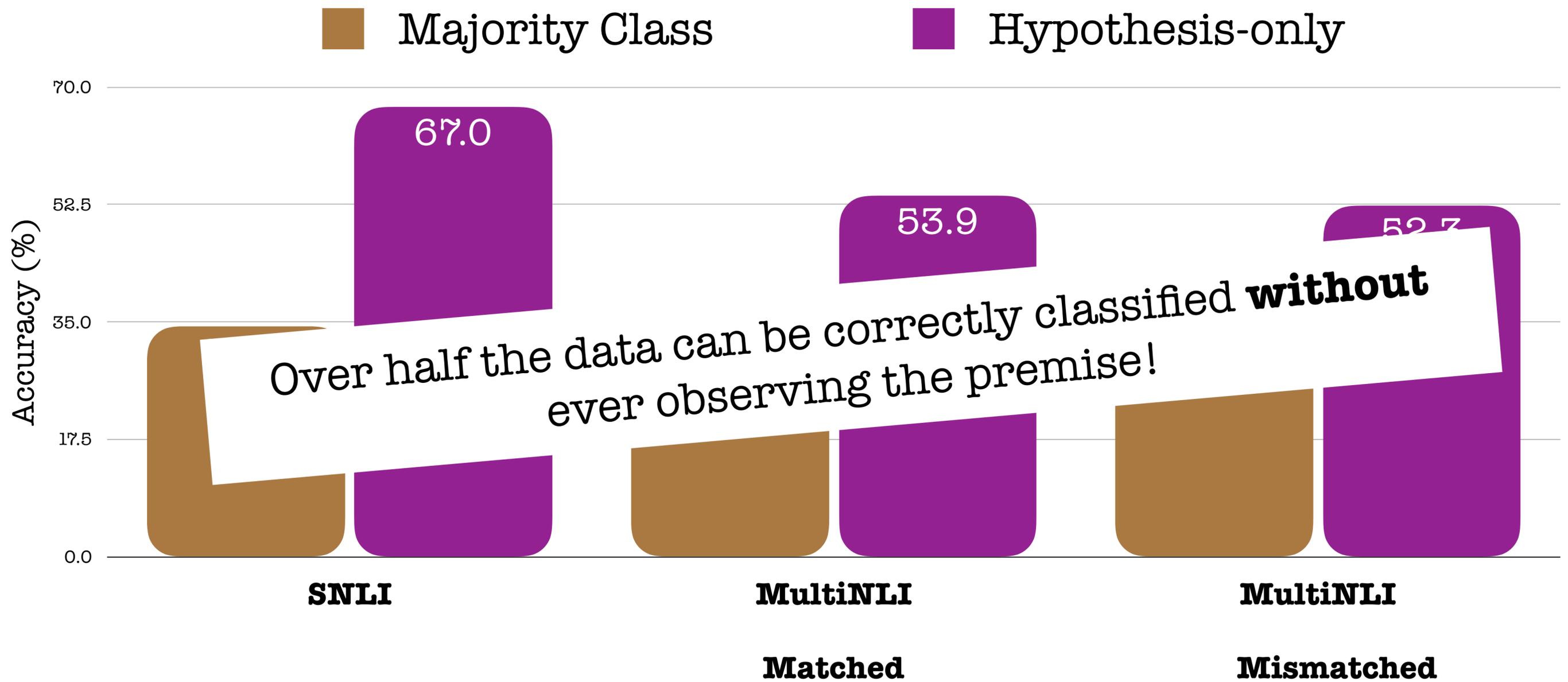
The little boy is diving off the diving board because he is an excellent swimmer.

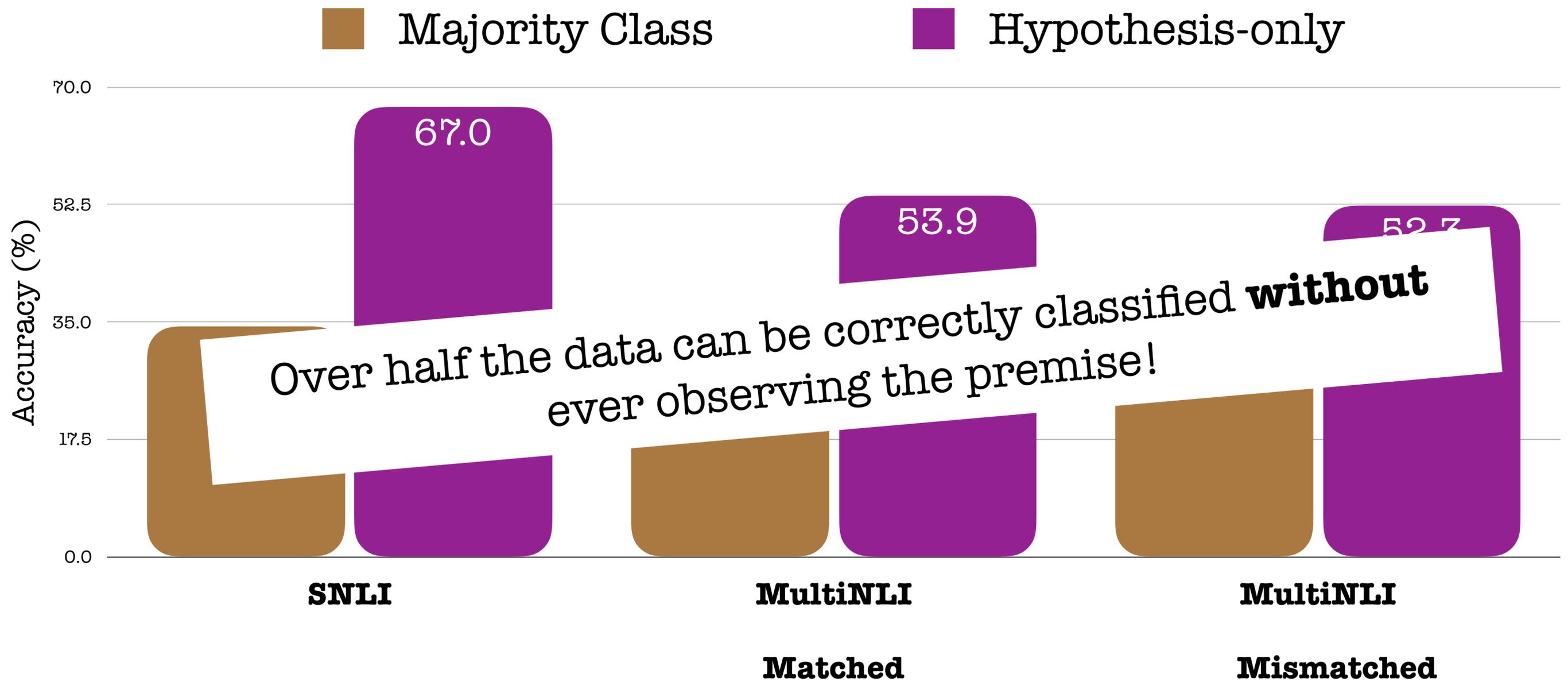
Given **no** premise, is a hypothesis true, false or neither?

- True → **Entailment**
- False → **Contradiction**
- Cannot Say → **Neutral**









Poliak et. al., 2018,
Glockner et. al., 2018

Digging Deeper



Digging Deeper



- **Annotation Artifacts:** Clues which give away the correct prediction without any reasoning.

Digging Deeper



- **Annotation Artifacts:** Clues which give away the correct prediction without any reasoning.
- Hypothesis-only artifacts are class-specific.

Digging Deeper



- **Annotation Artifacts:** Clues which give away the correct prediction without any reasoning.
- Hypothesis-only artifacts are class-specific.
- Word-class association via PPMI:

Digging Deeper



- **Annotation Artifacts:** Clues which give away the correct prediction without any reasoning.
- Hypothesis-only artifacts are class-specific.
- Word-class association via PPMI:

$$\max\left\{0, \log \frac{p(\mathbf{word}, \mathbf{class})}{p(\cdot, \mathbf{class})p(\mathbf{word}, \cdot)}\right\}$$

Digging Deeper



- **Annotation Artifacts:** Clues which give away the correct prediction without any reasoning.
- Hypothesis-only artifacts are class-specific.
- Word-class association via PPMI:

$$\max\left\{0, \log \frac{p(\mathbf{word}, \mathbf{class})}{p(., \mathbf{class})p(\mathbf{word}, .)}\right\}$$

N	C	E
OUTDOORS	NOBODY	TALL
LEAST	SLEEPING	FIRST
INSTRUMENT	No	COMPETITION
OUTSIDE	Tv	SAD
ANIMAL	CAT	FAVORITE

Entailment Artifacts

Entailment Artifacts



Some men and boys
are playing frisbee in
a grassy area.

Premise

Generalization

People play
frisbee **outdoors.**

Hypothesis

Entailment Artifacts



Some men and boys are playing frisbee in a grassy area.

Premise

Generalization

People play frisbee **outdoors**.

Hypothesis



A person in a red **shirt** is mowing the grass with a **green** riding mower.

Premise

Simplification

A person in red is cutting the grass on a riding mower.

Hypothesis

Neutral Artifacts

Neutral Artifacts



Modifiers

A man is doing work on a **black** Amtrak train.

Hypothesis

A middle-aged man works under the engine of a train on rail tracks.

Premise

Neutral Artifacts



Modifiers

A man is doing work on a **black** Amtrak train.

Hypothesis

A middle-aged man works under the engine of a train on rail tracks.

Premise



Purpose Clauses

A group of female athletes are huddled together and excited.

Premise

They are huddled together **because** they are working together.

Hypothesis

Contradiction Artifacts

Contradiction Artifacts



Negation

Nobody
wears a cap.

Hypothesis

Older man with white hair and a red cap painting the golden gate bridge on the shore with the golden gate bridge in the distance.

Premise

Contradiction Artifacts



Negation

Nobody
wears a cap.

Hypothesis

Older man with white hair and a red cap painting the golden gate bridge on the shore with the golden gate bridge in the distance.

Premise



Cats!

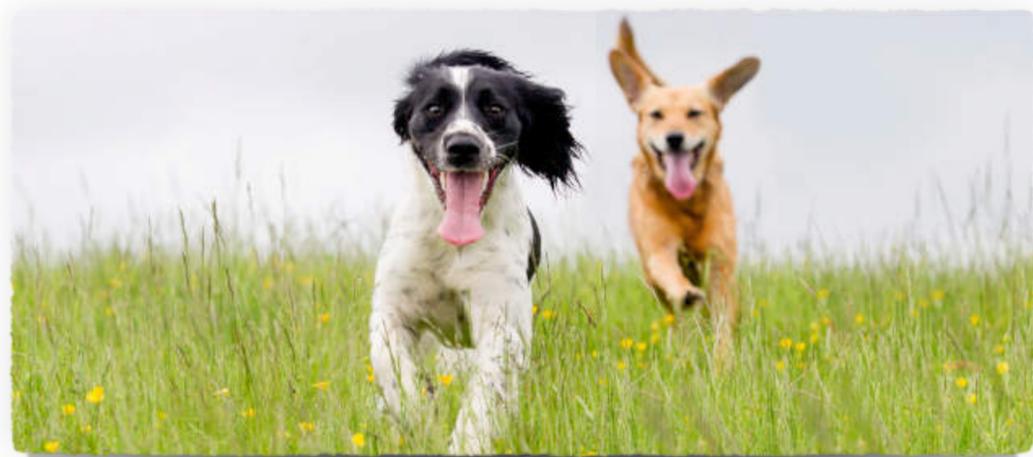
Three **cats** race on a track.

Hypothesis

Three dogs racing on racetrack.

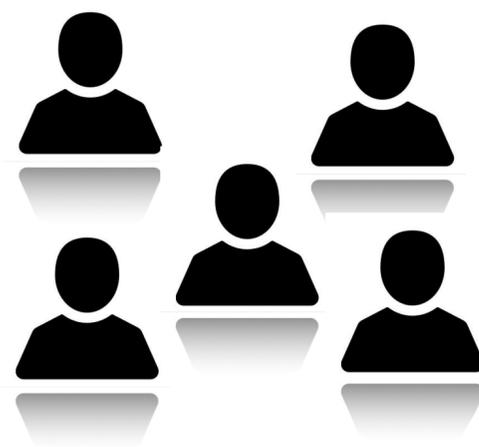
Premise

A possible explanation



Two dogs are running through a field.

Premise



Entailment

Neutral

Contradiction

There are **animals** outdoors.

Some puppies are running **to catch a stick**.

The pets are sitting on a couch.

A possible explanation



Two dogs are running through a field.

Premise



Entailment

Neutral

Contradiction

There are **animals** outdoors.

Some puppies are running **to catch a stick**.

The pets are sitting on a couch.

Are seed examples responsible?

Premise

A woman selling bamboo sticks talking to two men on a loading dock.

Entailment

There are at least three **people** on a loading dock.

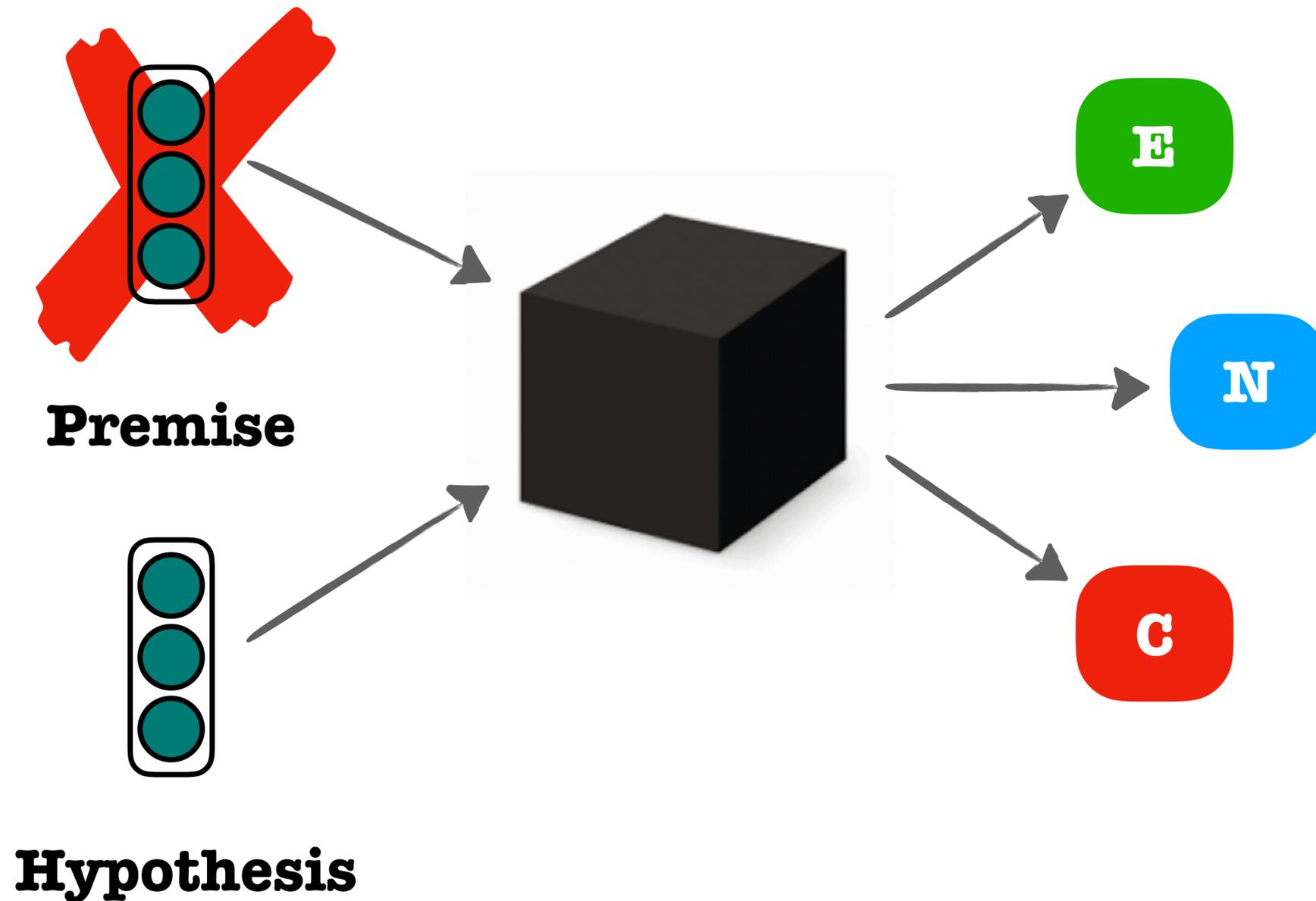
Neutral

A woman is selling bamboo sticks **to** help provide for her family

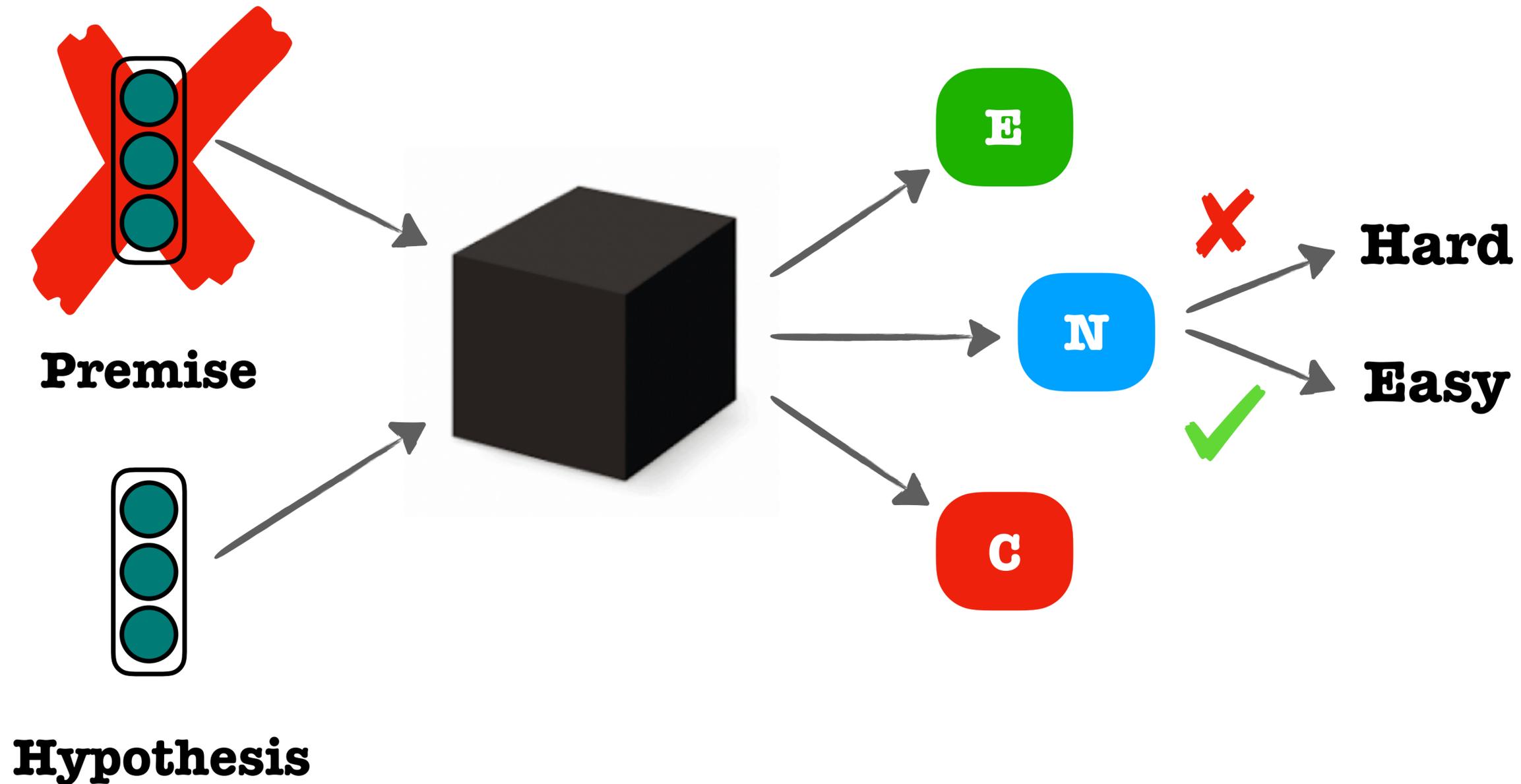
Contradiction

A woman is **not** taking money for any of her sticks.

Identifying examples with artifacts

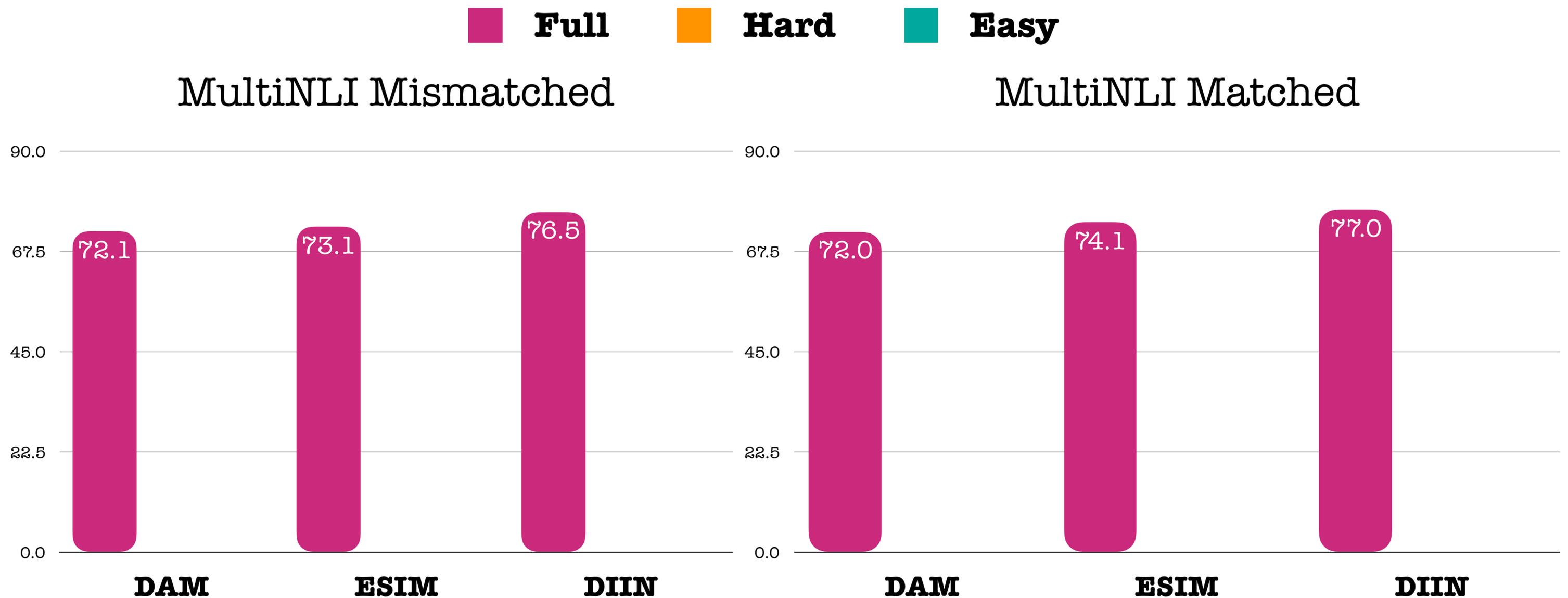


Identifying examples with artifacts



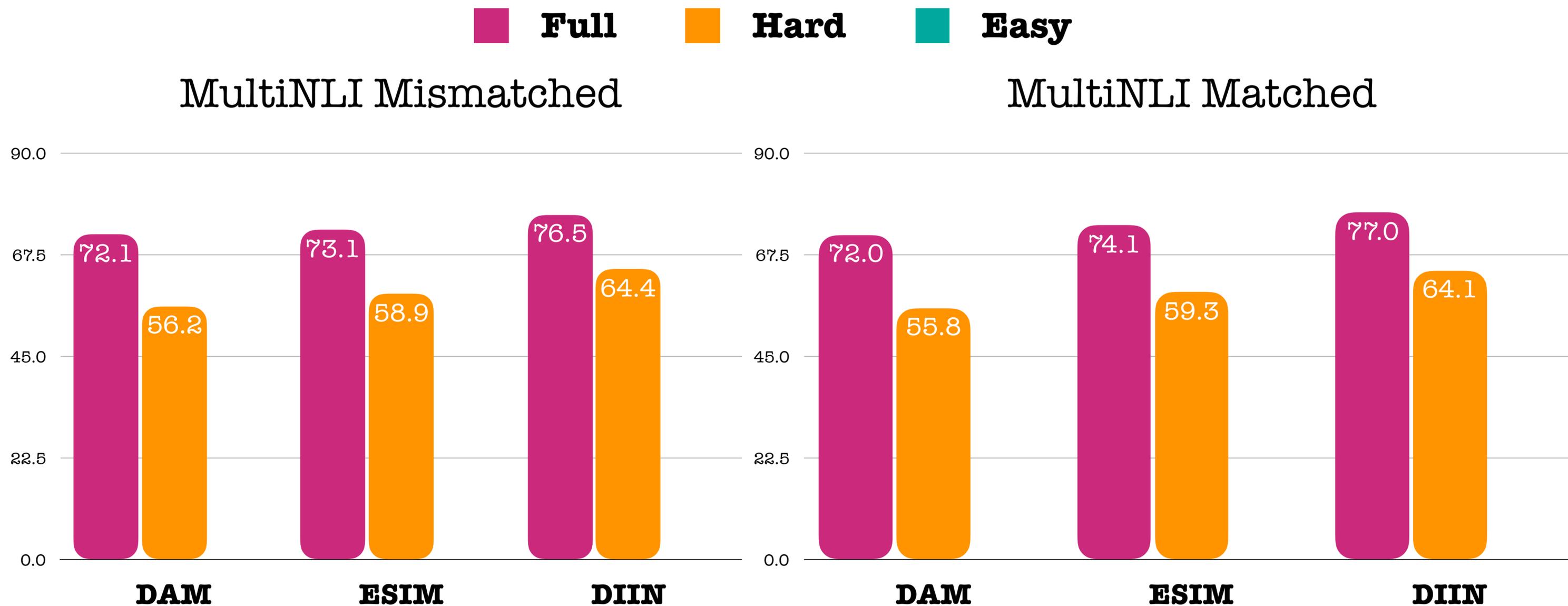
Revisiting NLI models

Revisiting NLI models



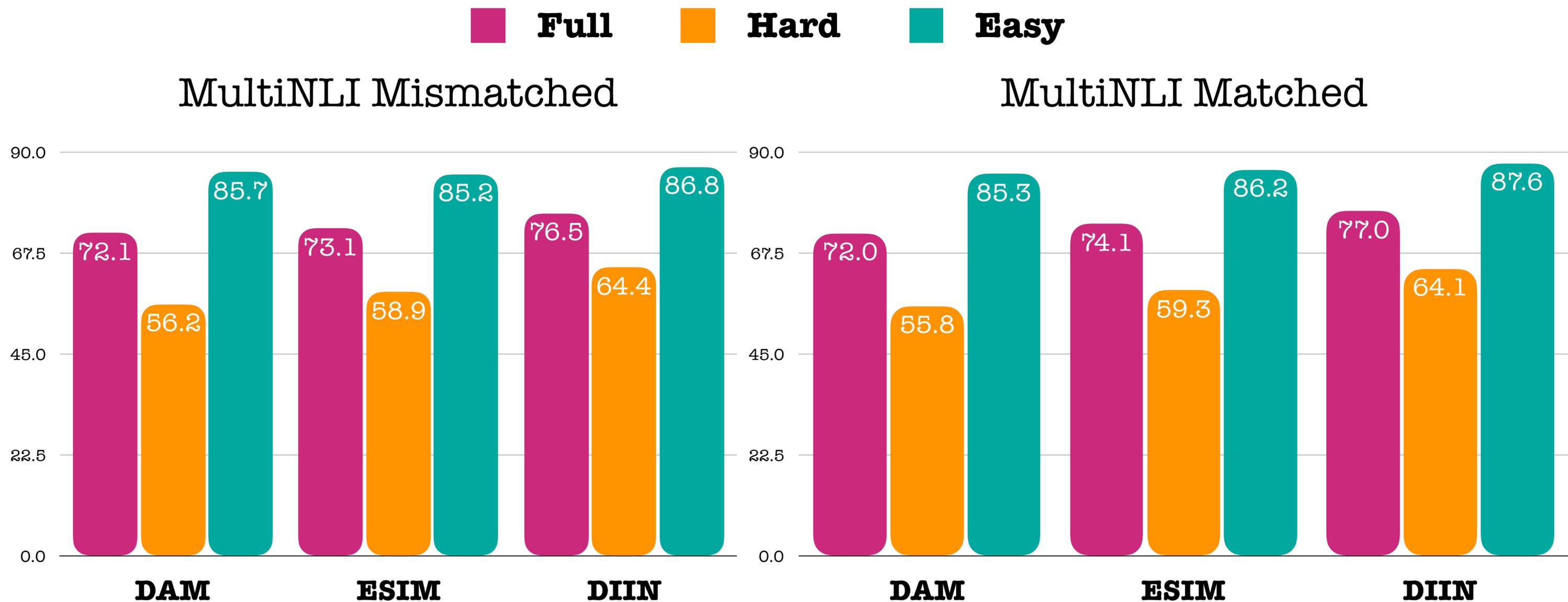
DAM - Decomposable Attention Model [Parikh et. al. 2016]
ESIM - Enhanced Sequential Inference Model [Chen et. al., 2017]
DIIN - Densely Interactive Inference Network [Gong et. al. 2018]

Revisiting NLI models



DAM - Decomposable Attention Model [Parikh et. al. 2016]
ESIM - Enhanced Sequential Inference Model [Chen et. al., 2017]
DIIN - Densely Interactive Inference Network [Gong et. al. 2018]

Revisiting NLI models

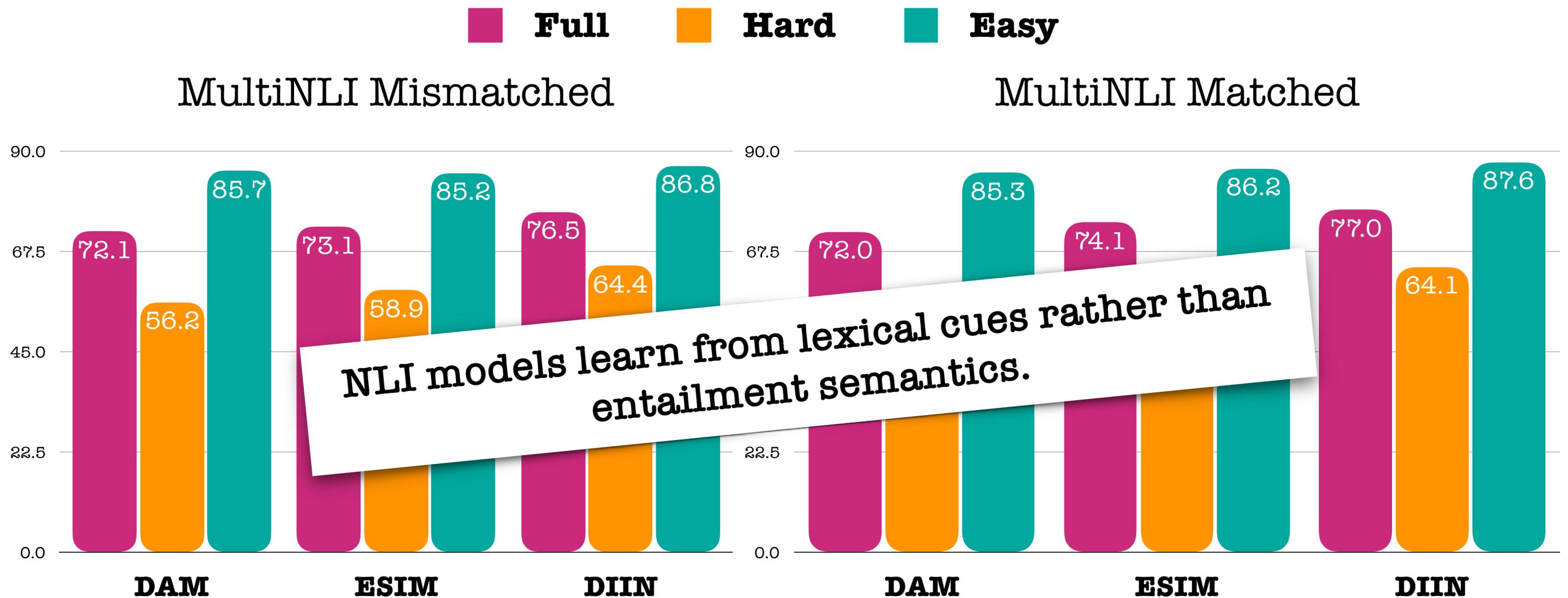


DAM - Decomposable Attention Model [Parikh et. al. 2016]

ESIM - Enhanced Sequential Inference Model [Chen et. al., 2017]

DIIN - Densely Interactive Inference Network [Gong et. al. 2018]

Revisiting NLI models



DAM - Decomposable Attention Model [Parikh et. al. 2016]
ESIM - Enhanced Sequential Inference Model [Chen et. al., 2017]
DIIN - Densely Interactive Inference Network [Gong et. al. 2018]

Not unique to NLI...

Not unique to NLI...

The logo for SQuAD, featuring the text "SQuAD" in a white, sans-serif font on a red rectangular background.

The Stanford Question Answering Dataset

Jia & Liang et al., 2017

Not unique to NLI...

Story Cloze Test and ROCStories Corpora

Schwartz et al., 2017;
Cai et al., 2017

The logo for SQuAD (Stanford Question Answering Dataset) features the word "SQuAD" in a white, serif font against a red rectangular background.

The Stanford Question Answering Dataset

Jia & Liang et al., 2017

Not unique to NLI...

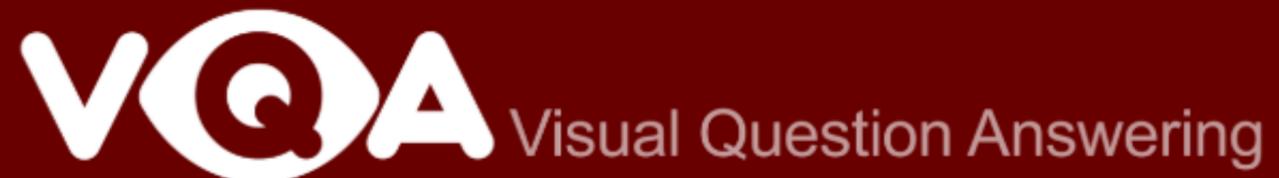
Story Cloze Test and ROCStories Corpora

Schwartz et al., 2017;
Cai et al., 2017

The logo for SQuAD (The Stanford Question Answering Dataset) features the word "SQuAD" in a white, sans-serif font with a thin outline, set against a solid red rectangular background.

The Stanford Question Answering Dataset

Jia & Liang et al., 2017

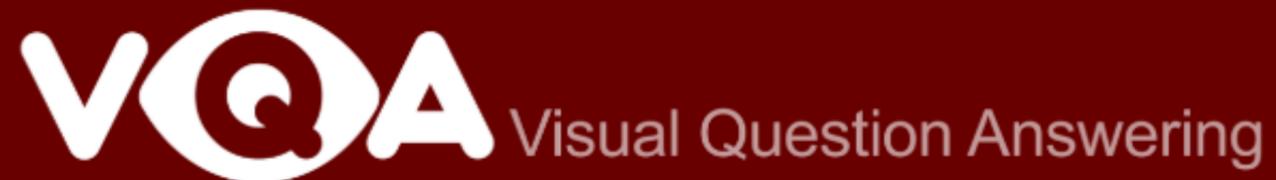
The logo for VQA (Visual Question Answering) features the letters "VQA" in a large, white, bold, sans-serif font. The letter "Q" is stylized with a white circle around its center. To the right of "VQA", the words "Visual Question Answering" are written in a smaller, white, sans-serif font, all set against a solid dark red rectangular background.

Jabri et al., 2016

Not unique to NLI...

Story Cloze Test and ROCStories Corpora

Schwartz et al., 2017;
Cai et al., 2017

The logo for Visual Question Answering (VQA) features the letters 'VQA' in a large, bold, white font. The 'Q' is stylized with a white circle around it. To the right of 'VQA', the words 'Visual Question Answering' are written in a smaller, white, sans-serif font. The entire logo is set against a dark red background.

VQA Visual Question Answering

Jabri et al., 2016

The logo for SQuAD (Stanford Question Answering Dataset) features the word 'SQuAD' in a large, white, sans-serif font. The 'Q' is stylized with a white circle around it. The logo is set against a red background.

SQuAD

The Stanford Question Answering Dataset

Jia & Liang et al., 2017

cnn_dailymail

Chen et al., 2017

Looking Ahead: Improved Data Collection



Looking Ahead: Improved Data Collection



- Partial input baselines. E.g.

Looking Ahead: Improved Data Collection



- Partial input baselines. E.g.
 - SWAG [Zellers et. al., 2018],
 - DROP [Dua et. al., 2019],
 - Diverse NLI [Poliak et. al., 2018]

Looking Ahead: Improved Data Collection



- Partial input baselines. E.g.
 - SWAG [Zellers et. al., 2018],
 - DROP [Dua et. al., 2019],
 - Diverse NLI [Poliak et. al., 2018]
- Alternatives to human elicitation for building datasets?

Looking Ahead: Improved Data Collection



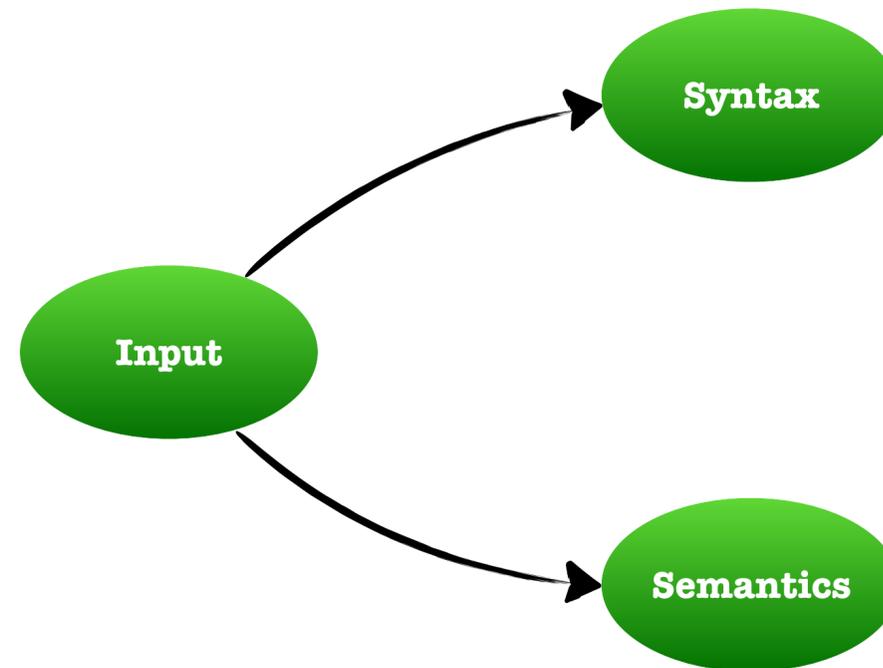
- Partial input baselines. E.g.
 - SWAG [Zellers et. al., 2018],
 - DROP [Dua et. al., 2019],
 - Diverse NLI [Poliak et. al., 2018]
- Alternatives to human elicitation for building datasets?



Looking Ahead: Improved Data Collection



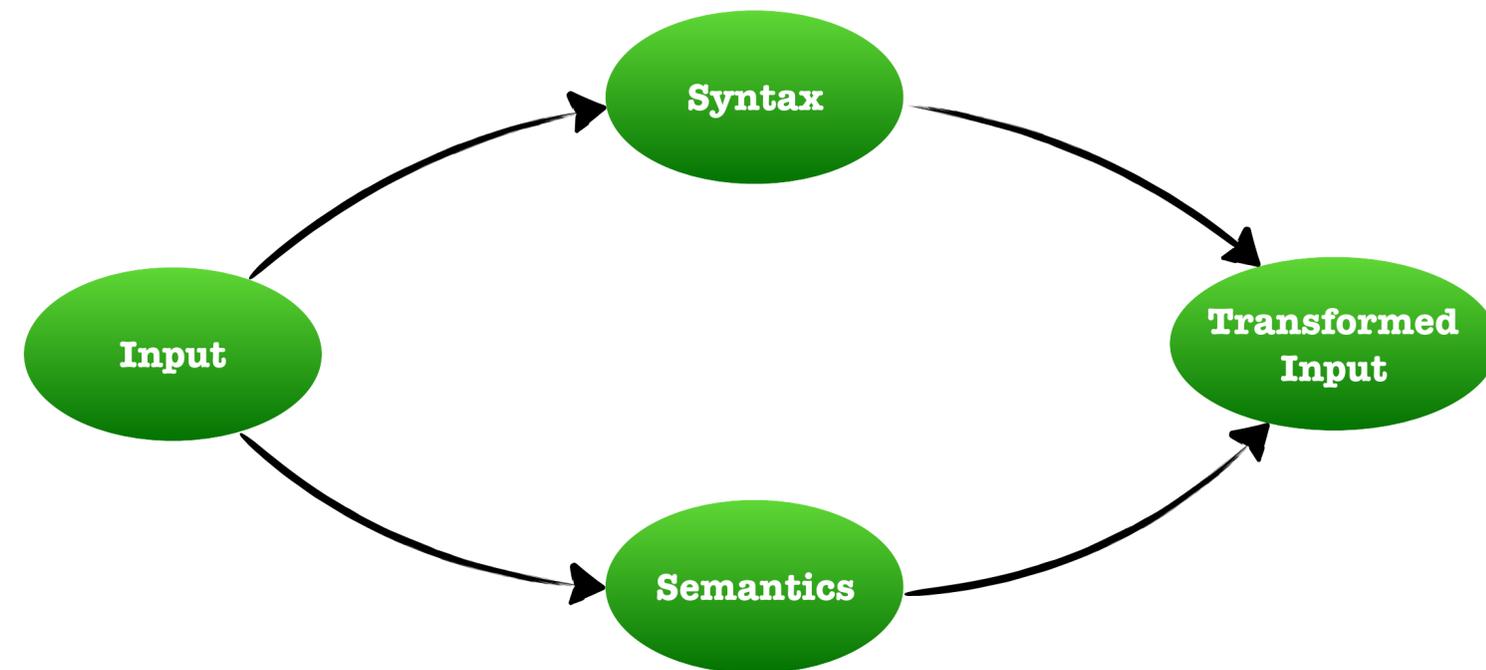
- Partial input baselines. E.g.
 - SWAG [Zellers et. al., 2018],
 - DROP [Dua et. al., 2019],
 - Diverse NLI [Poliak et. al., 2018]
- Alternatives to human elicitation for building datasets?



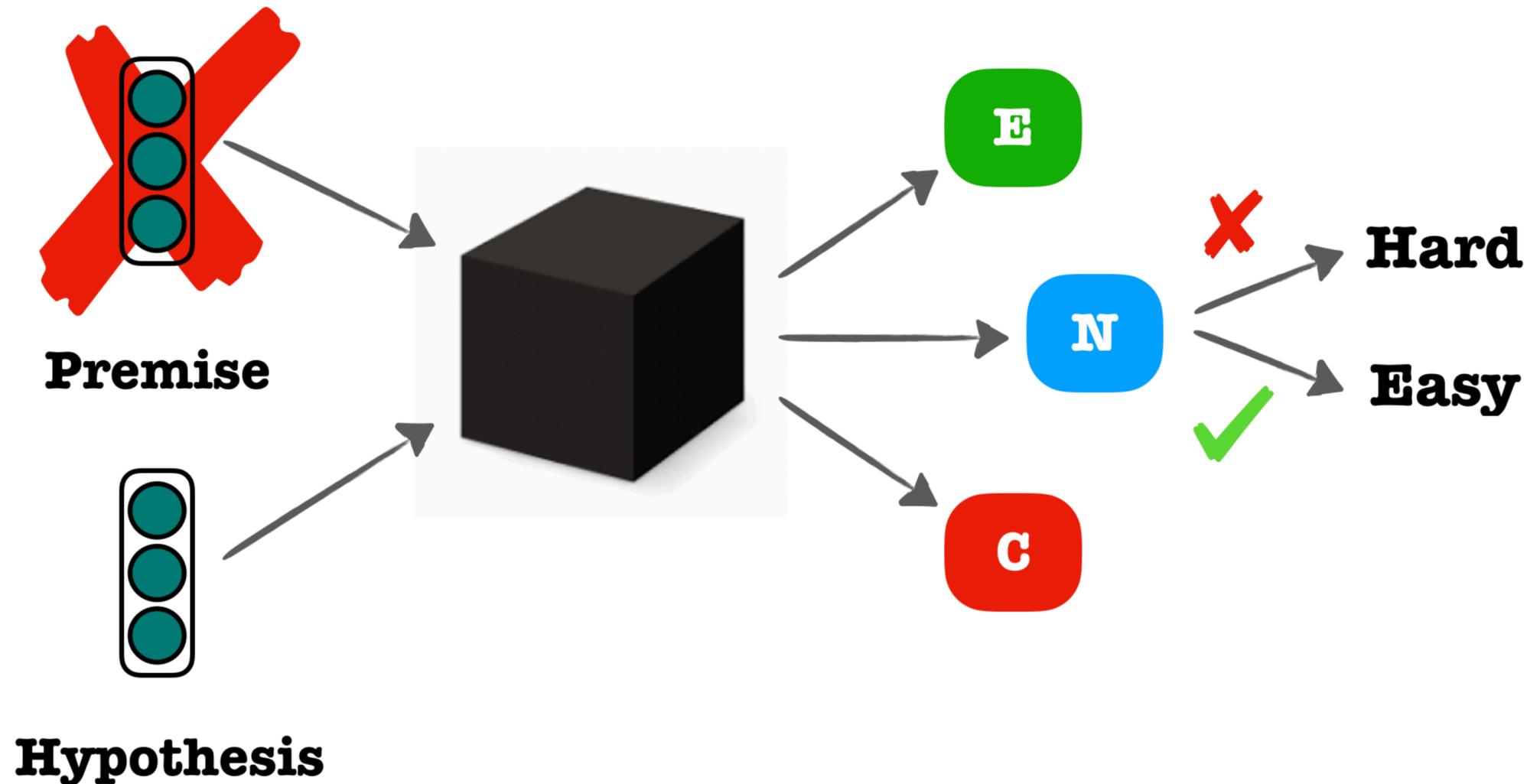
Looking Ahead: Improved Data Collection



- Partial input baselines. E.g.
 - SWAG [Zellers et. al., 2018],
 - DROP [Dua et. al., 2019],
 - Diverse NLI [Poliak et. al., 2018]
- Alternatives to human elicitation for building datasets?

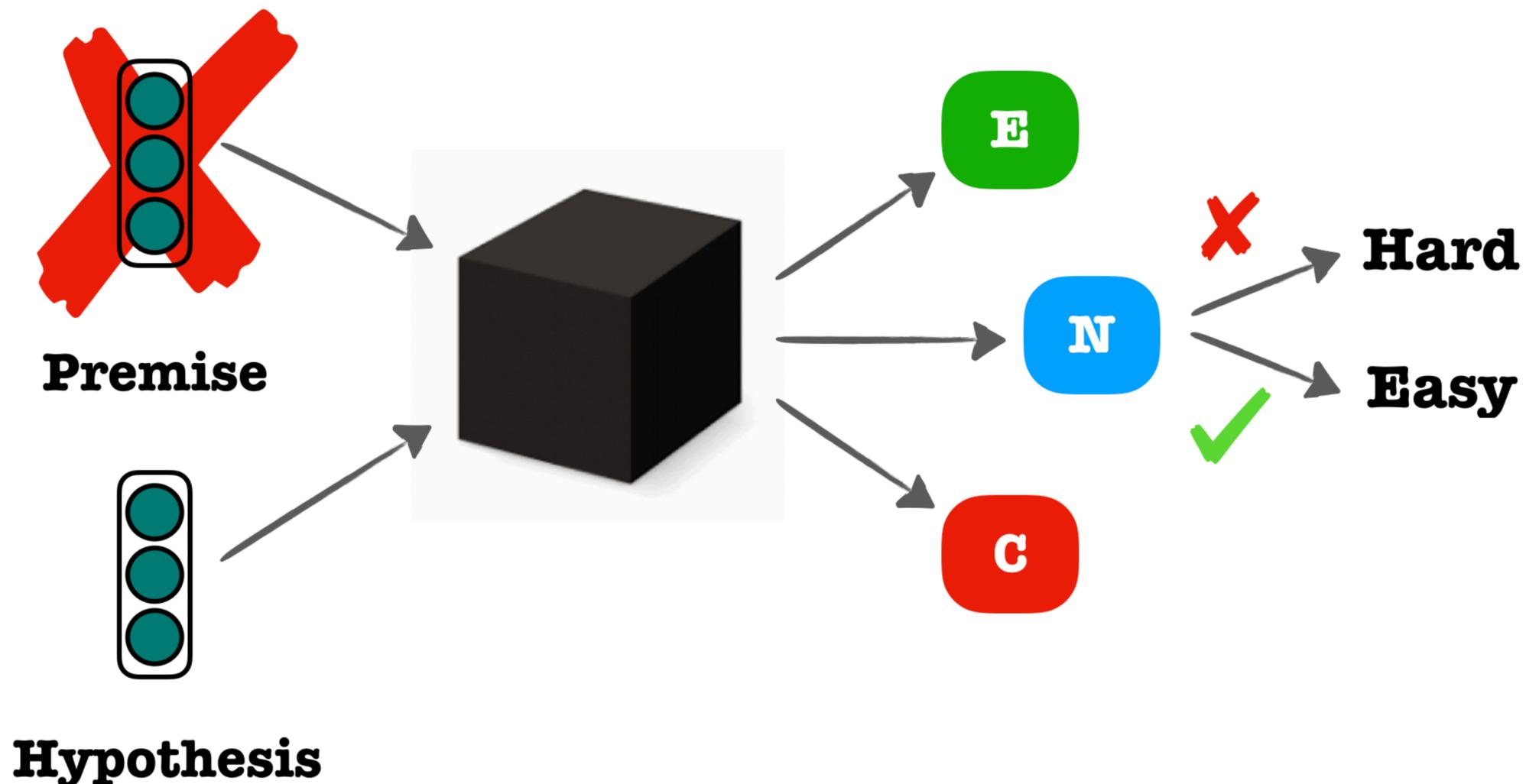


Filtering hypothesis-only artifacts



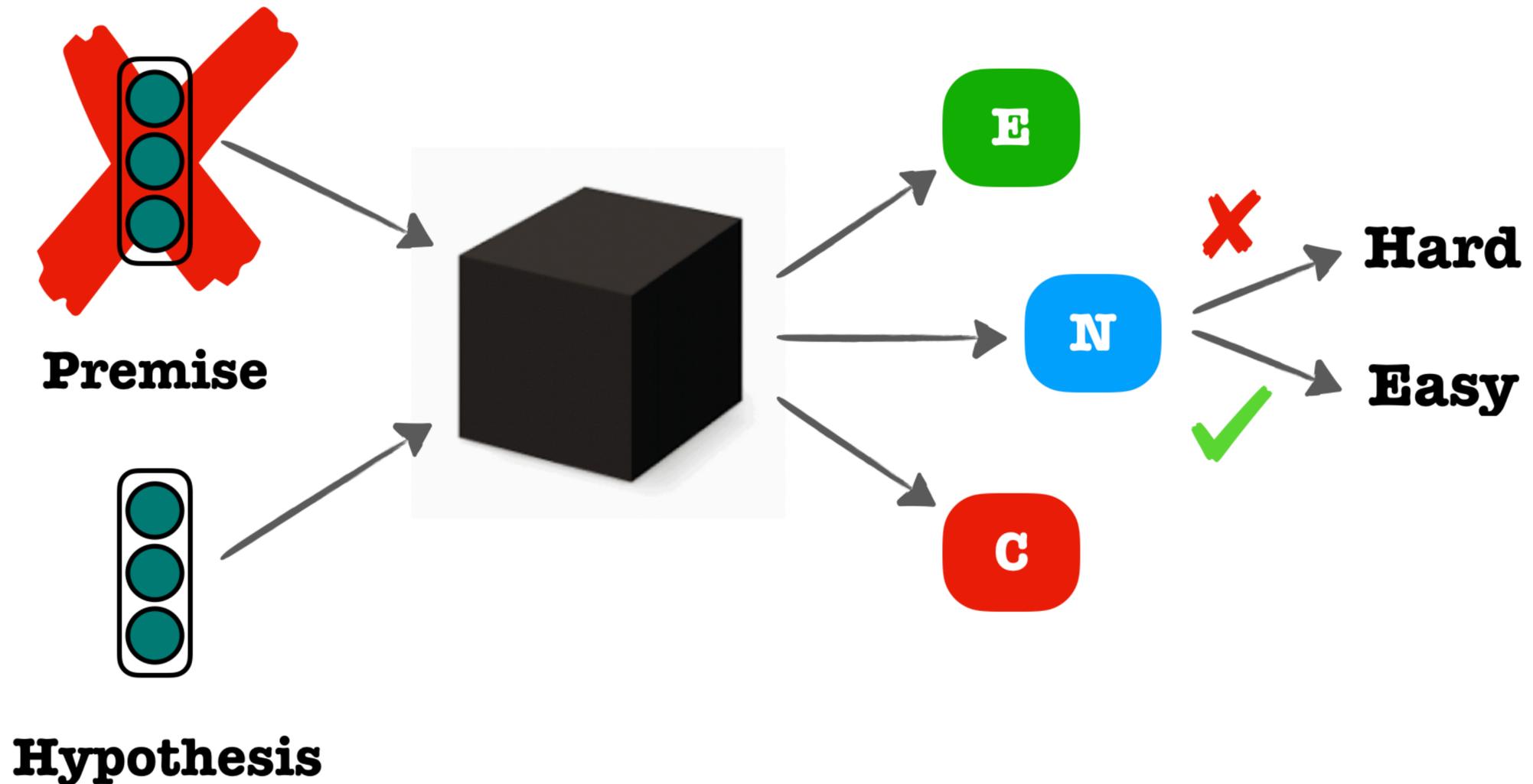
Filtering hypothesis-only artifacts

- Other kinds of artifacts. For e.g. shortening the premise.



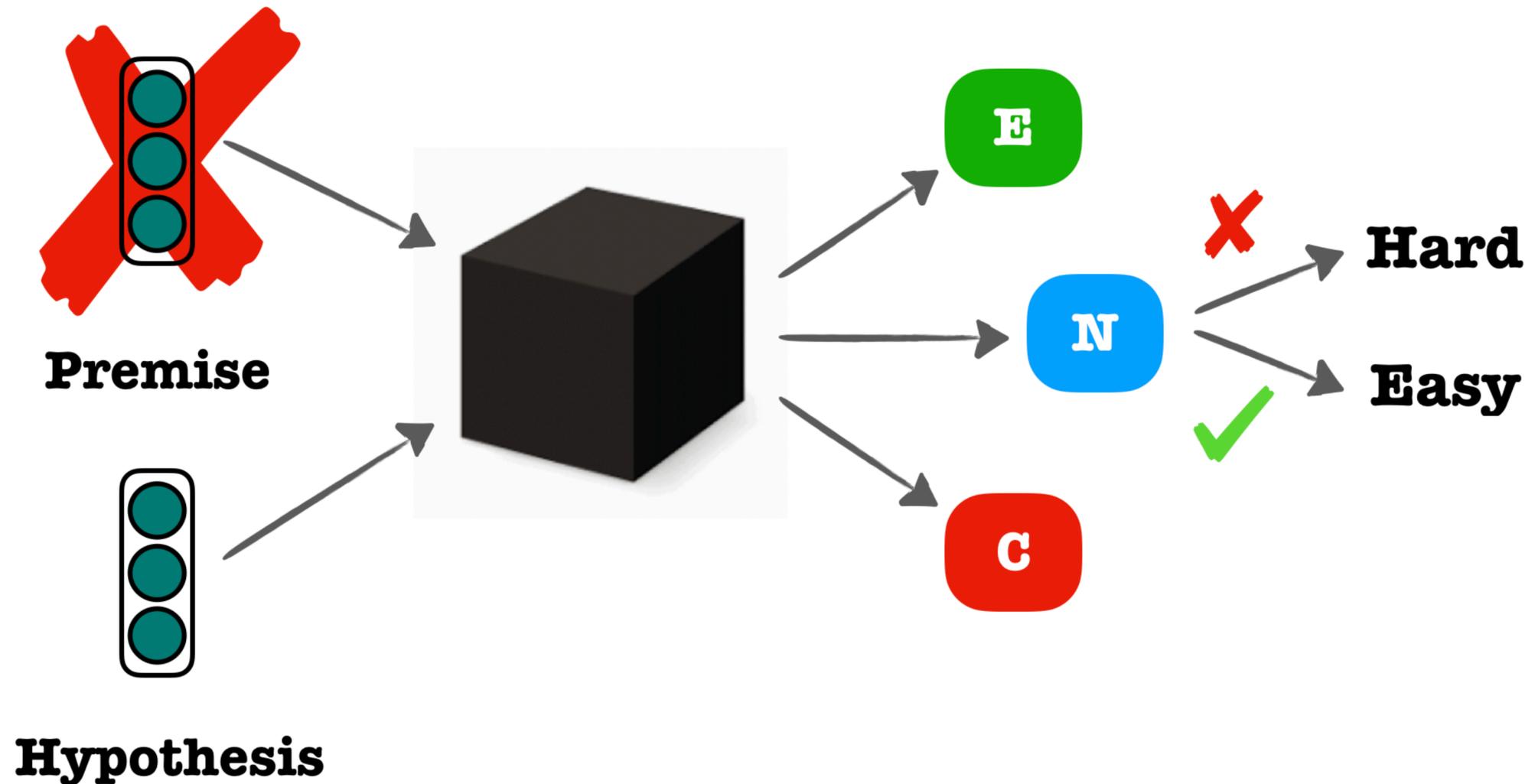
Filtering hypothesis-only artifacts

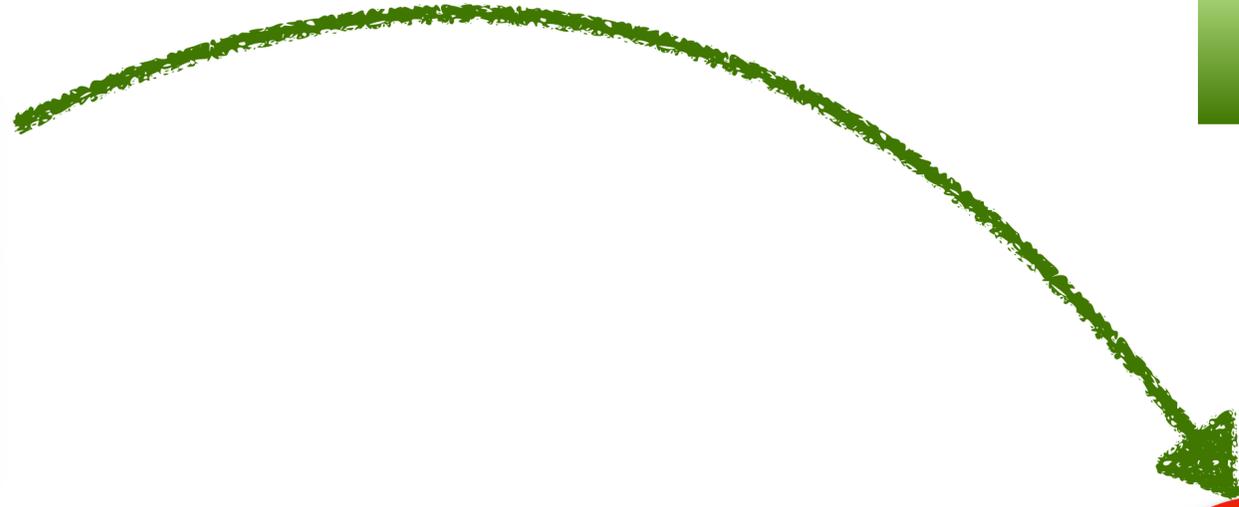
- Other kinds of artifacts. For e.g. shortening the premise.
- Examples with artifacts are still valid examples...



Filtering hypothesis-only artifacts

- Other kinds of artifacts. For e.g. shortening the premise.
- Examples with artifacts are still valid examples...
- Hard examples exhibit their own artifacts!





What do predictions tell us about the data?



Annotation Artifacts Abound!



Question #2



What do predictions tell us about the data?

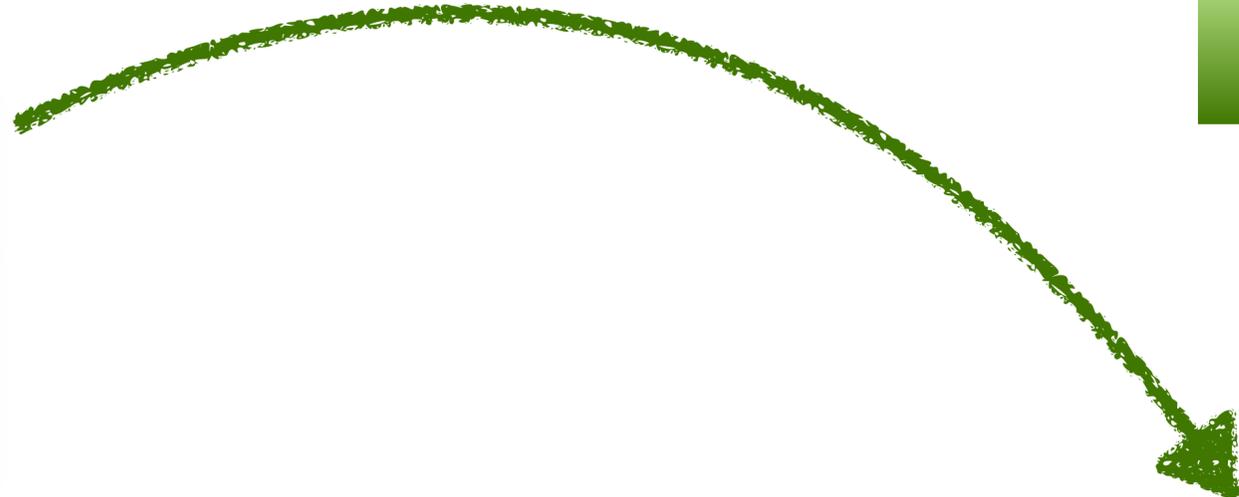


Annotation Artifacts Abound!

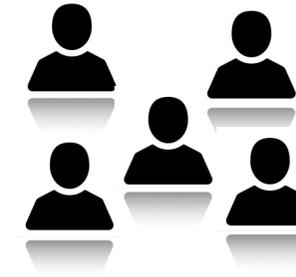


How can we use this information to spruce up our datasets?

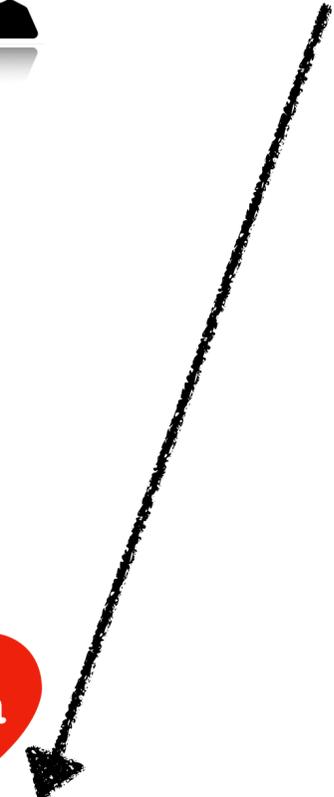
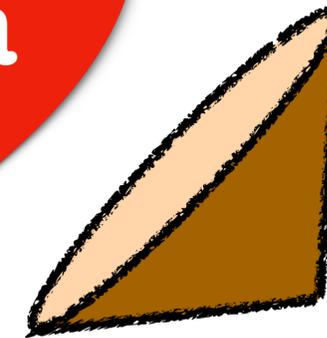
Question #2



What do predictions tell us about the data?

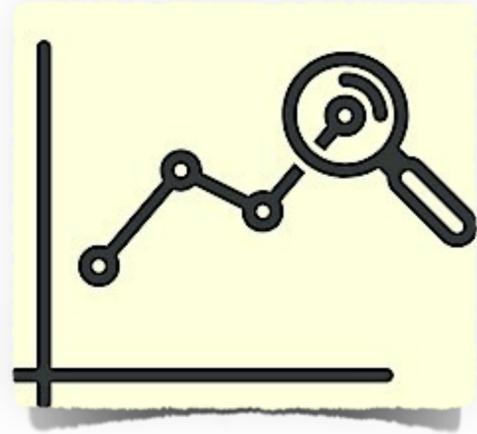


Annotation Artifacts Abound!



How can we use this information to spruce up our datasets?

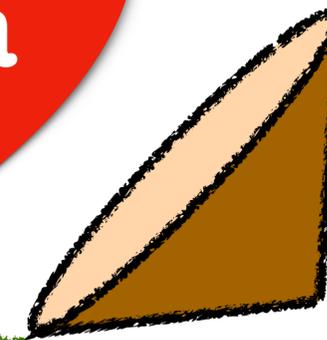
Question #2



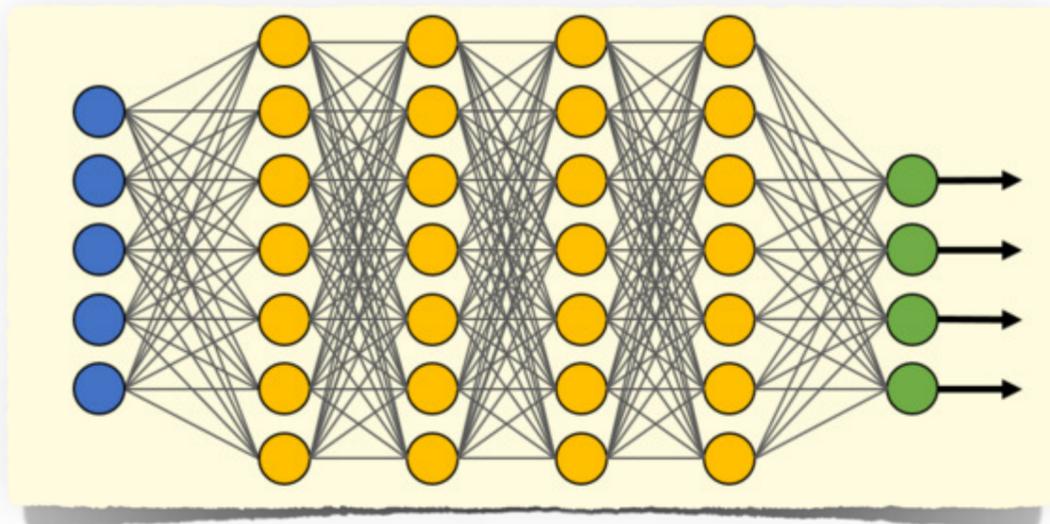
What do predictions tell us about the data?



Annotation Artifacts Abound!



How can we use this information to spruce up our datasets?



Goal

Goal

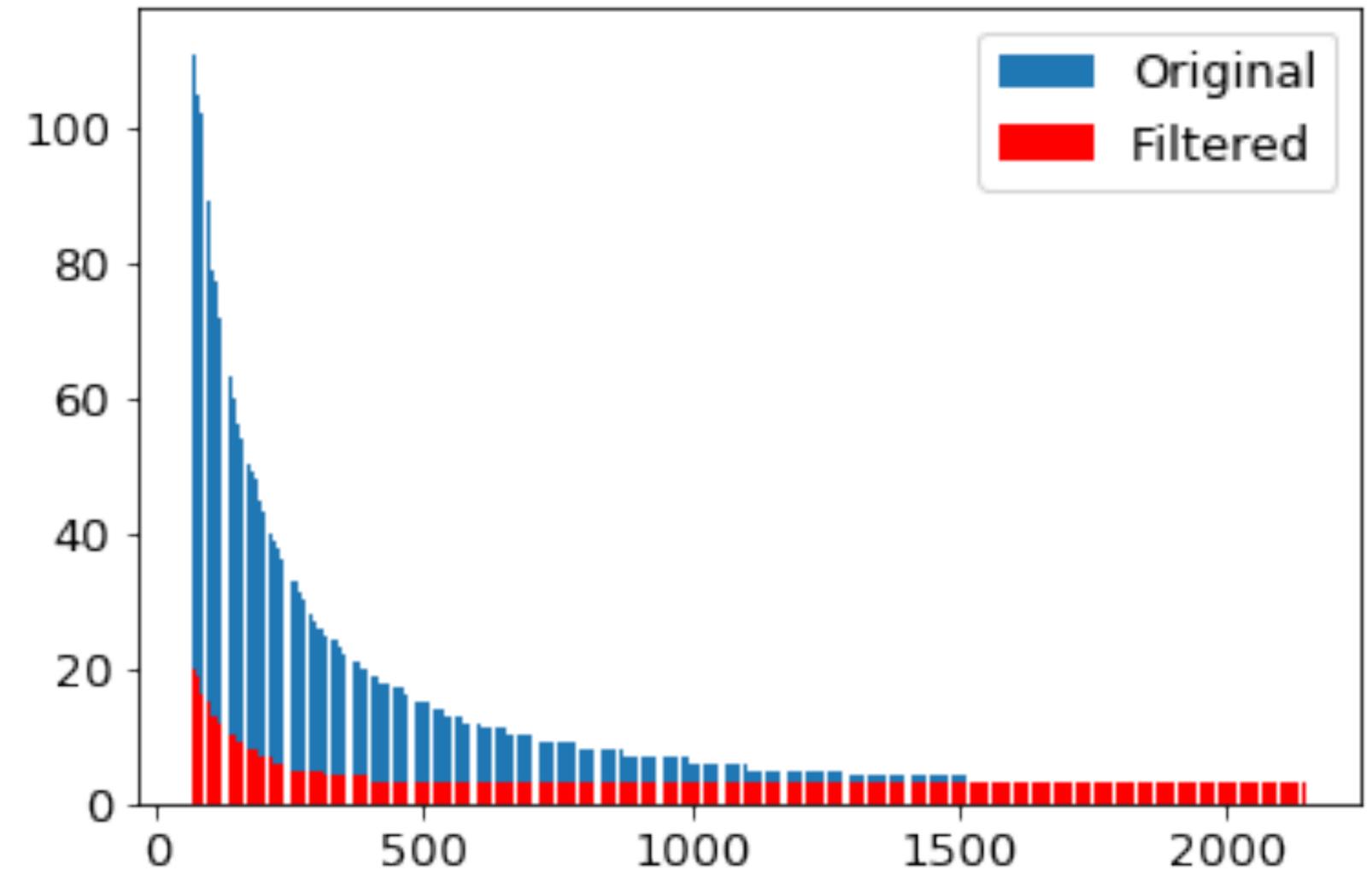
- **Balance out** the occurrence of different phenomena in the dataset
 - Filter out a majority of the samples which exhibit artifacts.

Goal

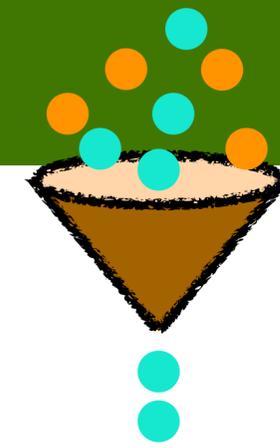
- **Balance out** the occurrence of different phenomena in the dataset
 - Filter out a majority of the samples which exhibit artifacts.

Goal

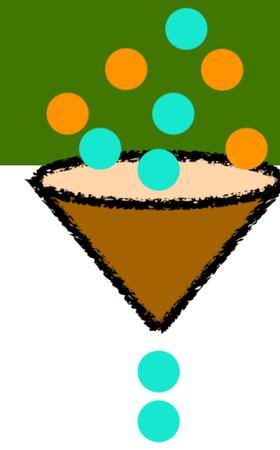
- **Balance out** the occurrence of different phenomena in the dataset
 - Filter out a majority of the samples which exhibit artifacts.
- Avoid **head phenomena redundancy**



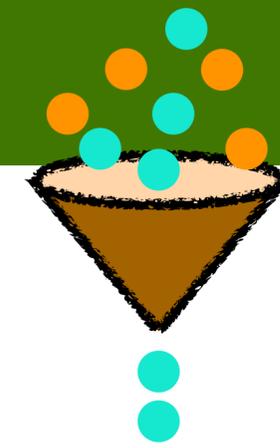
Insights



Insights

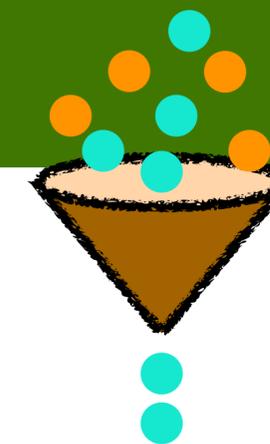


Insights



- i. The models that exploit artifacts, can be used to **detect artifacts!** Better than manual identification... (e.g. hypothesis-only artifacts)

Insights



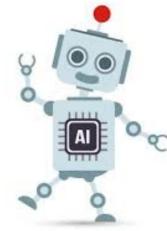
- i. The models that exploit artifacts, can be used to **detect artifacts!** Better than manual identification... (e.g. hypothesis-only artifacts)
- ii. Examples with artifacts can be classified correctly by **multiple** models.



How predictable is a sample?

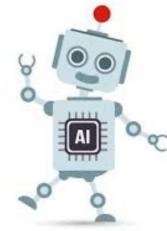
How predictable is a sample?

- Can it be predicted by a simple model?



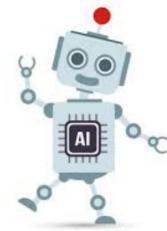
How predictable is a sample?

- Can it be predicted by a simple model?
- How much training does it take to predict it?



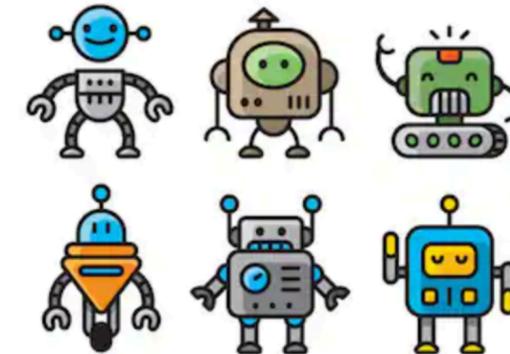
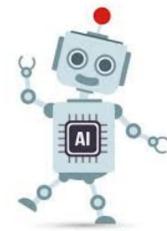
How predictable is a sample?

- Can it be predicted by a simple model?
- How much training does it take to predict it?
- How confident is the model?



How predictable is a sample?

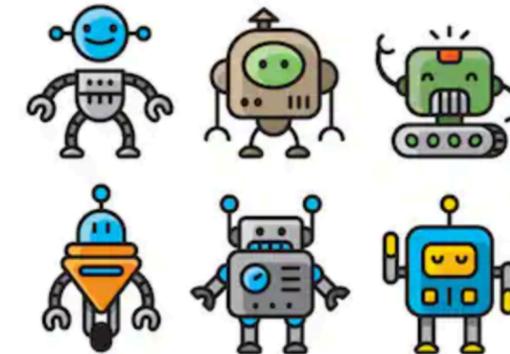
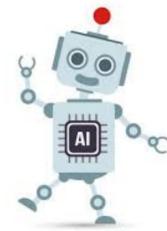
- Can it be predicted by a simple model?
- How much training does it take to predict it?
- How confident is the model?
- Can it be predicted by **several** simple models?



Ensembles of
Linear Classifiers

How predictable is a sample?

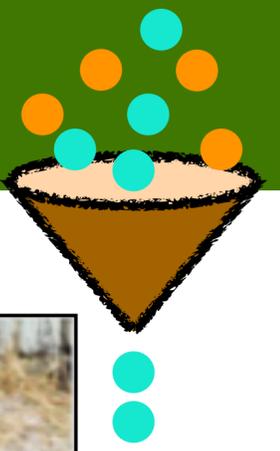
- Can it be predicted by a simple model?
- How much training does it take to predict it?
- How confident is the model?
- Can it be predicted by **several** simple models?



Predictability Score

**Ensembles of
Linear Classifiers**

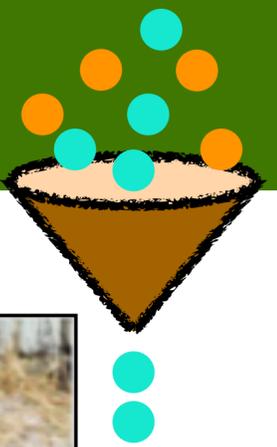
Predictability Scores



Ribiero et al., 2017



Predictability Scores



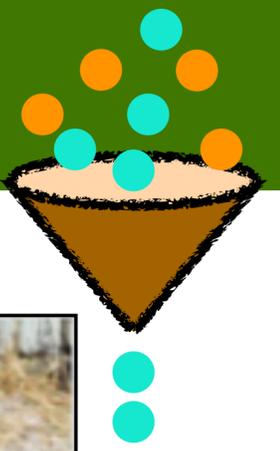
Ribiero et al., 2017



Iteration 1

0.99	0.93	0.97	0.87	0.89	0.78
------	------	------	------	------	------

Predictability Scores



Ribiero et al., 2017



Iteration 1

0.99

0.93

0.97

0.87

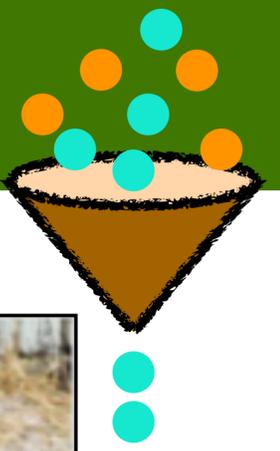
0.89

0.78

Iteration 2

0.99	0.93	0.97	0.87	0.89	0.78

Predictability Scores



Ribiero et al., 2017



Iteration 1

0.99

0.93

0.97

0.87

0.89

0.78

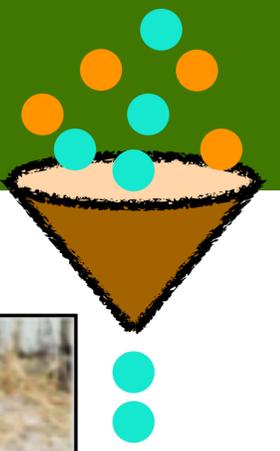
Iteration 2

0.83

0.99

0.54

Predictability Scores



Ribiero et al., 2017



Iteration 1

0.99

0.93

0.97

0.87

0.89

0.78

Iterative and Greedy!

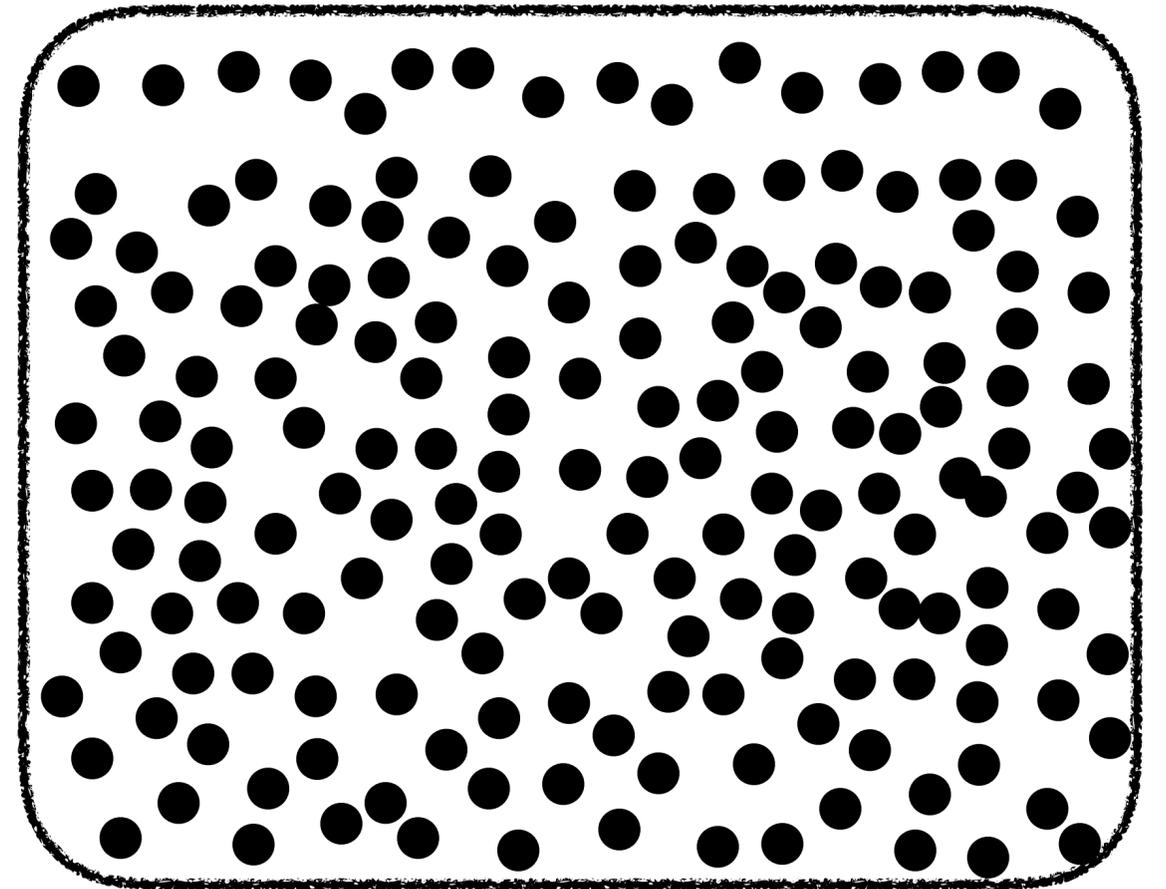
Iteration 2

0.83

0.99

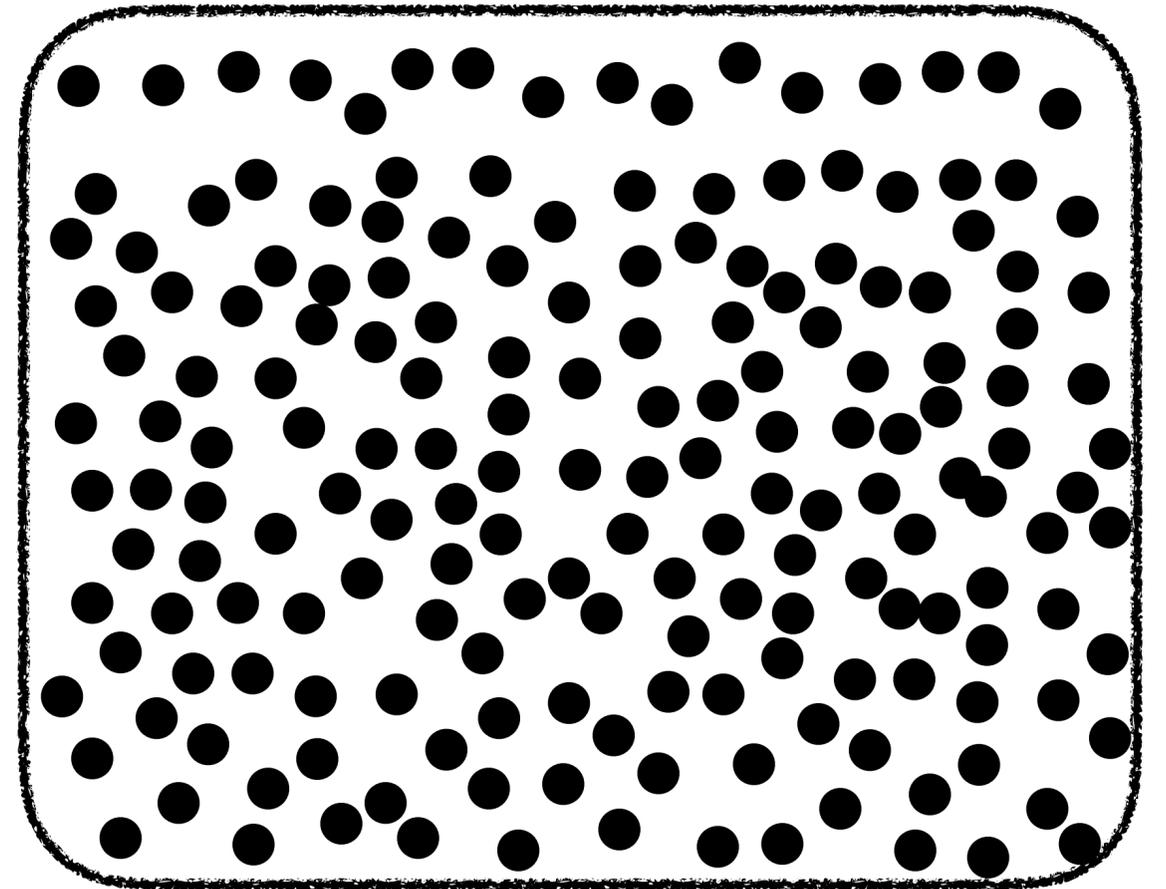
0.54

Algorithm



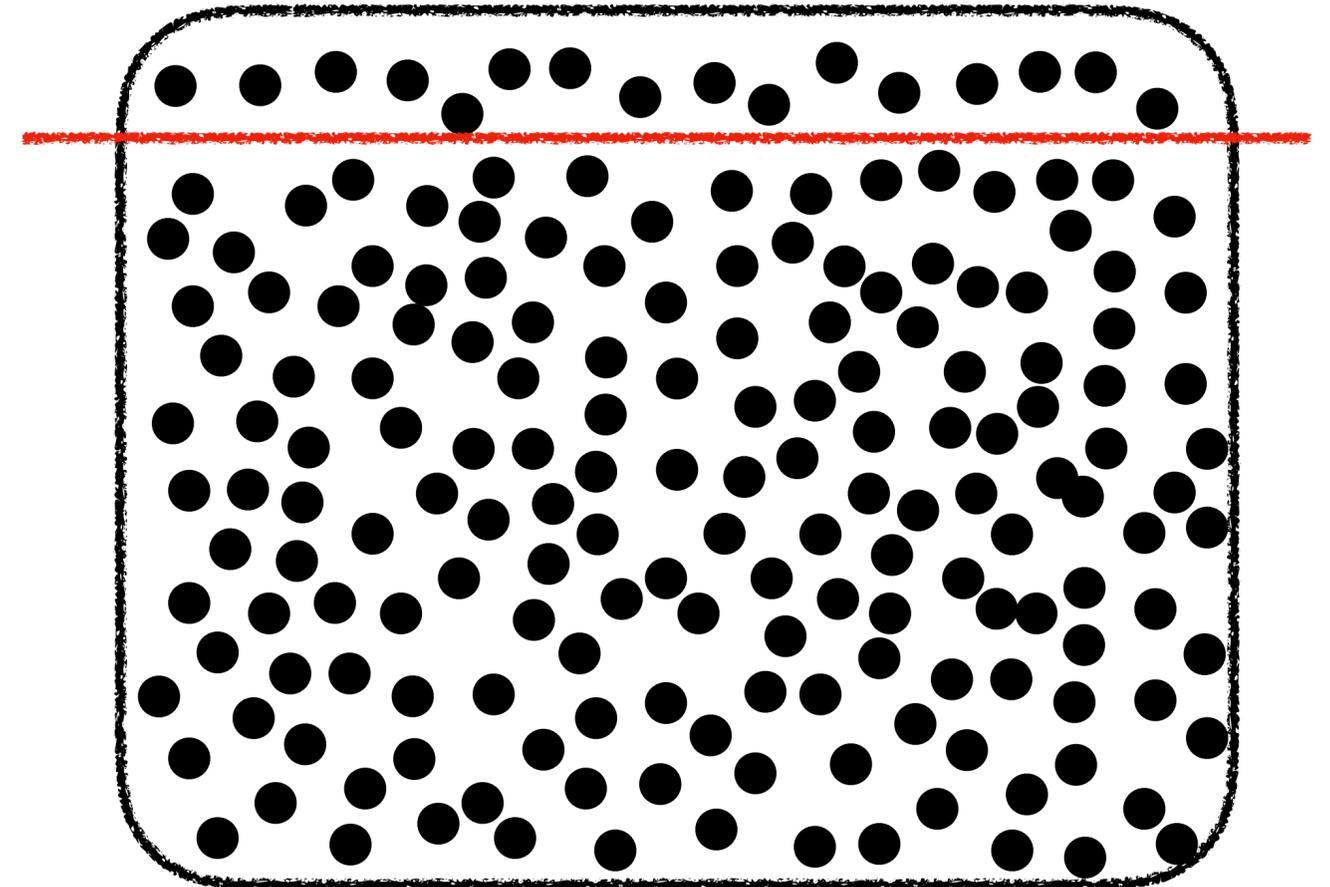
Algorithm

- Start with an initial feature representation, ϕ



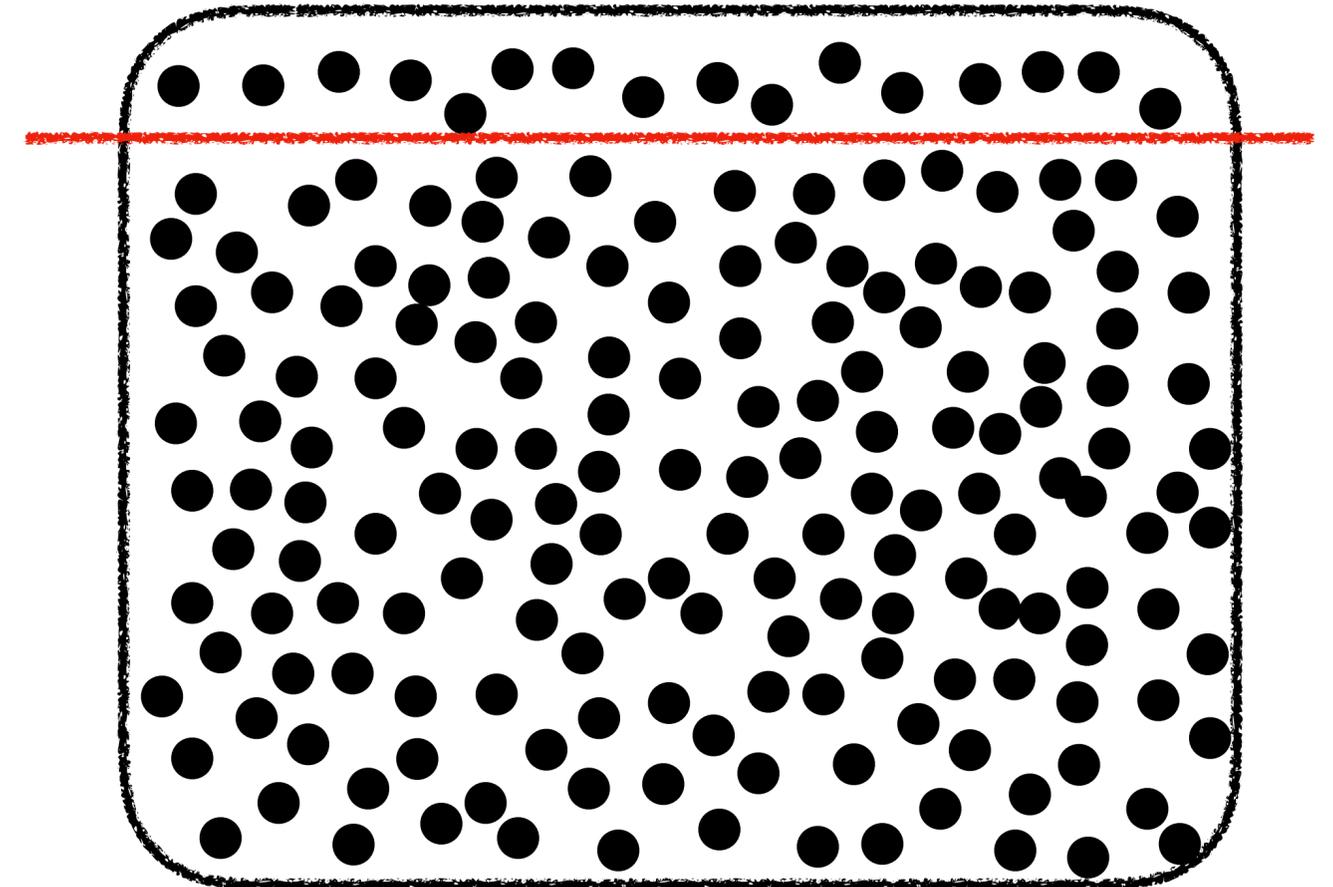
Algorithm

- Start with an initial feature representation, ϕ



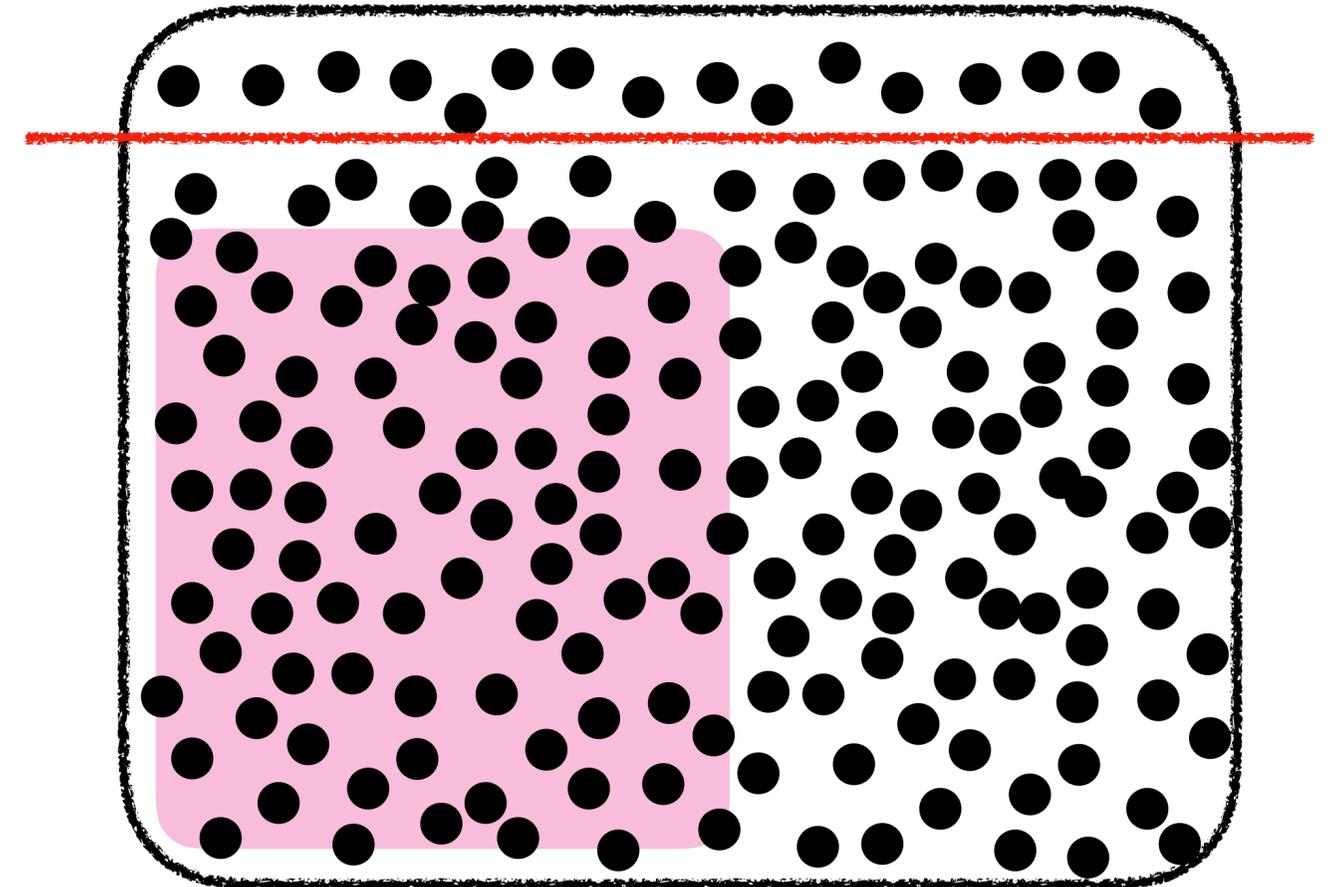
Algorithm

- Start with an initial feature representation, ϕ
- Train multiple models on random partitions of the remaining data.



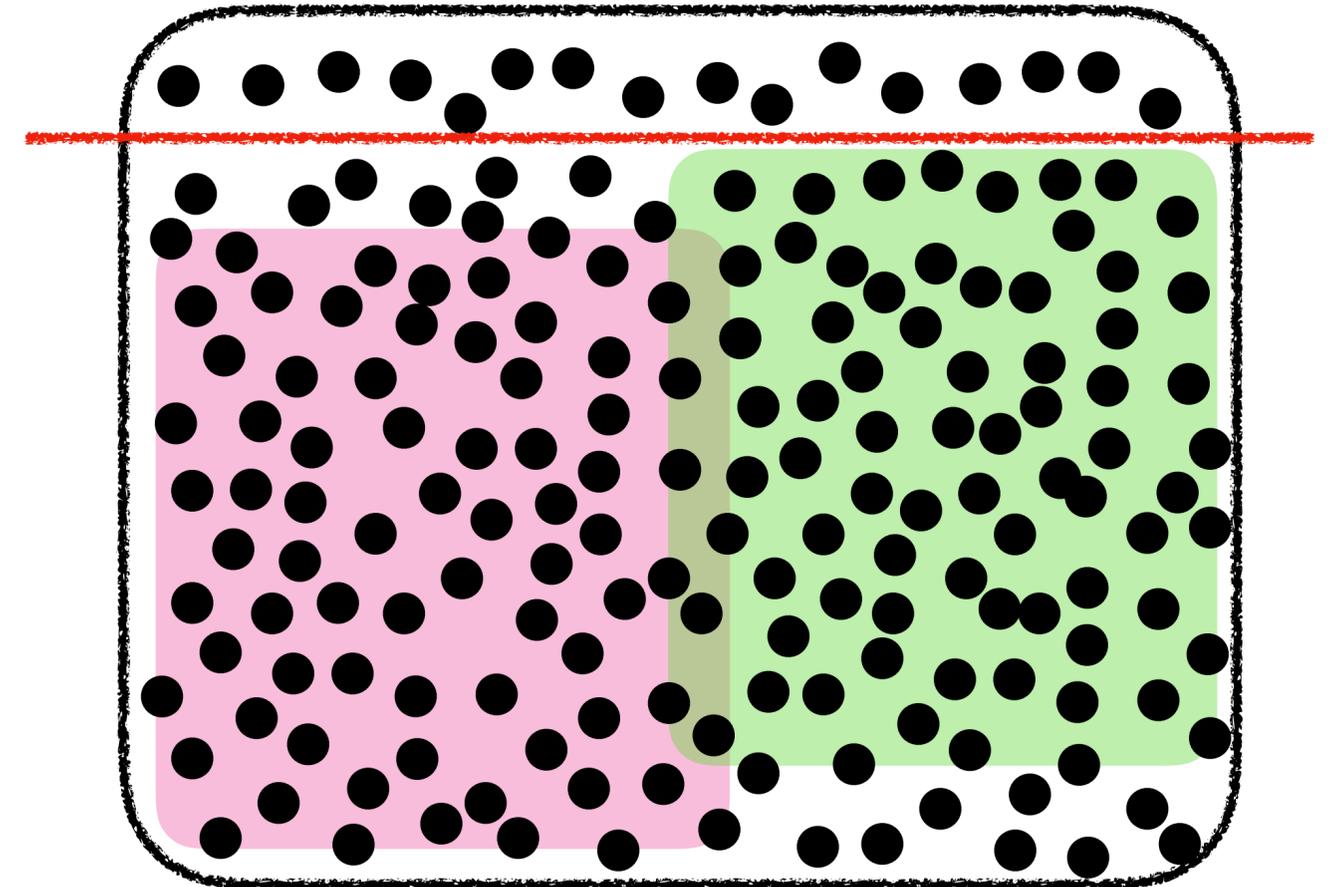
Algorithm

- Start with an initial feature representation, ϕ
- Train multiple models on random partitions of the remaining data.



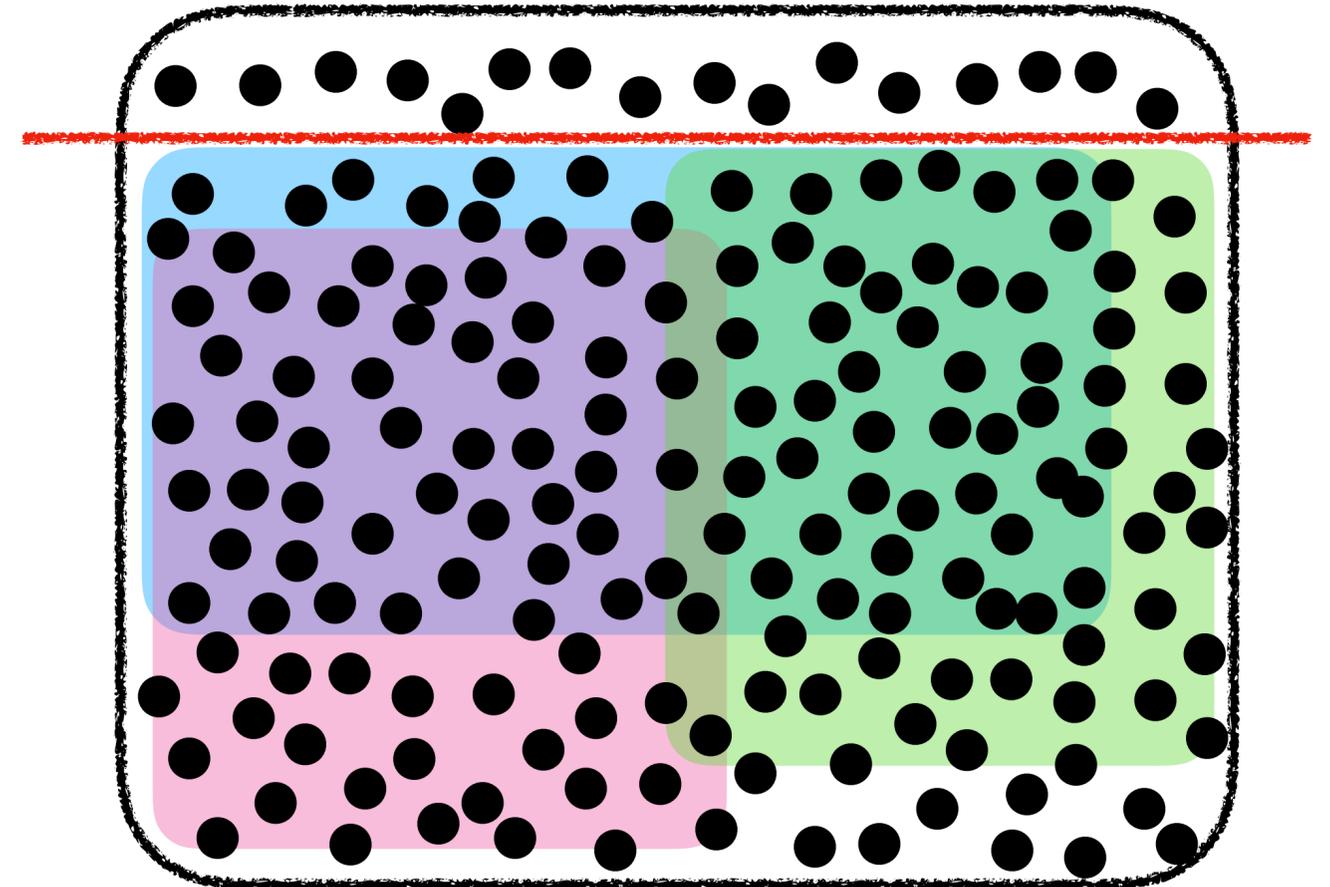
Algorithm

- Start with an initial feature representation, ϕ
- Train multiple models on random partitions of the remaining data.



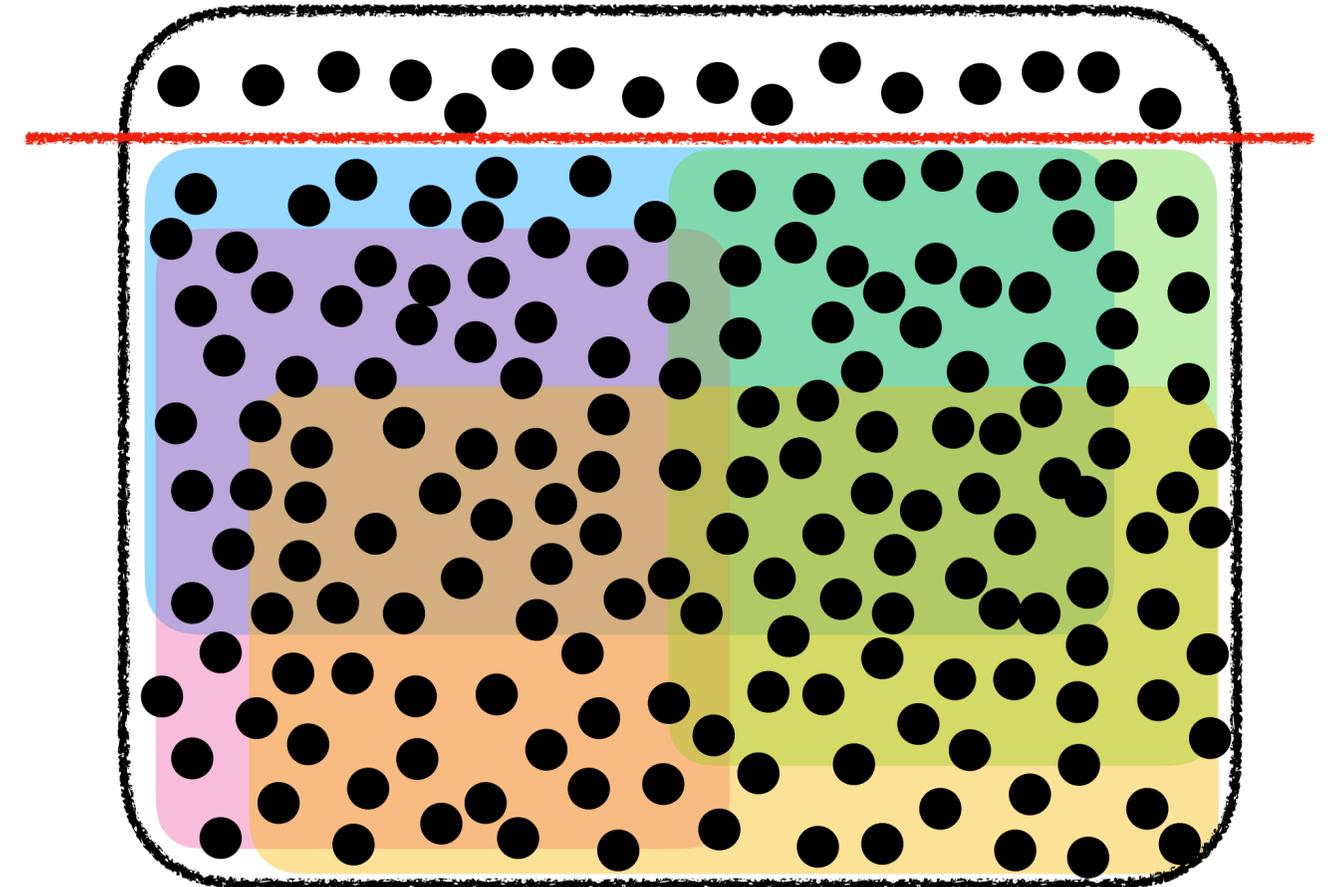
Algorithm

- Start with an initial feature representation, ϕ
- Train multiple models on random partitions of the remaining data.



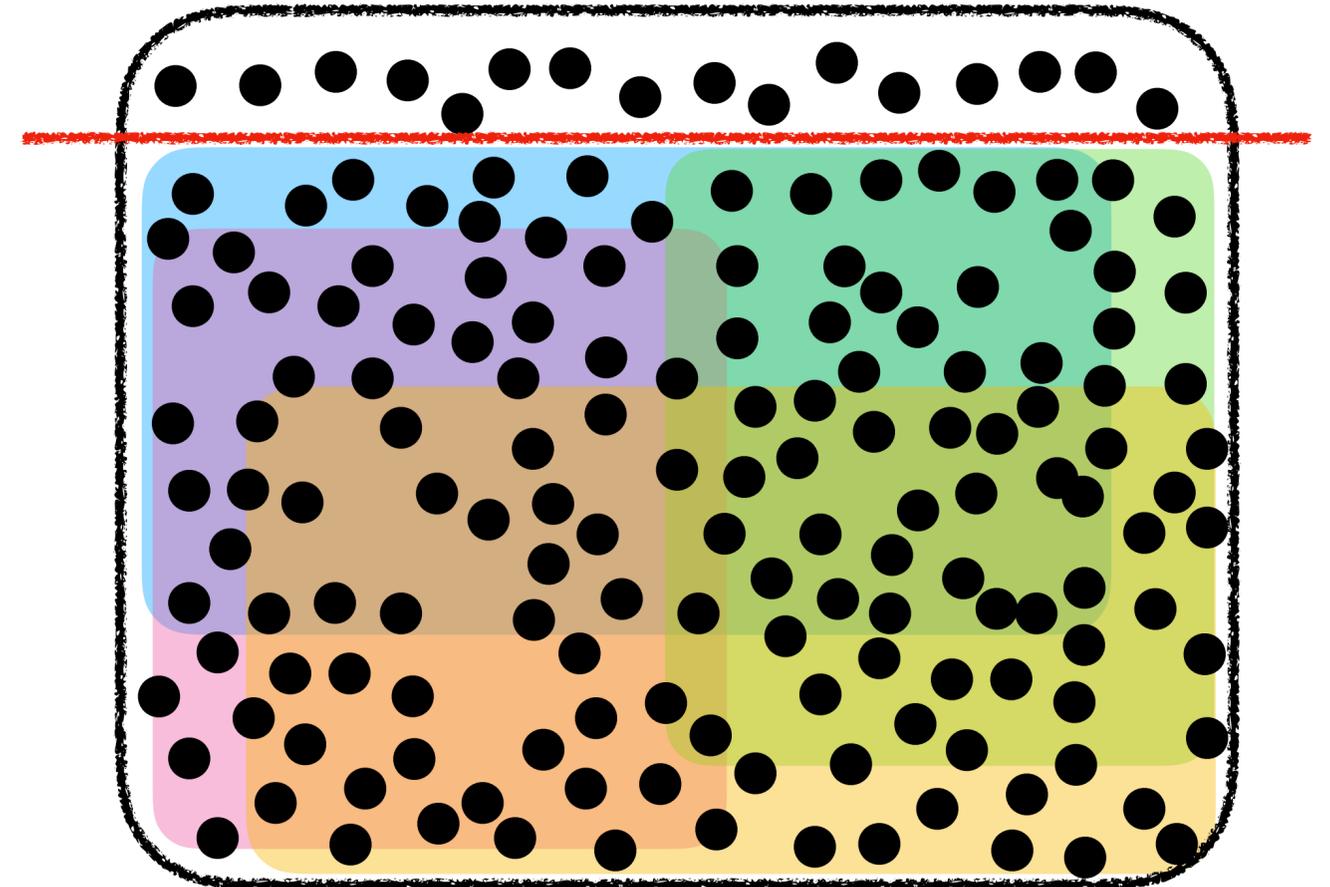
Algorithm

- Start with an initial feature representation, ϕ
- Train multiple models on random partitions of the remaining data.



Algorithm

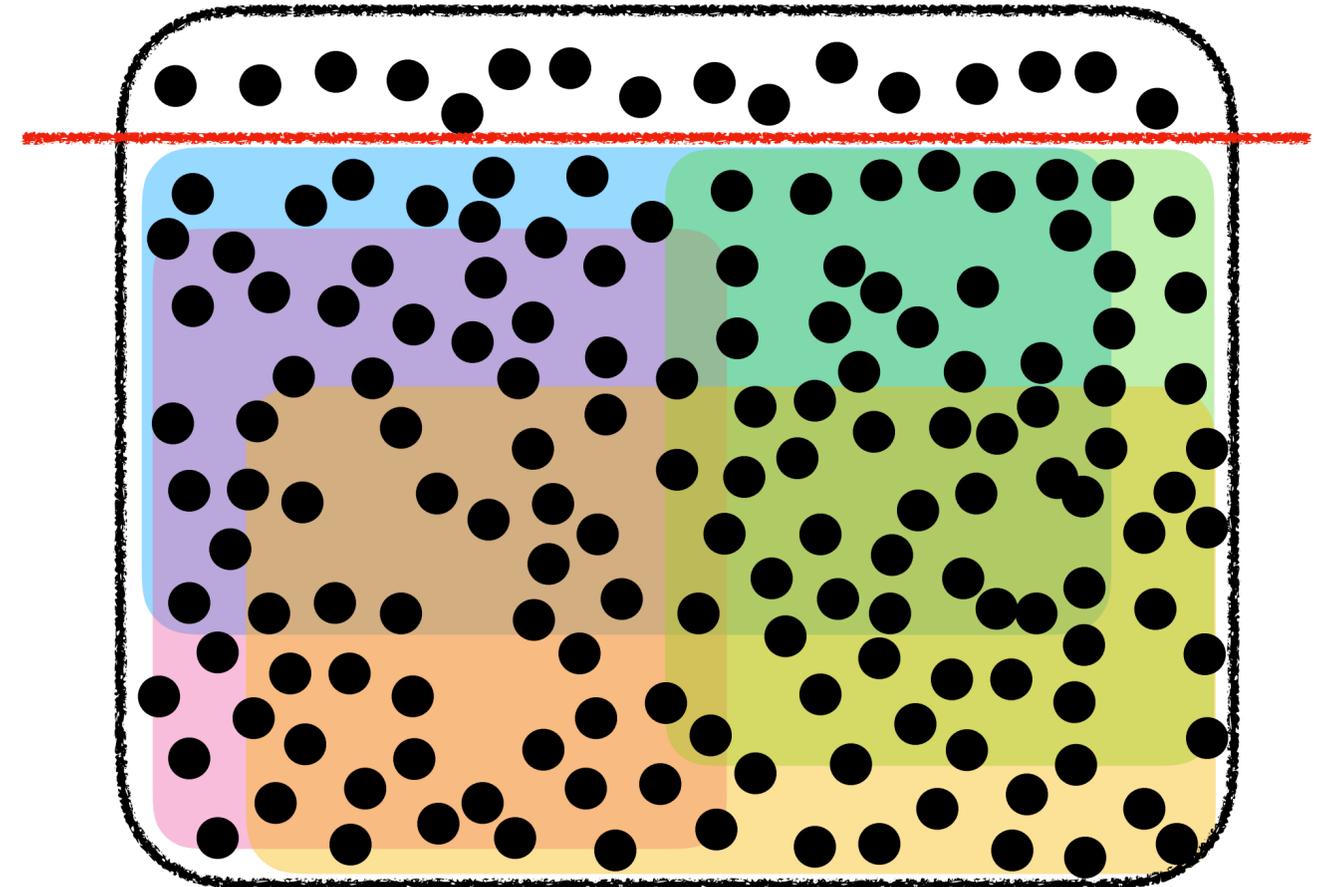
- Start with an initial feature representation, ϕ
- Train multiple models on random partitions of the remaining data.
- Discard the top-k examples which are correctly identified by most models, iteratively, till the ensemble is no longer confident.



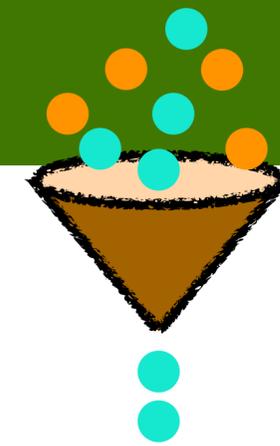
Algorithm

Precursor: Zellers et al., 2018; 2019

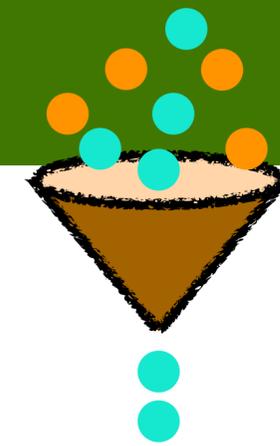
- Start with an initial feature representation, ϕ
- Train multiple models on random partitions of the remaining data.
- Discard the top-k examples which are correctly identified by most models, iteratively, till the ensemble is no longer confident.
- Lightweight Adversarial Filtering (**AFLite**):
 - linear models
 - fixed feature representation, ϕ .



Follow-up questions



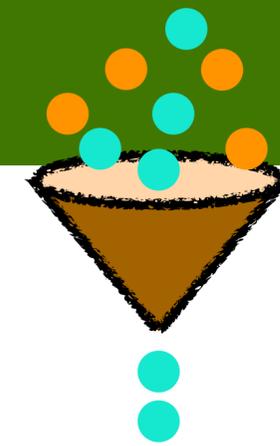
Follow-up questions



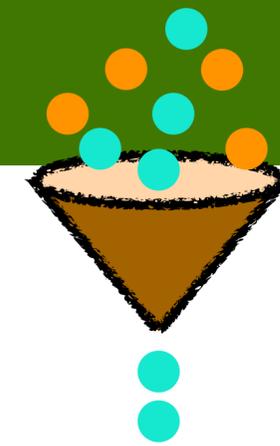
- Is AFLite optimal?
- No! It's a greedy procedure

Follow-up questions

- Is AFLite optimal?
 - No! It's a greedy procedure
- Is there an optimal variant?
 - Yes! But intractable – subset selection problem.



Follow-up questions



- Is AFLite optimal?
 - No! It's a greedy procedure
- Is there an optimal variant?
 - Yes! But intractable – subset selection problem.
- How would it work?
 - Find the smallest subset, any train-test split of which achieves high accuracy, but does not generalize to held-out data outside the subset.
 - Mini-max optimization problem

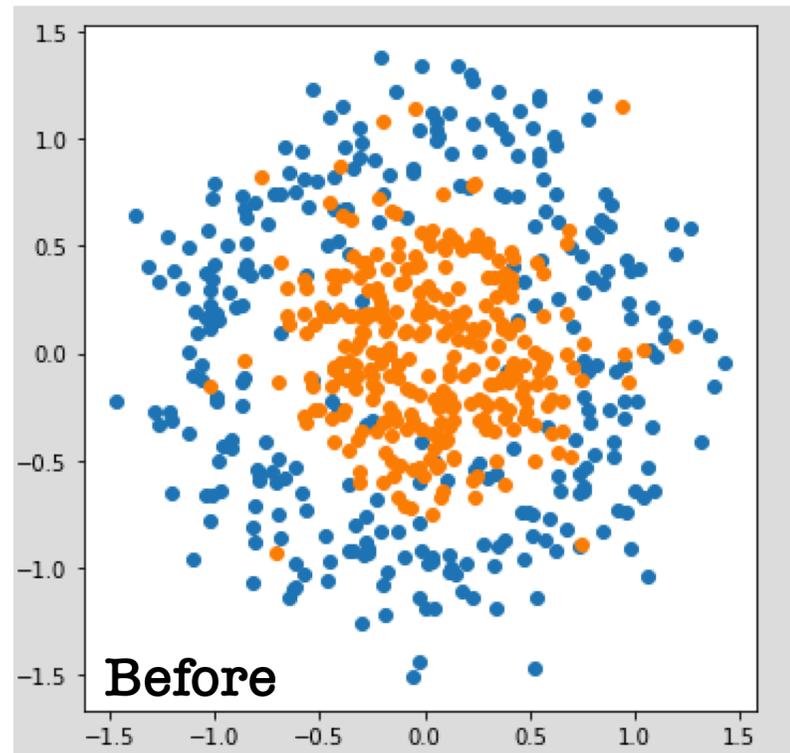


Evaluation Setting

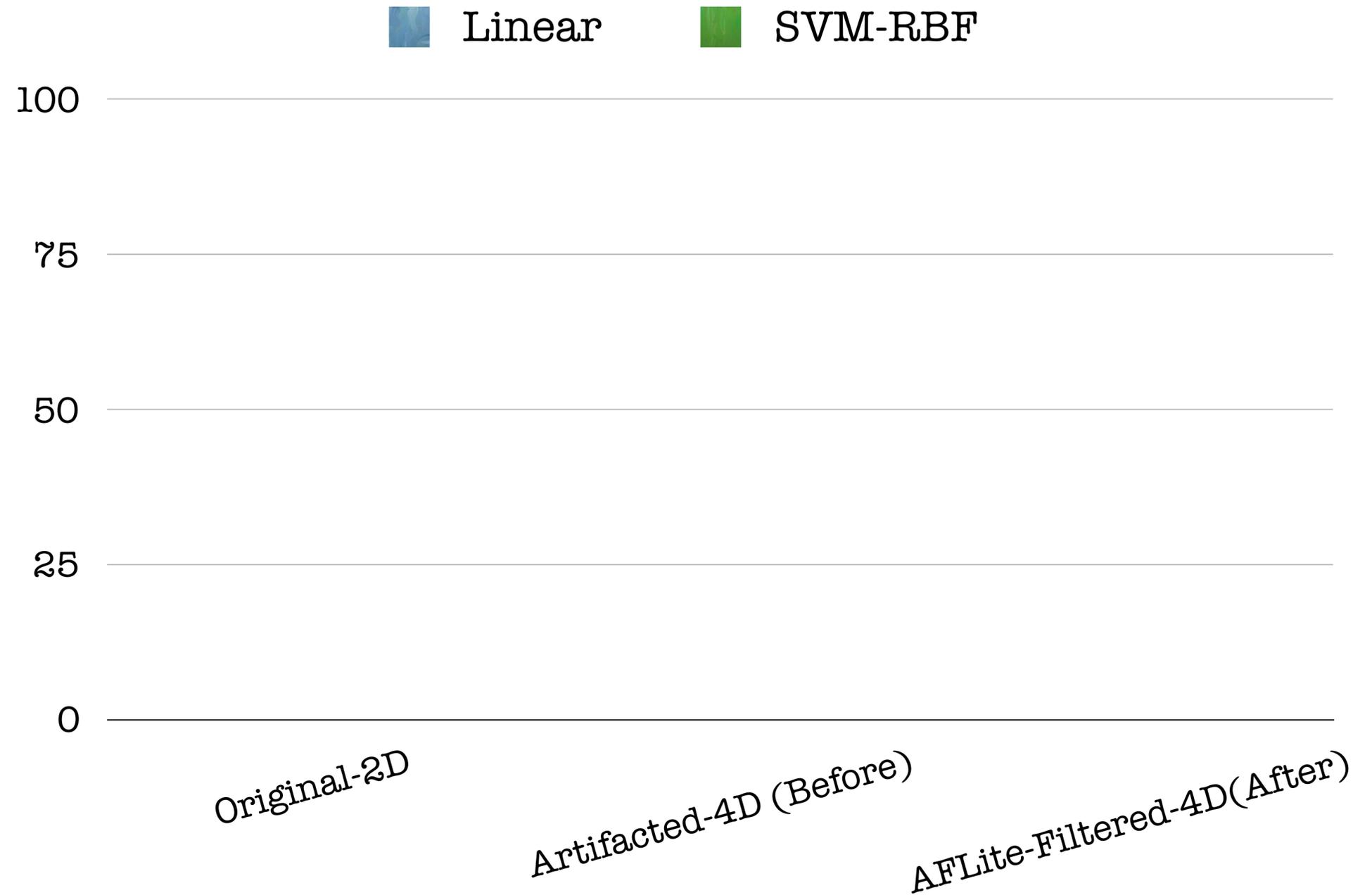
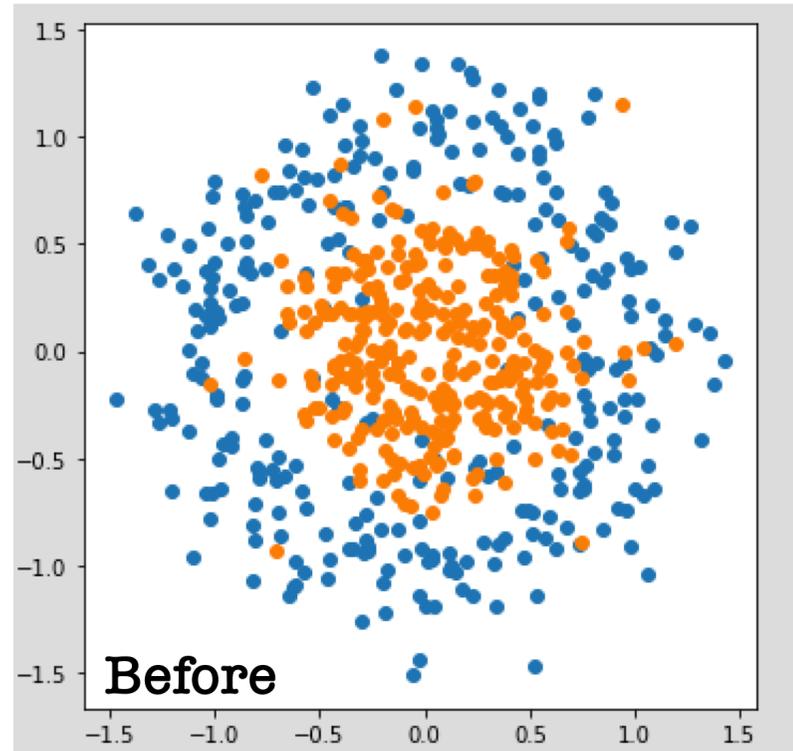
- Extrinsic Evaluation:
 - Model performance on test before / after filtering.
 - Training data also changes to account for distributional differences
- Intrinsic Evaluation:
 - Filtered dataset properties

	Unfiltered Train	Filtered Train
Unfiltered Test	✓	
Filtered Test		✓

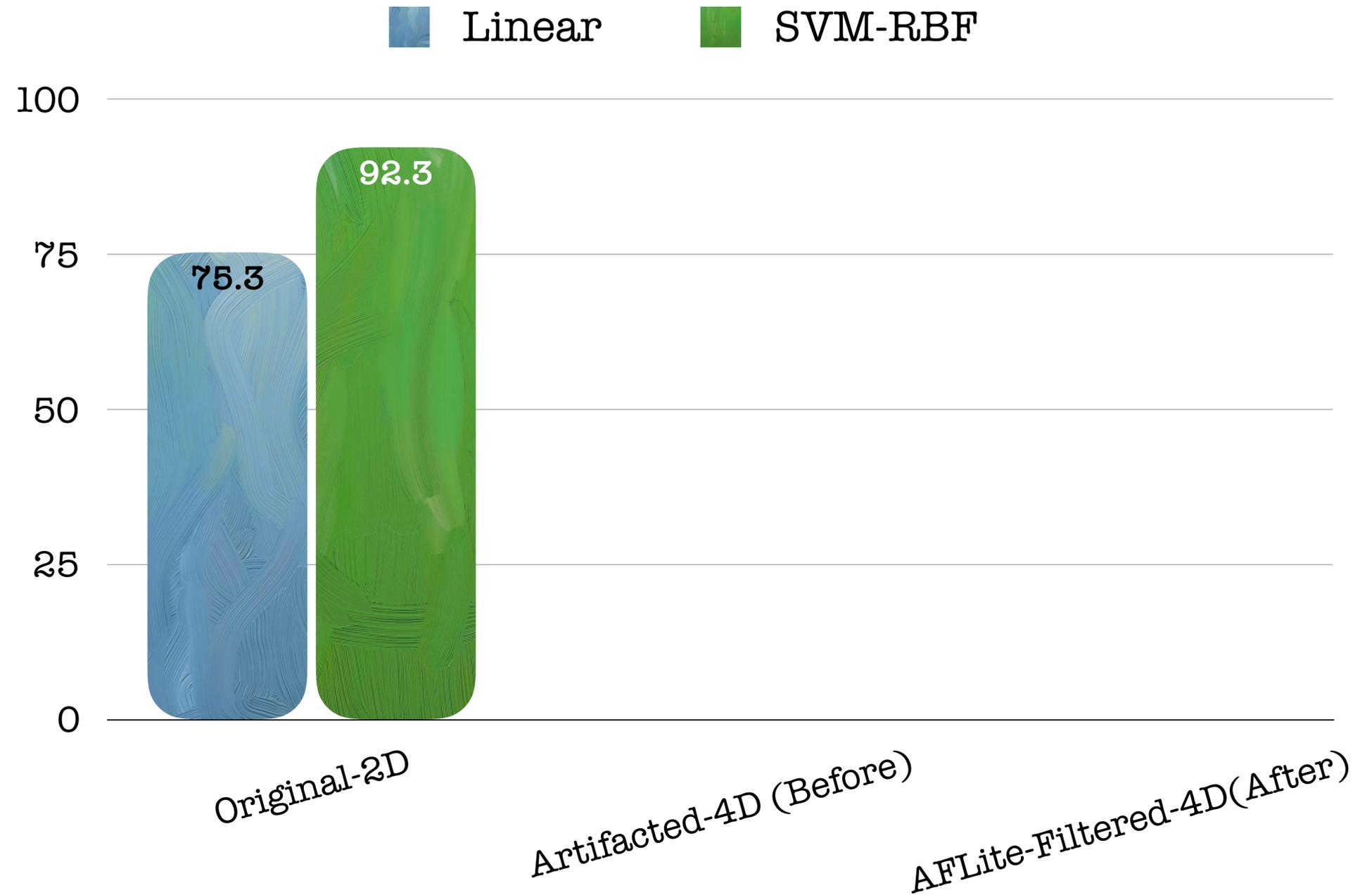
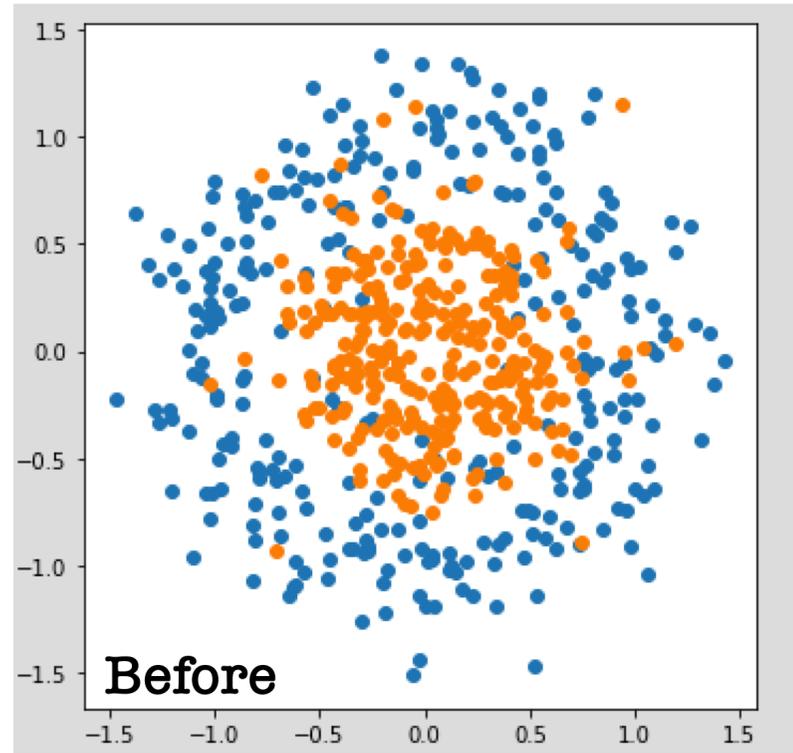
Task 0: Synthetic



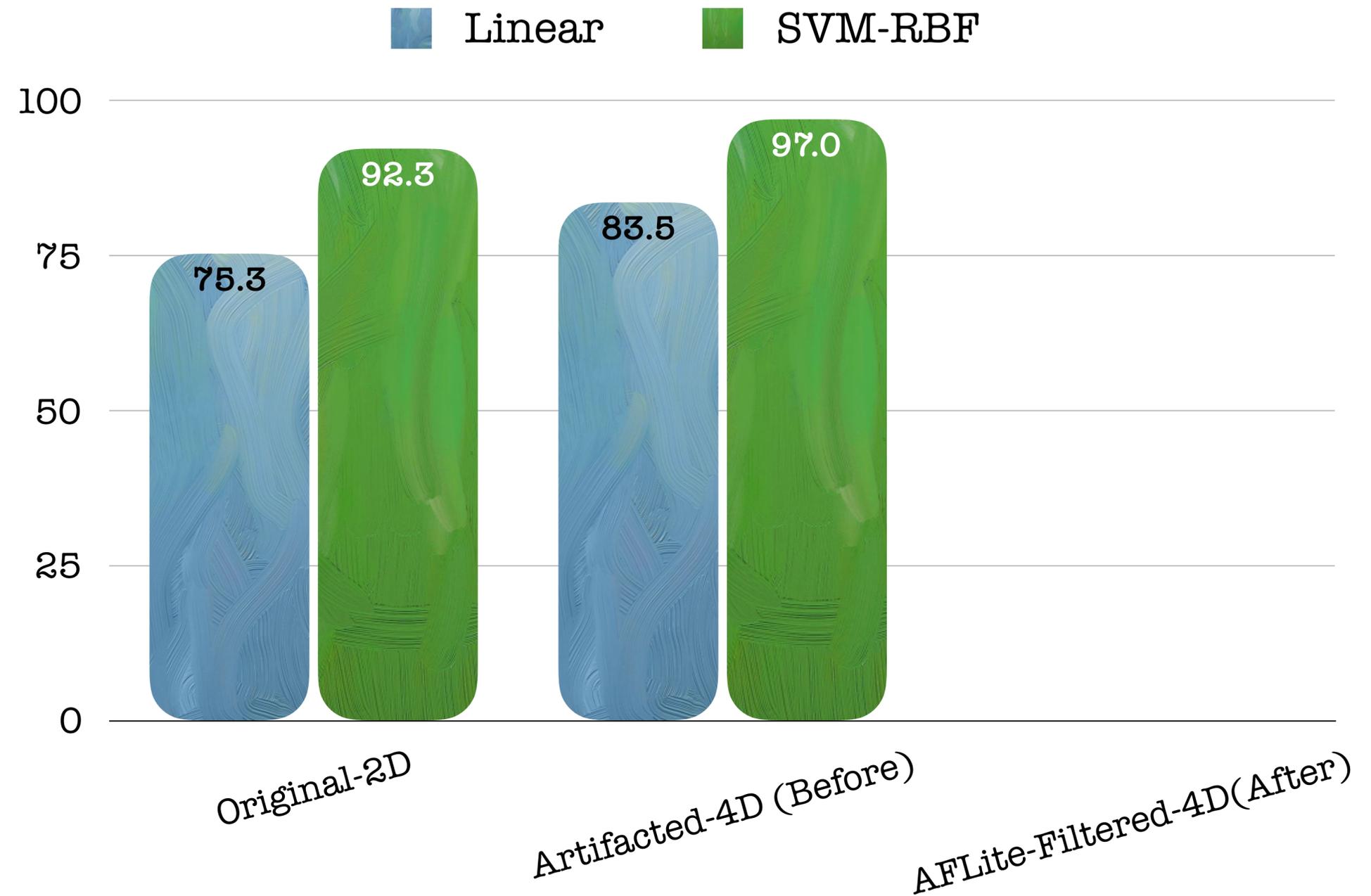
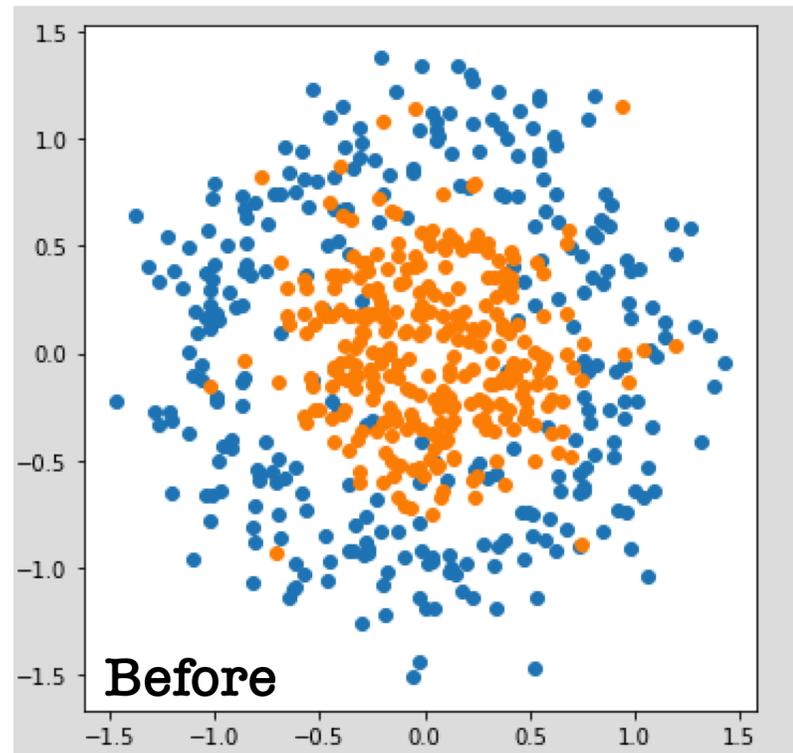
Task 0: Synthetic



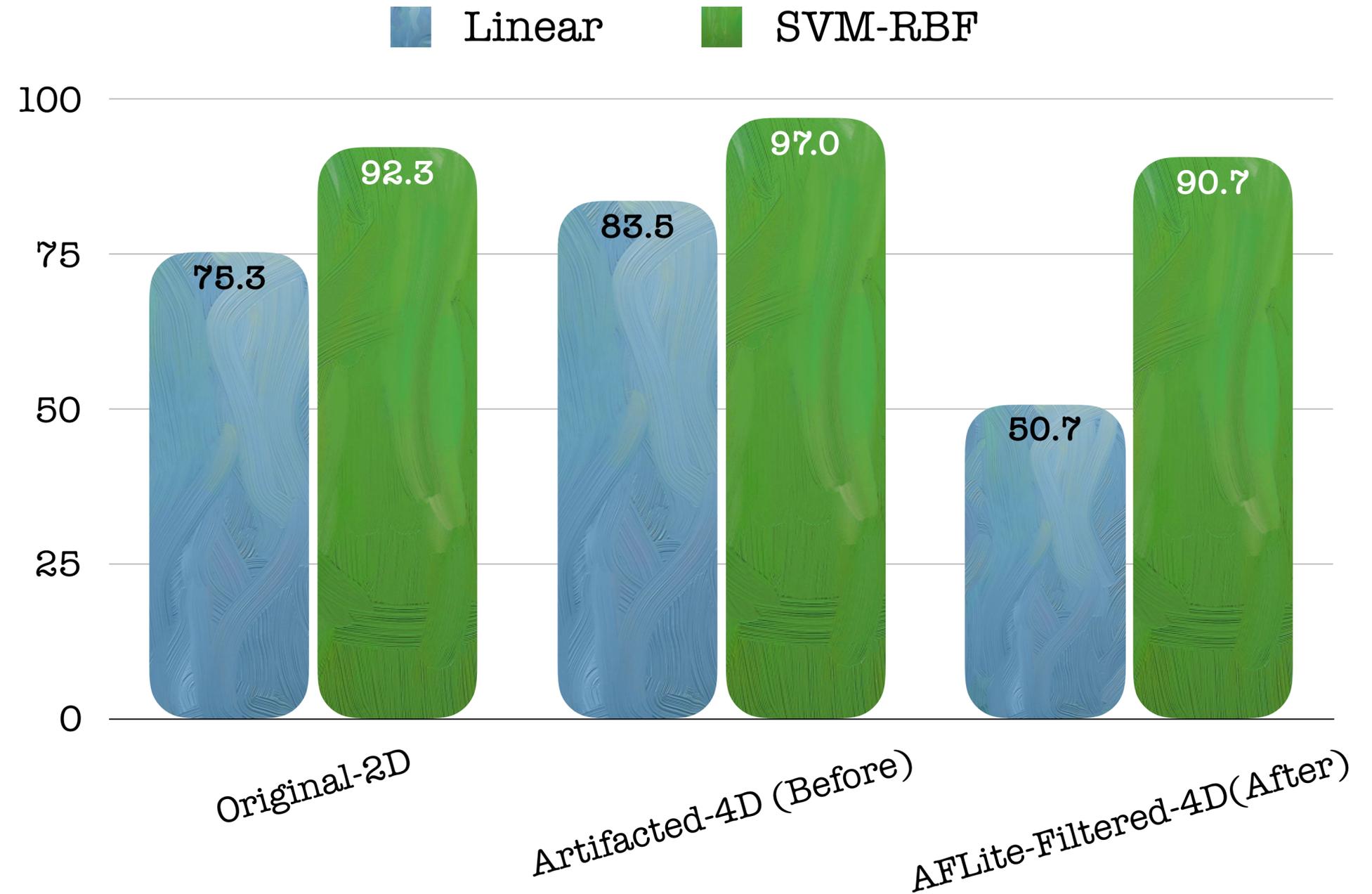
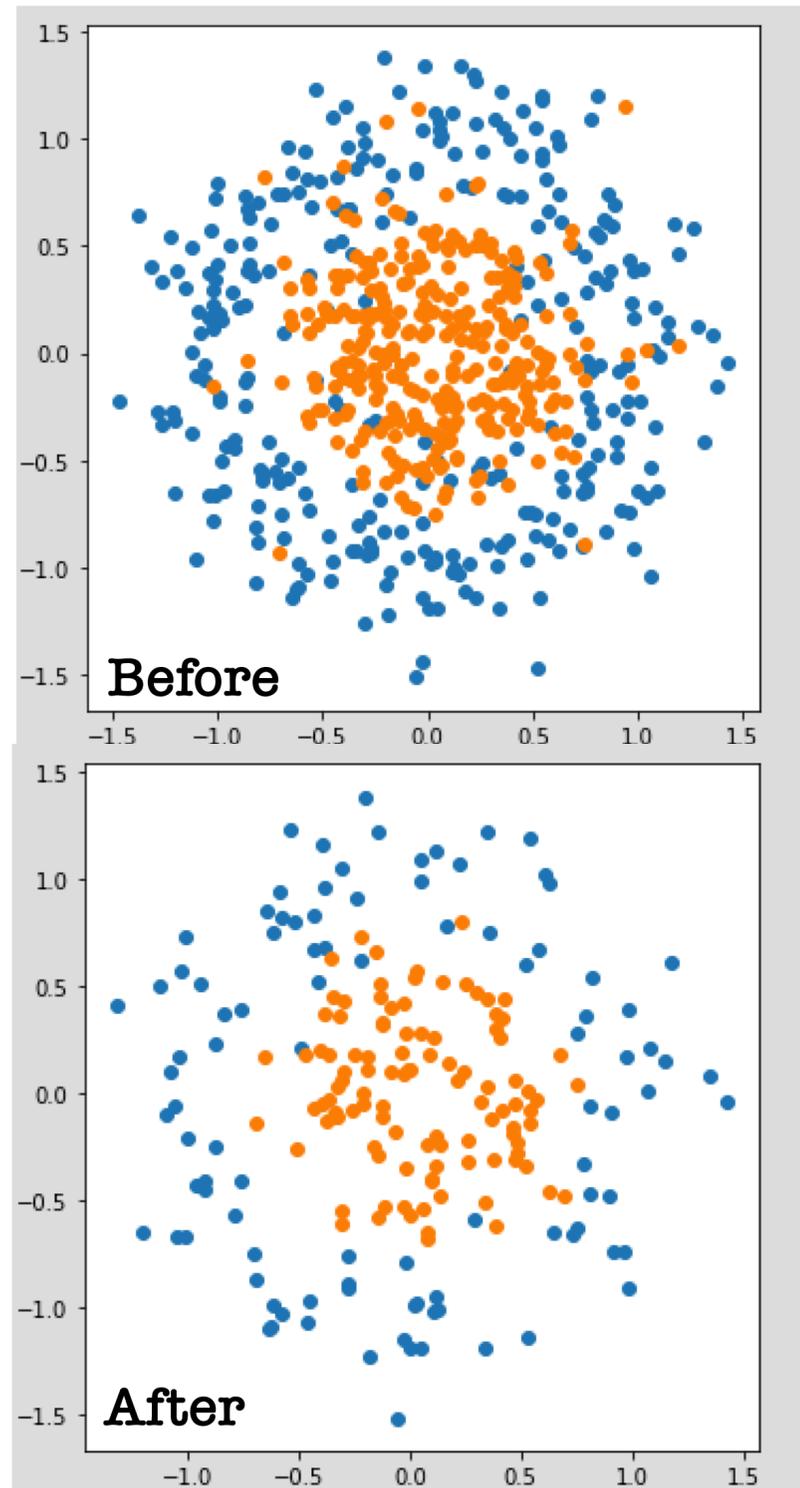
Task 0: Synthetic



Task 0: Synthetic

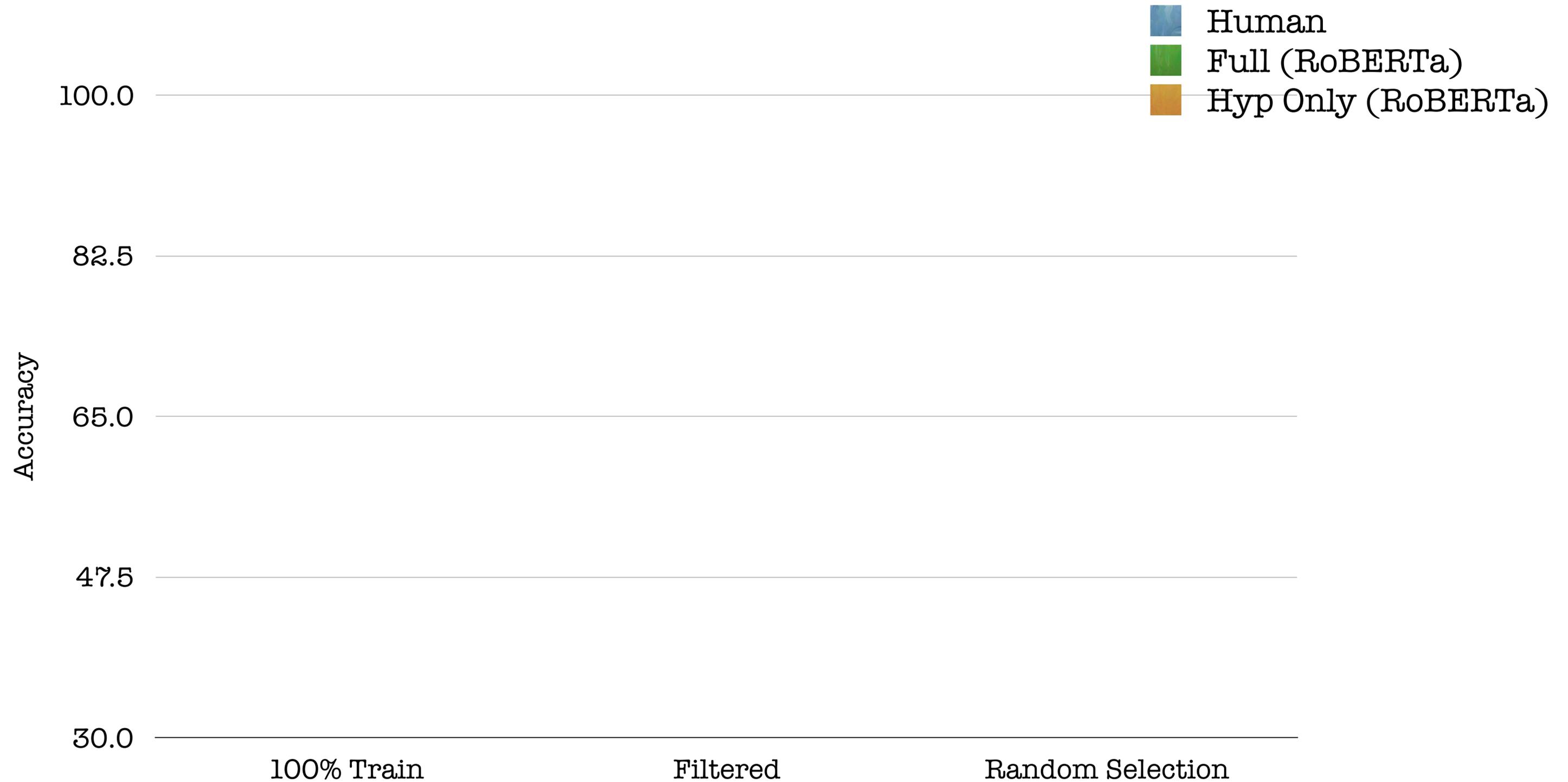


Task 0: Synthetic

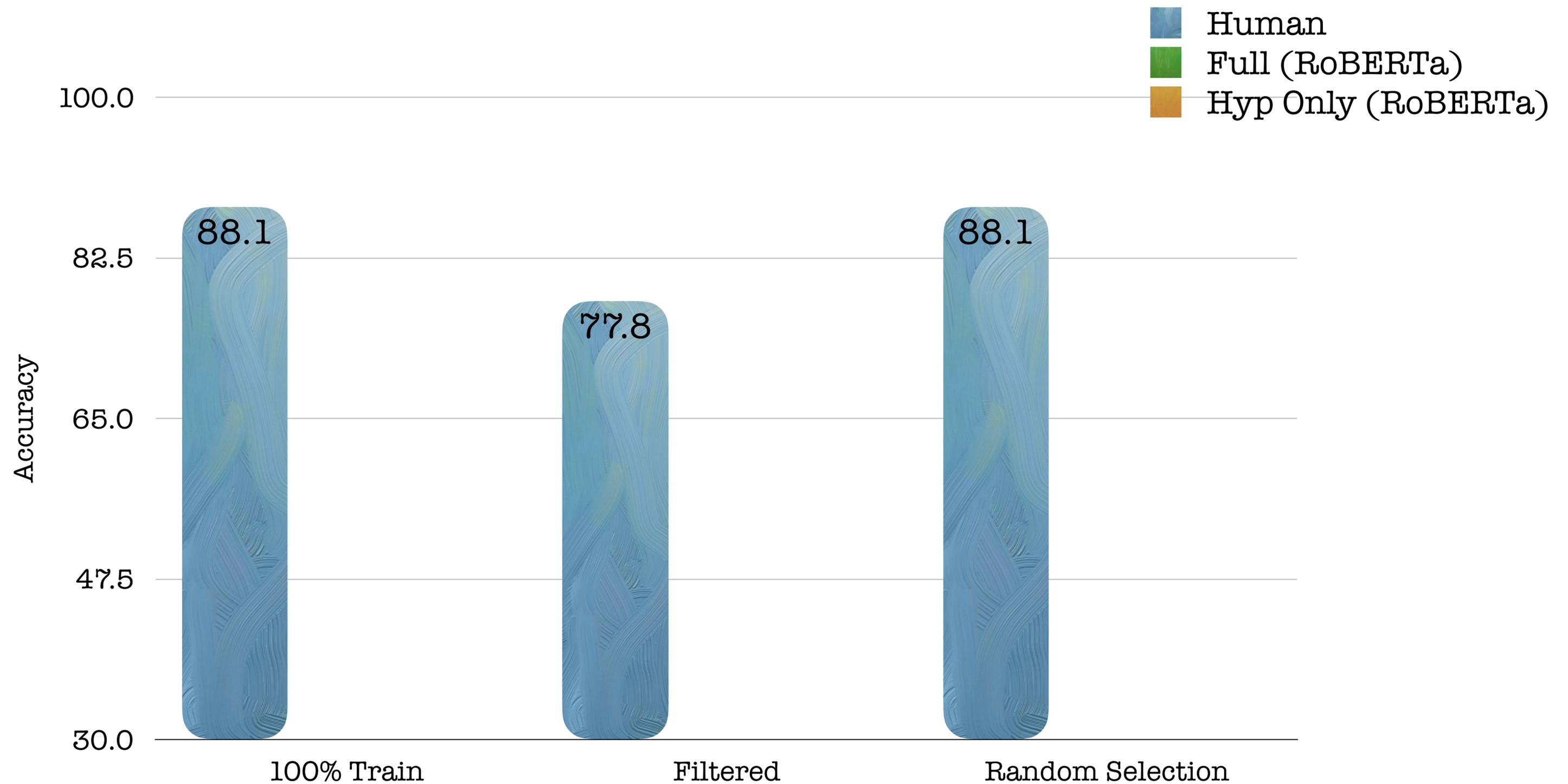


Task 1: SNLI performance

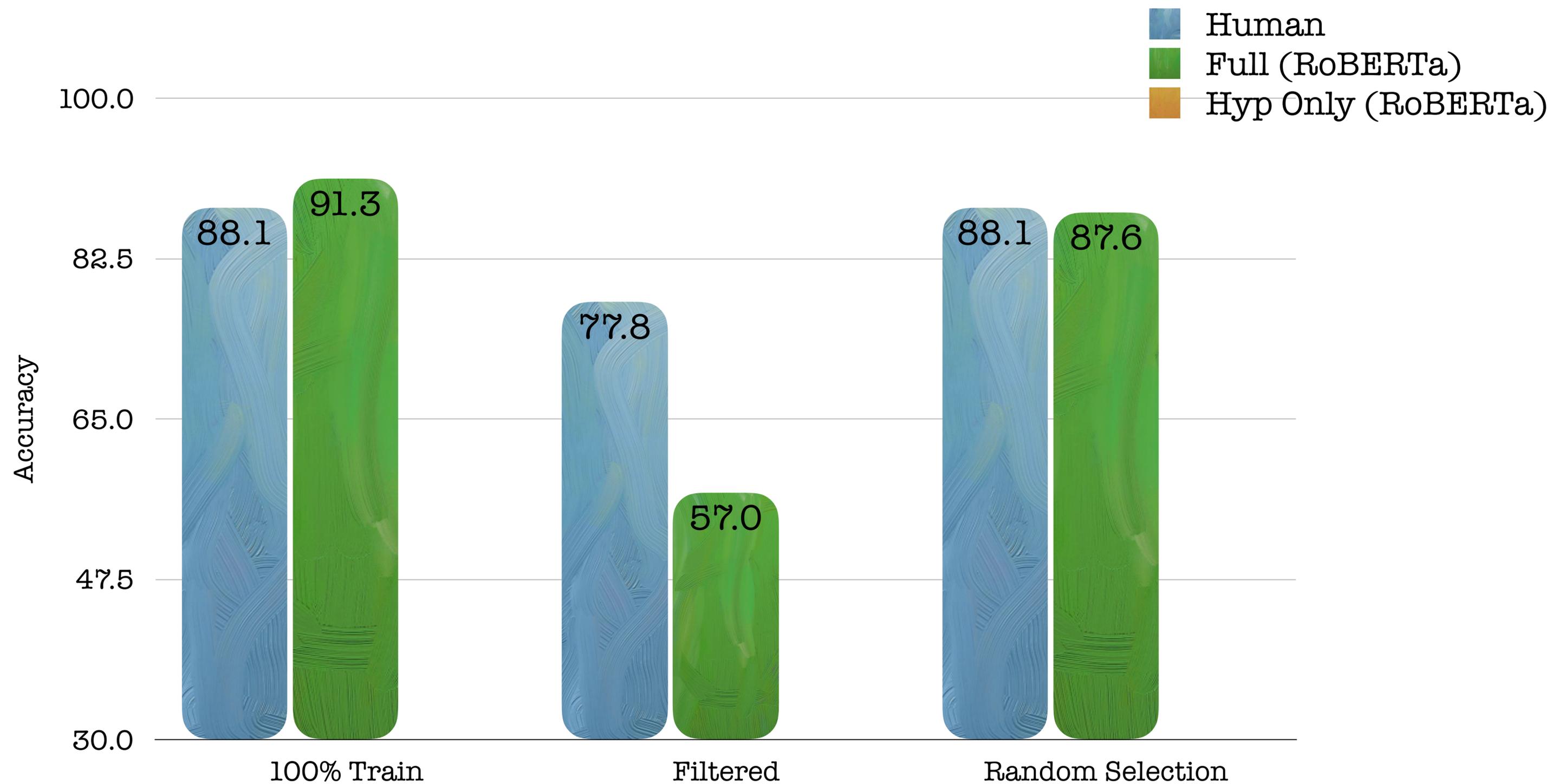
Task 1: SNLI performance



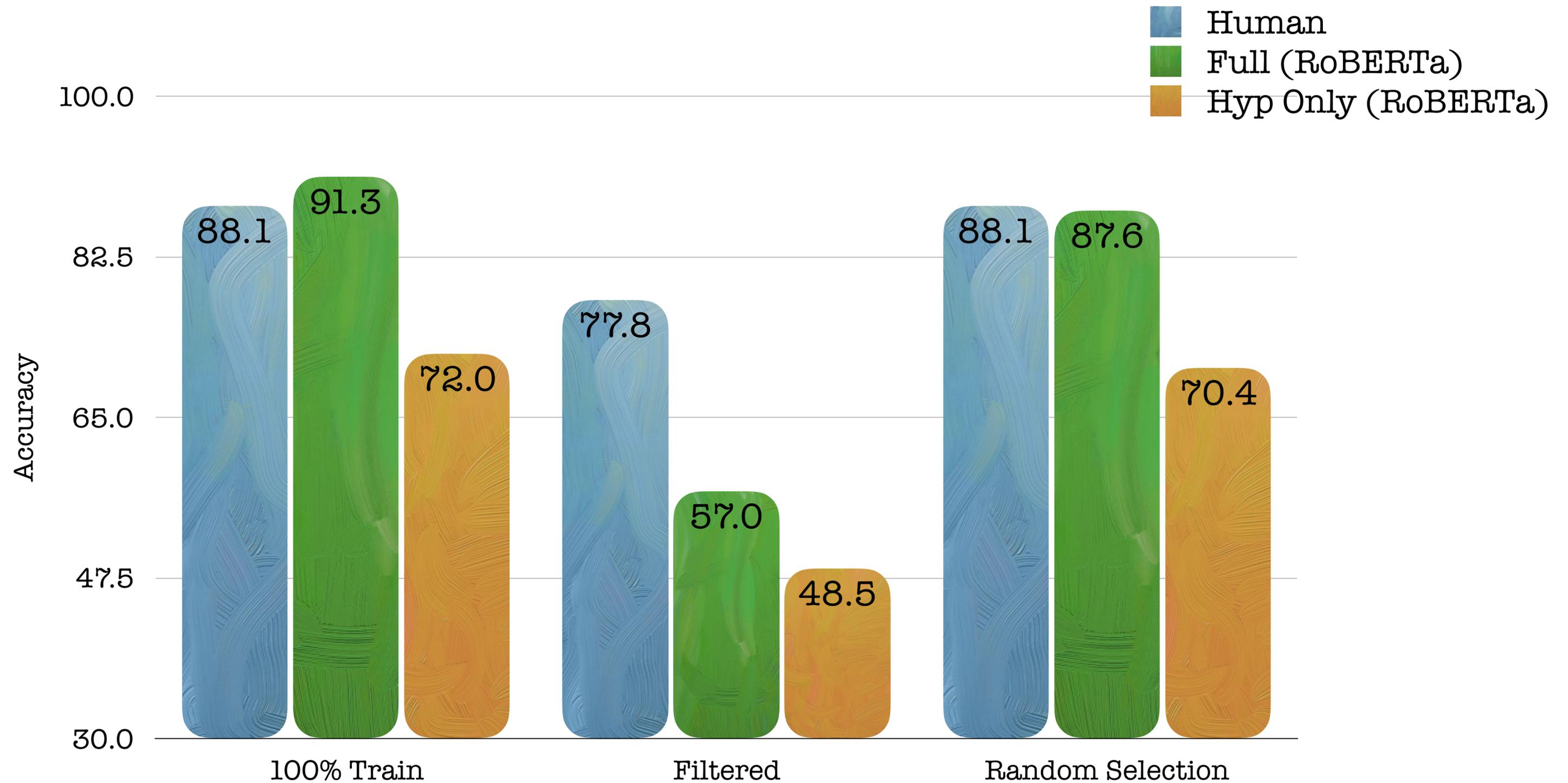
Task 1: SNLI performance



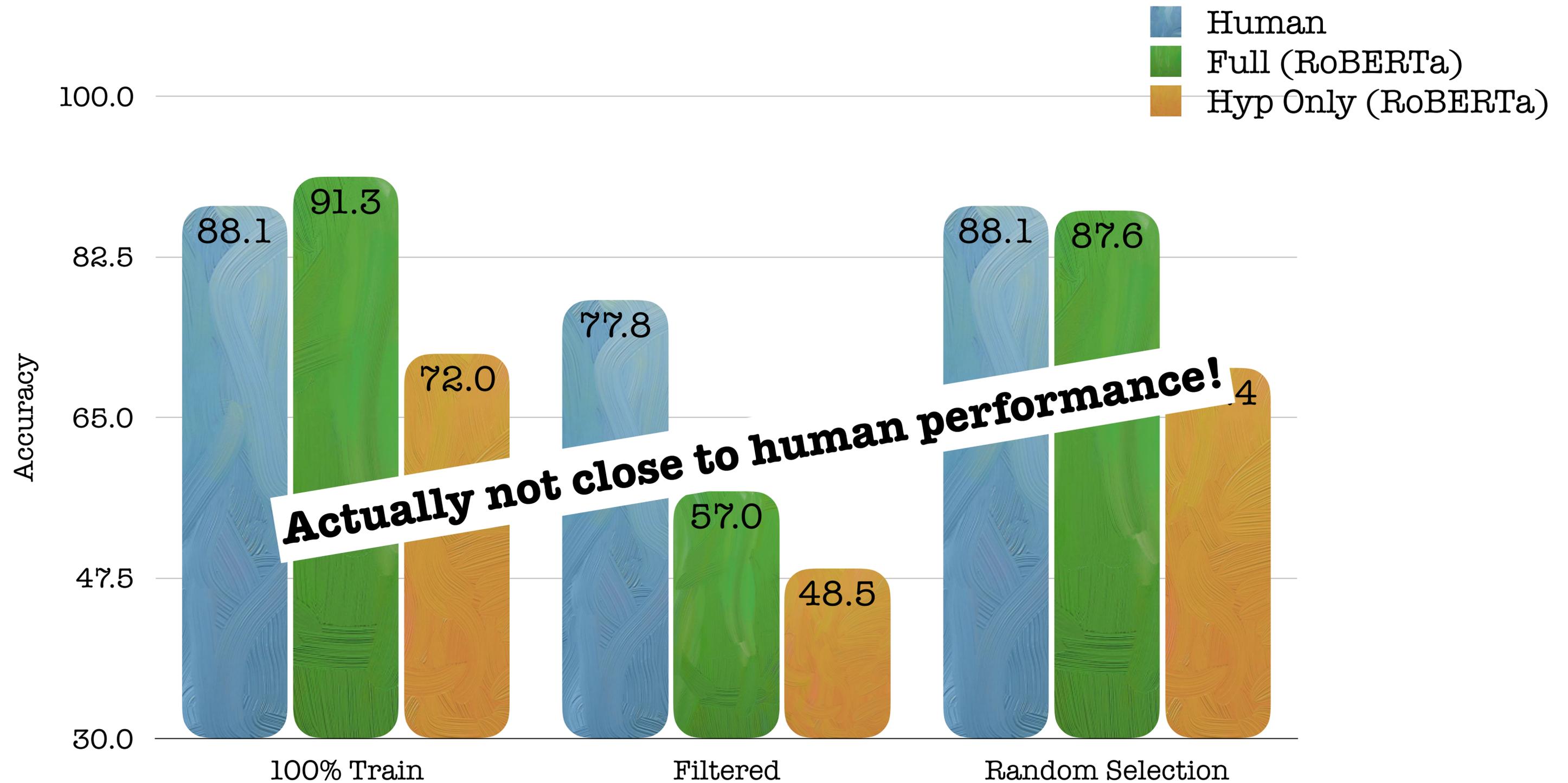
Task 1: SNLI performance



Task 1: SNLI performance

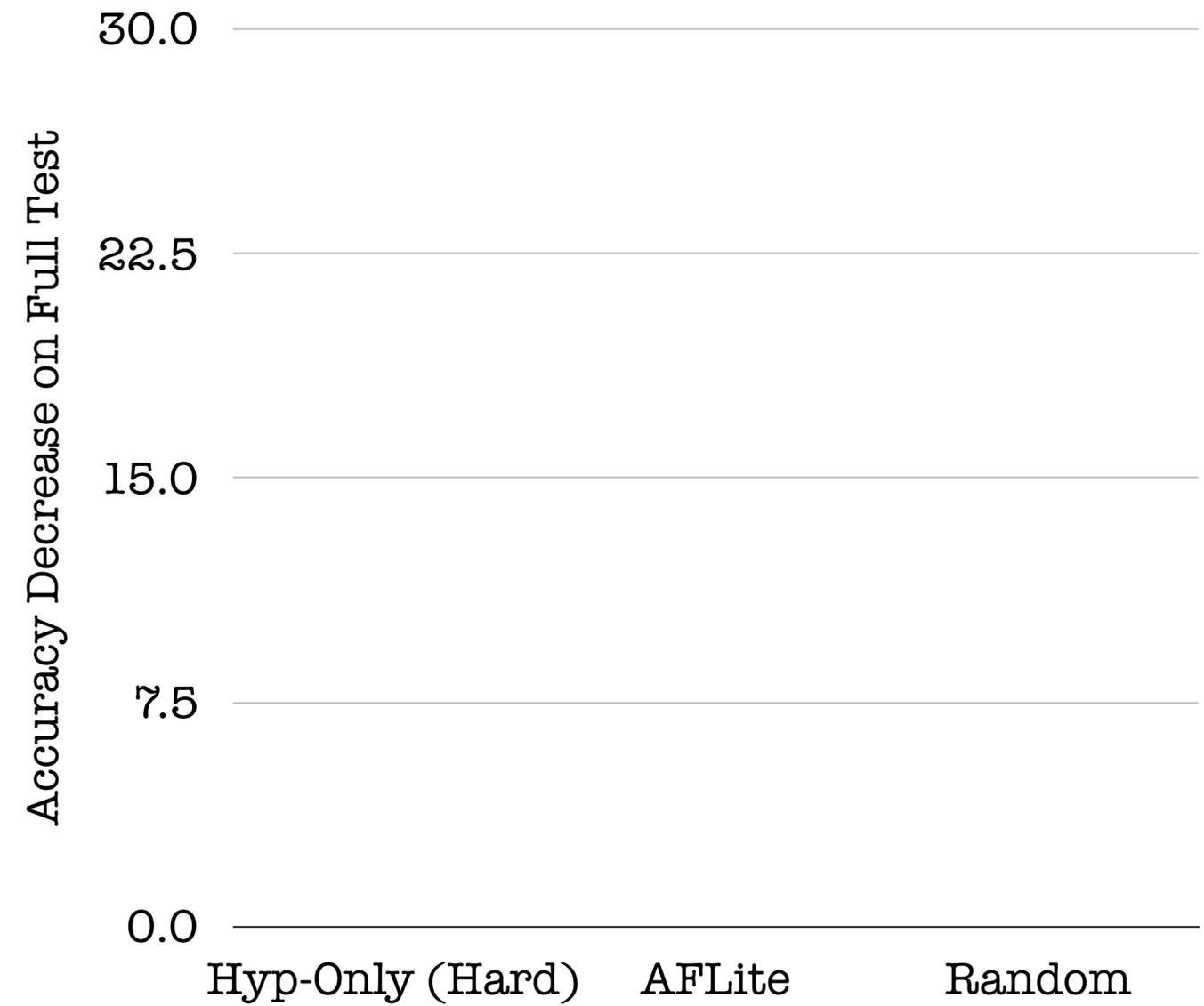


Task 1: SNLI performance

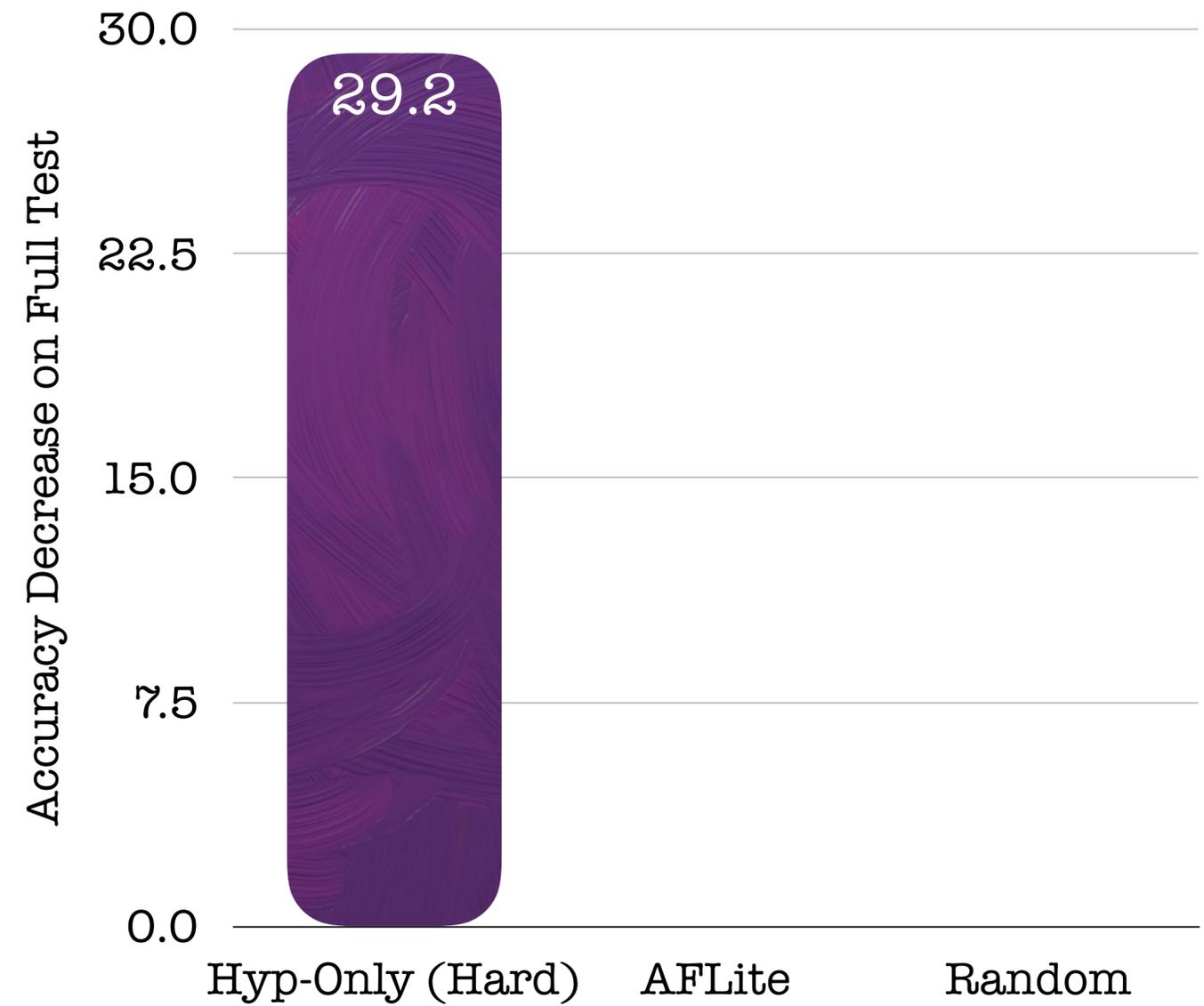


Comparison with Hyp-Only Filtering

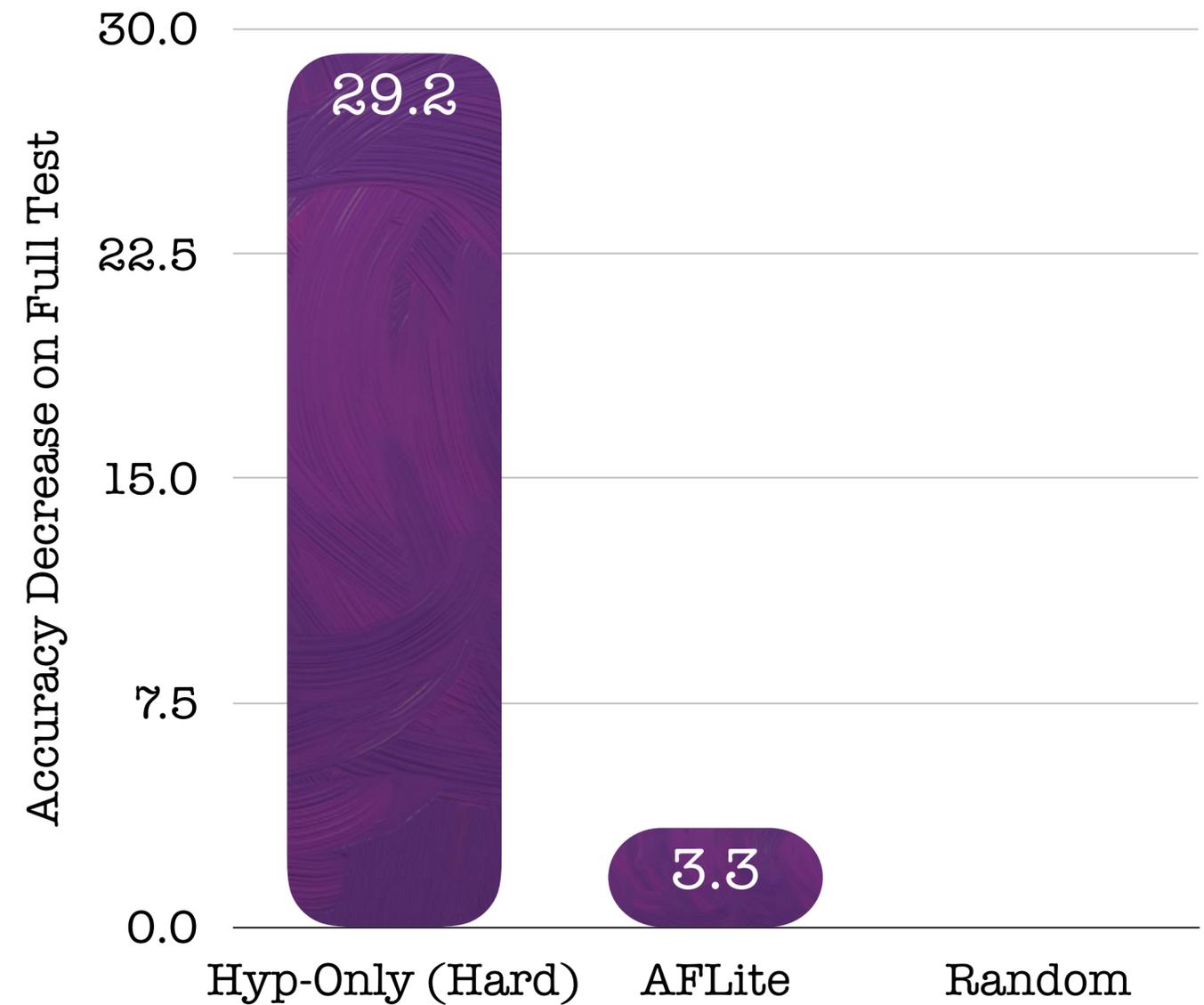
Comparison with Hyp-Only Filtering



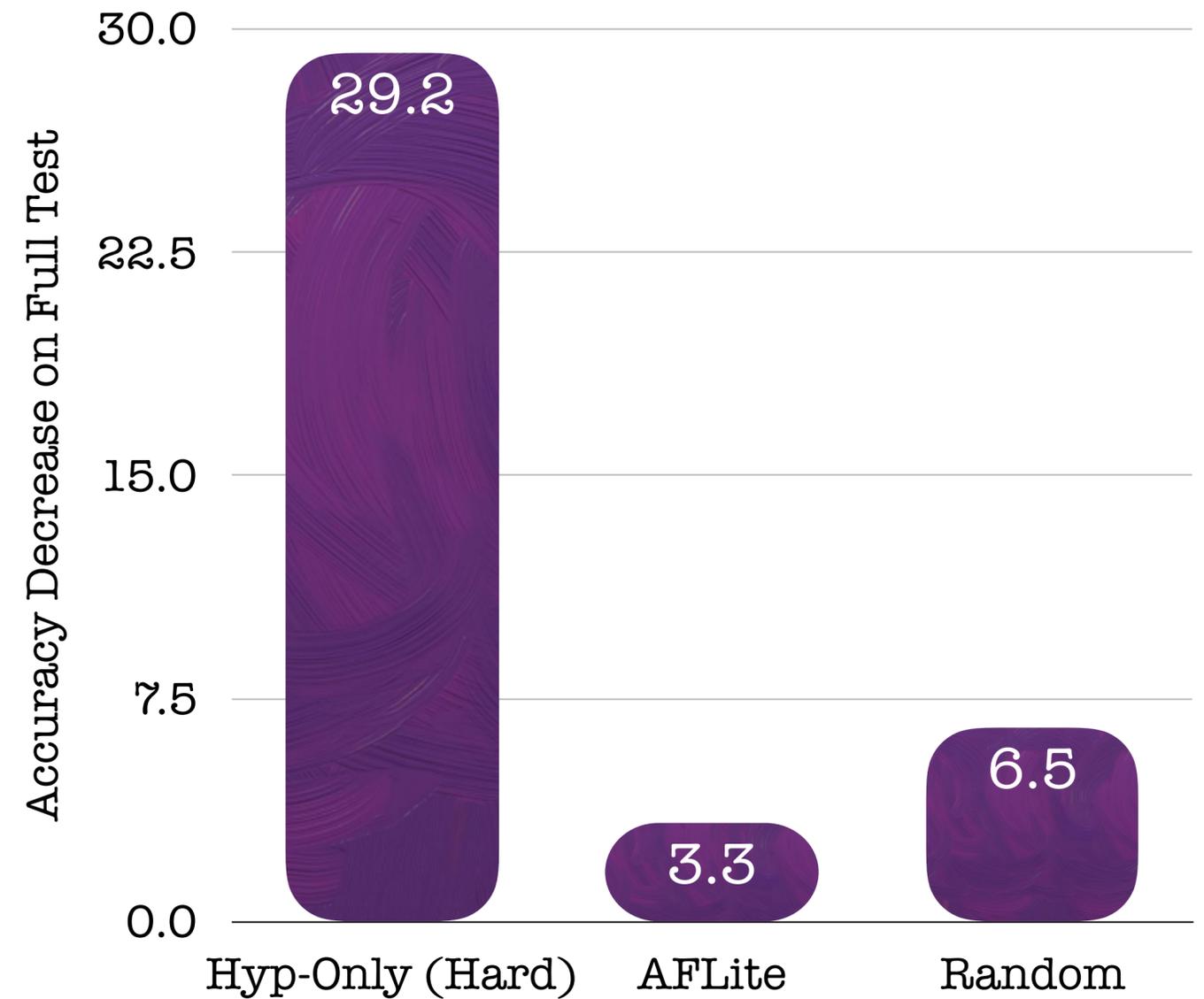
Comparison with Hyp-Only Filtering



Comparison with Hyp-Only Filtering

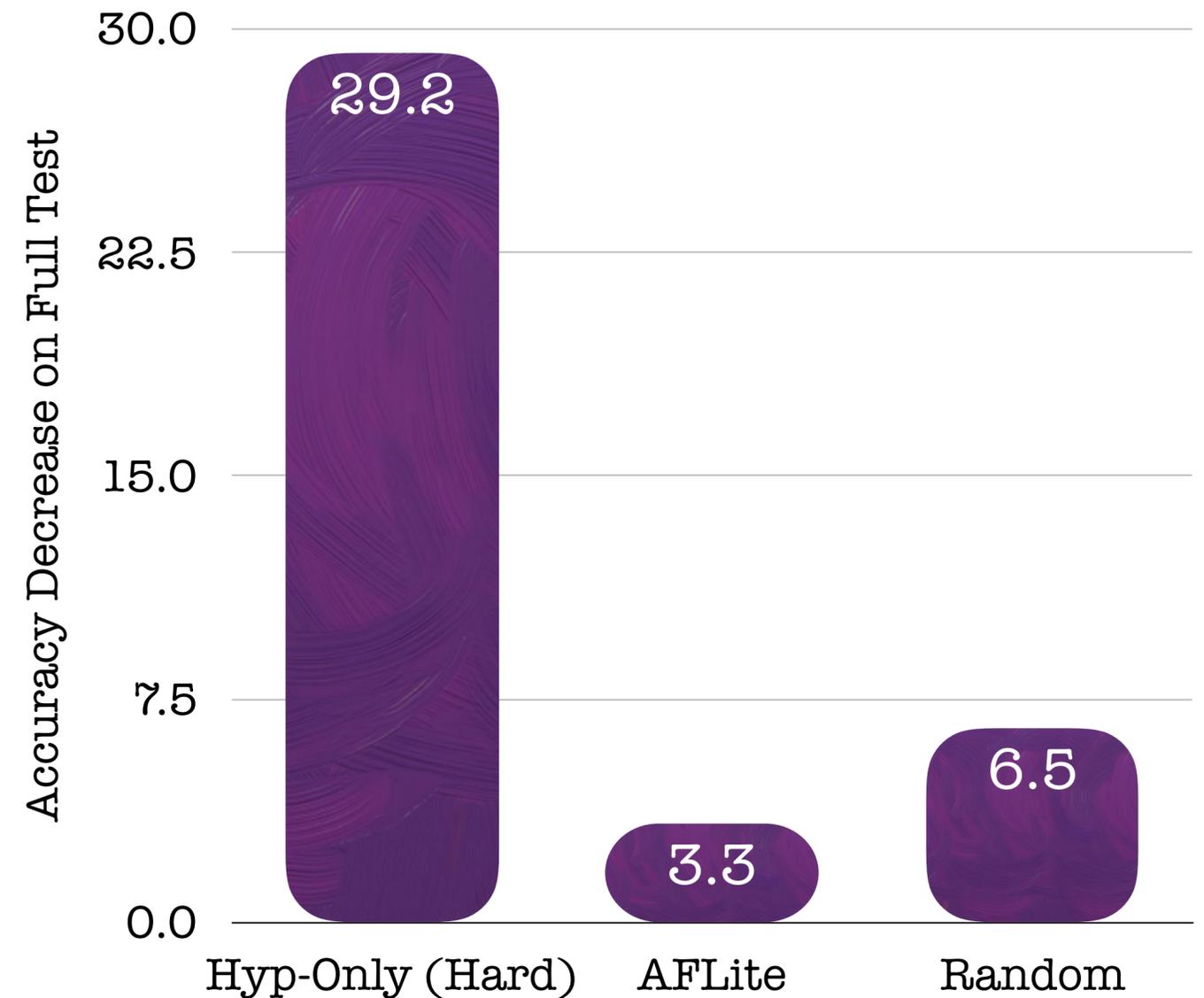


Comparison with Hyp-Only Filtering



Comparison with Hyp-Only Filtering

- AFLite retains generalizability to many examples.
- Manually detecting artifacts can only get rid of some
- Manual filtering for **balancing** artifacts might not be effective.



Digging Deeper Again



- Different word-class associations.

E

N

C

TRANSPORTATION

NEARBY

NOBODY

LEAST

PRETENDING

STEALING

AWAKE

BACKYARD

ALIENS

SPORT

DALMATION

SANDBOX

MULTIPLE

FARMER

STOLE

Digging Deeper Again



- Different word-class associations.

E

N

C

TRANSPORTATION

NEARBY

NOBODY

LEAST

PRETENDING

STEALING

AWAKE

BACKYARD

ALIENS

SPORT

DALMATION

SANDBOX

MULTIPLE

FARMER

STOLE

E

N

C

OUTDOORS

TALL

NOBODY

LEAST

FIRST

SLEEPING

INSTRUMENT

COMPETITION

No

OUTSIDE

SAD

Tv

ANIMAL

FAVORITE

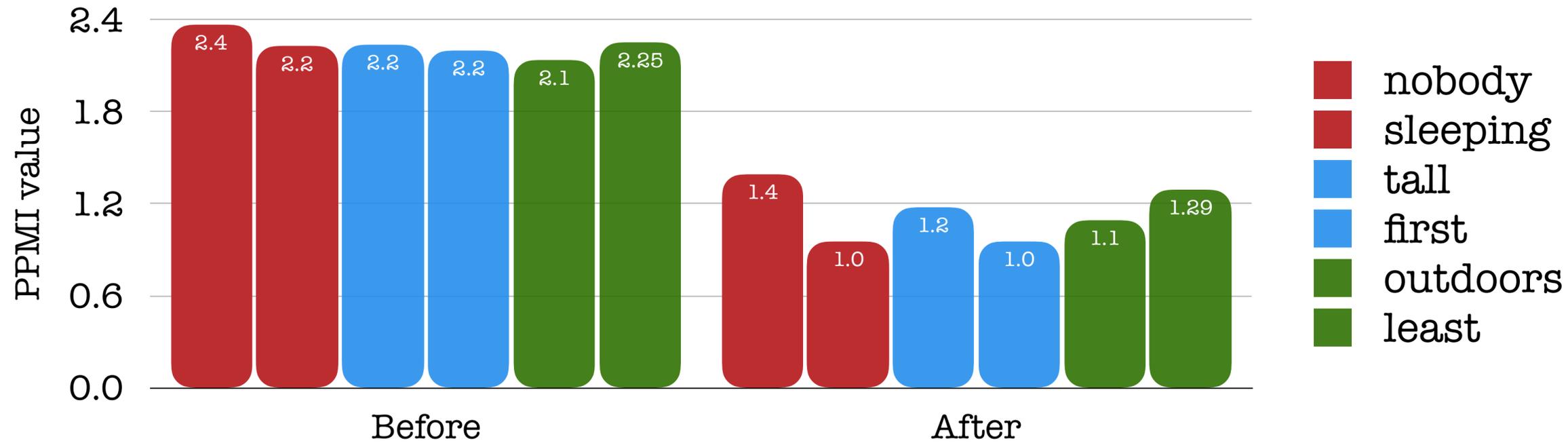
CAT

Effect 1: Word-Class Association

- Overall word-class association decreases

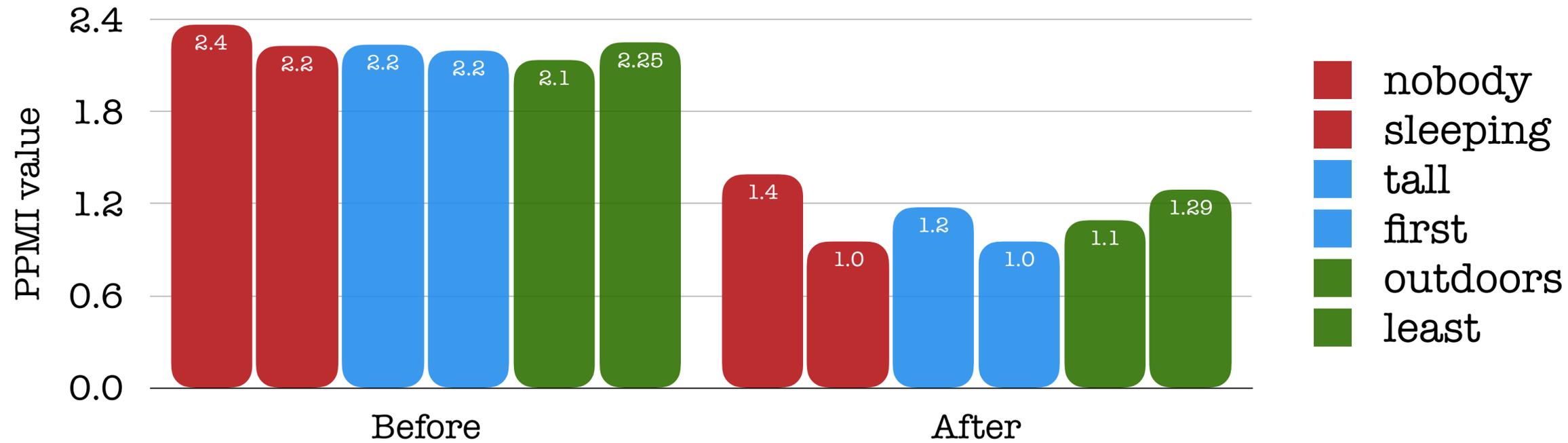
Effect 1: Word-Class Association

- Overall word-class association decreases



Effect 1: Word-Class Association

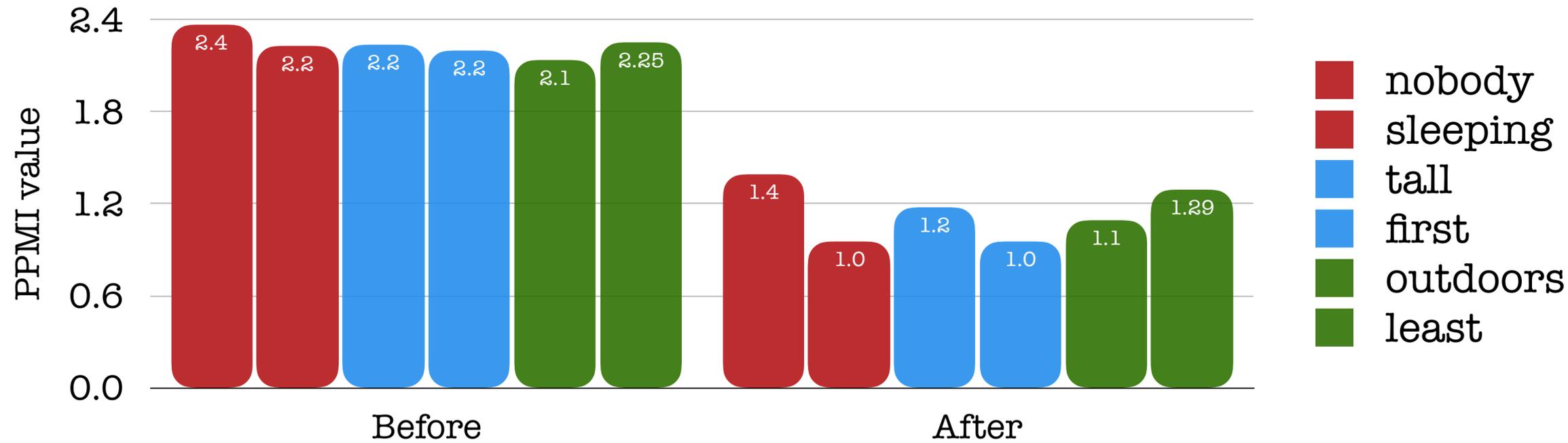
- Overall word-class association decreases



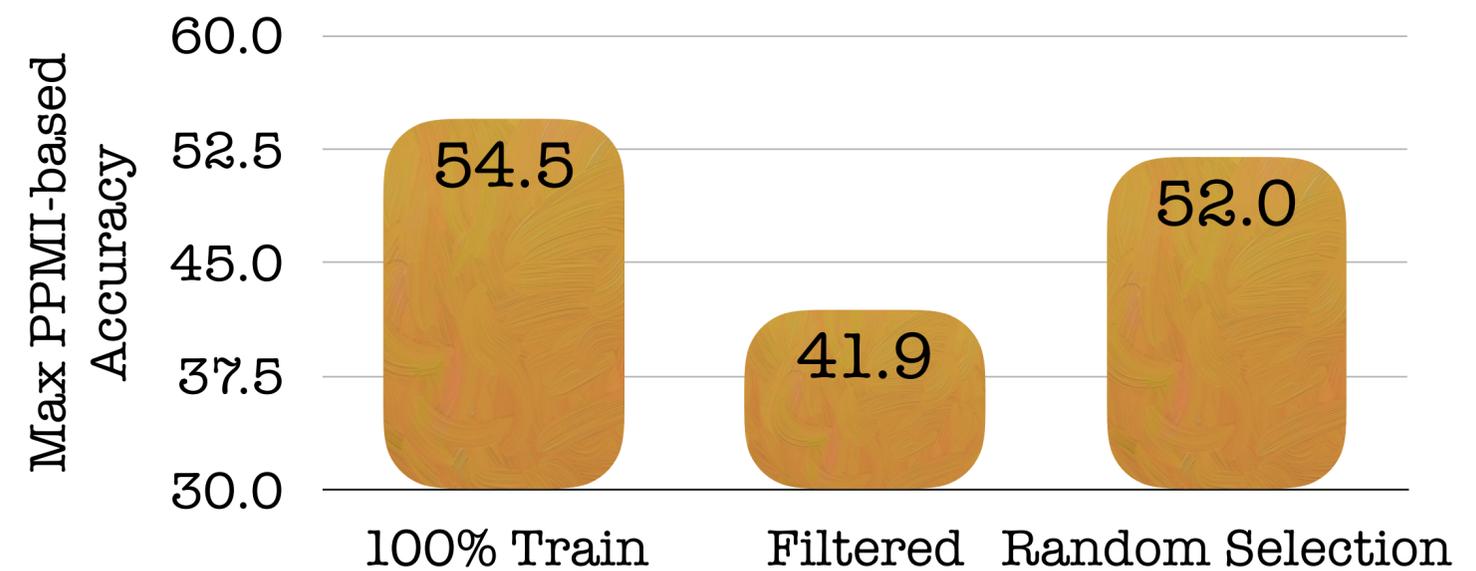
- If word association was the only indicator

Effect 1: Word-Class Association

- Overall word-class association decreases



- If word association was the only indicator

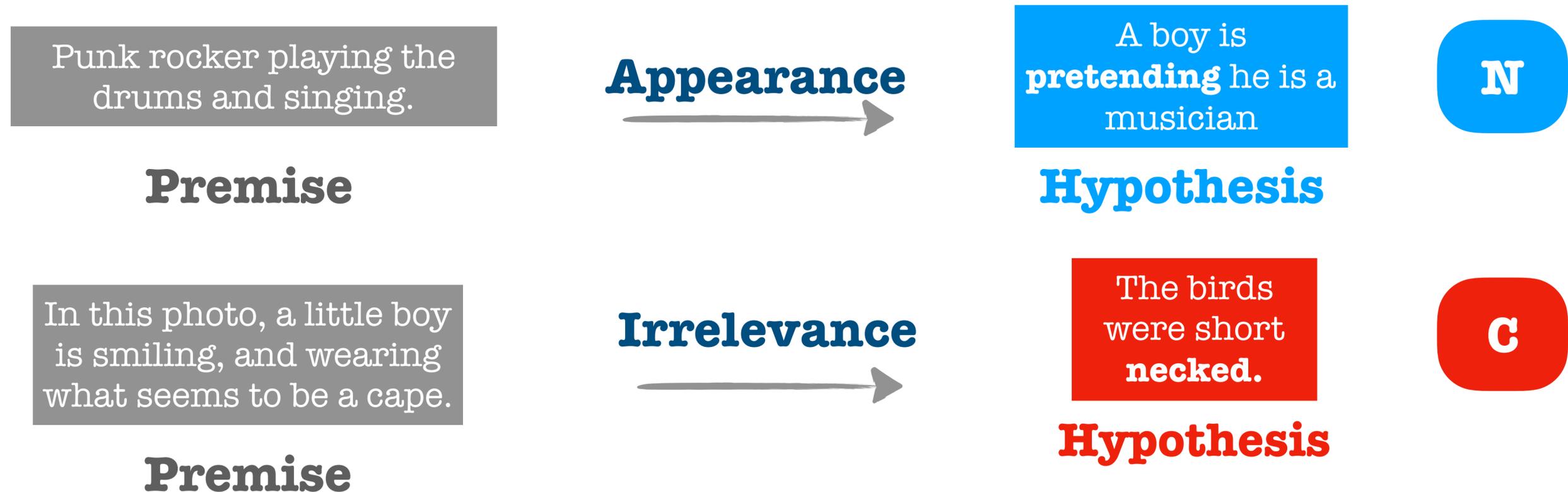


Effect 2: Newer, less common artifacts

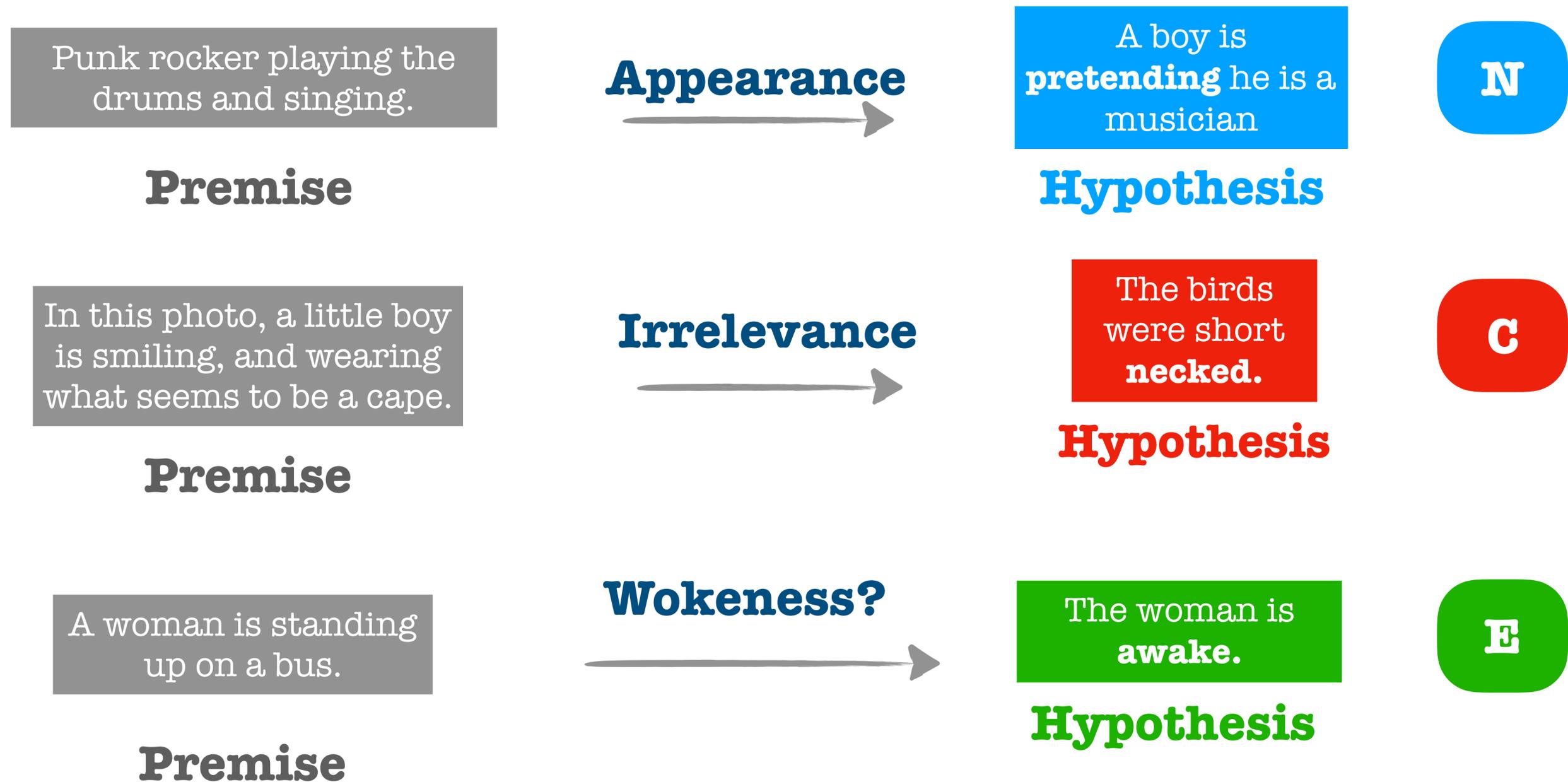
Effect 2: Newer, less common artifacts



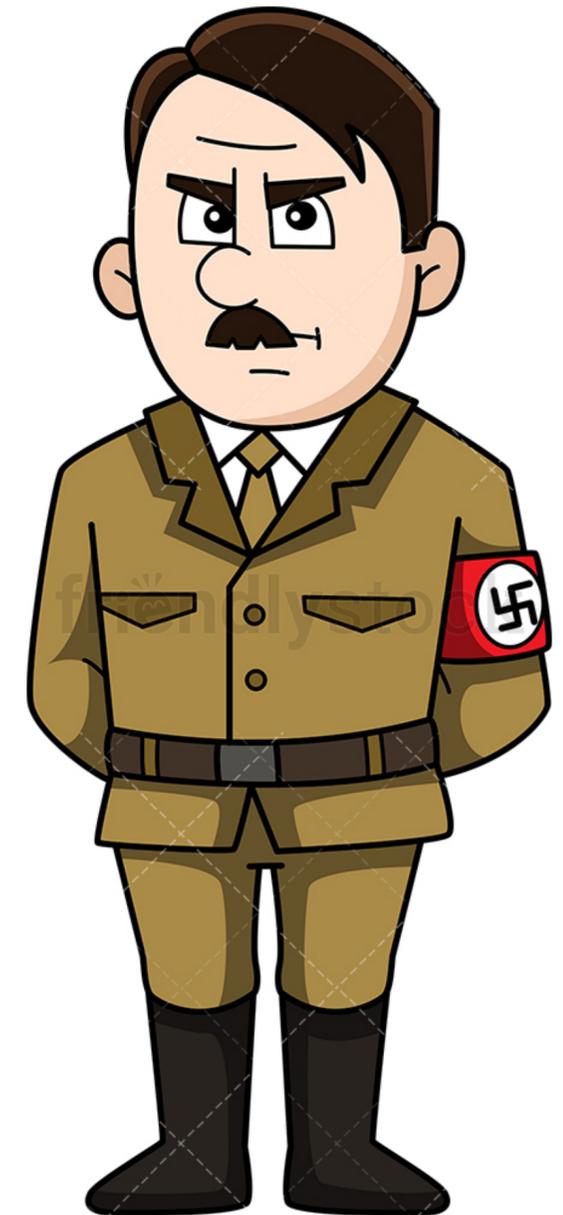
Effect 2: Newer, less common artifacts



Effect 2: Newer, less common artifacts



Effect 3: Ambiguous Artifacts



Effect 3: Ambiguous Artifacts

Two men sit casually in folding chairs, gesturing, and speaking to one another.

Premise

Neutral

Hitler and Stalin share tea and crumpets.

Hypothesis

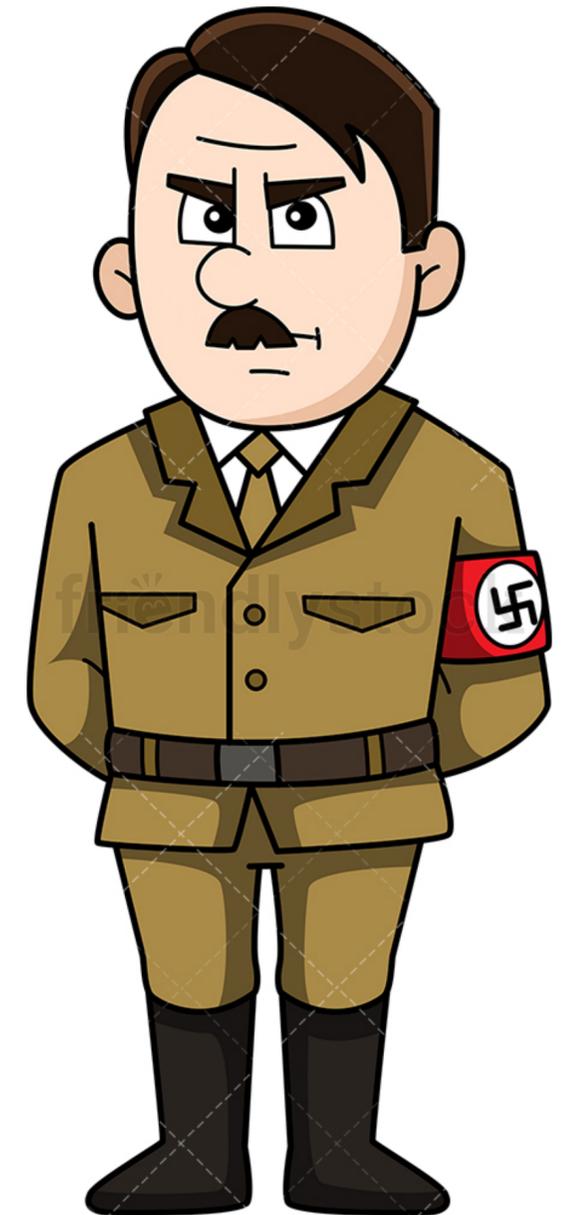
A man in a tie is holding a microphone while the people around him are cheering.

Premise

Contradiction

Hitler gives a speech.

Hypothesis

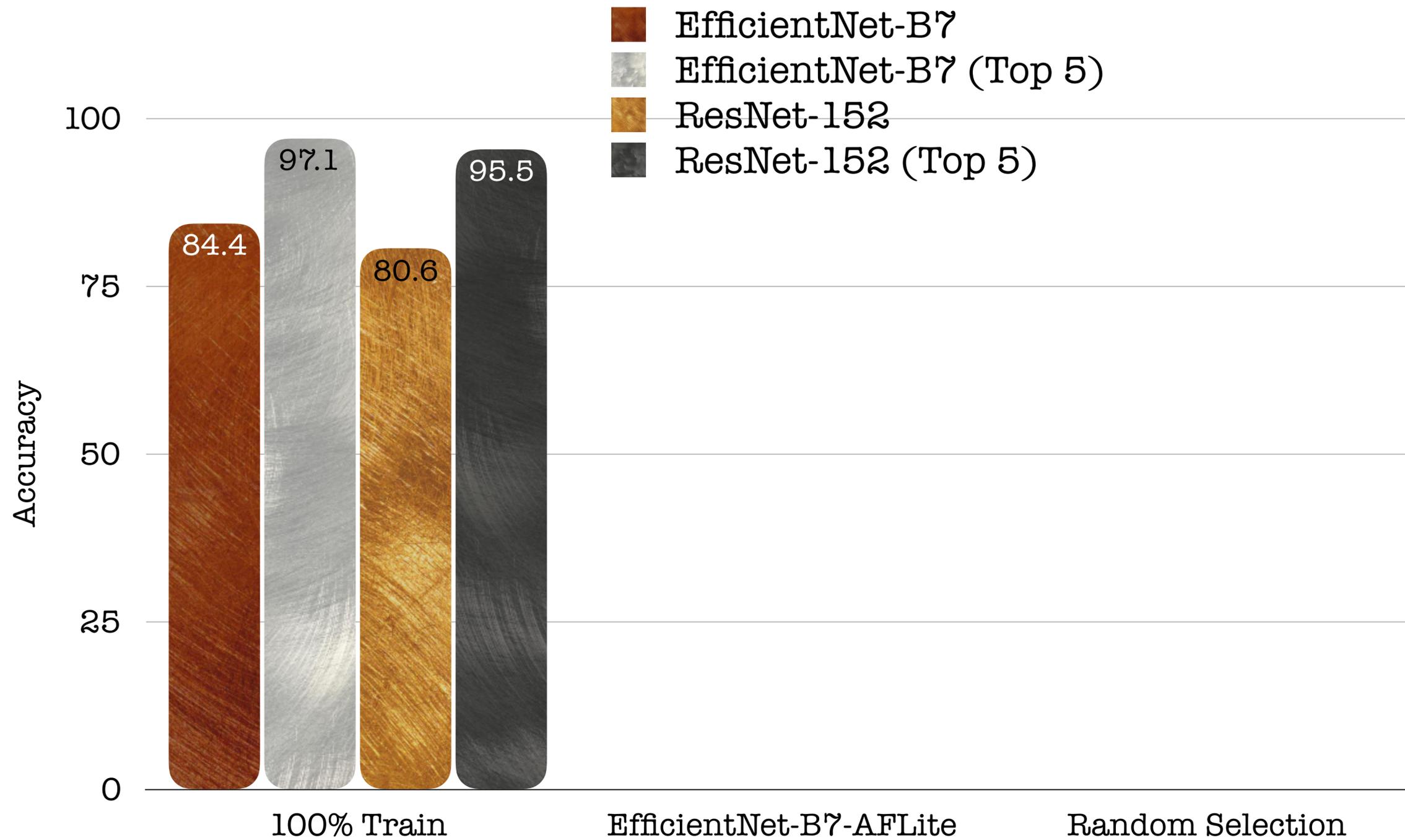


Task 4: Image Classification

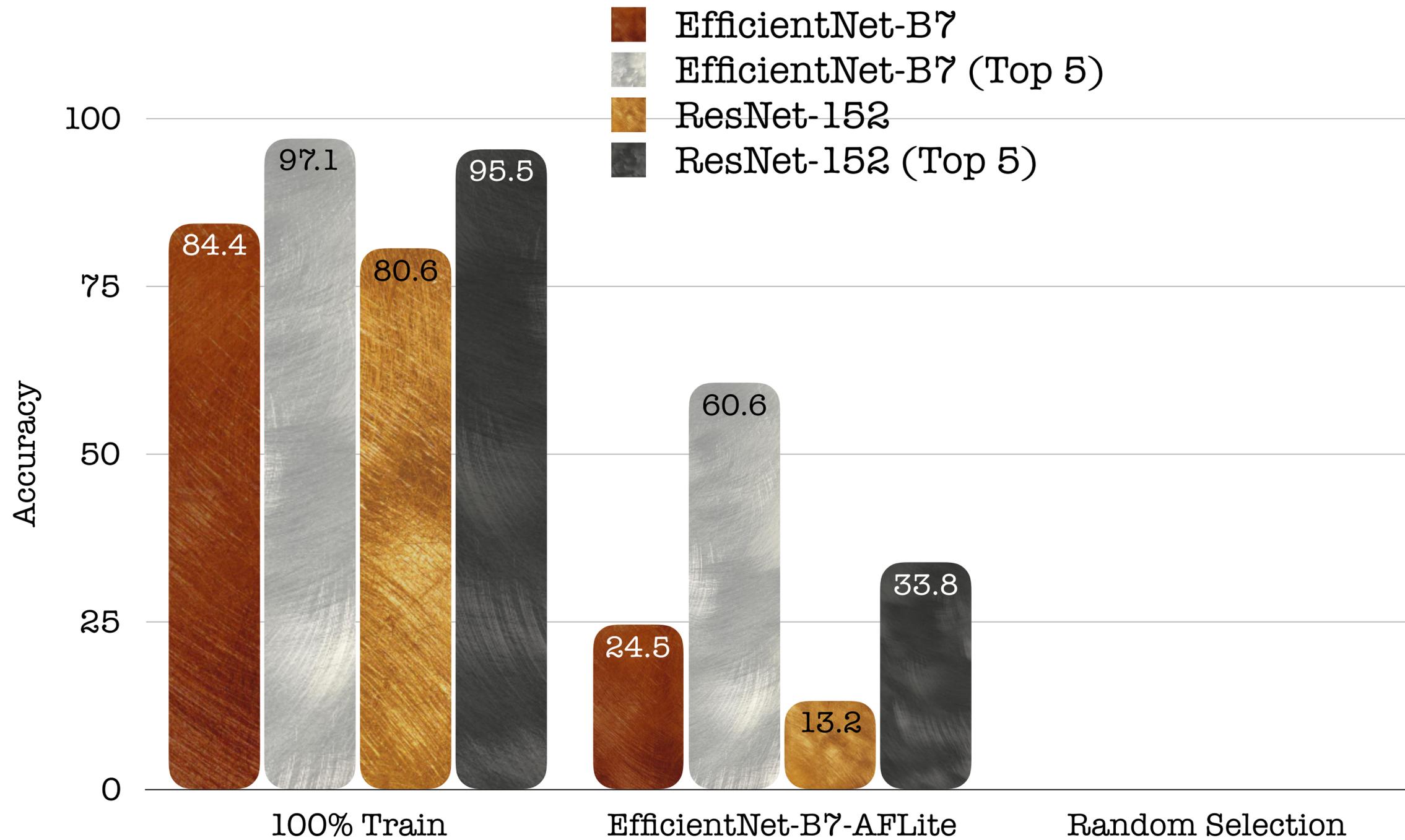
Task 4: Image Classification



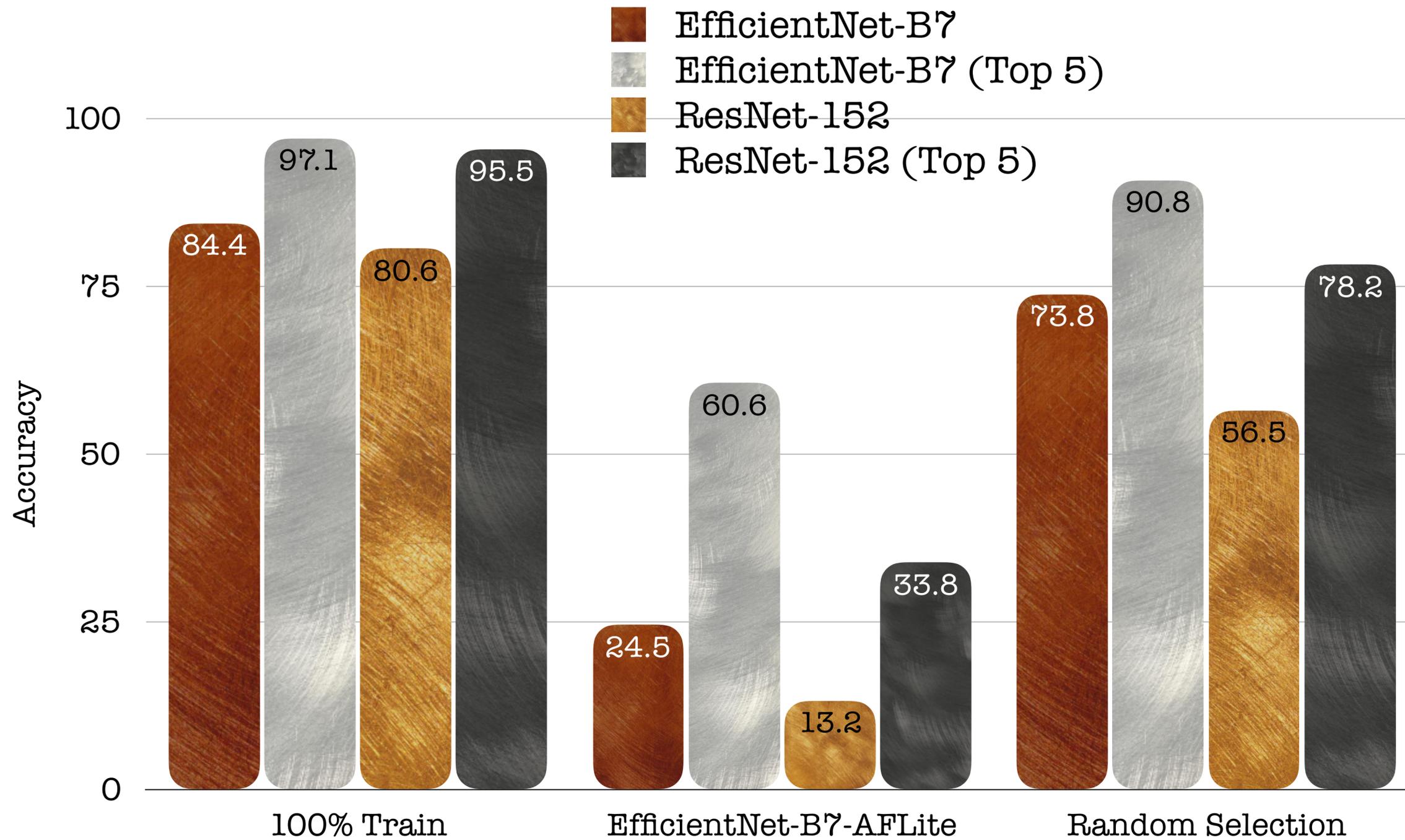
Task 4: Image Classification



Task 4: Image Classification



Task 4: Image Classification



A nearest neighbor perspective

monarch chosen by AFLite



1.0	0.1	0.1	0.1	0.1
0.1	1.0	0.1	0.0	0.1
0.1	0.1	1.0	0.4	0.8
0.1	0.0	0.4	1.0	0.4
0.1	0.1	0.8	0.4	1.0

monarch excluded by AFLite



1.0	0.7	0.7	0.8	0.7
0.7	1.0	0.9	0.8	0.7
0.7	0.9	1.0	0.8	0.7
0.8	0.8	0.8	1.0	0.7
0.7	0.7	0.7	0.7	1.0

chickadee chosen by AFLite



1.0	0.1	0.1	0.2	0.1
0.1	1.0	0.3	0.1	0.1
0.1	0.3	1.0	0.0	0.1
0.2	0.1	0.0	1.0	0.6
0.1	0.1	0.1	0.6	1.0

chickadee excluded by AFLite



1.0	0.5	0.5	0.4	0.5
0.5	1.0	0.8	0.8	0.6
0.5	0.8	1.0	0.7	0.6
0.4	0.8	0.7	1.0	0.6
0.5	0.6	0.6	0.6	1.0



A nearest neighbor perspective

monarch chosen by AFLite



1.0	0.1	0.1	0.1	0.1
0.1	1.0	0.1	0.0	0.1
0.1	0.1	1.0	0.4	0.8
0.1	0.0	0.4	1.0	
0.1	0.1	0.8		

monarch excluded by AFLite



1.0			0.8	0.7
			0.8	0.7
		1.0	0.8	0.7
	0.8	0.8	1.0	0.7
0.7	0.7	0.7	0.7	1.0

chickadee chosen



1.0	0.1	0.1	0.2	0.1
0.1	1.0	0.3	0.1	0.1
0.1	0.3	1.0	0.0	0.1
0.2	0.1	0.0	1.0	0.6
0.1	0.1	0.1	0.6	1.0

chickadee excluded by AFLite



1.0	0.5	0.5	0.4	0.5
0.5	1.0	0.8	0.8	0.6
0.5	0.8	1.0	0.7	0.6
0.4	0.8	0.7	1.0	0.6
0.5	0.6	0.6	0.6	1.0

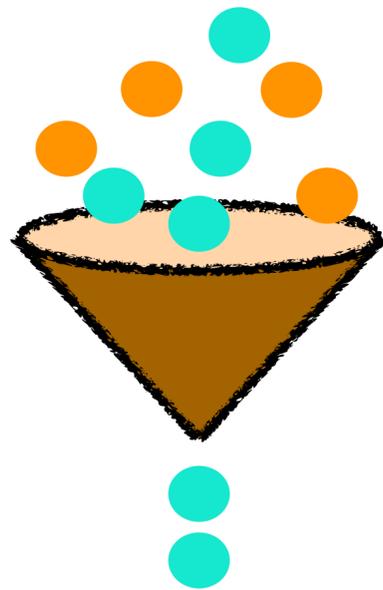
AFLite retains diversity of examples



Takeaways

Takeaways

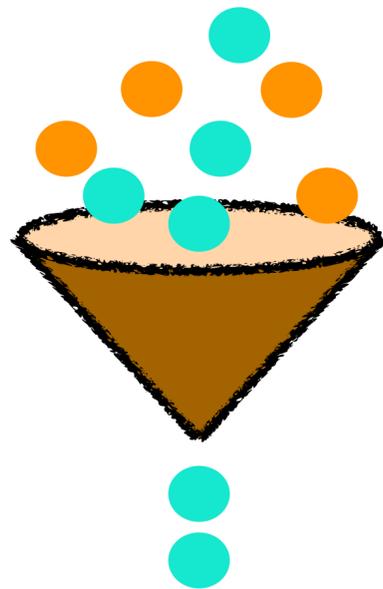
Performance on
some large-scale
datasets could be
misleading.



English-only #BenderRule

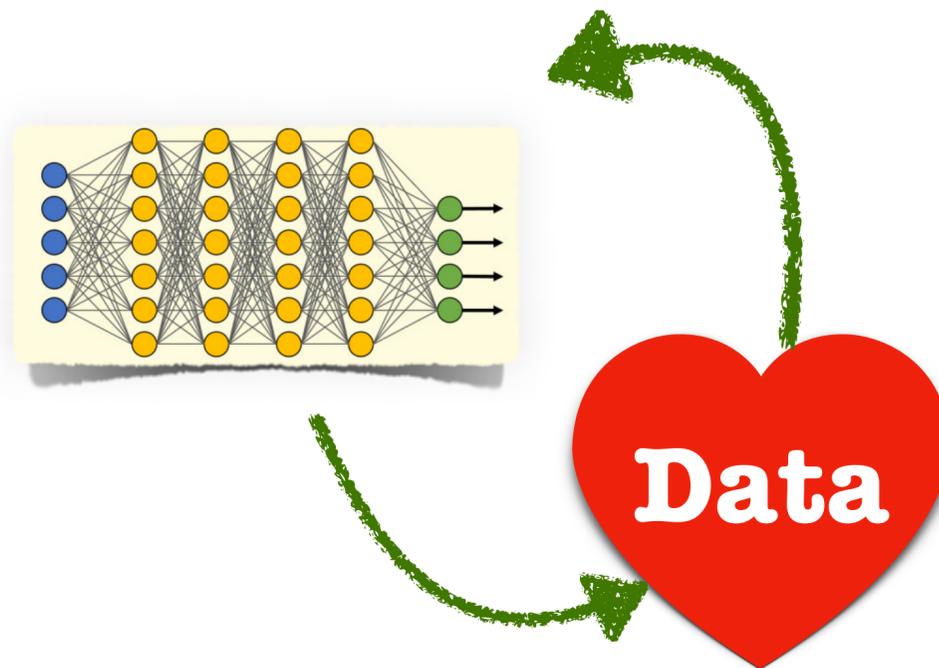
Takeaways

Performance on some large-scale datasets could be misleading.



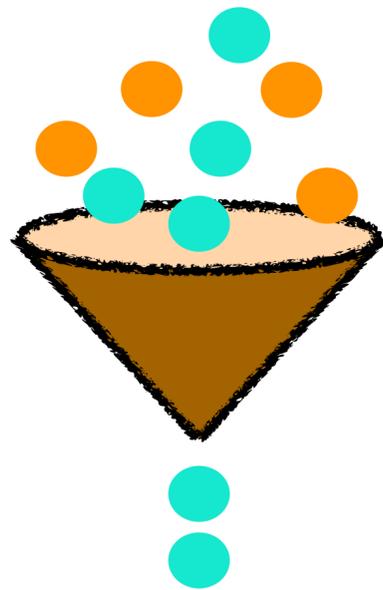
English-only #BenderRule

What makes a good feature representation for adversarial filtering? How to reduce model dependence?



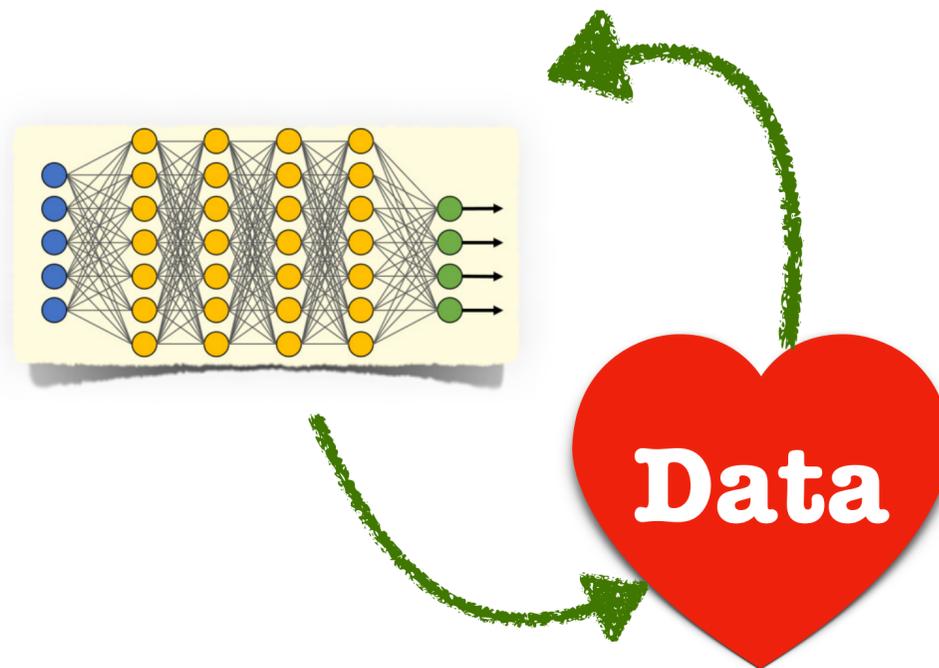
Takeaways

Performance on some large-scale datasets could be misleading.



English-only #BenderRule

What makes a good feature representation for adversarial filtering? How to reduce model dependence?

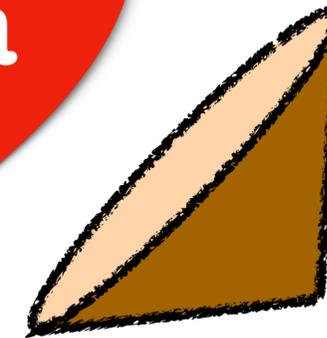
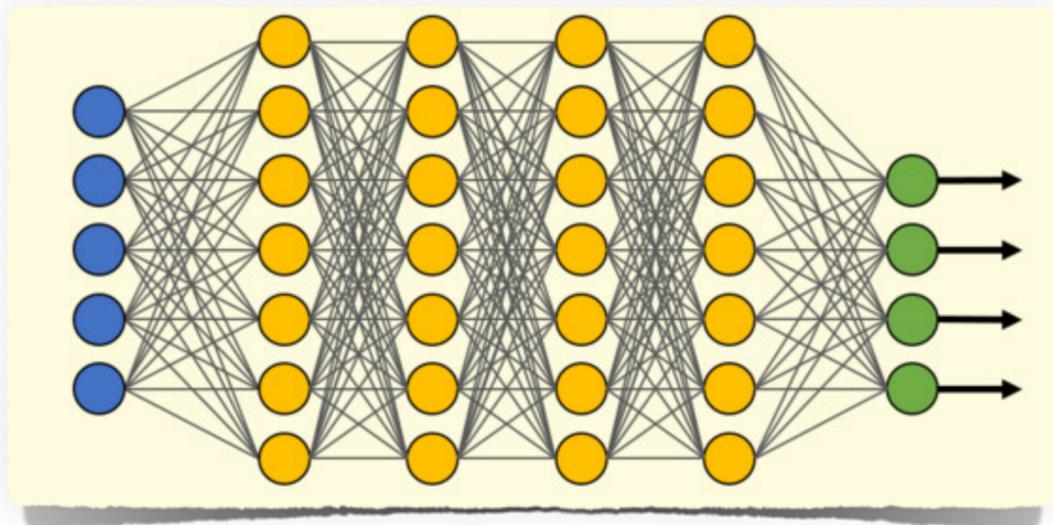


Looking forward: Can we reuse filtered out data?



In Summary

What do predictions tell us about the data?



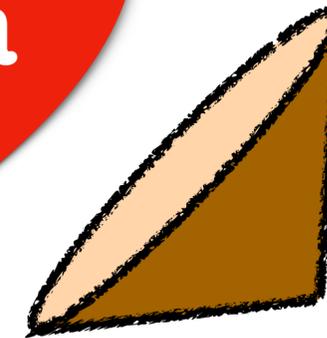
How can we use this information to spruce up our datasets?

In Summary

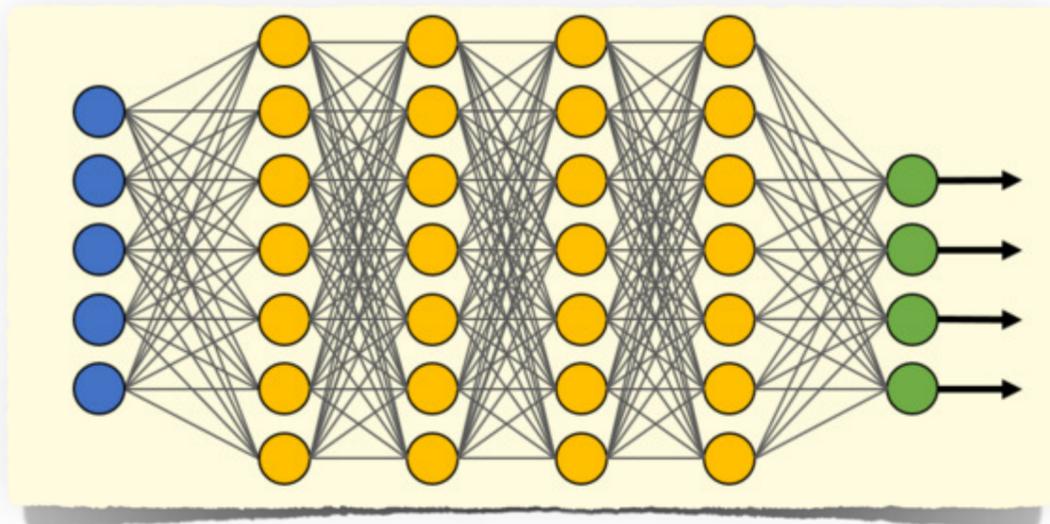
What do predictions tell us about the data?



**Annotation Artifacts
Abound!**



**How can we use this information to
spruce up our datasets?**

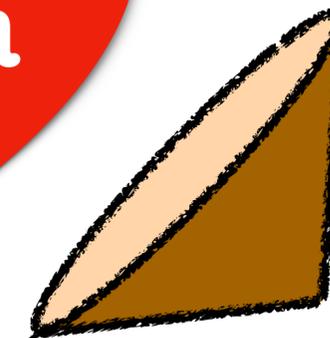


In Summary

What do predictions tell us about the data?



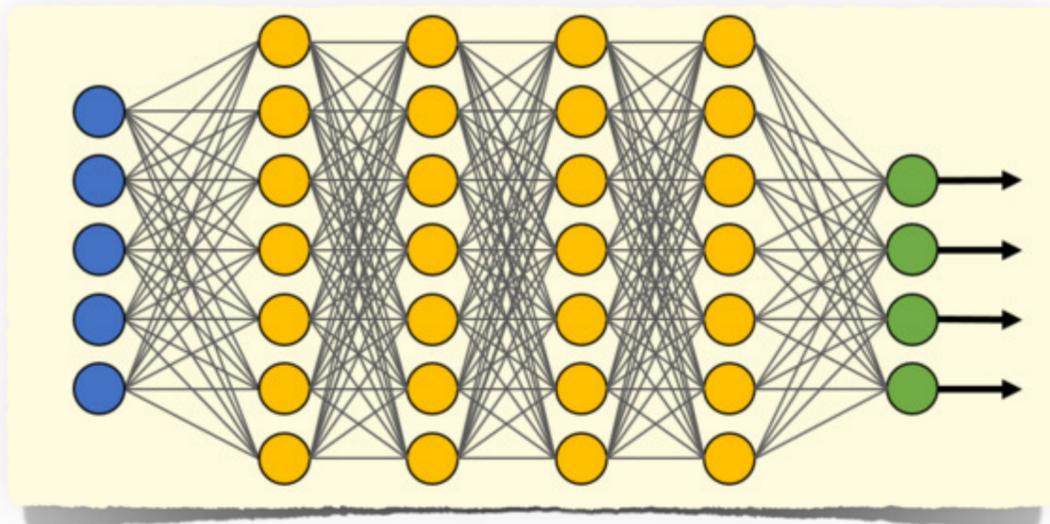
**Annotation Artifacts
Abound!**



**Lightweight Adversarial
Filtering**



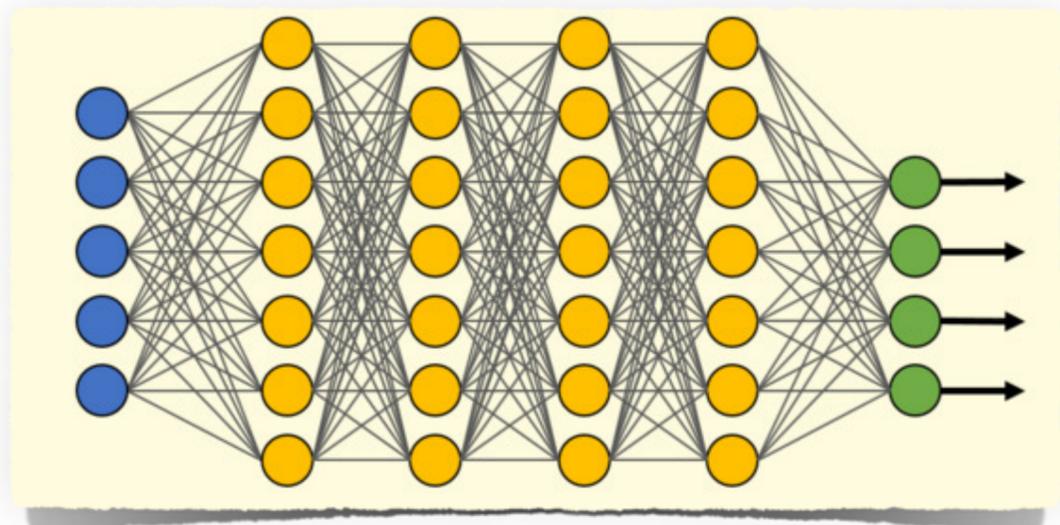
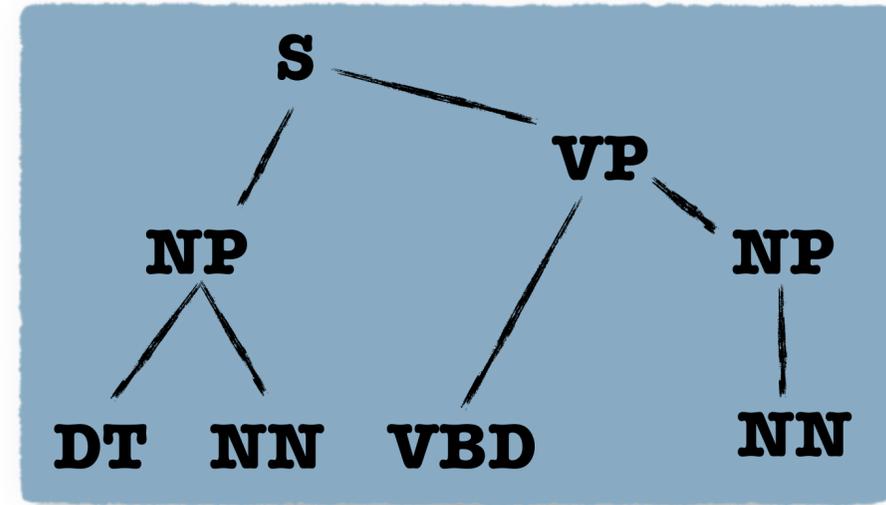
**How can we use this information to
spruce up our datasets?**



Machine Learning



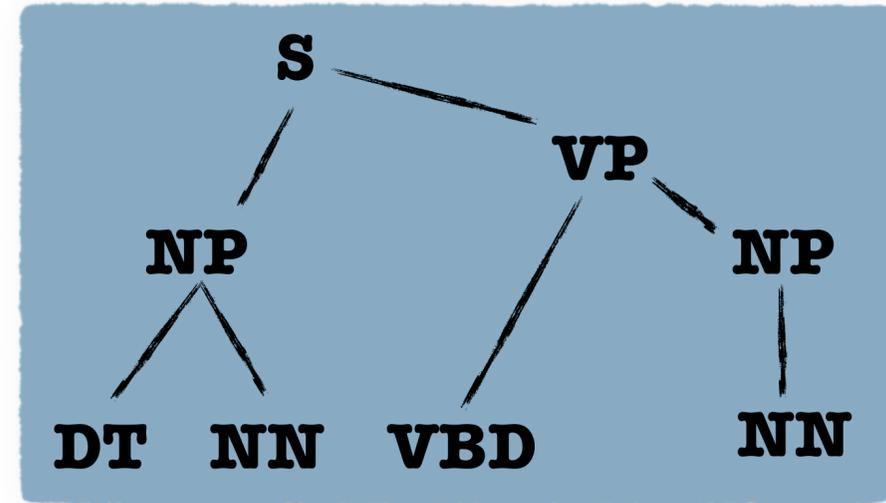
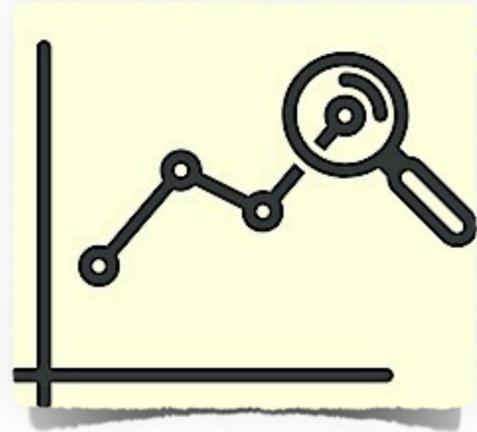
NLP



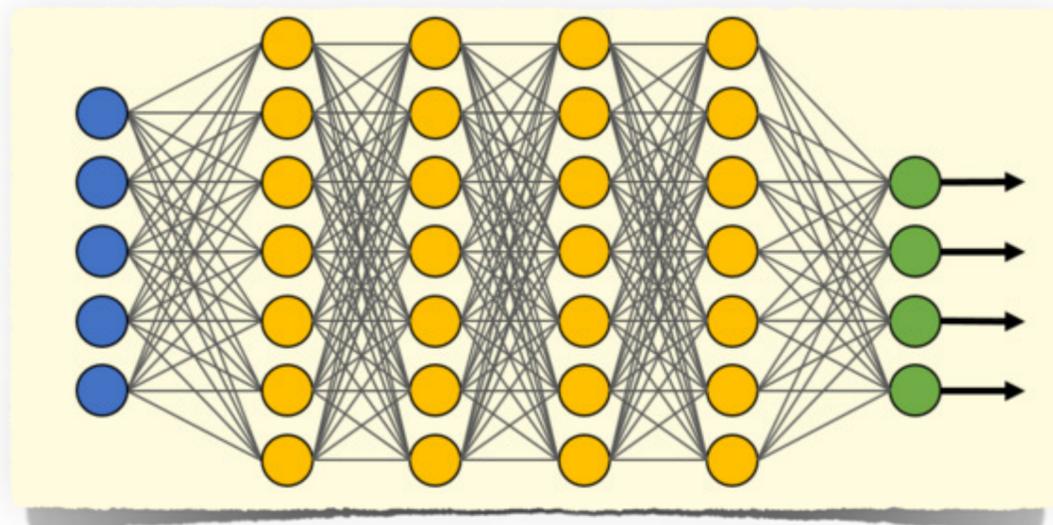
Linguistics

Machine Learning

What does our knowledge of language tell us about models?

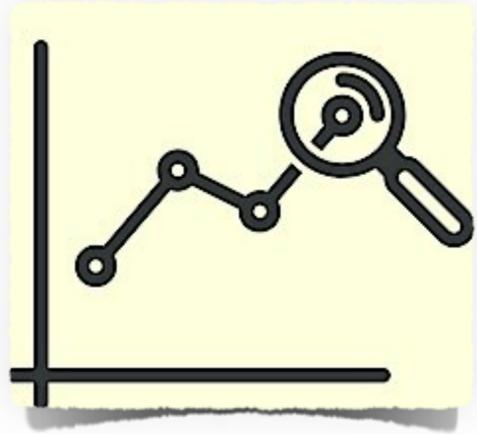


NLP

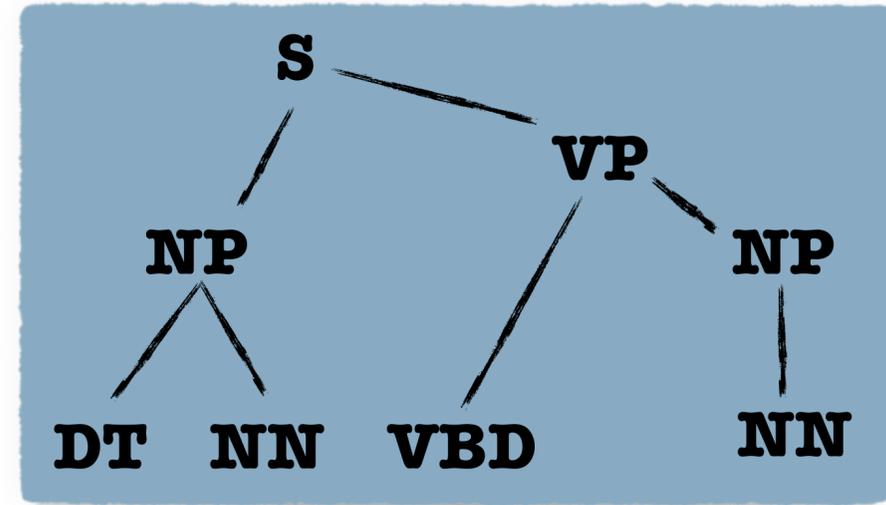


Linguistics

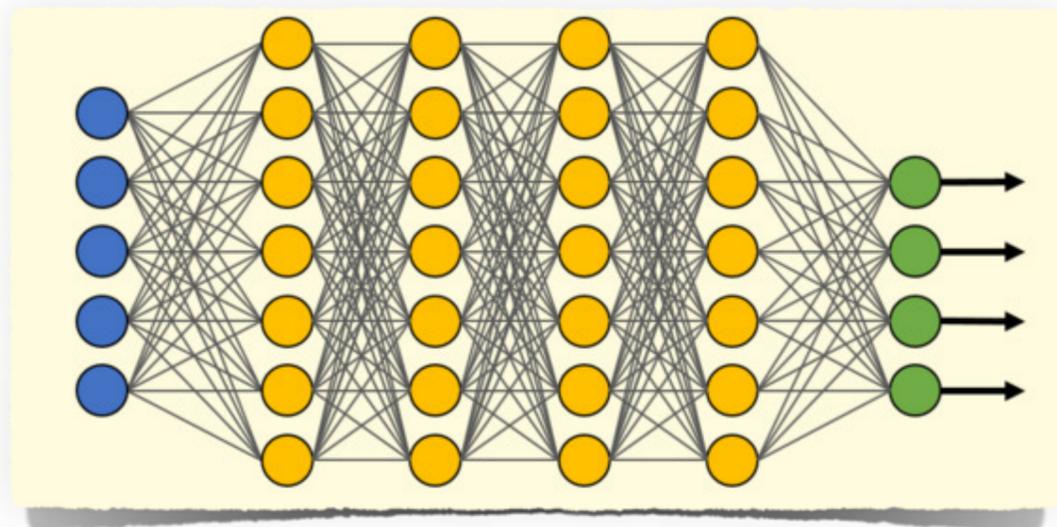
Machine Learning



NLP



Data

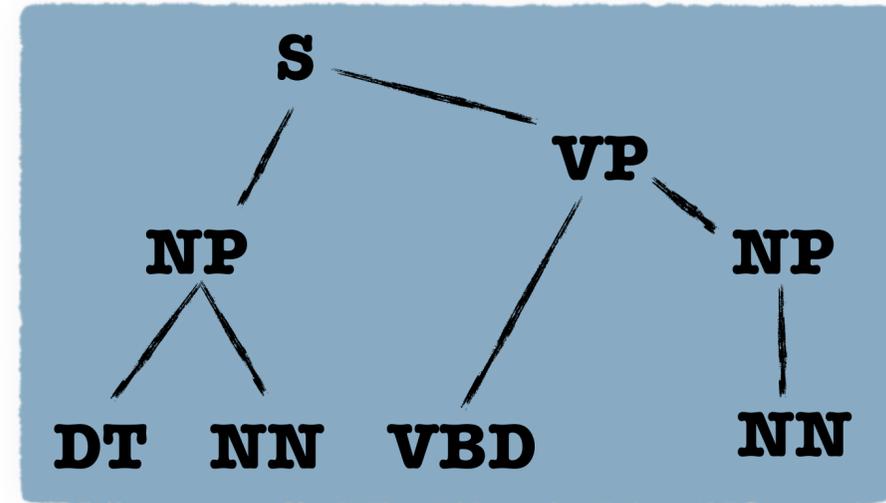


Linguistics

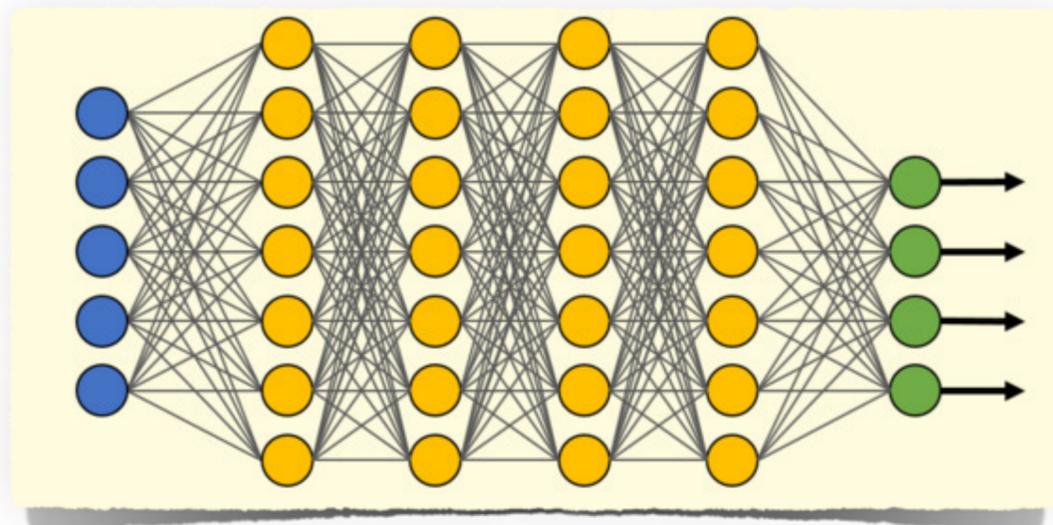
Machine Learning



NLP



Can we improve performance by accounting for linguistic structure?



Linguistics



Sam
Bowman



Suchin
Gururangan



Matt
Peters



Noah A.
Smith



Chandra
Bhagavatula



Ronan
LeBras



Ashish
Sabharwal



Rowan
Zellers



Yejin
Choi



Omer
Levy



Roy
Schwartz

Thanks!



swabhs.com



swabhs



swabhz