*Models-schmodels*: why you should care about **Data-Centric AI**

Marysia Winkels, PyData London 2022

# The way I learned
## data science

🎓 Study the algorithms



**PATTERN RECOGNITION**
**AND MACHINE LEARNING**
**CHRISTOPHER M. BISHOP**

GO
DATA
DRIVEN

## The way I learned

## data science

🎓 Study the algorithms

🛠️ Implement the algorithms

GO
DATA
DRIVEN

```python
def sigmoid(x):
    return 1 / (1 + np.exp(-x))

def sigmoid_derivative(x):
    return x * (1 - x)

# define layers
n_input = 2
n_hidden = 6
n_output = 1

# weight initialization
hidden_weights = np.random.uniform(size=(n_input, n_hidden))
output_weights = np.random.uniform(size=(n_hidden, n_output))


epochs = 10000
for _ in range(epochs):
    # Forward pass.
    hidden_layer = X @ hidden_weights
    hidden_activated = sigmoid(hidden_layer)

    output_layer = hidden_activated @ output_weights
    output_activated = sigmoid(output_layer)
    y_hat = output_activated

    # Backpropagation / error calculation
    error_output = y - y_hat
    delta_output = error_output * sigmoid_derivative(output_activated)

    error_hidden = delta_output @ output_weights.T
    delta_hidden = error_hidden * sigmoid_derivative(hidden_activated)

    # Update weights.
    output_weights += hidden_activated.T @ delta_output
    hidden_weights += X.T @ delta_hidden
```

# The way I learned

## data science

🎓 Study the algorithms

🛠️ Implement the algorithms

💪 Practice on toy datasets

GO
DATA
DRIVEN

kaggle

🏛️ GettingStarted Prediction Competition

**Titanic - Machine Learning from Disaster**

Start here! Predict survival on the Titanic and get familiar with ML basics

K  Kaggle · 14,027 teams · Ongoing

🏛️ GettingStarted Prediction Competition

**Spaceship Titanic**

Predict which passengers are transported to an alternate dimension

K  Kaggle · 2,217 teams · Ongoing

If machine learning is *20% modelling* and *80% data prep*...

GO
DATA
DRIVEN

**If machine learning is *20% modelling* and *80% data prep*...**

*....* why is *data prep* not taught?

# Data scientists treat datasets as static

👩‍🏫 What they learn in courses

# Data scientists treat datasets as static

👩‍🏫 What they learn in courses

🎲 It's what most online competitions focus on

coursera
education for evervone

kaggle

GO
DATA
DRIVEN

# Data scientists treat datasets as static

👩‍🏫 What they learn in courses

🎲 It's what most online competitions focus on

🎓 Because that's what they do in academia

**Data scientists treat datasets as static**

👩‍🏫 What they learn in courses

🎲 It's what most online competitions focus on

🎓 Because that's what they do in academia

🛠️ It's what most tools are being built for
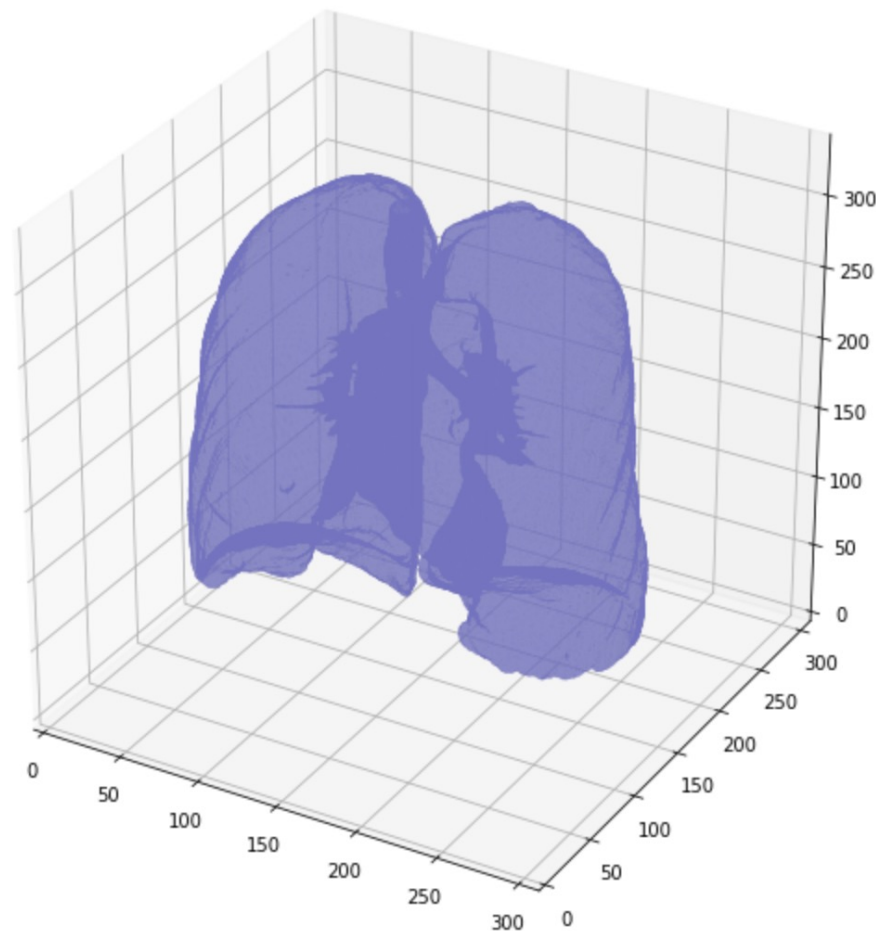
coursera
education for evervone

kaggle

IM**A**GENET

scikit learn

GO
DATA
DRIVEN

# But datasets should not be *static*

GO
DATA
DRIVEN

**Example**: Data Science Bowl 2017

**Example**: Data Science Bowl 2017

ESTIMATE PROBABILITY CANCER → 0.86

"For this solution, engineering the train set was an essential – if not *the* most essential – part."
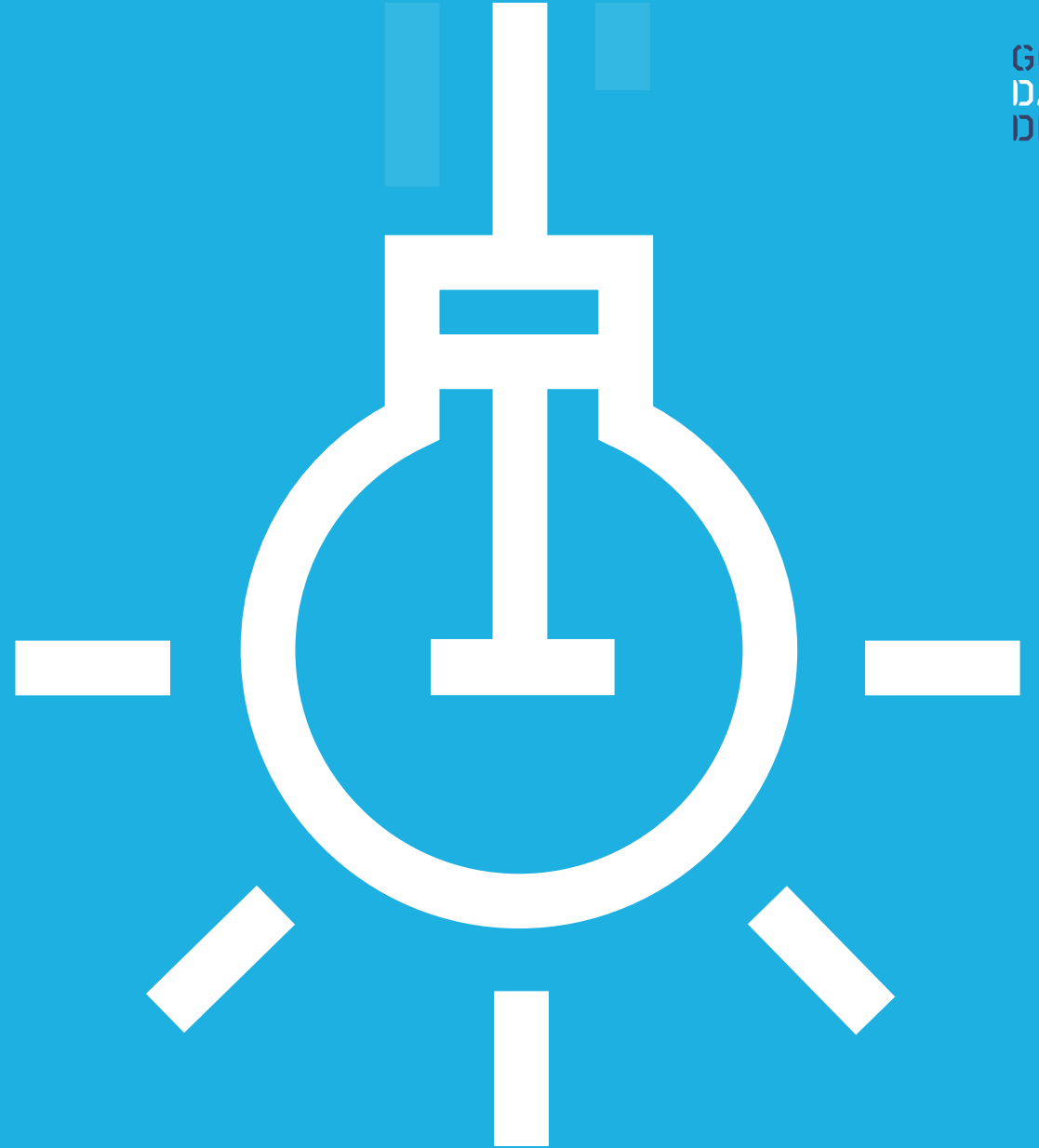
– Julian de Wit, 2$^{nd}$ place

GO
DATA
DRIVEN

**Example**: Data Science Bowl 2017

# Data-Centric AI competition

*__Data-centric AI__ is the discipline of systematically engineering the data used to build an __AI__ system.*

# Data-Centric AI competition

- Model is fixed (ResNet50)

GO
DATA
DRIVEN

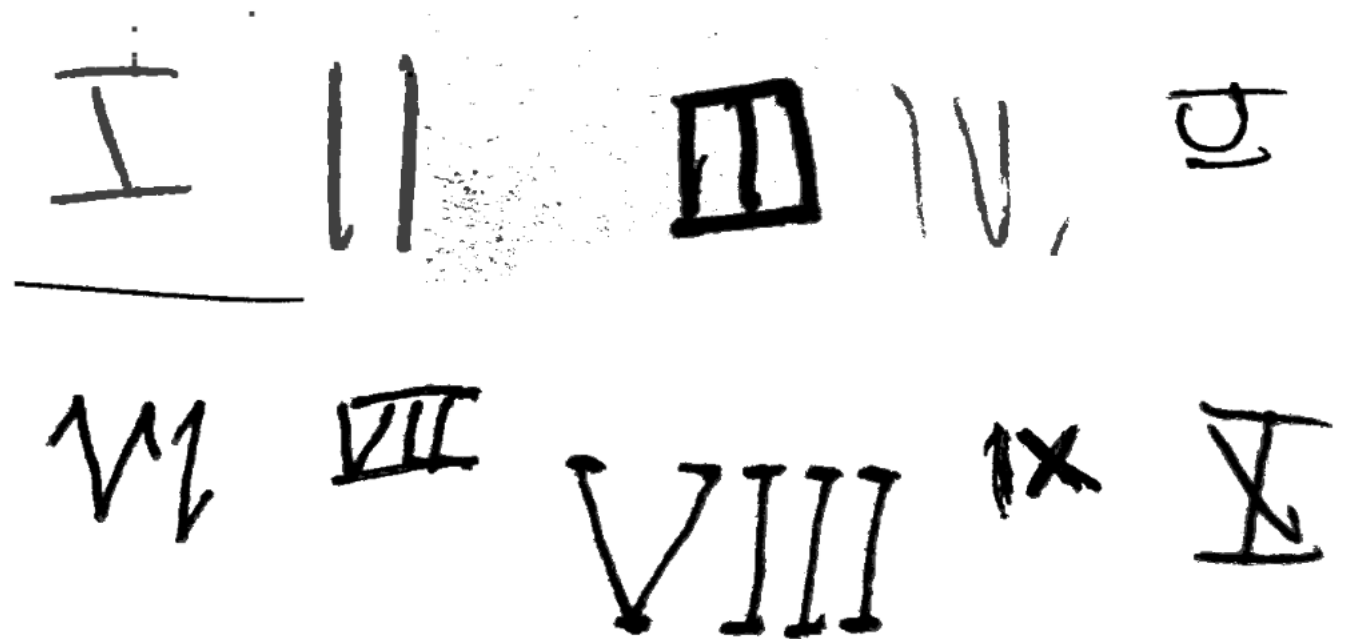# Data-Centric AI competition

- Model is fixed (ResNet50)

- Roman numerals from 1 to 10
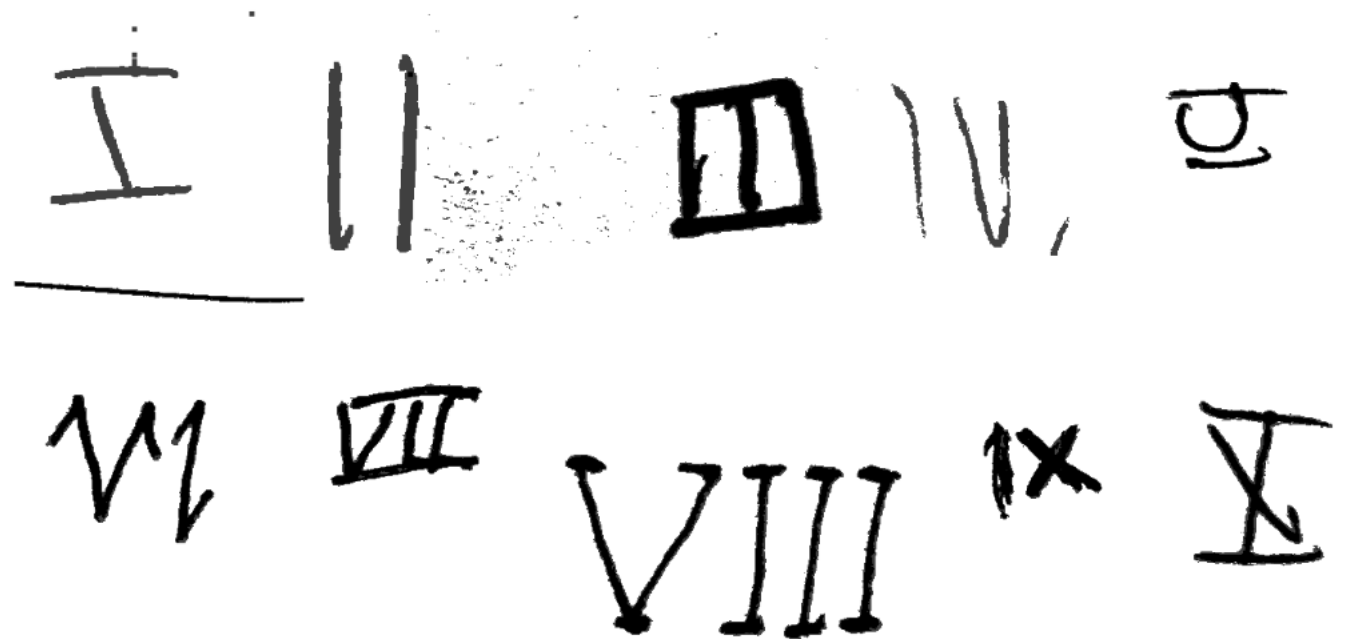
DeepLearning.AI | LANDING AI

# Data-Centric AI Competition

Join the data-centric AI movement!

Click here to enter the contest!

# Data-Centric AI competition

- Model is fixed (ResNet50)

- Roman numerals from 1 to 10

- 3K images in a train/validation set split

## Data-Centric AI Competition

Join the data-centric AI movement!

Click here to enter the contest!

GO
DATA
DRIVEN

# Data-Centric AI competition

- Model is fixed (ResNet50)

- Roman numerals from 1 to 10

- 3K images in a train/validation set split
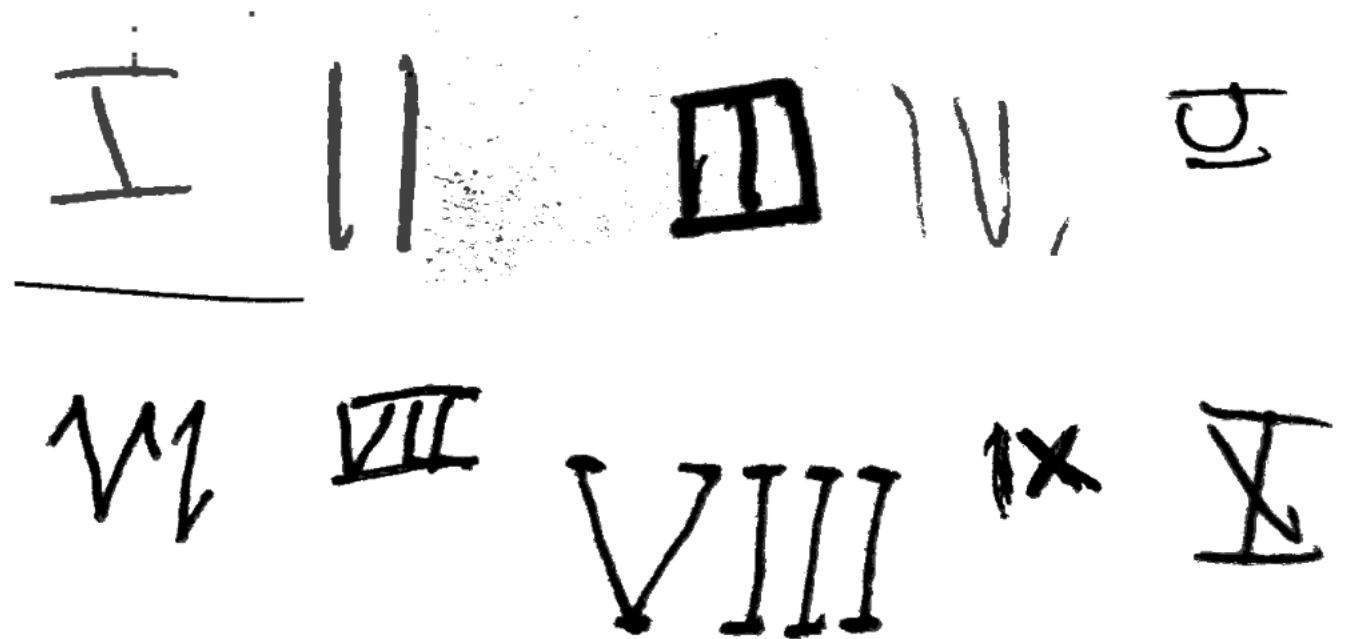
- Labelbook with to-be-expected examples of each class

GO
DATA
DRIVEN

# Data-Centric AI Competition

Join the data-centric AI movement!

Click here to enter the contest!

# Data-Centric AI competition

- Model is fixed (ResNet50)

- Roman numerals from 1 to 10

- 3K images in a train/validation set split

- Labelbook with to-be-expected examples of each class

**TASK**

**Enhance the dataset to a maximum of 10K images that maximizes the model accuracy on a hidden test set**

GO
DATA
DRIVEN

# Data-Centric AI Competition

Join the data-centric AI movement!

Click here to enter the contest!

# roman-numerals-labeling-public

File   Edit   View   Insert   Format   Data   Tools   Add-ons   Help

100%   $   %   .0   .00   123 ▾   Default (Ari... ▾   10 ▾   **B**   *I*   S̶   A   ...

C2   |   fx   |   i

| | image | fname | label | subset | img_url | sim1 | sim2 |
|---|---|---|---|---|---|---|---|
| 1 | | | | | | | |
| 2 |  | ac604594-ce5d- | i | train | https://workshee | 2.971025 | -0.068182476 |
| 3 |  | ac19def6-ce5d-1 | i | train | https://workshee | 3.1976612 | -0.041533813 |
| 4 |  | ac4fe7b2-ce5d-1 | i | train | https://workshee | 3.134209 | -0.02839761 |

# Improve label quality

1. Get predictions from baseline model

# Improve label quality

1. Get predictions from baseline model

2. Focus on discrepancies between the model and the ground truth

# Improve label quality

1. Get predictions from baseline model

2. Focus on discrepancies between the
   model and the ground truth

3. Individually annotate and create
   annotator consensus

| img_url | img_full | annotator 1 | annotator 2 | marysia :) | agreement | override |
|---------|----------|-------------|-------------|------------|-----------|----------|
| | | 4 | x | x | 0 | 3 |
| | | 2 | x | x | 0 | 3 |
| | | 5 | x | x | 0 | 9 |
| | | 3 | x | 3 | 0 | 7 |
| | | x | x | 2 | 0 | x |
| | | 6 | 7 | 7 | 0 | 7 |
| | | 6 | x | x | 0 | 6 |

GO
DATA
DRIVEN

# Lessons learnt

❌ Some data points simply needed to be removed

GO DATA DRIVEN

# Lessons learnt

❌    Some data points simply

needed to be removed

🙅    Lack of consensus between

annotators was often about

the same classes

| | | 2 | | | 6 | | 2 |
|---|---|---|---|---|---|---|---|
| | | 4 | | | 2 | | 2 |
| | | 7 | | | 7 | | 3 |
| | | 7 | | | 3 | | 7 |

GO
DATA
DRIVEN

# Lessons learnt

❌ Some data points simply needed to be removed

🤷‍♀️ Lack of consensus between annotators was often about the same classes

✍️ Different styles of writing

# Lessons learnt

❌    Some data points simply needed to be removed

🤷    Lack of consensus between annotators was often about the same classes

✍️    Different styles of writing

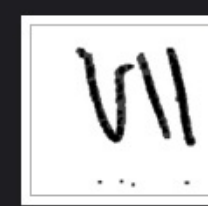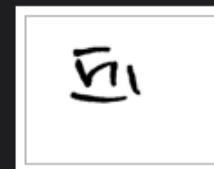⁉️    Difference in train/validation..??

set?!

GO
DATA
DRIVEN

## Sample from training set
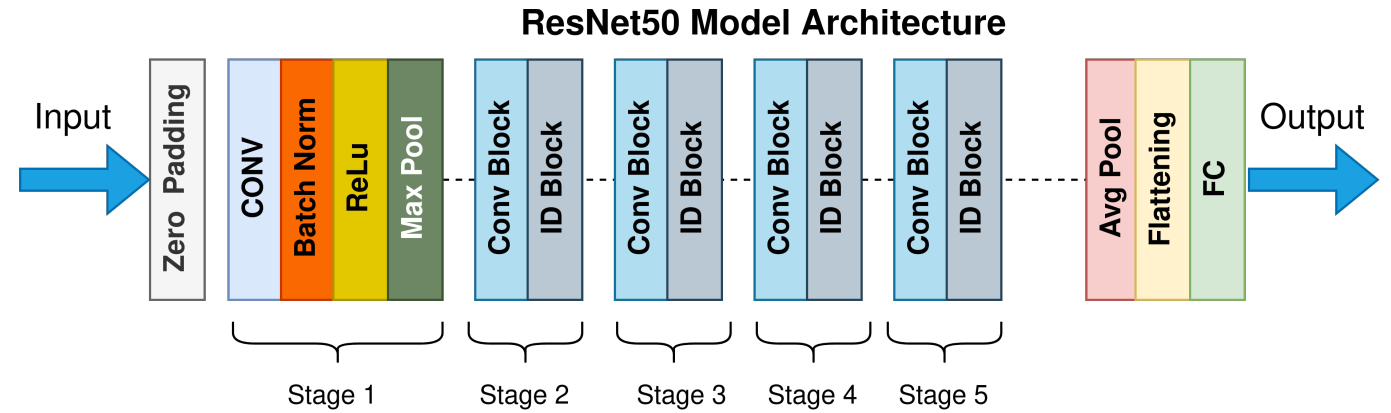


## Sample from validation set

# Visualize the data

1. Pass all the data through the network to obtain the embeddings



**ResNet50 Model Architecture**

# Visualize the data

1. Pass all the data through the network to obtain the embeddings



ResNet50 Model Architecture

Input → Zero Padding → [CONV | Batch Norm | ReLu | Max Pool] (Stage 1) → [Conv Block | ID Block] (Stage 2) → [Conv Block | ID Block] (Stage 3) → [Conv Block | ID Block] (Stage 4) → [Conv Block | ID Block] (Stage 5) → [Avg Pool | Flattening | FC] → Output

ResNet50 Model Architecture

Input → Zero Padding → [CONV | Batch Norm | ReLu | Max Pool] (Stage 1) → [Conv Block | ID Block] (Stage 2) → [Conv Block | ID Block] (Stage 3) → [Conv Block | ID Block] (Stage 4) → [Conv Block | ID Block] (Stage 5) → ~~[Avg Pool | Flattening | FC]~~ → Output
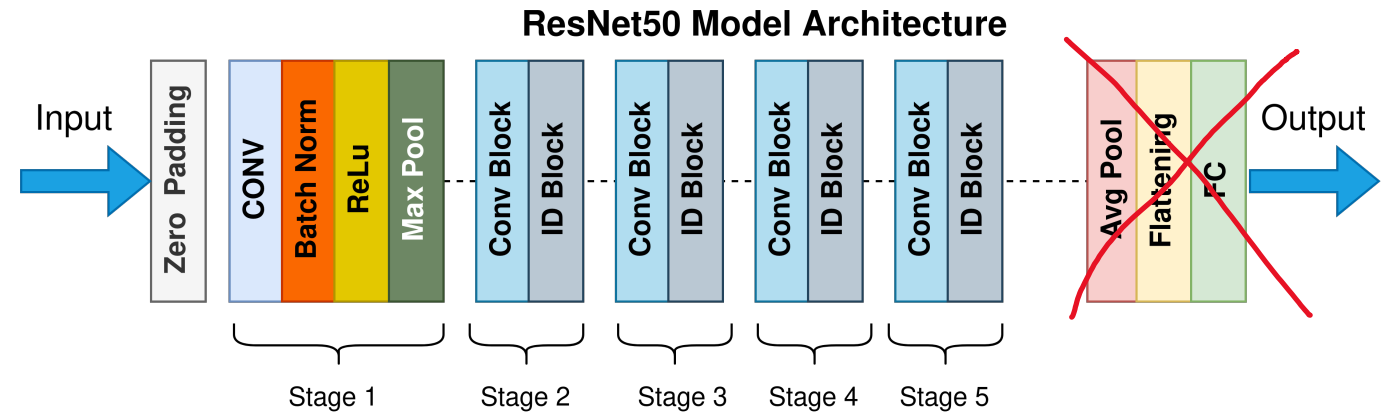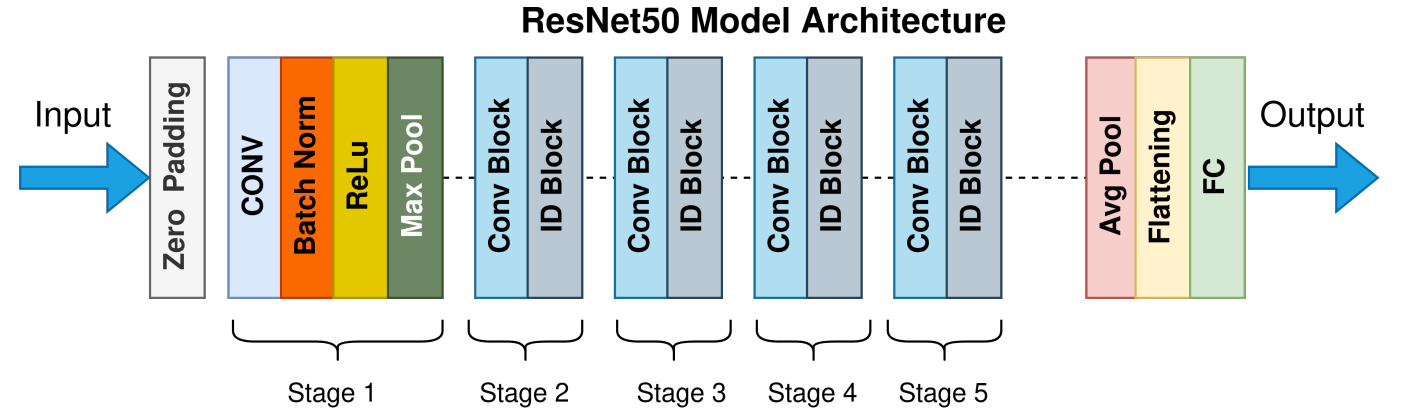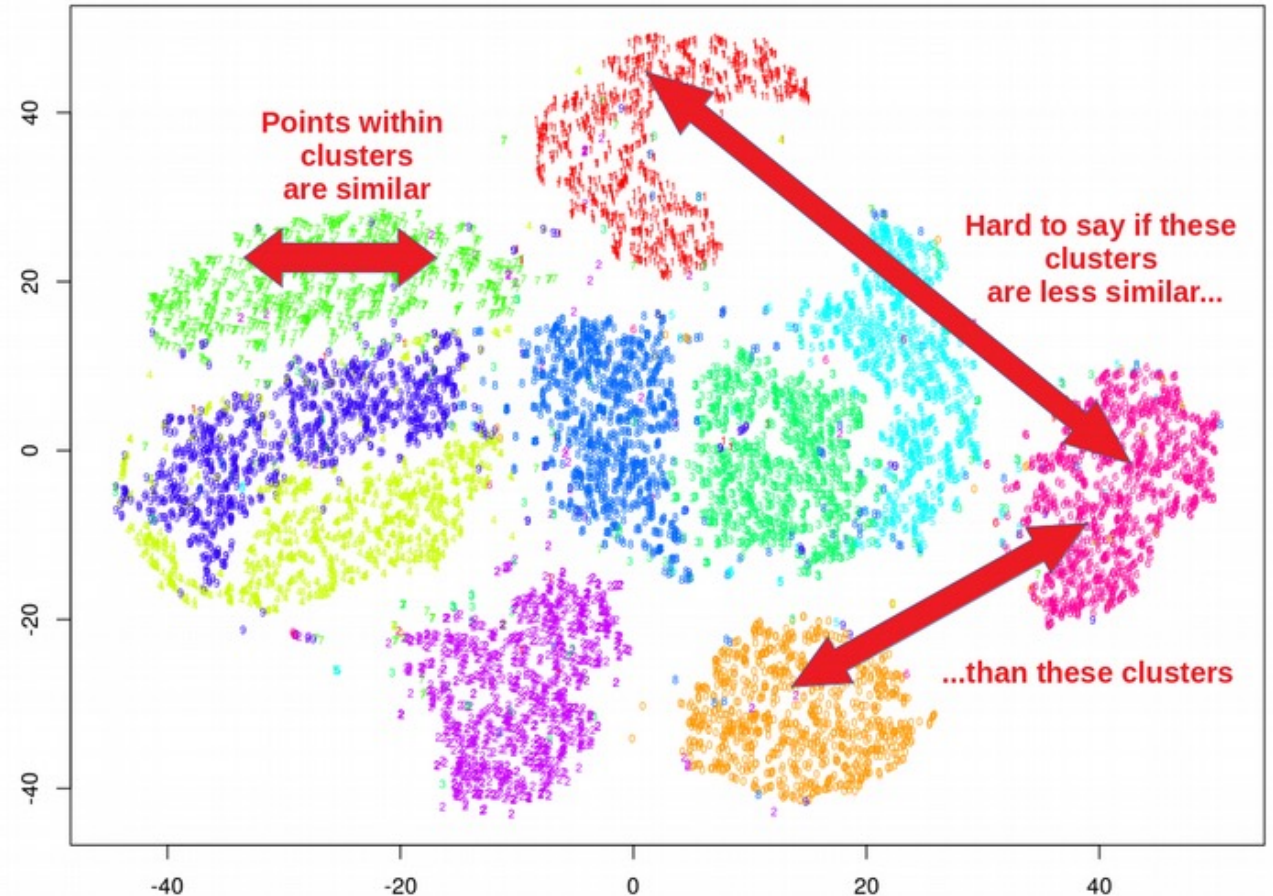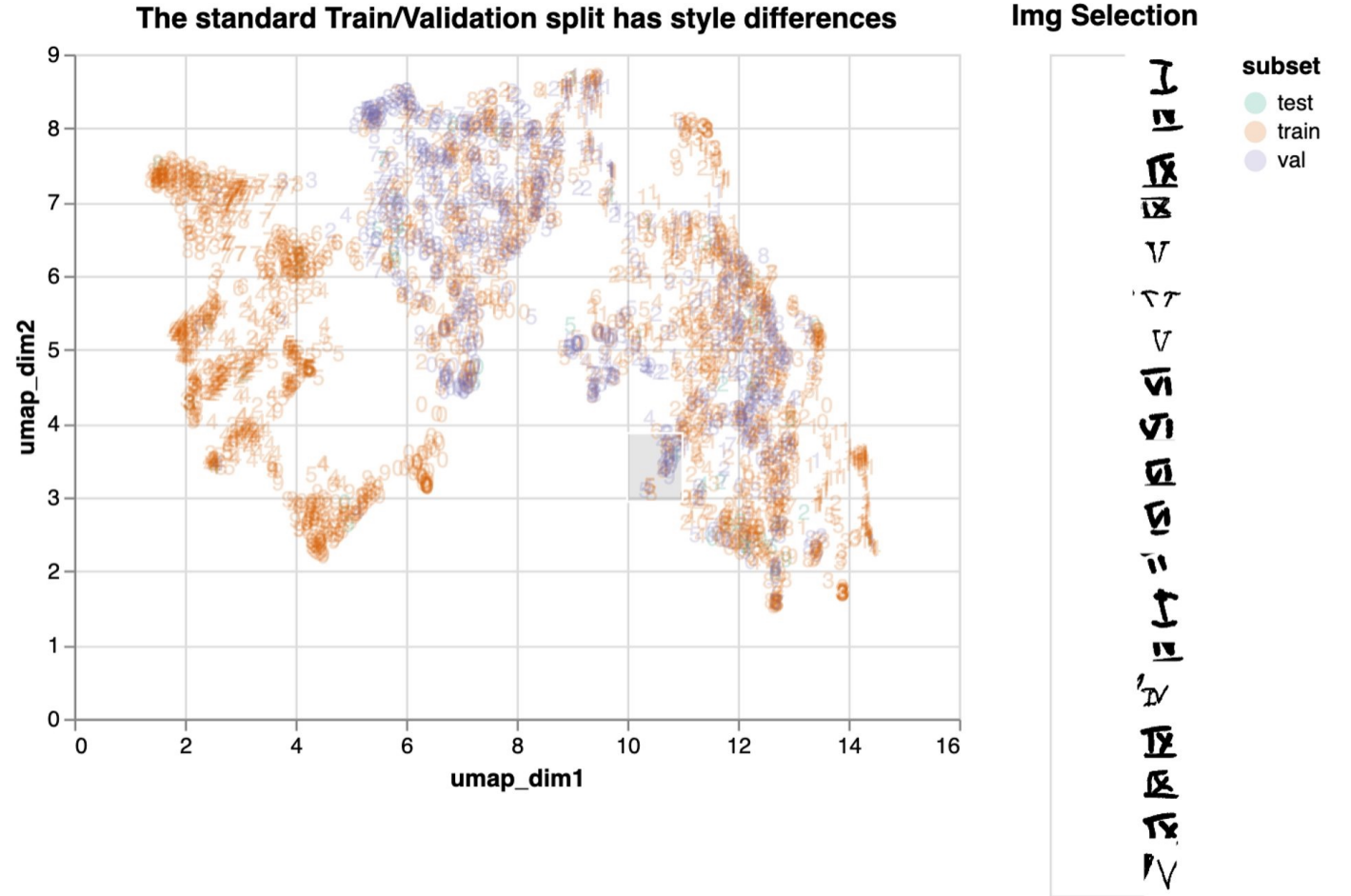
# Visualize the data

1. Pass all the data through the network to obtain the embeddings
2. Perform UMAP

# Visualize the data

1. Pass all the data through the network to obtain the embeddings
2. Perform UMAP
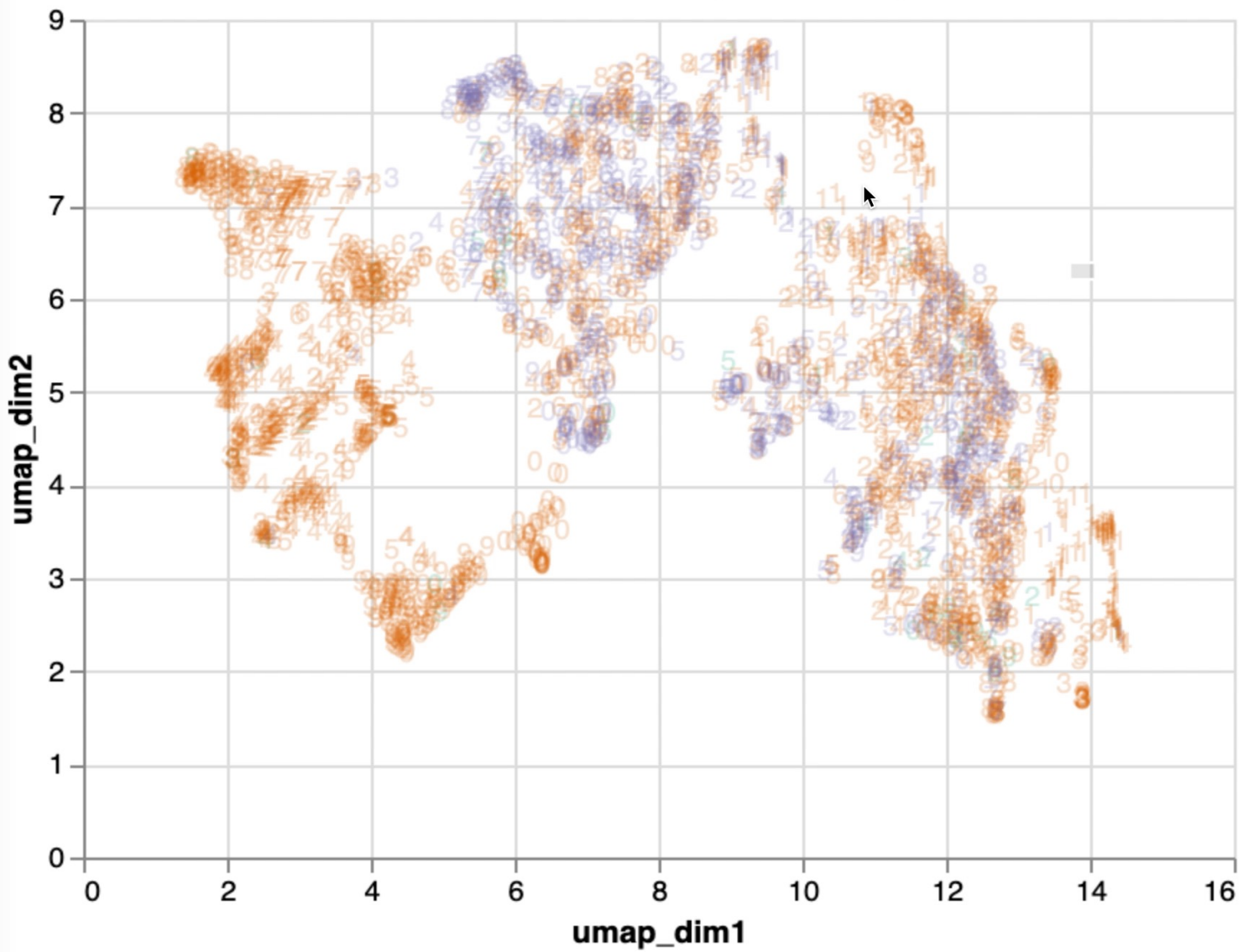3. Visualize using interactive library

The standard Train/Validation split has style differences

Img Selection

subset
- test
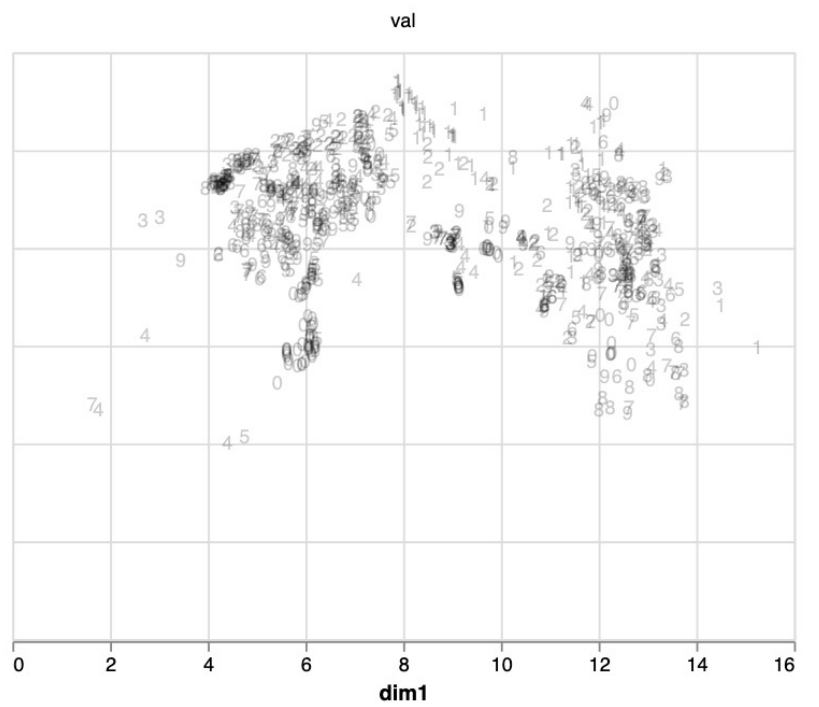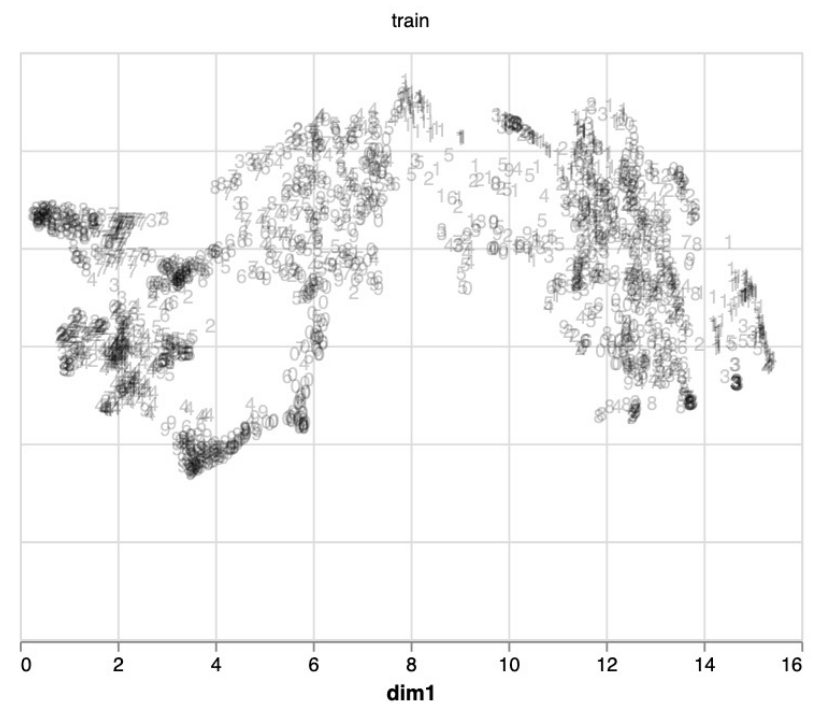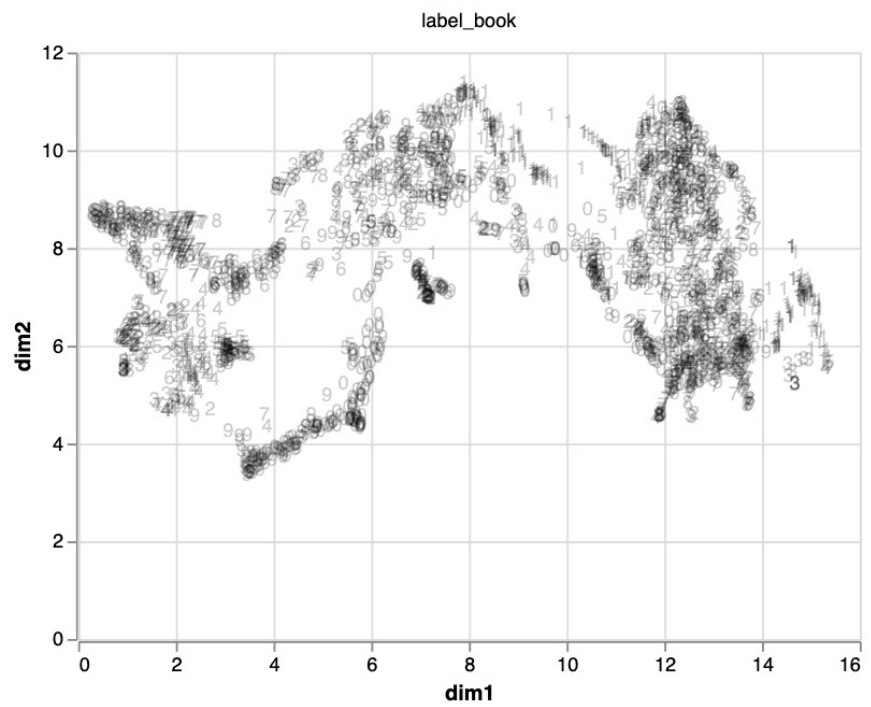- train
- val

**Label book / Train / Validation / style differences.**

# Augment your dataset

🎨 Transform existing data points to create augmented versions

# Augment your dataset

🎨 Transform existing data points to create augmented versions

🧐 Create counterfactuals

# Augment your dataset

1. Transform existing data points to create augmented versions

2. Create counterfactuals

# DCAI competition

*What did others do?*

# Synaptic-AnN

**Best performance**

1. Manual data cleaning

2. Manual data generation

3. Auto data generation

4. Distribution and style

   replication

5. Filtering by vote



Figure 2: Style replication applied on class I of the label book — images bordered in blue are the original label book images.



bottom-left     bottom-right     centre     top-left     top-right

Read about it at
www.deeplearning.ai/data-centric-ai-competition-synaptic-ann/

# Innotescus

**Best performance**

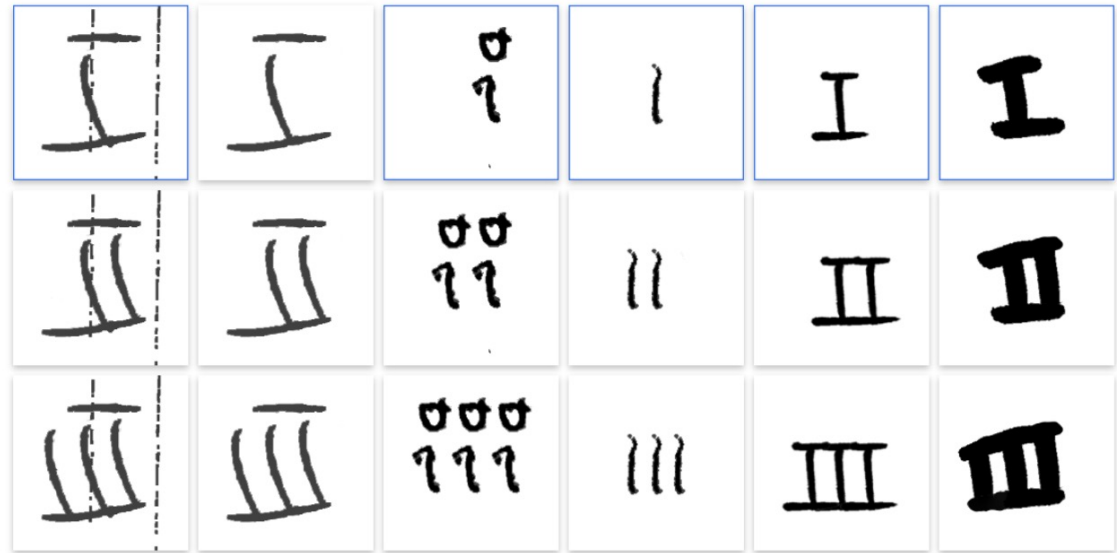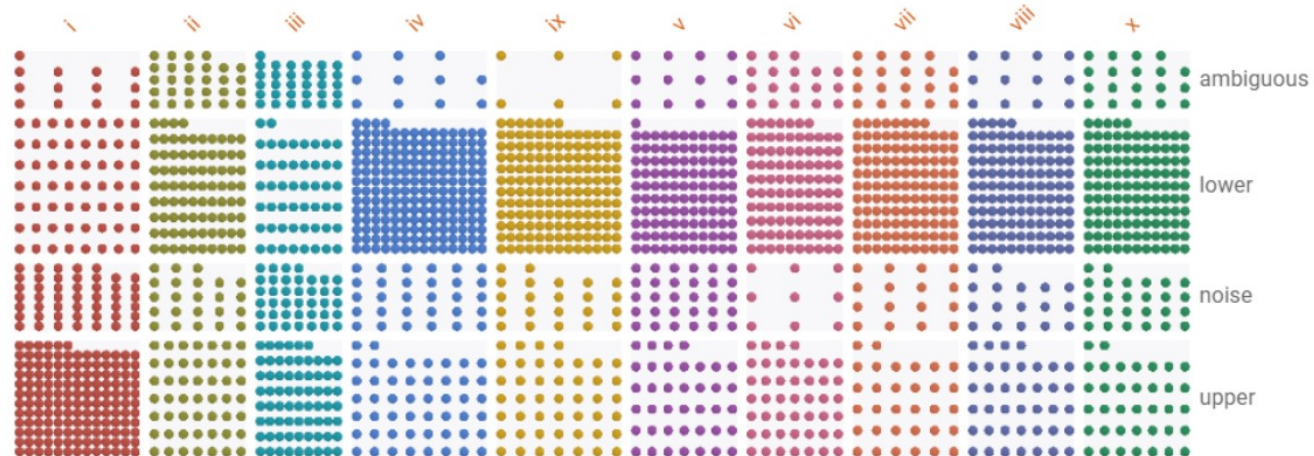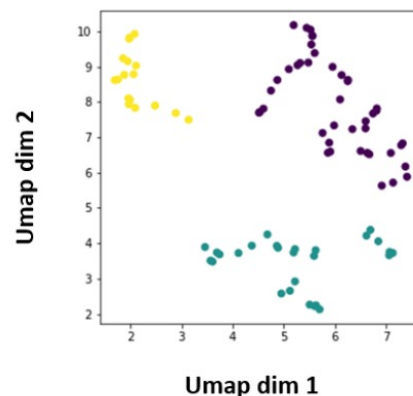1. Data cleaning

2. Rebalancing train & test dataset

3. Rebalancing subclasses using embeddings

4. Rebalancing edge cases with hard examples



Imbalance between lower and uppercase numerals
(Innotescus chart)



Read about it at
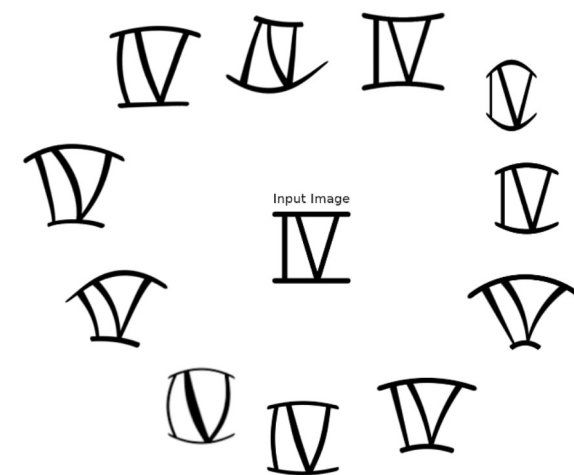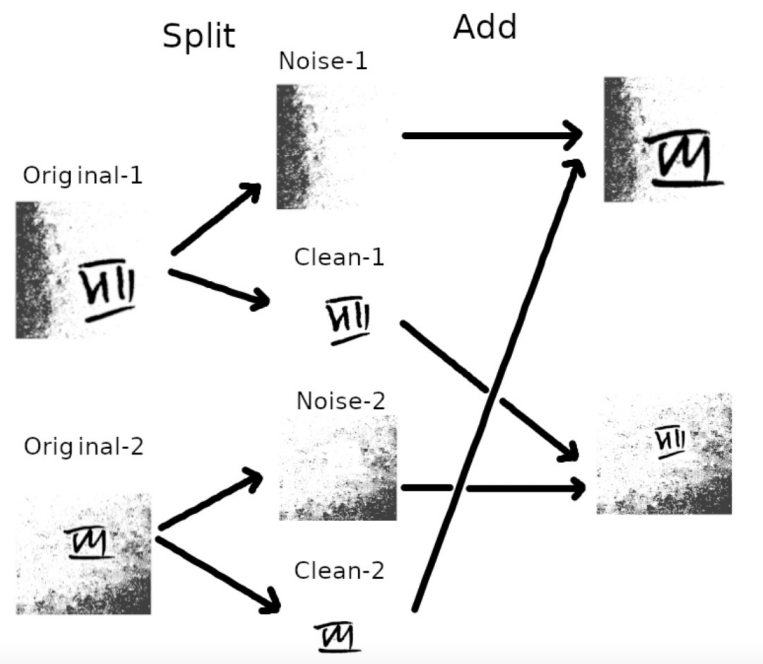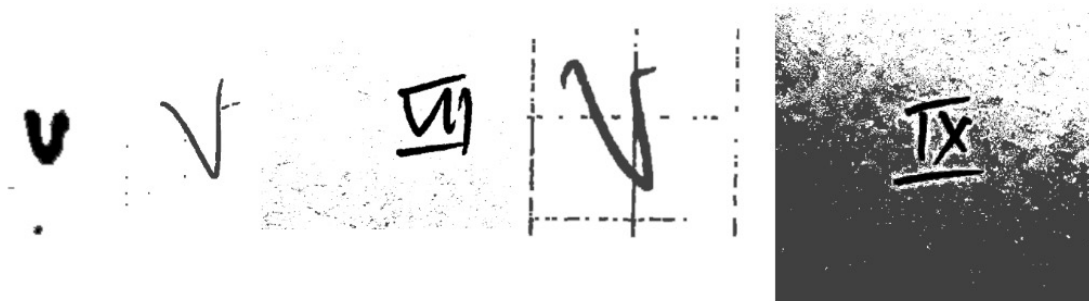https://www.deeplearning.ai/data-centric-ai-competition-innotescus/
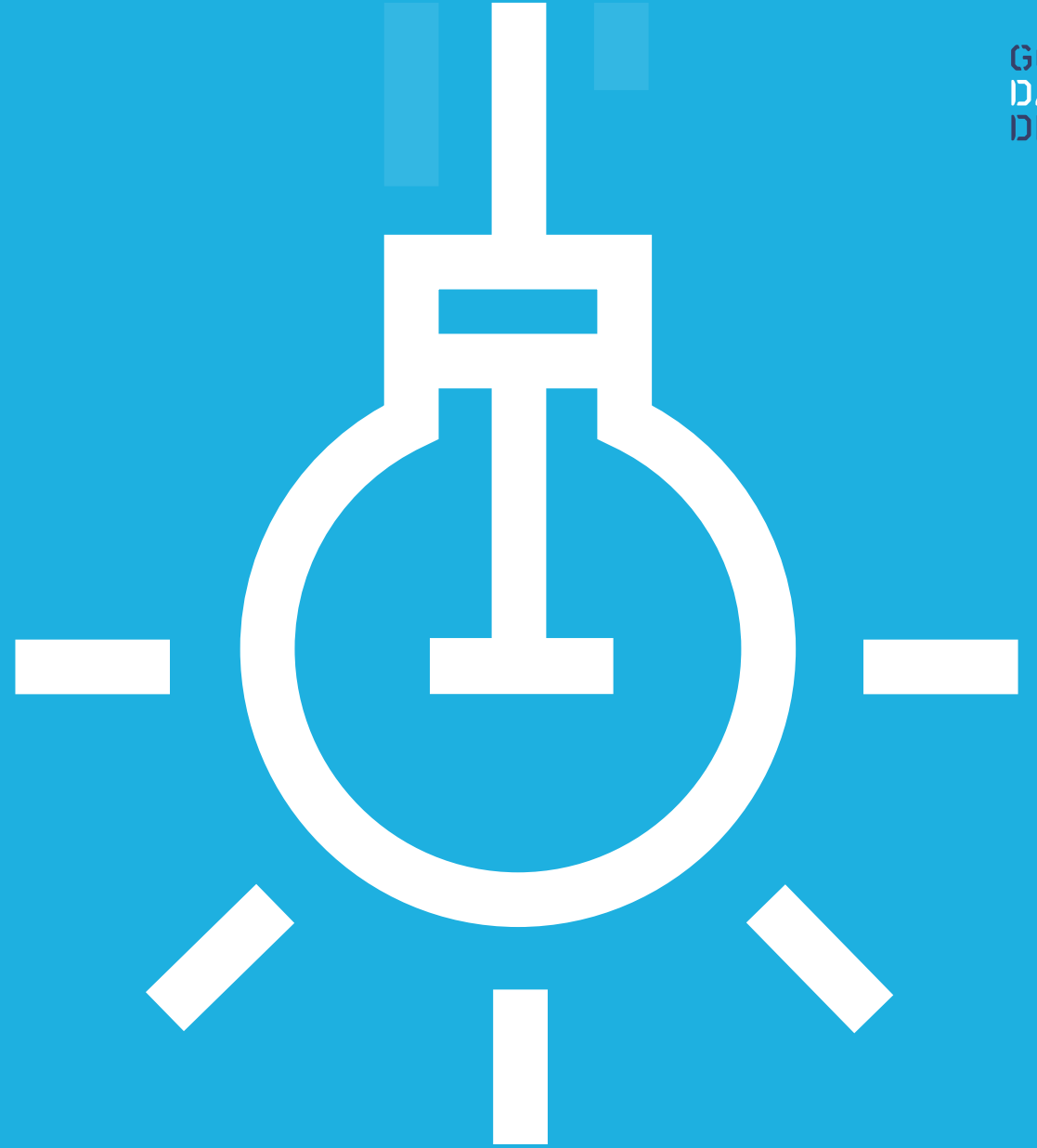
# Divakar Roy

**Best performance**

1. Separate noise

2. Camera distortion onto skewed grid

3. Overlay into canvas

4. Data quality assessment and cleaning up

The *why* of Data-Centric AI

# The *why* of data-centric AI

🧱 Foundation models & transfer learning

*The competitive advantage of data scientists*

*lies in everything that surrounds the model.*

# The *why* of data-centric AI

🧱 Foundation models & transfer learning

💪 Improve performance

# Improving the code vs. the data

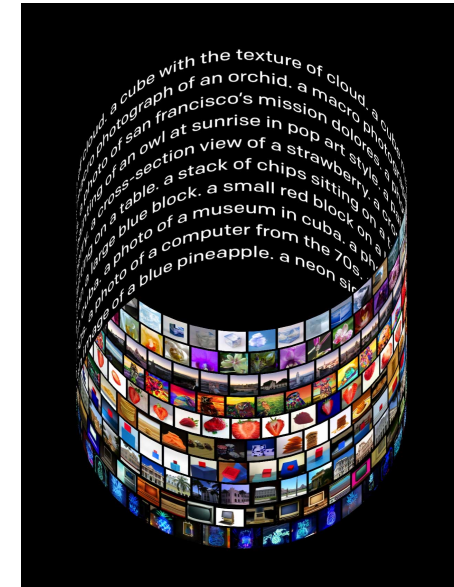|  | Steel defect detection | Solar panel | Surface inspection |
|---|---|---|---|
| Baseline | 76.2% | 75.68% | 85.05% |
| Model-centric | +0%<br>(76.2%) | +0.04%<br>(75.72%) | +0.00%<br>(85.05%) |
| Data-centric | +16.9%<br>(93.1%) | +3.06%<br>(78.74%) | +0.4%<br>(85.45%) |

GO
DATA
DRIVEN

# The *why* of data-centric AI

🧱 Foundation models & transfer learning

💪 Improve performance



Credit: Landing.ai

# The *why* of data-centric AI

🧱 Foundation models & transfer learning

💪 Improve performance

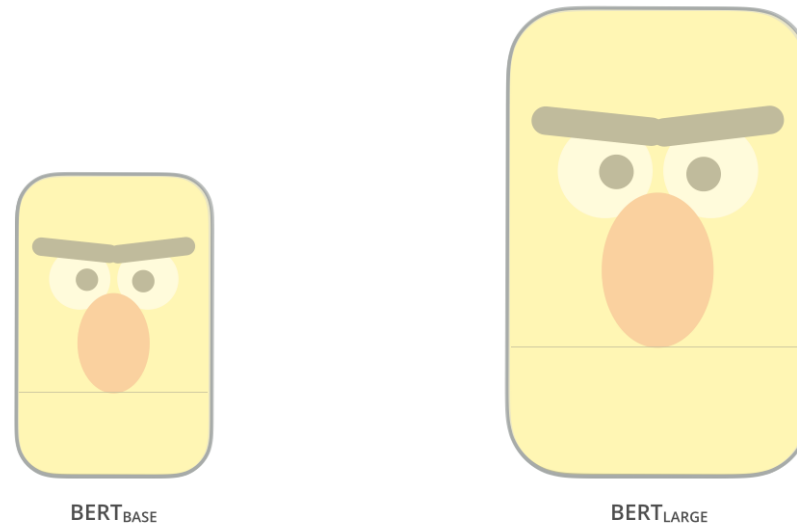*Big data problems can also small data problems*

GO
DATA
DRIVEN

# The *why* of data-centric AI

🧱 Foundation models & transfer learning

💪 Improve performance

🤝 Enables better collaboration

## Nodule type classification



Solid

nodule

Part-solid

nodule

Ground-glass

nodule

GO
DATA
DRIVEN

# The *why* of data-centric AI

🧱 Foundation models & transfer learning

💪 Improve performance

🤝 Enables better collaboration

GO
DATA
DRIVEN

# The *why* of data-centric AI

🧱 Foundation models & transfer learning

💪 Improve performance

🤝 Enables better collaboration

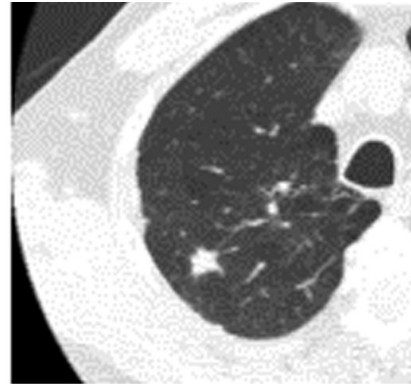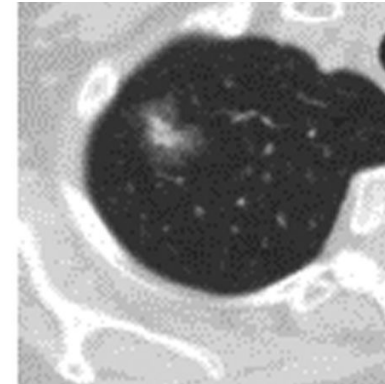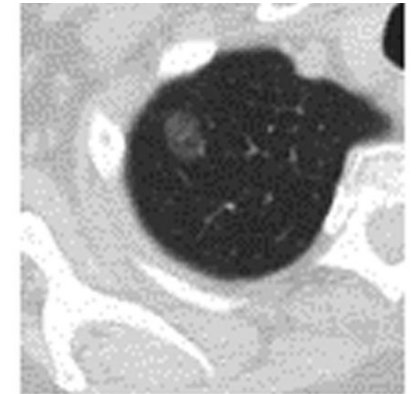# The *why* of data-centric AI

🧱 Foundation models & transfer learning

💪 Improve performance

🤝 Enables better collaboration

Data centric scientist



i solved the medical problem by looping over a bunch of settings

Let's ask why these image were labeled as "good"!
Can they gather more of those examples? Can we help?

uh oh.

HOSPITAL

ivory tower

← distance →

"The focus has to shift from **big data** to **good data**. Having 50 thoughtfully engineered examples can be sufficient to explain to the neural network what you want it to learn."

*- Andrew Ng*

# Data quality

- Consistent data labels

# Data quality

- Consistent data labels



Credit: Michael Bernstein

# Data quality

- **Consistent** data labels

- **Complete** and **representative** data

## Train set

# Data quality

- Consistent data labels

- Complete and representative data

**Train set**



**Real world example**



GO
DATA
DRIVEN

# Data quality

- Consistent data labels

- Complete and representative data

- Unbiased data

```python
from transformers import pipeline

unmasker = pipeline("fill-mask", model="bert-base-uncased")
result = unmasker("This man works as a [MASK].")
print([r["token_str"] for r in result])

result = unmasker("This woman works as a [MASK].")
print([r["token_str"] for r in result])
```
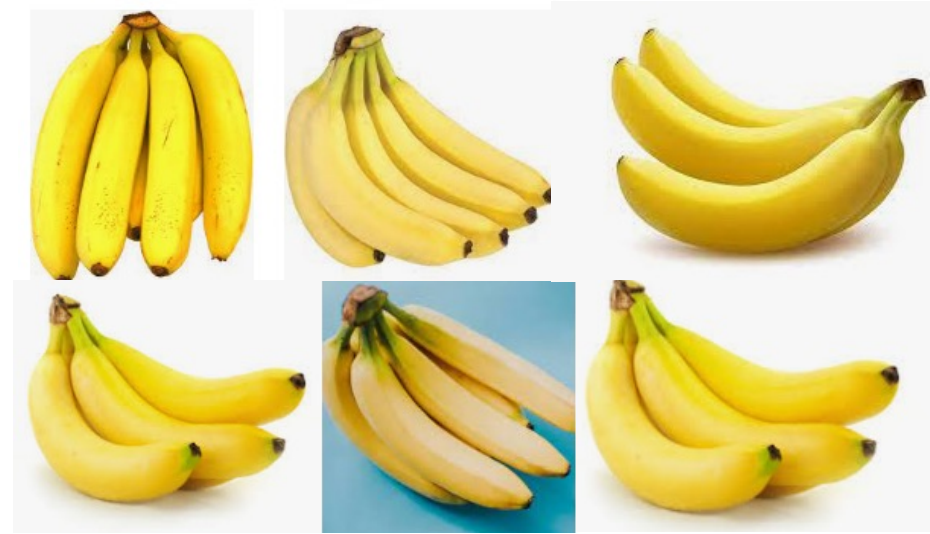
```
['lawyer', 'carpenter', 'doctor', 'waiter', 'mechanic']
['nurse', 'waitress', 'teacher', 'maid', 'prostitute']
```

GO
DATA
DRIVEN

## Data quality

- Consistent data labels

- Complete and representative data

- Unbiased data
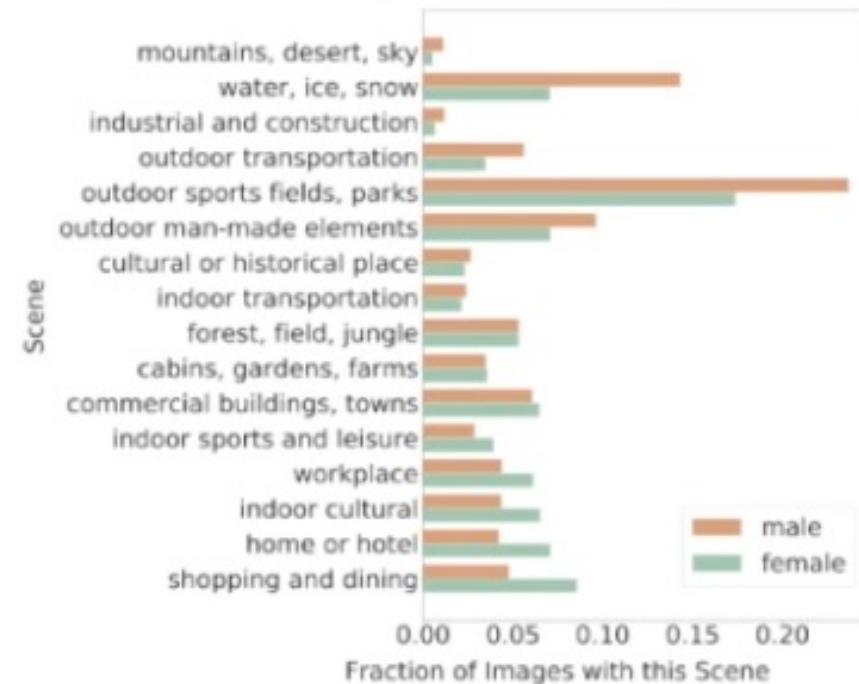


REVISE: A tool for measuring and mitigating bias in visual datasets

Angelina Wang, Arvind Narayanan and Olga Russakovsky
ECCV 2020
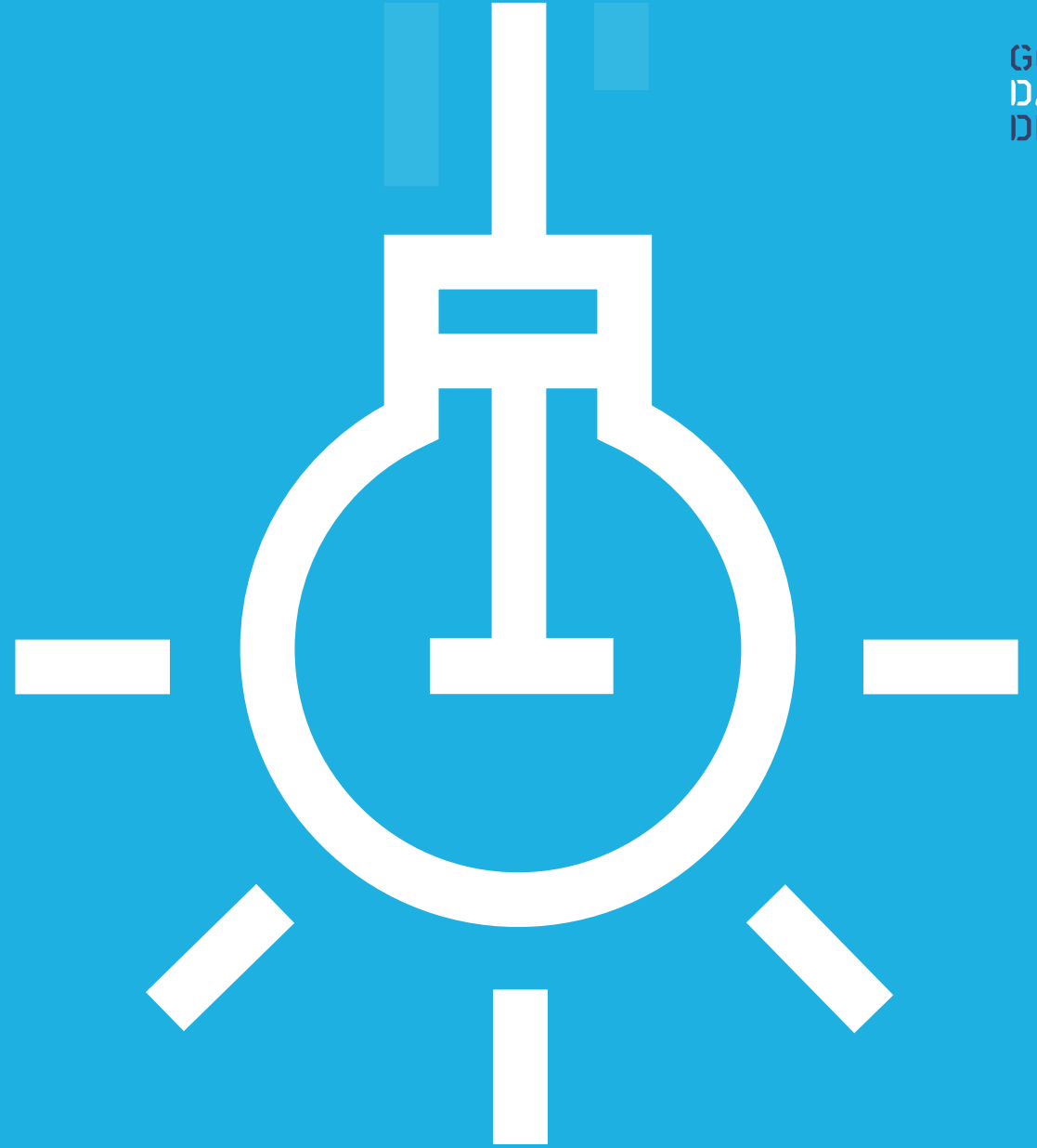https://github.com/princetonvisualai/revise-tool

Images: COCO dataset [Lin et al. ECCV'14]
Annotations: (1) inferred gender [Zhao et al. EMNLP'17],
(2) predicted scenes with the Places network [Zhou et al. TPAMI'7]
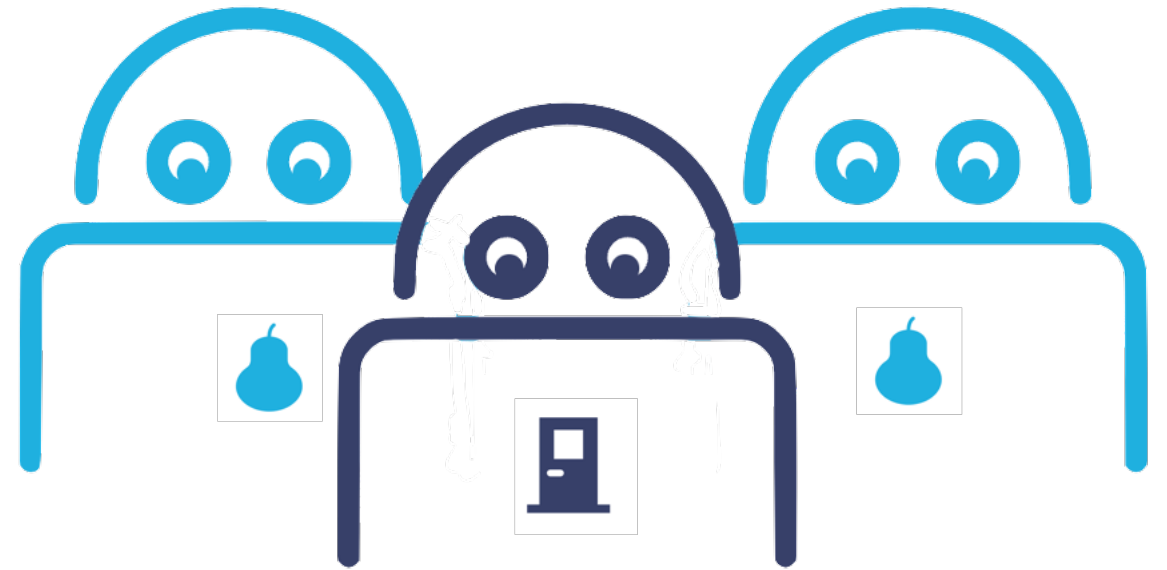
How do we get to Data-Centric AI?

# Development of Data-Centric AI

1. Individuals take an interest

## Development of Data-Centric AI

1. **Individuals** take an interest

2. **Adopted** by many as best practice

***Data-centric AI** is the discipline of systematically engineering the data used to build an **AI** system.*

GO
DATA
DRIVEN

# Development of Data-Centric AI

1. Individuals take an interest

2. Adopted by many as best practice

3. Systematic tools are developed

PyHard: a novel tool for generating hardness embeddings to support data-centric analysis

AutoAugment:
Learning Augmentation Strategies from Data

CircleNLU: A Tool for building Data-Driven Natural Language Understanding System

REVISE: A tool for measuring and mitigating bias in visual datasets

YMIR: A Rapid Data-centric Development Platform for Vision Applications

Augment & Valuate : A Data Enhancement Pipeline for Data-Centric AI

# Development of Data-Centric AI

1. **Individuals** take an interest

2. **Adopted** by many as best practice
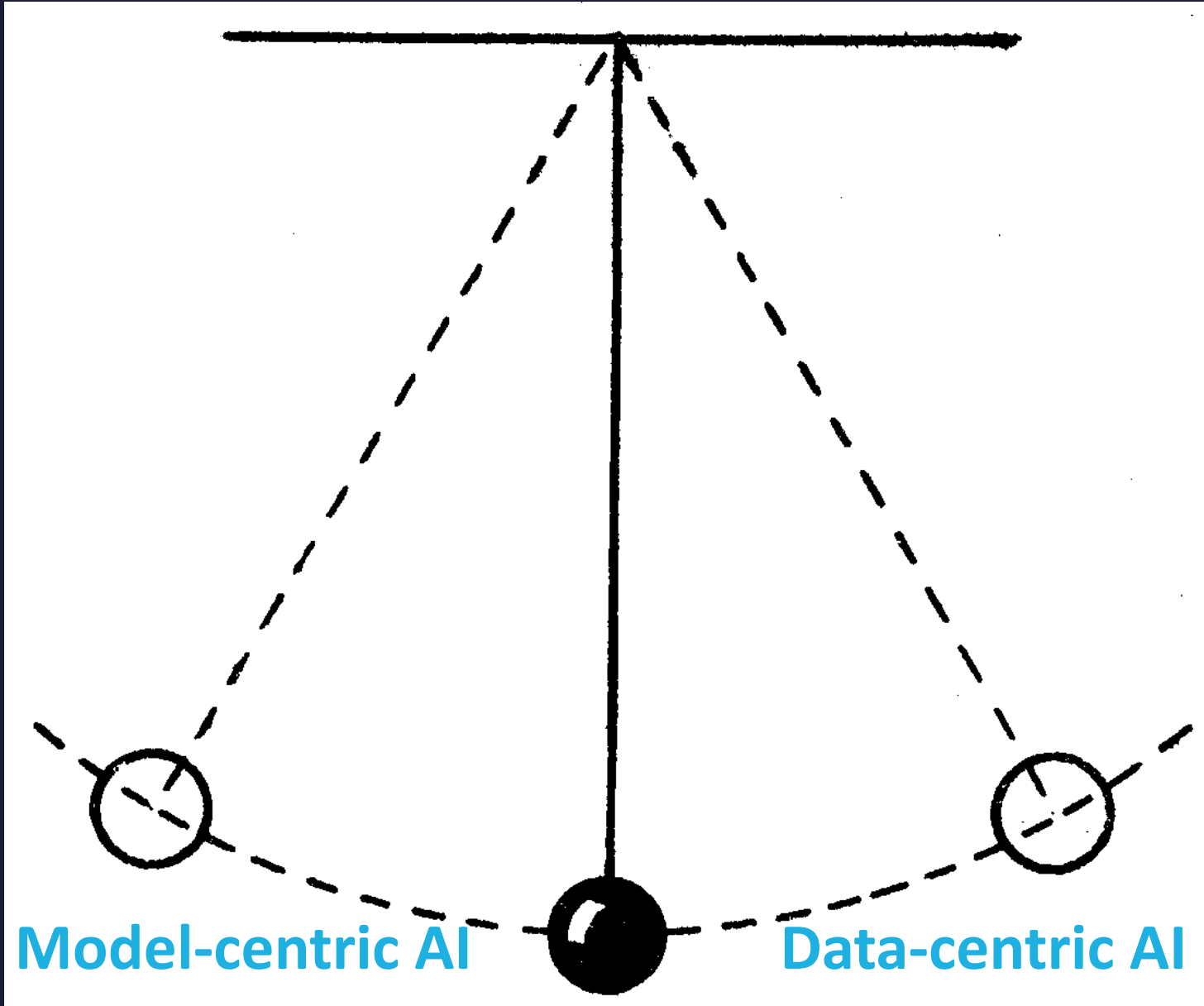
3. **Systematic** tools are developed

# "But I *like* building models!"

# Thank you!