# ABOUT ALECTIO
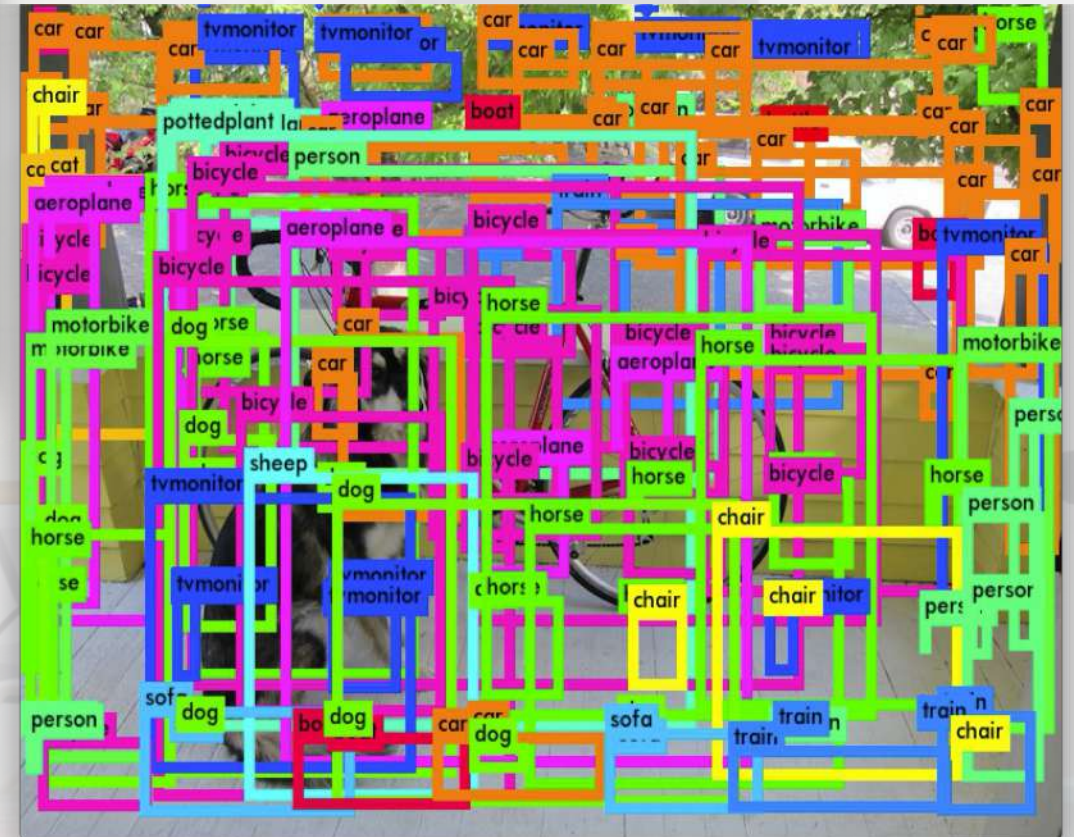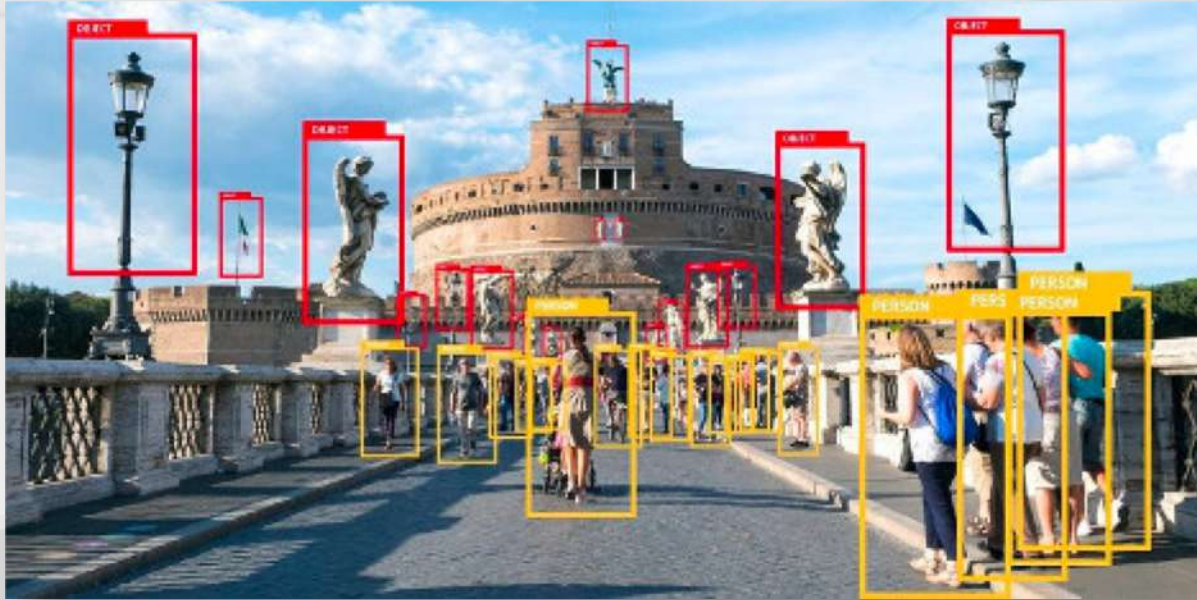
- The first ML company dedicated to **Data Curation**

- Founded in 2019

- Mission:

  Empower ML experts to build, train and retrain models with less data, and hence less resources.

# OUTLINE

- **The Big Data Labeling Crisis**

- **Understanding Class Separation**

- **Not All Data is Created Equal**

- **How to Best Spend your Labeling Budget**

- **Results and Conclusions**

# BIG DATA LABELING CRISIS

# OUR 'TOY' CASE STUDY: CIFAR-10

## The Data

- **CIFAR-10**
- **10 classes of everyday "objects"**
- **50,000 training samples**
- **10,000 testing samples**

## The Model

- **Small CNN**
  - **7 layers**
  - **309,290 total parameters**
  - **308,394 trainable parameters**
  - **896 non-trainable parameters**

# BASELINE RESULTS

## Results

- **Baseline accuracy**: **89%** **(across all classes)**



Confusion Marix with 50K Data Samples & 0% Noise

|  | airplane | automobile | bird | cat | deer | dog | frog | horse | ship | truck |
|---|---|---|---|---|---|---|---|---|---|---|
| **airplane** | 897 | 2 | 14 | 11 | 3 | 5 | 4 | 8 | 21 | 9 |
| **automobile** | 13 | 966 | 0 | 7 | 1 | 5 | 2 | 1 | 19 | 41 |
| **bird** | 25 | 0 | 889 | 42 | 36 | 29 | 17 | 10 | 3 | 2 |
| **cat** | 4 | 0 | 15 | 706 | 13 | 77 | 5 | 7 | 2 | 0 |
| **deer** | 3 | 0 | 14 | 32 | 875 | 17 | 0 | 12 | 0 | 0 |
| **dog** | 0 | 0 | 8 | 75 | 6 | 790 | 0 | 9 | 0 | 0 |
| **frog** | 1 | 1 | 41 | 82 | 38 | 32 | 963 | 5 | 3 | 3 |
| **horse** | 6 | 0 | 10 | 19 | 26 | 35 | 3 | 944 | 1 | 1 |
| **ship** | 29 | 3 | 4 | 6 | 1 | 1 | 4 | 1 | 928 | 8 |
| **truck** | 22 | 28 | 5 | 20 | 1 | 9 | 2 | 3 | 23 | 936 |

## Results

- **Baseline accuracy**: **89%** (across all classes)



Confusion Marix with 50K Data Samples & 0% Noise

|  | airplane | automobile | bird | cat | deer | dog | frog | horse | ship | truck |
|---|---|---|---|---|---|---|---|---|---|---|
| airplane | 897 | 2 | 14 | 11 | 3 | 5 | 4 | 8 | 21 | 9 |
| automobile | 13 | 966 | 0 | 7 | 1 | 5 | 2 | 1 | 19 | 41 |
| bird | 25 | 0 | 889 | 42 | 36 | 29 | 17 | 10 | 3 | 2 |
| cat | 4 | 0 | 15 | 706 | 13 | 77 | 5 | 7 | 2 | 0 |
| deer | 3 | 0 | 14 | 32 | 875 | 17 | 0 | 12 | 0 | 0 |
| dog | 0 | 0 | 8 | 75 | 6 | 790 | 0 | 9 | 0 | 0 |
| frog | 1 | 1 | 41 | 82 | 38 | 32 | 963 | 5 | 3 | 3 |
| horse | 6 | 0 | 10 | 19 | 26 | 35 | 3 | 944 | 1 | 1 |
| ship | 29 | 3 | 4 | 6 | 1 | 1 | 4 | 1 | 928 | 8 |
| truck | 22 | 28 | 5 | 20 | 1 | 9 | 2 | 3 | 23 | 936 |

*Ground truth* →

# BASELINE RESULTS

## Results

- **Baseline accuracy**: **89%** (**across all classes**)



Confusion Marix with 50K Data Samples & 0% Noise

*Predictions*

## Results

- **Baseline accuracy**: **89%** (**across all classes**)

**True Positive Rate**

**(x 1000)**



Confusion Marix with 50K Data Samples & 0% Noise

| | airplane | automobile | bird | cat | deer | dog | frog | horse | ship | truck |
|---|---|---|---|---|---|---|---|---|---|---|
| **airplane** | 897 | 2 | 14 | 11 | 3 | 5 | 4 | 8 | 21 | 9 |
| **automobile** | 13 | 966 | 0 | 7 | 1 | 5 | 2 | 1 | 19 | 41 |
| **bird** | 25 | 0 | 889 | 42 | 36 | 29 | 17 | 10 | 3 | 2 |
| **cat** | 4 | 0 | 15 | 706 | 13 | 77 | 5 | 7 | 2 | 0 |
| **deer** | 3 | 0 | 14 | 32 | 875 | 17 | 0 | 12 | 0 | 0 |
| **dog** | 0 | 0 | 8 | 75 | 6 | 790 | 0 | 9 | 0 | 0 |
| **frog** | 1 | 1 | 41 | 82 | 38 | 32 | 963 | 5 | 3 | 3 |
| **horse** | 6 | 0 | 10 | 19 | 26 | 35 | 3 | 944 | 1 | 1 |
| **ship** | 29 | 3 | 4 | 6 | 1 | 1 | 4 | 1 | 928 | 8 |
| **truck** | 22 | 28 | 5 | 20 | 1 | 9 | 2 | 3 | 23 | 936 |

# BASELINE RESULTS

## Results

- **Baseline accuracy**: **89%** (across all classes)

*False Positive Rate*



Confusion Marix with 50K Data Samples & 0% Noise

|  | airplane | automobile | bird | cat | deer | dog | frog | horse | ship | truck |
|---|---|---|---|---|---|---|---|---|---|---|
| airplane | 897 | 2 | 14 | 11 | 3 | 5 | 4 | 8 | 21 | 9 |
| automobile | 13 | 966 | 0 | 7 | 1 | 5 | 2 | 1 | 19 | 41 |
| bird | 25 | 0 | 889 | 42 | 36 | 29 | 17 | 10 | 3 | 2 |
| cat | 4 | 0 | 15 | 706 | 13 | 77 | 5 | 7 | 2 | 0 |
| deer | 3 | 0 | 14 | 32 | 875 | 17 | 0 | 12 | 0 | 0 |
| dog | 0 | 0 | 8 | 75 | 6 | 790 | 0 | 9 | 0 | 0 |
| frog | 1 | 1 | 41 | 82 | 38 | 32 | 963 | 5 | 3 | 3 |
| horse | 6 | 0 | 10 | 19 | 26 | 35 | 3 | 944 | 1 | 1 |
| ship | 29 | 3 | 4 | 6 | 1 | 1 | 4 | 1 | 928 | 8 |
| truck | 22 | 28 | 5 | 20 | 1 | 9 | 2 | 3 | 23 | 936 |

## Results

- **Baseline accuracy**: **89%** (**across all classes**)

**False Negative Rate**



Confusion Marix with 50K Data Samples & 0% Noise

# BASELINE RESULTS

## Results

- **Baseline accuracy**: **89%** (**across all classes**)



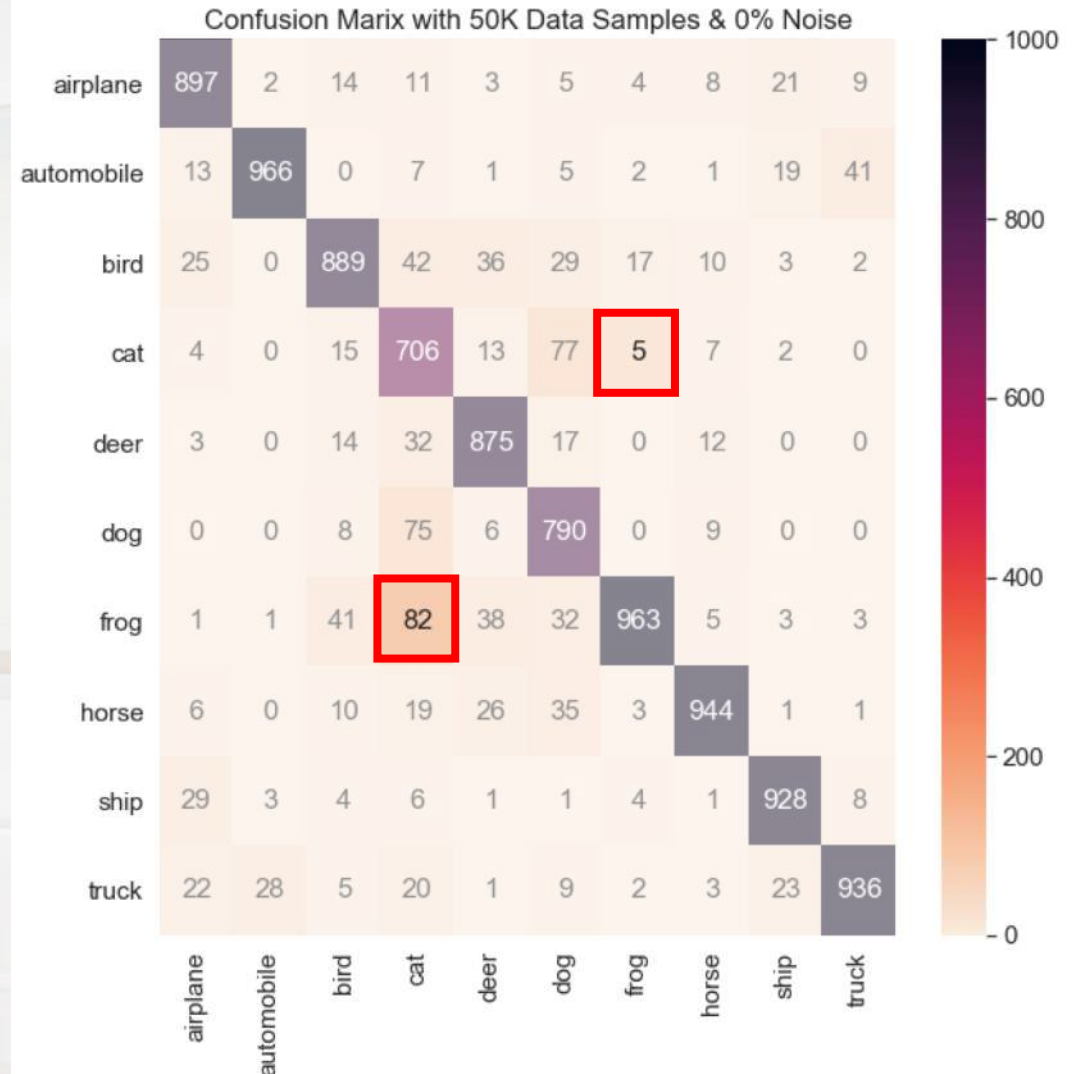Confusion Marix with 50K Data Samples & 0% Noise

*True Negative Rate*

# BASELINE RESULTS

## Results

- **Baseline accuracy**: **89% (across all classes)**
- **Accuracy varies dramatically across classes**



Confusion Marix with 50K Data Samples & 0% Noise

# BASELINE RESULTS

## Results

- **Baseline accuracy**: **89%** (across all classes)
- **Accuracy varies dramatically across classes**

**More details...**

- **Lowest accuracy for class 'cat' and 'dog'**



Confusion Marix with 50K Data Samples & 0% Noise

# BASELINE RESULTS

## Results

- **Baseline accuracy**: **89% (across all classes)**
- **Accuracy varies dramatically across classes**

**More details...**

- **Lowest accuracy for class 'cat' and 'dog'**
- **Class 'bird' has a fairly high accuracy**



Confusion Marix with 50K Data Samples & 0% Noise

# BASELINE RESULTS

## Results

- **Baseline accuracy**: **89%** (**across all classes**)
- **Accuracy varies dramatically across classes**

**More details…**

- **Lowest accuracy for class 'cat' and 'dog'**

- **Class 'bird' has a fairly high accuracy**

- **Higher confusion for 'cat' ➔ 'frog' and for 'cat' ➔ 'dog'**



Confusion Marix with 50K Data Samples & 0% Noise

|  | airplane | automobile | bird | cat | deer | dog | frog | horse | ship | truck |
|---|---|---|---|---|---|---|---|---|---|---|
| airplane | 897 | 2 | 14 | 11 | 3 | 5 | 4 | 8 | 21 | 9 |
| automobile | 13 | 966 | 0 | 7 | 1 | 5 | 2 | 1 | 19 | 41 |
| bird | 25 | 0 | 889 | 42 | 36 | 29 | 17 | 10 | 3 | 2 |
| cat | 4 | 0 | 15 | 706 | 13 | 77 | 5 | 7 | 2 | 0 |
| deer | 3 | 0 | 14 | 32 | 875 | 17 | 0 | 12 | 0 | 0 |
| dog | 0 | 0 | 8 | 75 | 6 | 790 | 0 | 9 | 0 | 0 |
| frog | 1 | 1 | 41 | 82 | 38 | 32 | 963 | 5 | 3 | 3 |
| horse | 6 | 0 | 10 | 19 | 26 | 35 | 3 | 944 | 1 | 1 |
| ship | 29 | 3 | 4 | 6 | 1 | 1 | 4 | 1 | 928 | 8 |
| truck | 22 | 28 | 5 | 20 | 1 | 9 | 2 | 3 | 23 | 936 |

# BASELINE RESULTS

## Results

- **Baseline accuracy**: **89%** (**across all classes**)
- **Accuracy varies dramatically across classes**

**More details...**

- **Lowest accuracy for class 'cat' and 'dog'**
- **Class 'bird' has a fairly high accuracy**
- **Higher confusion for 'cat' ➔ 'frog' and for 'cat' ➔ 'dog'**
- **As easy to mistake a cat for a dog, than a dog for a cat**



Confusion Marix with 50K Data Samples & 0% Noise

|  | airplane | automobile | bird | cat | deer | dog | frog | horse | ship | truck |
|---|---|---|---|---|---|---|---|---|---|---|
| airplane | 897 | 2 | 14 | 11 | 3 | 5 | 4 | 8 | 21 | 9 |
| automobile | 13 | 966 | 0 | 7 | 1 | 5 | 2 | 1 | 19 | 41 |
| bird | 25 | 0 | 889 | 42 | 36 | 29 | 17 | 10 | 3 | 2 |
| cat | 4 | 0 | 15 | 706 | 13 | 77 | 5 | 7 | 2 | 0 |
| deer | 3 | 0 | 14 | 32 | 875 | 17 | 0 | 12 | 0 | 0 |
| dog | 0 | 0 | 8 | 75 | 6 | 790 | 0 | 9 | 0 | 0 |
| frog | 1 | 1 | 41 | 82 | 38 | 32 | 963 | 5 | 3 | 3 |
| horse | 6 | 0 | 10 | 19 | 26 | 35 | 3 | 944 | 1 | 1 |
| ship | 29 | 3 | 4 | 6 | 1 | 1 | 4 | 1 | 928 | 8 |
| truck | 22 | 28 | 5 | 20 | 1 | 9 | 2 | 3 | 23 | 936 |

# BASELINE RESULTS

## Results

- **Baseline accuracy**: **89%** (**across all classes**)
- **Accuracy varies dramatically across classes**

**More details...**

- **Lowest accuracy for class 'cat' and 'dog'**

- **Class 'bird' has a fairly high accuracy**

- **Higher confusion for 'cat' ➔ 'frog' and for 'cat' ➔ 'dog'**

- **As easy to mistake a cat for a dog, than a dog for a cat**

- **Easier to mistake a cat for a frog, than a frog for a cat**

- **Confusion is NOT SYMMETRICAL across classes**



Confusion Marix with 50K Data Samples & 0% Noise

|  | airplane | automobile | bird | cat | deer | dog | frog | horse | ship | truck |
|---|---|---|---|---|---|---|---|---|---|---|
| airplane | 897 | 2 | 14 | 11 | 3 | 5 | 4 | 8 | 21 | 9 |
| automobile | 13 | 966 | 0 | 7 | 1 | 5 | 2 | 1 | 19 | 41 |
| bird | 25 | 0 | 889 | 42 | 36 | 29 | 17 | 10 | 3 | 2 |
| cat | 4 | 0 | 15 | 706 | 13 | 77 | 5 | 7 | 2 | 0 |
| deer | 3 | 0 | 14 | 32 | 875 | 17 | 0 | 12 | 0 | 0 |
| dog | 0 | 0 | 8 | 75 | 6 | 790 | 0 | 9 | 0 | 0 |
| frog | 1 | 1 | 41 | 82 | 38 | 32 | 963 | 5 | 3 | 3 |
| horse | 6 | 0 | 10 | 19 | 26 | 35 | 3 | 944 | 1 | 1 |
| ship | 29 | 3 | 4 | 6 | 1 | 1 | 4 | 1 | 928 | 8 |
| truck | 22 | 28 | 5 | 20 | 1 | 9 | 2 | 3 | 23 | 936 |

# EXPERIMENT #1: LABELING POLLUTION

## Goal:

**Study impact of noise in labeling process on model performance**

## Protocol:

- **We randomly shuffle the labels within the selected subset**

- **We select n% of the 50,000 records (full dataset)**
  - **Those records are chosen randomly, with no distinction of the class**

- **We repeat the same experiment 5 times for each amount to eliminate noisy results**
  - **Different levels of noise of data might lead to different results**
  - **We chose 5 times because of compute power limitations**

- **We observe the accuracy and the confusion matrix**

# EXPERIMENT #1: LABELING POLLUTION



Confusion Matrix with 5% Noise

**Average Confusion Matrix with 5% noisy labels**

# EXPERIMENT #1: LABELING POLLUTION



Confusion Marix with 10% Noise

**Average Confusion Matrix with 10% noisy labels**

# EXPERIMENT #1: LABELING POLLUTION



Confusion Marix with 15% Noise

**Average Confusion Matrix with 15% noisy labels**

# EXPERIMENT #1: LABELING POLLUTION



Confusion Matrix with 20% Noise

**Average Confusion Matrix with 20% noisy labels**

# EXPERIMENT #1: LABELING POLLUTION



Confusion Matrix with 25% Noise

**Average Confusion Matrix with 25% noisy labels**

# EXPERIMENT #1: LABELING POLLUTION



Confusion Matrix with 30% Noise

**Average Confusion Matrix with 30% noisy labels**

# EXPERIMENT #1: LABELING POLLUTION



Confusion Matrix with 5% Noise

Confusion Matrix Difference from Baseline with 5% Noise

# EXPERIMENT #1: LABELING POLLUTION

# EXPERIMENT #1: LABELING POLLUTION

# EXPERIMENT #1: LABELING POLLUTION



Confusion {dog ➔ cat}
vs. labeling noise level

Confusion {cat ➔ frog}
vs. labeling noise level

# EXPERIMENT #1: LABELING POLLUTION

## Results

- **Accuracy seems to drop linearly with the amount of noise in the labels**



Noise vs Accuracy

# EXPERIMENT #2: DATA VOLUME REDUCTION

## Goal:

**Study impact of size of training set on model performance**

## Protocol:

- **We increase the size of the training set from 5,000 records (10%) to 50,000 records (full dataset)**
  - **Those records are chosen randomly**

- **We repeat the same experiment 5 times for each amount to eliminate noisy results**
  - **Different subsets of data might lead to different results**
  - **We chose 5 times because of compute power limitations**

- **We report the accuracy and the confusion matrix**

# EXPERIMENT #2: DATA VOLUME REDUCTION



Confusion Marix with 5K Data Samples

**Average Confusion Matrix with size 5k samples**

# EXPERIMENT #2: DATA VOLUME REDUCTION



Confusion Marix with 10K Data Samples

**Average Confusion Matrix with size 10k samples**

Confusion Marix with 20K Data Samples

**Average Confusion Matrix with size 20k samples**

# EXPERIMENT #2: DATA VOLUME REDUCTION



Confusion Marix with 30K Data Samples

**Average Confusion Matrix with size 30k samples**

# EXPERIMENT #2: DATA VOLUME REDUCTION


Confusion Matrix with 40K Data Samples

**Average Confusion Matrix with size 40k samples**
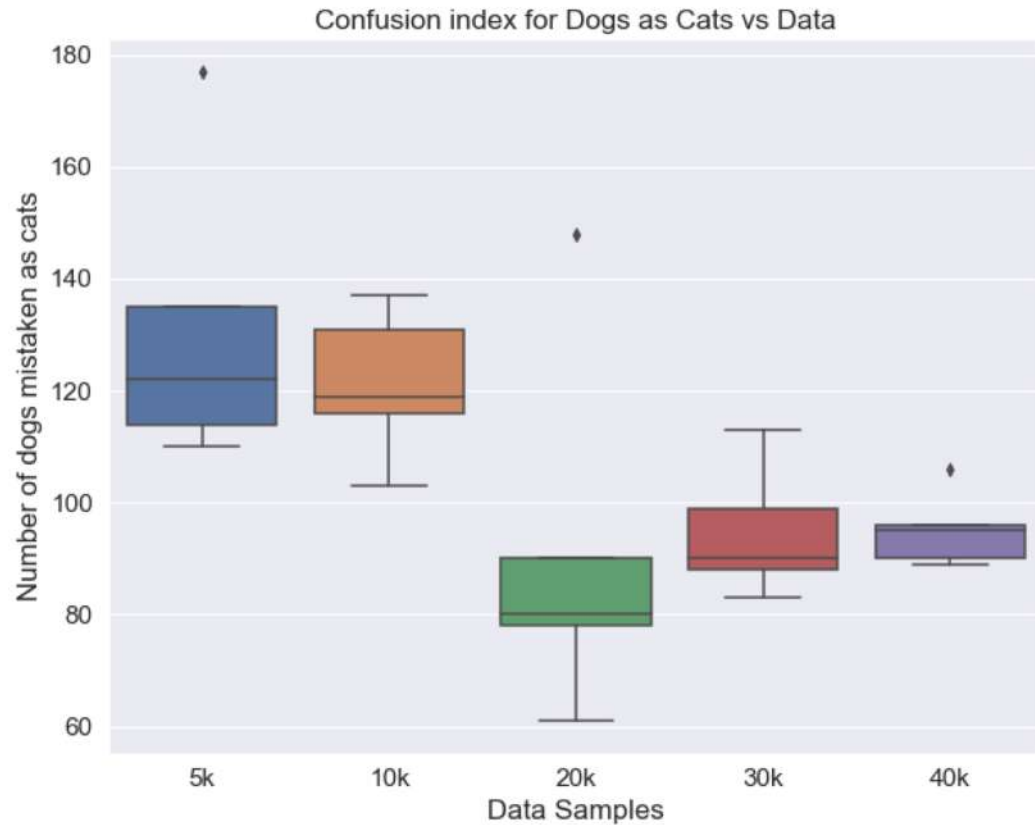
# EXPERIMENT #2: DATA VOLUME REDUCTION



Confusion Matrix with 10K Data Samples

Confusion Matrix Difference from Baseline with 10K Data Samples

# EXPERIMENT #2: DATA VOLUME REDUCTION



Confusion Marix with 20K Data Samples

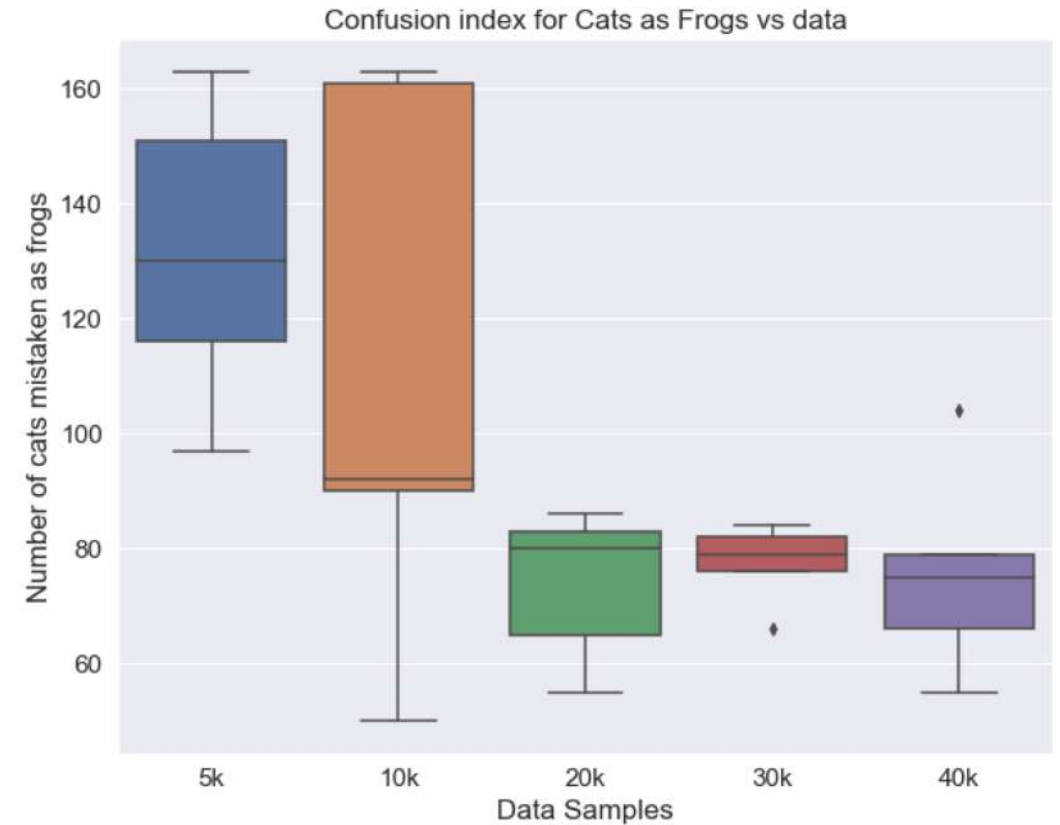Confusion Matrix Difference from Baseline with 20K Data Samples

Confusion Marix with 30K Data Samples

Confusion Matrix Difference from Baseline with 30K Data Samples

# EXPERIMENT #2: DATA VOLUME REDUCTION



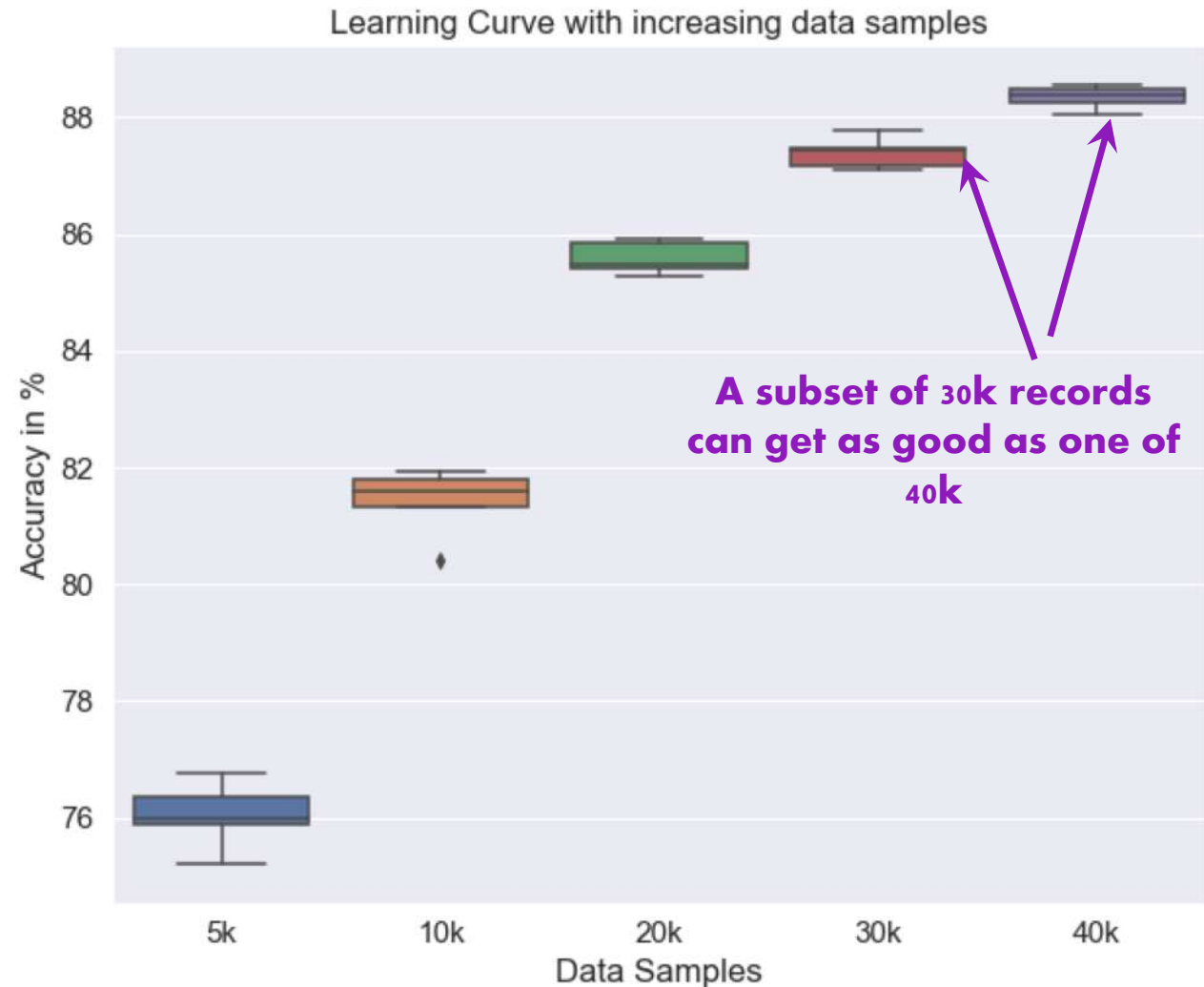**Confusion {dog ➔ cat}
vs. volume of training data**

**Confusion {cat ➔ frog}
vs. volume of training data**

## Results

- **Accuracy gets asymptotically better with more data**

- **10k gets us more than 90% of the way there**

- **20k (less than half) gets us 95+% of the way there**

- **The best sample of size 30k gets similar accuracy to the worse one with size 40k**



Learning Curve with increasing data samples

**A subset of 30k records can get as good as one of 40k**

TIME TO DRAW SOME (REAL) CONCLUSIONS

# DISCUSSION: ARE ALL CLASSES EQUALLY IMPACTED?

## A FEW CONCLUSIONS...

**Lowest accuracy** → **Highest accuracy**

| Baseline | | 30% Labeling Noise | | 5K Data Samples | |
|---|---|---|---|---|---|
| Cat | 706 | Cat | 608 | Cat | 520 |
| Dog | 790 | Dog | 698 | Bird | 594 |
| Deer | 875 | Bird | 744 | Dog | 599 |
| Bird | 889 | Deer | 811 | Deer | 716 |
| Airplane | 897 | Airplane | 837 | Airplane | 789 |
| Ship | 928 | Horse | 877 | Horse | 830 |
| Truck | 936 | Ship | 897 | Ship | 853 |
| Horse | 944 | Automobile | 922 | Truck | 887 |
| Frog | 963 | Truck | 928 | Automobile | 905 |
| Automobile | 966 | Frog | 951 | Frog | 914 |

# DISCUSSION: ARE ALL CLASSES EQUALLY IMPACTED?

## A FEW CONCLUSIONS…

- **'Cat'** is the least accurate class even with labeling noise and data quantity

Lowest accuracy → Highest accuracy

| Baseline | | 30% Labeling Noise | | 5K Data Samples | |
|---|---|---|---|---|---|
| Cat | 706 | Cat | 608 | Cat | 520 |
| Dog | 790 | Dog | 698 | Bird | 594 |
| Deer | 875 | Bird | 744 | Dog | 599 |
| Bird | 889 | Deer | 811 | Deer | 716 |
| Airplane | 897 | Airplane | 837 | Airplane | 789 |
| Ship | 928 | Horse | 877 | Horse | 830 |
| Truck | 936 | Ship | 897 | Ship | 853 |
| Horse | 944 | Automobile | 922 | Truck | 887 |
| Frog | 963 | Truck | 928 | Automobile | 905 |
| Automobile | 966 | Frog | 951 | Frog | 914 |

# DISCUSSION: ARE ALL CLASSES EQUALLY IMPACTED?

## A FEW CONCLUSIONS…

- **'Cat'** is the least accurate class even with labeling noise and data quantity

- **'Bird'** class relative performance decreases with labeling noise and volume reduction

Lowest accuracy → Highest accuracy

| Baseline | | 30% Labeling Noise | | 5K Data Samples | |
|---|---|---|---|---|---|
| Cat | 706 | Cat | 608 | Cat | 520 |
| Dog | 790 | Dog | 698 | Bird | 594 |
| Deer | 875 | Bird | 744 | Dog | 599 |
| Bird | 889 | Deer | 811 | Deer | 716 |
| Airplane | 897 | Airplane | 837 | Airplane | 789 |
| Ship | 928 | Horse | 877 | Horse | 830 |
| Truck | 936 | Ship | 897 | Ship | 853 |
| Horse | 944 | Automobile | 922 | Truck | 887 |
| Frog | 963 | Truck | 928 | Automobile | 905 |
| Automobile | 966 | Frog | 951 | Frog | 914 |

# DISCUSSION: ARE ALL CLASSES EQUALLY IMPACTED?

## A FEW CONCLUSIONS...

- **'Cat'** is the least accurate class even with labeling noise and data quantity

- **'Bird'** class relative performance decreases with labeling noise and volume reduction

- **'Frog'** class stays stable with noise induction as well as data volume reduction

**Lowest accuracy** → **Highest accuracy**

| Baseline | | 30% Labeling Noise | | 5K Data Samples | |
|---|---|---|---|---|---|
| Cat | 706 | Cat | 608 | Cat | 520 |
| Dog | 790 | Dog | 698 | Bird | 594 |
| Deer | 875 | Bird | 744 | Dog | 599 |
| Bird | 889 | Deer | 811 | Deer | 716 |
| Airplane | 897 | Airplane | 837 | Airplane | 789 |
| Ship | 928 | Horse | 877 | Horse | 830 |
| Truck | 936 | Ship | 897 | Ship | 853 |
| Horse | 944 | Automobile | 922 | Truck | 887 |
| Frog | 963 | Truck | 928 | Automobile | 905 |
| Automobile | 966 | Frog | 951 | Frog | 914 |

# DISCUSSION: ARE ALL CLASSES EQUALLY IMPACTED?

## INTUITION: 'BIRD' CLASS VARIANCE

# DISCUSSION: ARE ALL CLASSES EQUALLY IMPACTED?

## Most Sensitive Class – 'Bird'

### Results with labeling pollution



### Results with data volume reduction

# DISCUSSION: IS IT THE MODEL OR THE DATA?

| Model | Epochs | Batch Size | Accuracy |
|---|---|---|---|
| Custom (Keras with TF backend) | 125 | 64 | 88.94 |
| LeNet (Pytorch) | 125 | 64 | 66.6 |
| ResNet18 (Pytorch) | 25 | 64 | 88.29 |
| UnResNet18 (Pytorch) | 25 | 64 | 85.77 |
| GoogLeNet (Pytorch) | 25 | 64 | 88.6 |

**Same Experiments, Different Models**

# DISCUSSION: IS IT THE MODEL OR THE DATA?



**Same Experiments, Different Models**

# DISCUSSION: VOLUME REDUCTION VS. LABELING NOISE



20% **Labeling Noise Induction**

20% **Data Volume Reduction**

# DISCUSSION: VOLUME REDUCTION VS. LABELING NOISE

**CLEAN LABELS
FULL TRAINING SET**

100% "GOOD" DATA

**20% VOLUME REDUCTION**

80% "GOOD" DATA

**20% LABELING POLLUTION**

80% "GOOD" DATA

20% "BAD" DATA

2 COMBINED EFFECTS TO DECOUPLE

Impact Index for Data Volume Reduction (X)



Impact Index for Labeling Noise Induction (Y)

# DISCUSSION: VOLUME REDUCTION VS. LABELING NOISE

LET'S SAVE
   SOME MONEY!

# TOWARDS A SMART LABELING STRATEGY

# TOWARDS A SMART LABELING STRATEGY

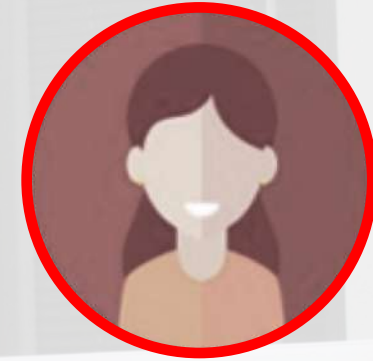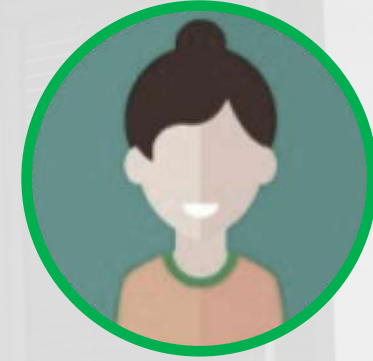# TOWARDS A SMART LABELING STRATEGY

TOWARDS A SMART LABELING STRATEGY

# TOWARDS A SMART LABELING STRATEGY

# TOWARDS A SMART LABELING STRATEGY



## SUPERVISED LEARNING

- All data is labeled
- No. of annotations is predetermined
- No. of annotations is arbitrary

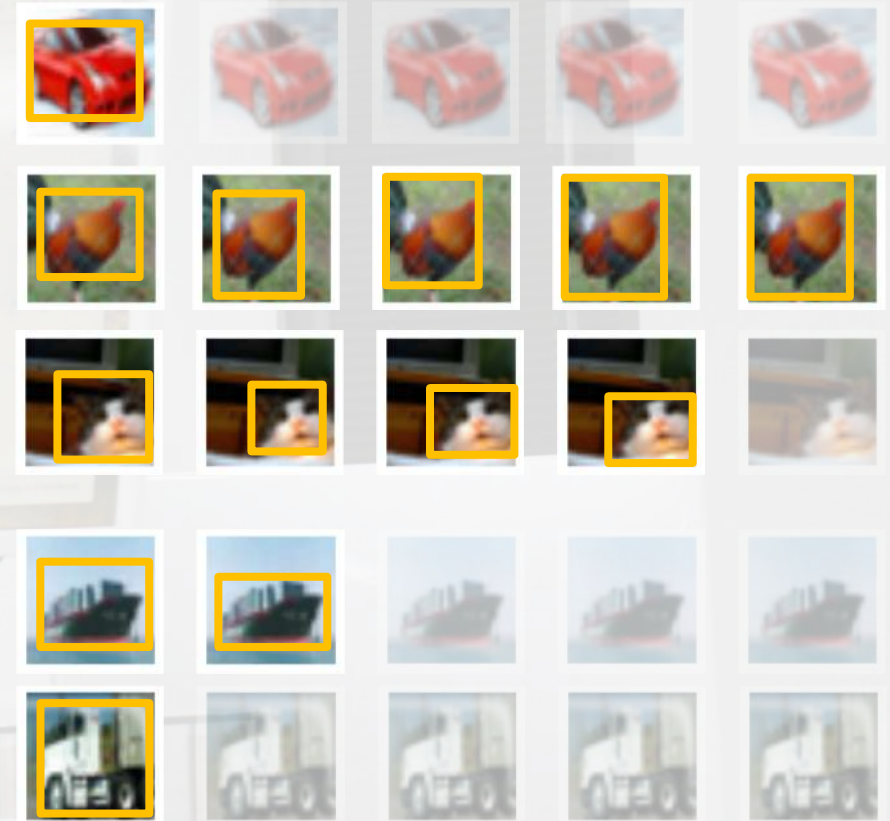# TOWARDS A SMART LABELING STRATEGY

## ACTIVE LEARNING

# TOWARDS A SMART LABELING STRATEGY

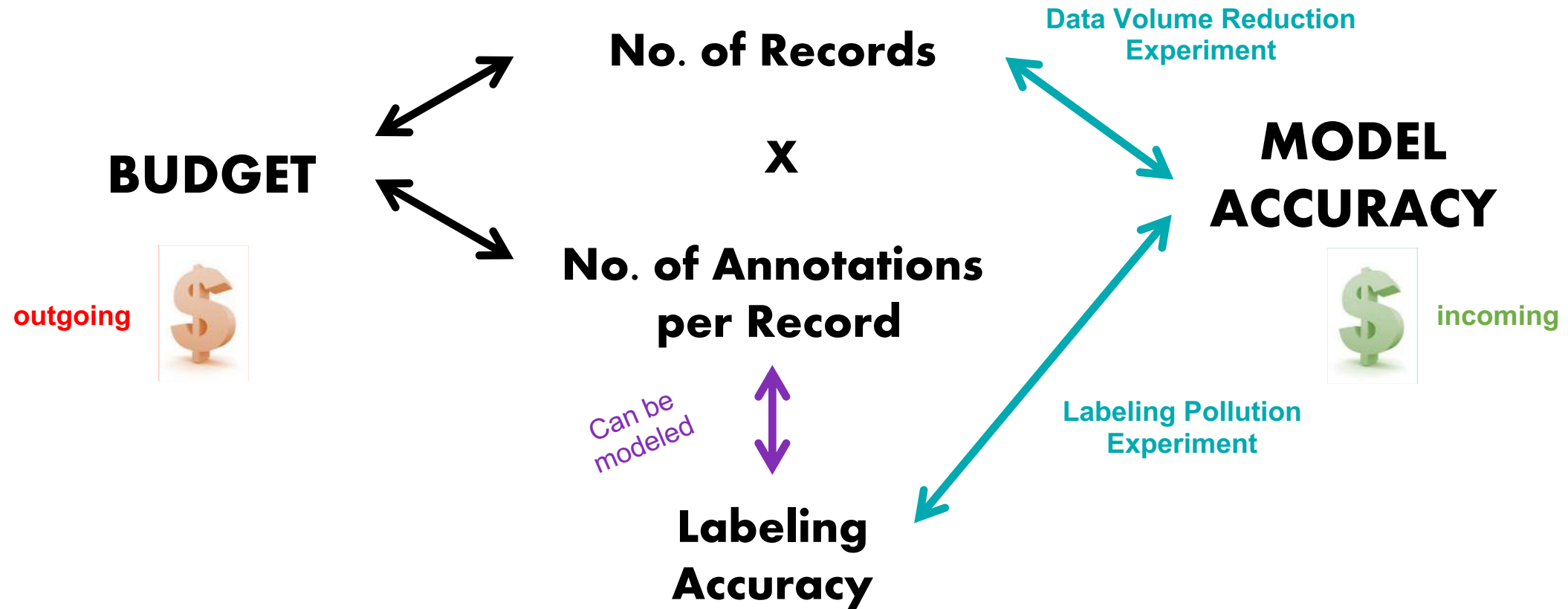**No. of Records**

**X**

**No. of Annotations per Record**

**BUDGET**

outgoing

**Data Volume Reduction Experiment**

**MODEL ACCURACY**

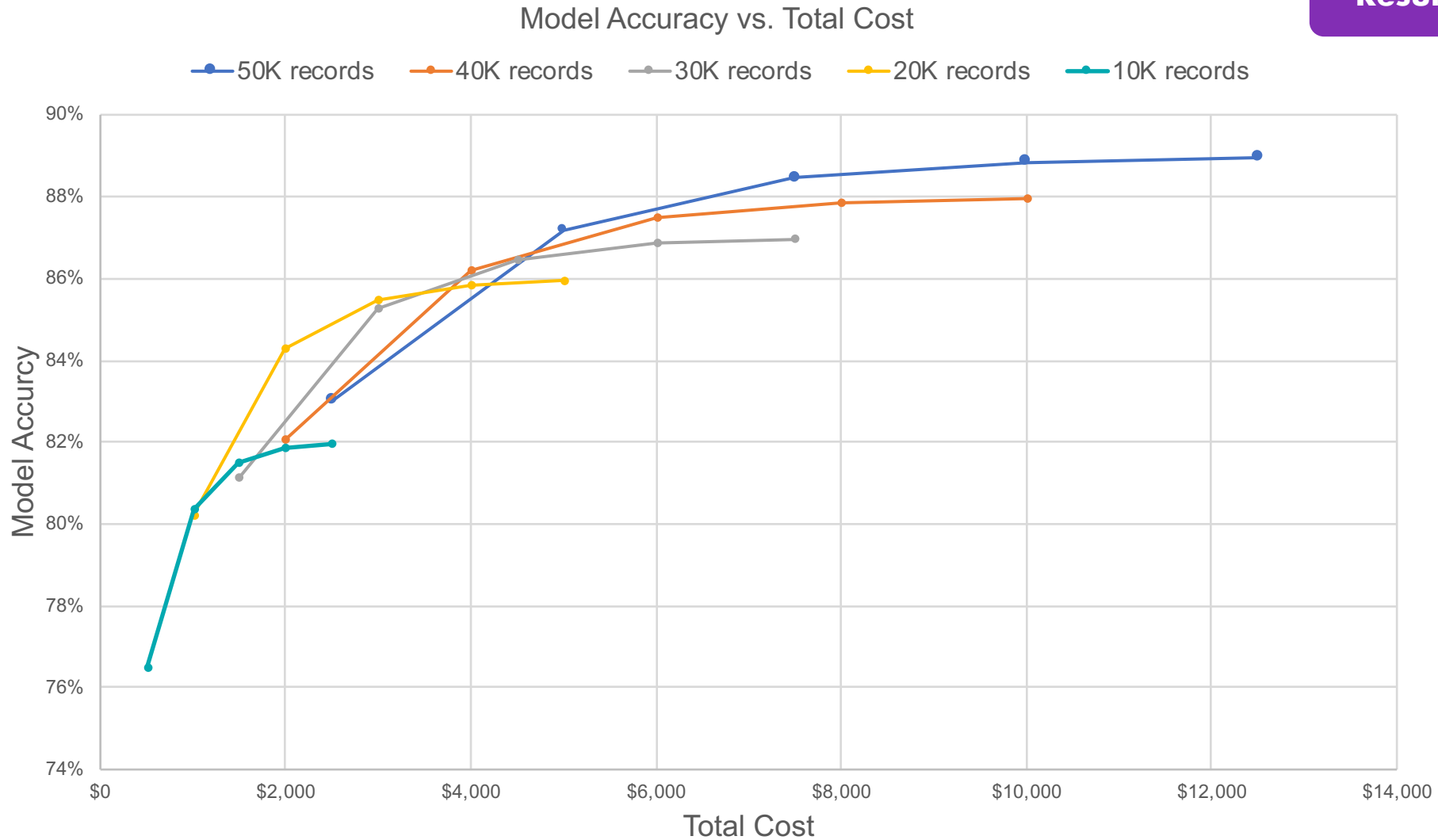incoming

Can be modeled

**Labeling Pollution Experiment**

**Labeling Accuracy**

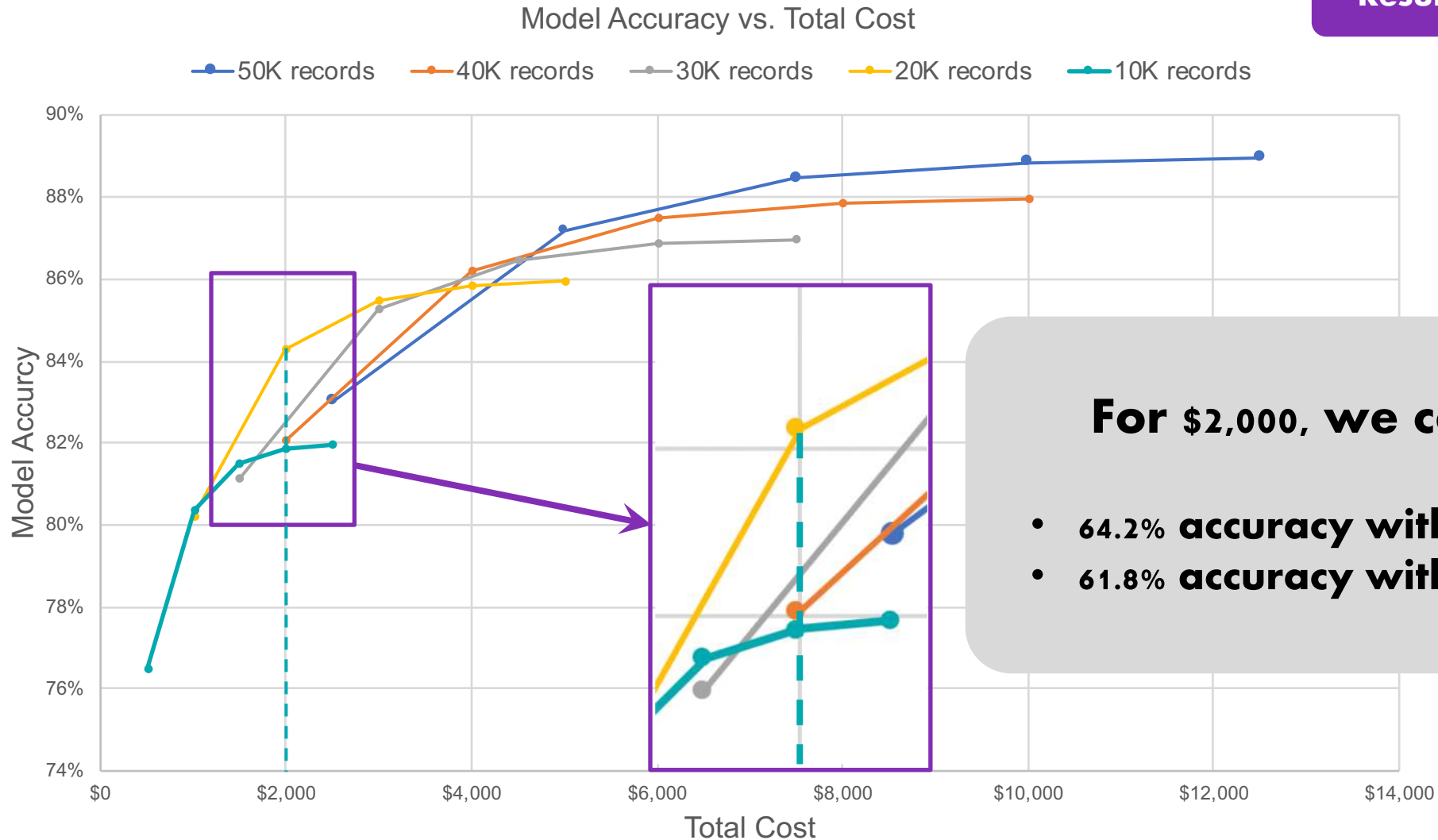# TOWARDS A SMART LABELING STRATEGY



Results on CIFAR-10 Study

# TOWARDS A SMART LABELING STRATEGY



Results on CIFAR-10 Study

Model Accuracy vs. Total Cost

50K records — 40K records — 30K records — 20K records — 10K records
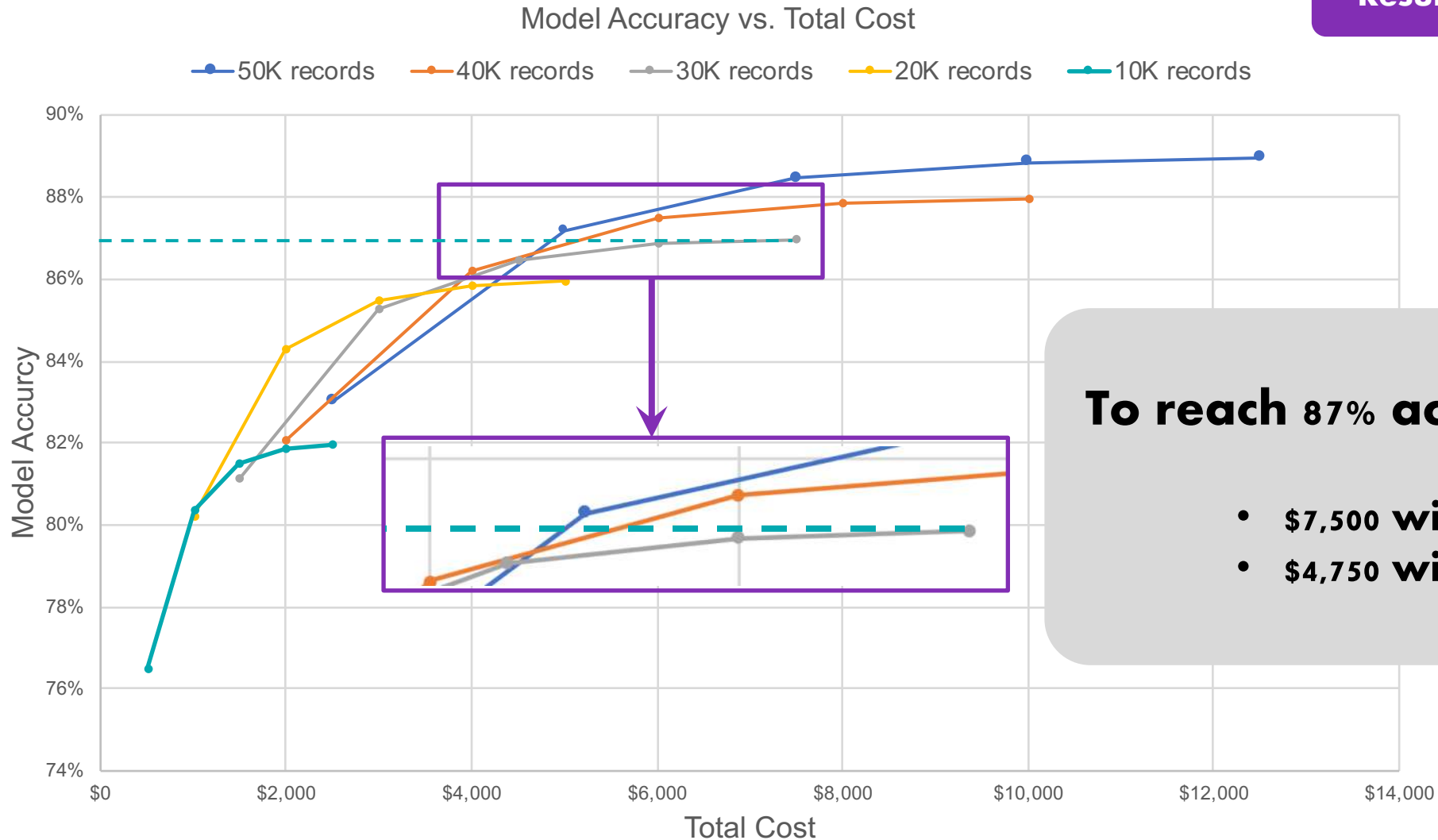
For $2,000, we can get:

- 64.2% accuracy with strategy 4
- 61.8% accuracy with strategy 1

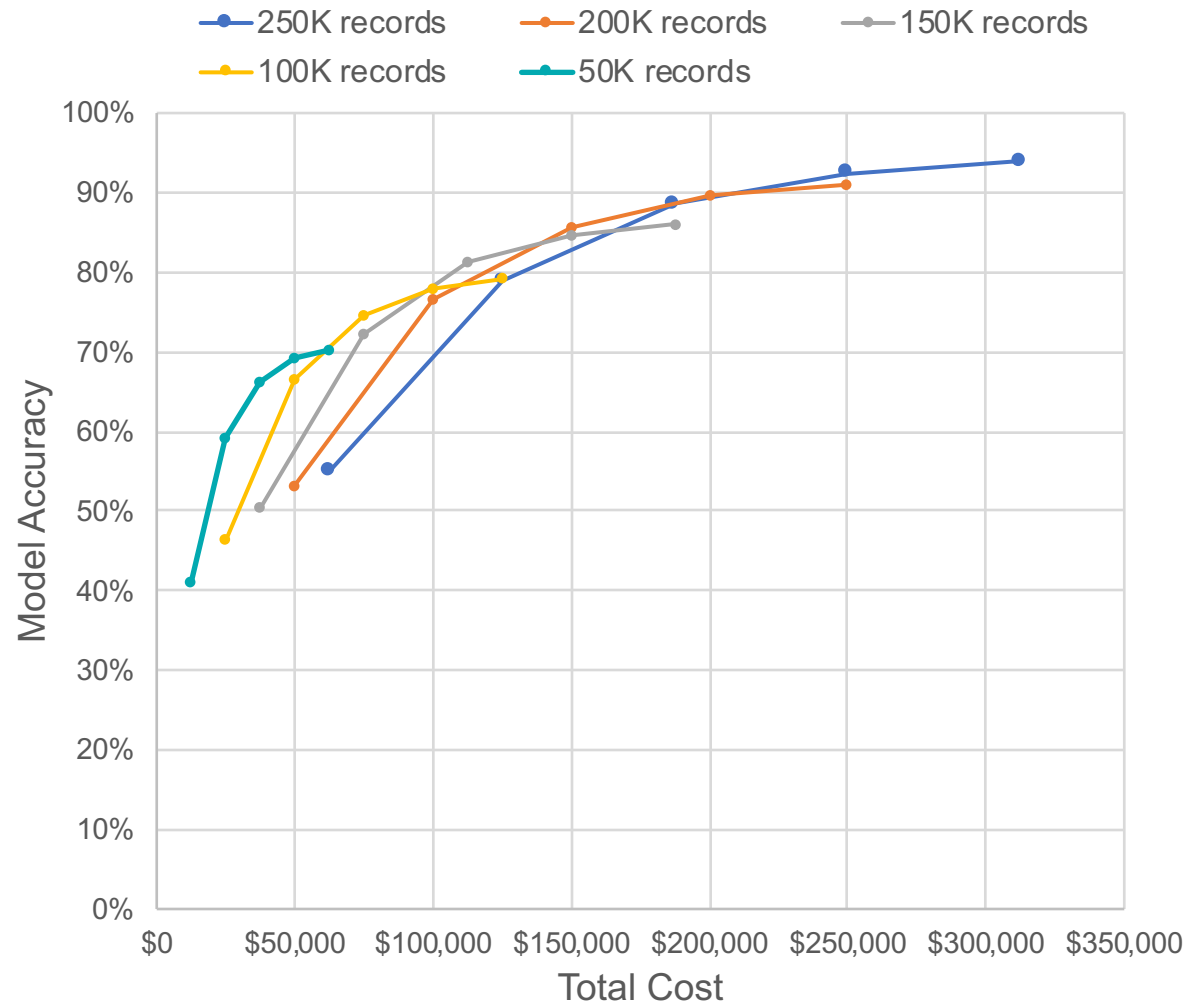# TOWARDS A SMART LABELING STRATEGY



Results on CIFAR-10 Study

Model Accuracy vs. Total Cost

50K records — 40K records — 30K records — 20K records — 10K records

To reach 87% accuracy we need:

- $7,500 **with strategy** 3
- $4,750 **with strategy** 1

# TOWARDS A SMART LABELING STRATEGY



Model Accuracy vs. Total Cost

More Realistic Use Case
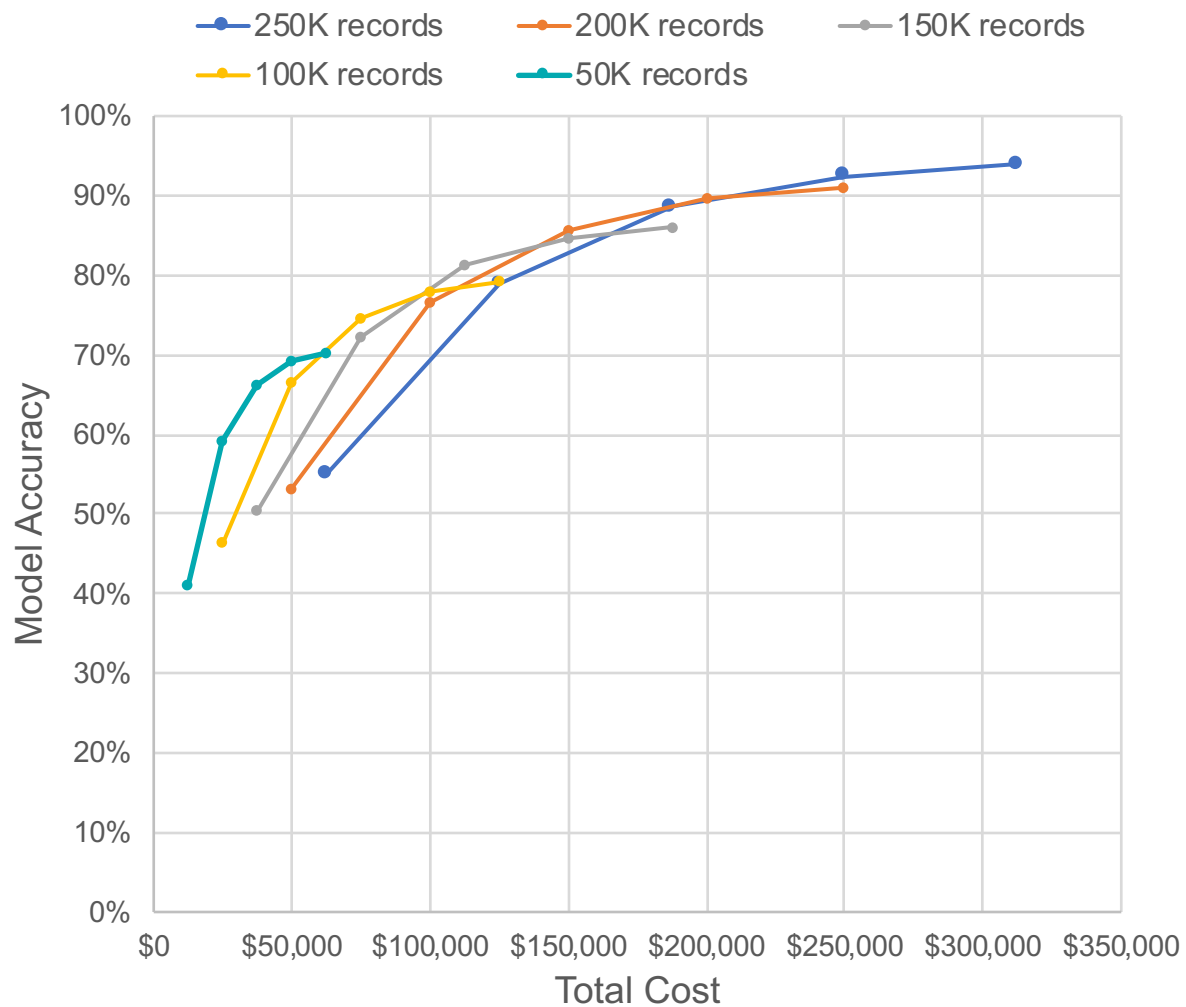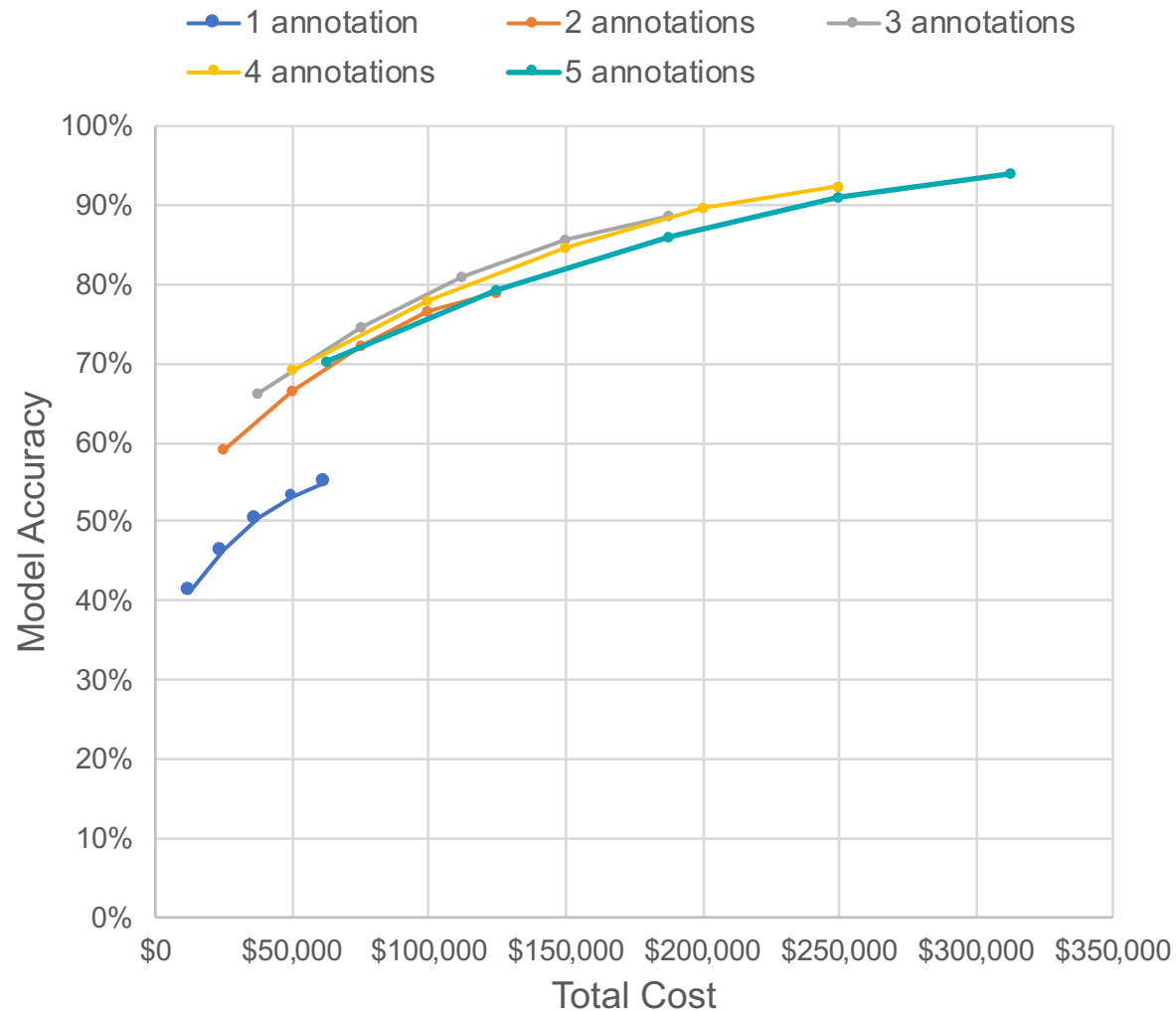
# TOWARDS A SMART LABELING STRATEGY

**More Realistic Use Case**



Model Accuracy vs. Total Cost

Model Accuracy vs. Total Cost

# TOWARDS A SMART LABELING STRATEGY

## NOT COVERED IN THIS TALK:

- **Sensitivity by cluster (instead of class)**

- **Combining data usefulness with difficulty to label**

- **Combining with AL: "non-binary" Active Learning**

# CONCLUSIONS

- **Class sensitivity is inerrant to the data**
  - **Not all data requires as much labeling care**
  - **Better models can't solve everything...**

- **"Compensating" for bad labels**
  - **Is more or less difficult depending on the class**
  - **Might not be possible as all**

- **Smarter labeling strategies are needed**
  - **Saving $$ on labeling doesn't necessarily imply labeling less data**
  - **Local optimization is coming (record level labeling recommendations)**
  - **Bring the area of non-binary Active Learning**

**THANK YOU!**

www.alectio.com

Follow us on LinkedIn & Medium