



**CORNELL
TECH**

Fall 2019

Deep Learning Clinic

Lecture 4 - Data

Jin Sun

10/1/2019

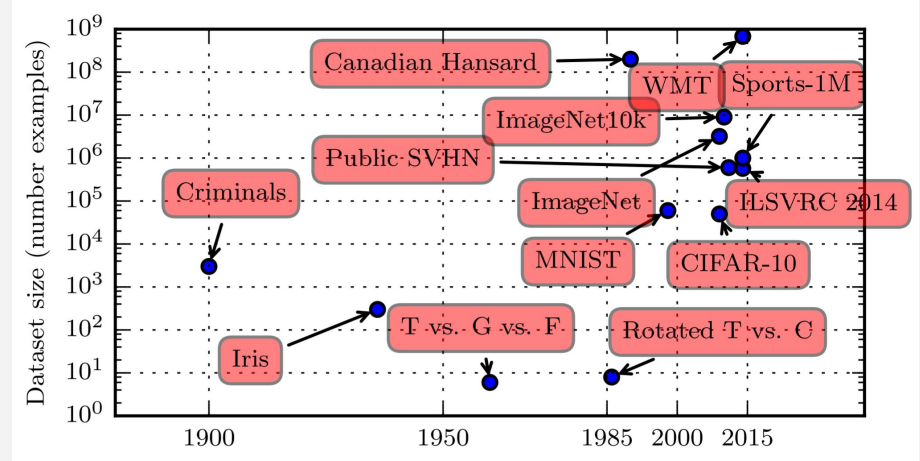
Summary

- **Overview**
- Existing Dataset
- Build A Dataset
- Data Loading and Processing in PyTorch
- Amazon MTurk Tutorial

Main Reasons Behind Deep Learning's Success



Hardware



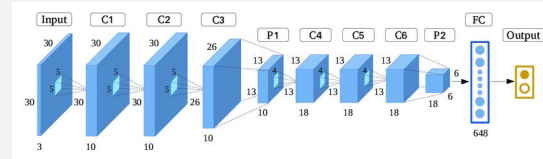
Data

Data in Machine Learning/Deep Learning Pipeline



Training
Data

Learning
Algo



Testing
Data

Model



Importance of Dataset

Datasets are used as benchmarks to compare learning systems

IMAGENET Large Scale Visual Recognition Challenge 2017 (ILSVRC2017)

[DET](#) [LOC](#) [VID](#) [Team information](#)

Legend:

Yellow background = winner in this task according to this metric; authors are willing to reveal the method

White background = authors are willing to reveal the method

Grey background = authors chose not to reveal the method

Italics = authors requested entry not participate in competition

Object detection (DET)^[top]

Task 1a: Object detection with provided training data

Ordered by number of categories won

Team name	Entry description	Number of object categories won	mean AP
BDAT	submission4	85	0.731392
BDAT	submission3	65	0.732227
BDAT	submission2	30	0.723712
DeepView(ETRI)	Ensemble_A	10	0.593084
NUS-Qihoo_DPNs (DET)	Ensemble of DPN models	9	0.656932
KAISTNIA_ETRI	Ensemble Model5	1	0.61022
KAISTNIA_ETRI	Ensemble Model4	0	0.609402
KAISTNIA_ETRI	Ensemble Model2	0	0.608299
KAISTNIA_ETRI	Ensemble Model1	0	0.608278
KAISTNIA_ETRI	Ensemble Model3	0	0.60631
DeepView(ETRI)	Single model A using ResNet for detection	0	0.587519

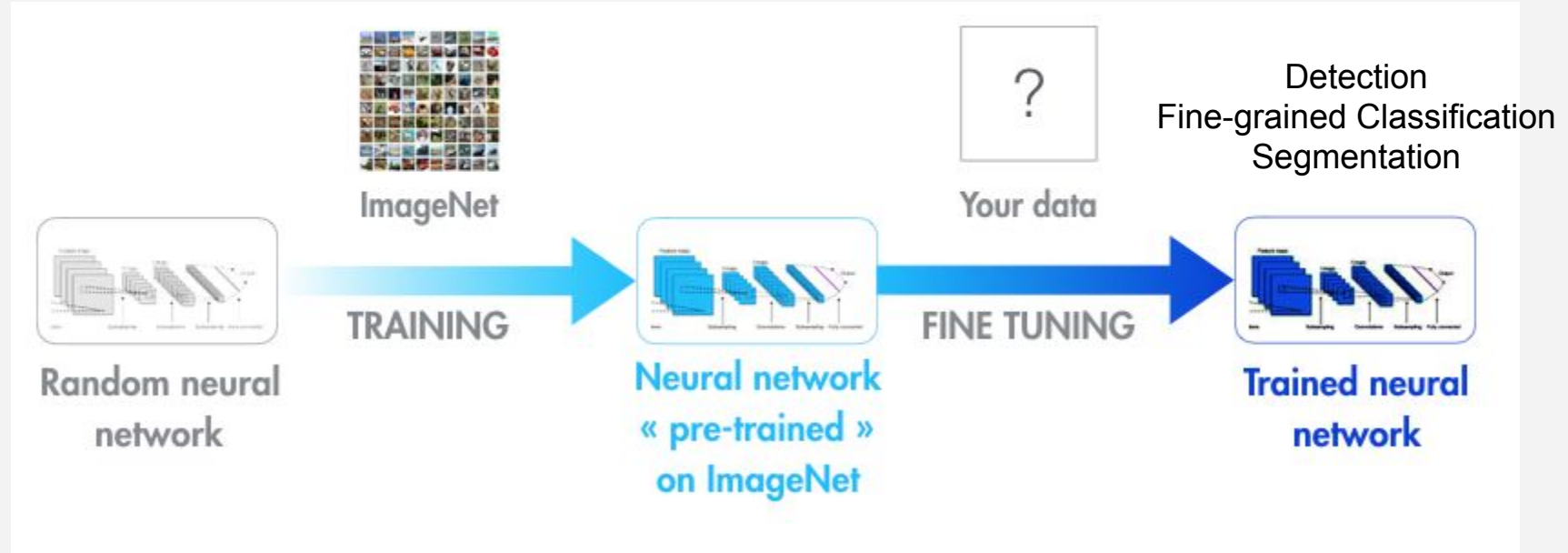
Importance of Dataset

Datasets are used as benchmarks to compare learning systems

Russian-English			
#	score	range	system
1	0.583	1	AFRL-PE
2	0.299	2	ONLINE-B
3	0.190	3-5	ONLINE-A
	0.178	3-5	PROMT-HYBRID
	0.123	4-7	PROMT-RULE
	0.104	5-8	UEDIN-PHRASE
	0.069	5-8	Y-SDA
	0.066	5-8	ONLINE-G
4	-0.017	9	AFRL
5	-0.159	10	UEDIN-SYNTAX
6	-0.306	11	KAZNU
7	-0.487	12	RBMT1
8	-0.642	13	RBMT4
English-Russian			
#	score	range	system
1	0.575	1-2	PROMT-RULE
	0.547	1-2	ONLINE-B
2	0.426	3	PROMT-HYBRID
3	0.305	4-5	UEDIN-UNCNSTR
	0.231	4-5	ONLINE-G
4	0.089	6-7	ONLINE-A
	0.031	6-7	UEDIN-PHRASE
5	-0.920	8	RBMT4
6	-1.284	9	RBMT1
French-English			
#	score	range	system
1	0.608	1	UEDIN-PHRASE
2	0.479	2-4	KIT
	0.475	2-4	ONLINE-B
	0.428	2-4	STANFORD
3	0.331	5	ONLINE-A
4	-0.389	6	RBMT1
5	-0.648	7	RBMT4
6	-1.284	8	ONLINE-C
English-French			
#	score	range	system
1	0.327	1	ONLINE-B
2	0.232	2-4	UEDIN-PHRASE
	0.194	2-5	KIT
	0.185	2-5	MATRAN
	0.142	4-6	MATRAN-RULES
	0.120	4-6	ONLINE-A
3	0.003	7-9	UU-DOCENT
	-0.019	7-10	PROMT-HYBRID
	-0.033	7-10	UA
	-0.069	8-10	PROMT-RULE
	-0.215	11	RBMT1
4	-0.328	12	RBMT4
6	-0.540	13	ONLINE-C
Hindi-English			
#	score	range	system
1	1.326	1	ONLINE-B
2	0.559	2-3	ONLINE-A
	0.476	2-4	UEDIN-SYNTAX
	0.434	3-4	CMU
3	0.323	5	UEDIN-PHRASE
4	-0.198	6-7	AFRL
	-0.280	6-7	IIT-BOMBAY
5	-0.549	8	DCU-LINGO24
6	-2.092	9	IIIT-HYDERABAD
English-Hindi			
#	score	range	system
1	1.008	1	ONLINE-B
2	0.915	2	ONLINE-A
3	0.214	3	UEDIN-UNCNSTR
4	0.120	4-5	UEDIN-PHRASE
	0.054	4-5	CU-MOSES
5	-0.111	6-7	IIT-BOMBAY
	-0.142	6-7	IPN-UPV-CNTXT
6	-0.233	8-9	DCU-LINGO24
	-0.261	8-9	IPN-UPV-NODEV
7	-0.449	10-11	MANAWI-H1
	-0.494	10-11	MANAWI
8	-0.622	12	MANAWI-RMOOV

Importance of Dataset

Datasets are used to learn a general purpose prior

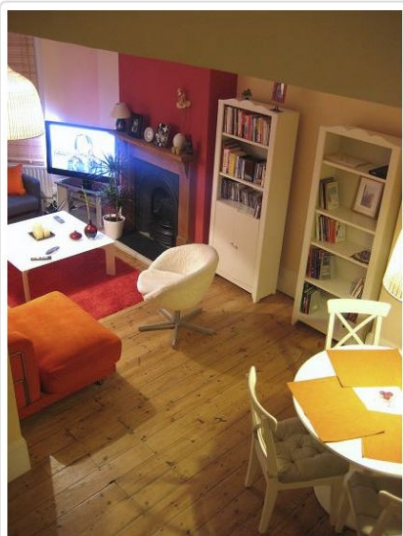


Importance of Dataset

New datasets inspire novel algorithms and research problems

Question : How many shelves?

Original Image | 8



Complementary Image | 11



Visual Q&A

Importance of Dataset

New datasets inspire novel algorithms and research problems



Recommendation Systems

Music, books, videos

Online shopping

Financial

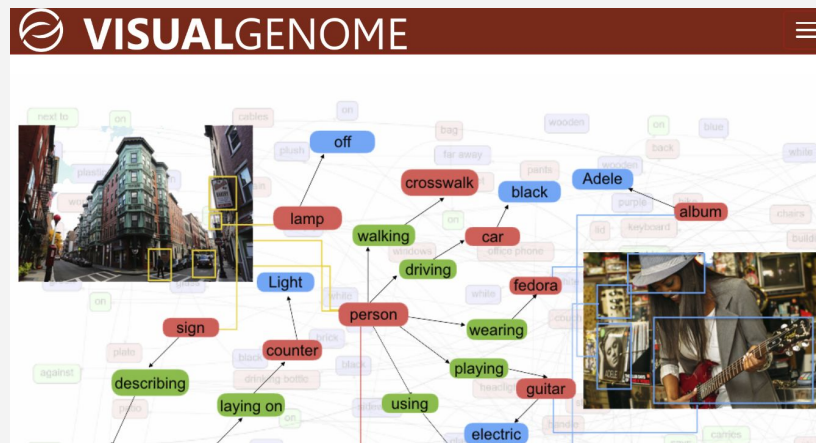
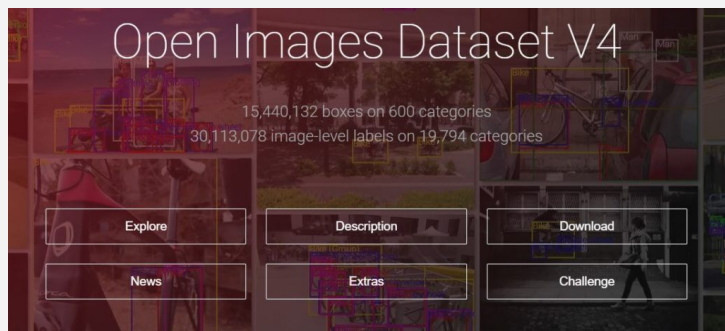
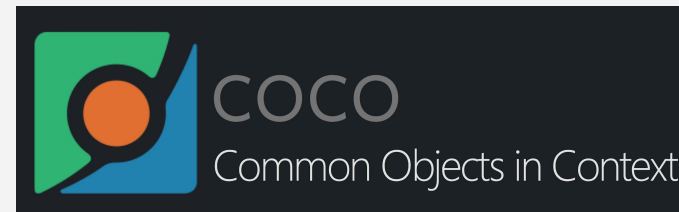
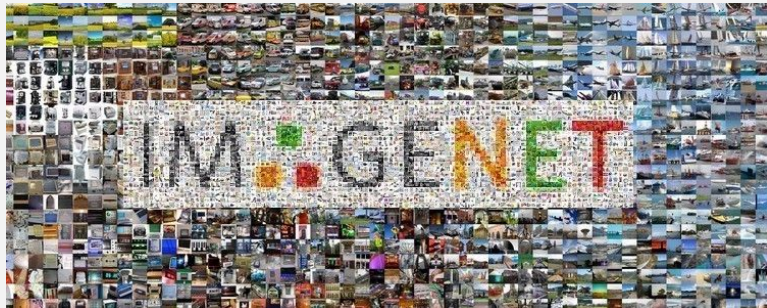
Online dating

...

Summary

- Overview
- **Existing Dataset**
- Build A Dataset
- Data Loading and Processing in PyTorch
- Amazon MTurk Tutorial

Existing Dataset - Vision



Existing Dataset - Natural Language

[IMDB Reviews](#), [Sentiment140](#) (Sentiment Analysis)

[1 Billion Word Language Model Benchmark](#) (Language Modeling)

[WordNet](#) (Database for English 'synsets')

[Google Books Ngrams](#)

Existing Dataset - Others

[HealthData.gov](#) (Health Care)

[OASIS brain images](#)

[Data.gov](#) (agriculture, climate, ecosystems, public safety...)

[Kaggle Dataset](#)

Summary

- Overview
- Existing Dataset
- **Build A Dataset**
 - Data Collection
 - Annotation
 - Verification
 - Tools
- Data Loading and Processing in PyTorch
- Amazon MTurk Tutorial

Why Build Your Own Dataset

Variation

Existing datasets do not contain enough variety.

E.g., non-traditional lighting and poses.

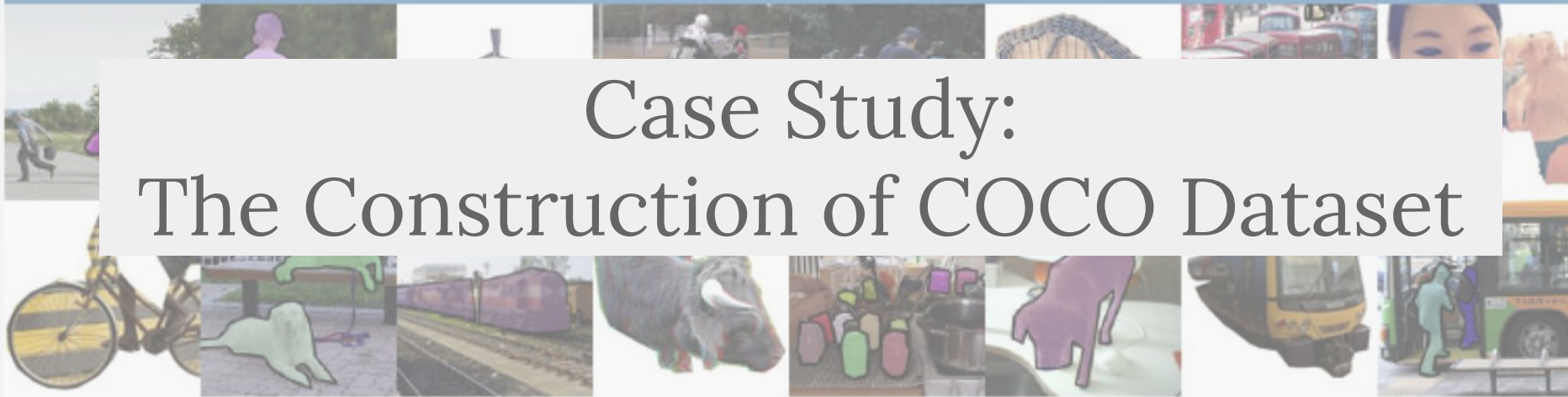
Annotation

Existing datasets do not provide the information you need.

E.g., no object segmentation masks in ImageNet.

Dataset examples

Case Study: The Construction of COCO Dataset



COCO Dataset Statistics

What is COCO?



COCO is a large-scale object detection, segmentation, and captioning dataset. COCO has several features:

- ✓ Object segmentation
- ✓ Recognition in context
- ✓ Superpixel stuff segmentation
- ✓ 330K images (>200K labeled)
- ✓ 1.5 million object instances
- ✓ 80 object categories
- ✓ 91 stuff categories
- ✓ 5 captions per image
- ✓ 250,000 people with keypoints

Collaborators

Tsung-Yi Lin Google Brain

Genevieve Patterson MSR, Trash TV

Matteo R. Ronchi Caltech

Yin Cui Cornell Tech

Michael Maire TTI-Chicago

Serge Belongie Cornell Tech

Lubomir Bourdev WaveOne, Inc.

Ross Girshick FAIR

James Hays Georgia Tech

Pietro Perona Caltech

Deva Ramanan CMU

Larry Zitnick FAIR

Piotr Dollár FAIR

Sponsors



CVDF



Microsoft

facebook



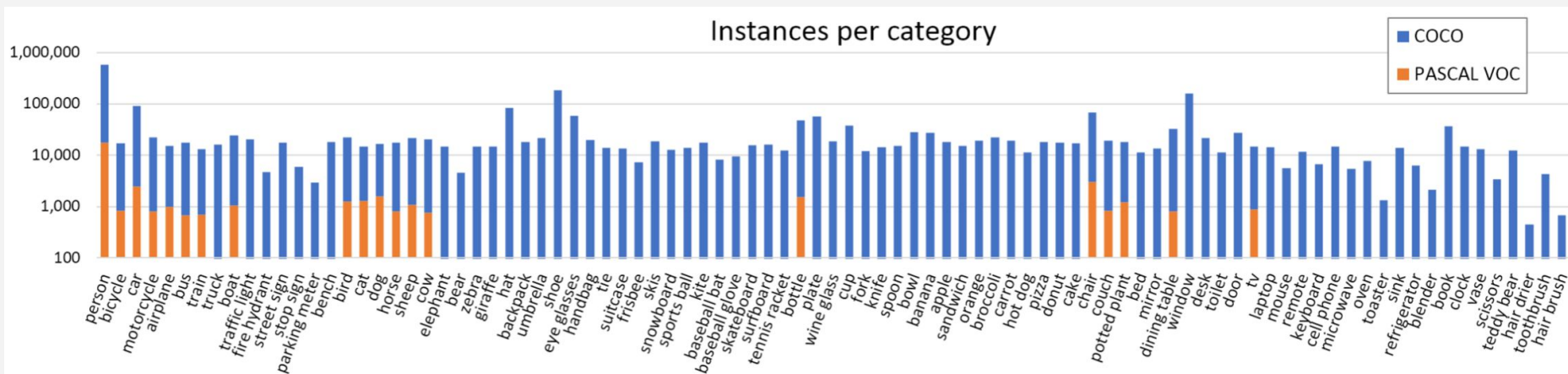
Mighty Ai

Data Collection

Identify Object Categories

PASCAL VOC + frequently used words for objects + survey on 4-8 years old children = 272 candidates

Voting to get final categories: 91.

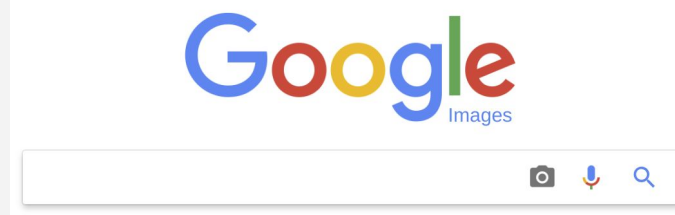


Data Collection

Collect Images For Each Object Category



Iconic Images



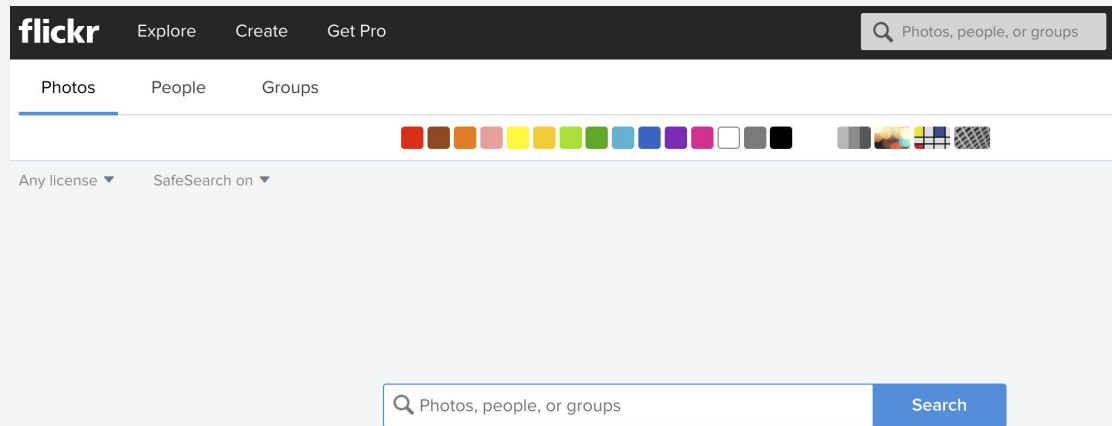
Data Collection

Collect Images For Each Object Category

328,000 images in total.



Non-Iconic Images



Data Annotation

How to label over 2.5 million object instances in 300K+ images?

Crowdsourcing.

Annotation Pipeline



(a) Category labeling



(b) Instance spotting



(c) Instance segmentation

Data Annotation

How to label over 2.5 million object instances in 300K+ images?

Crowdsourcing.

Annotation Pipeline



(a) Category labeling

8 Workers Per Image

~20k Worker Hours

Data Annotation

How to label over 2.5 million object instances in 300K+ images?

Crowdsourcing.

8 Workers Per Image

~10k Worker Hours



(b) Instance spotting

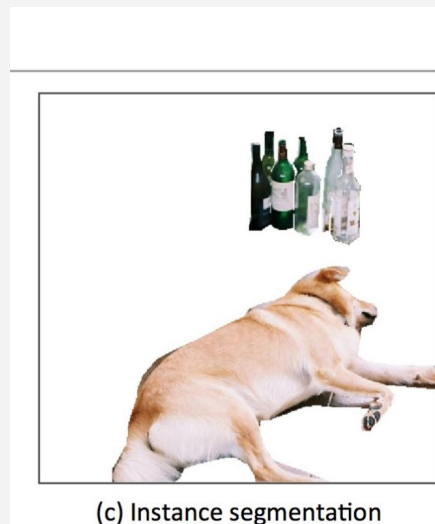
Data Annotation

How to label over 2.5 million object instances in 300K+ images?

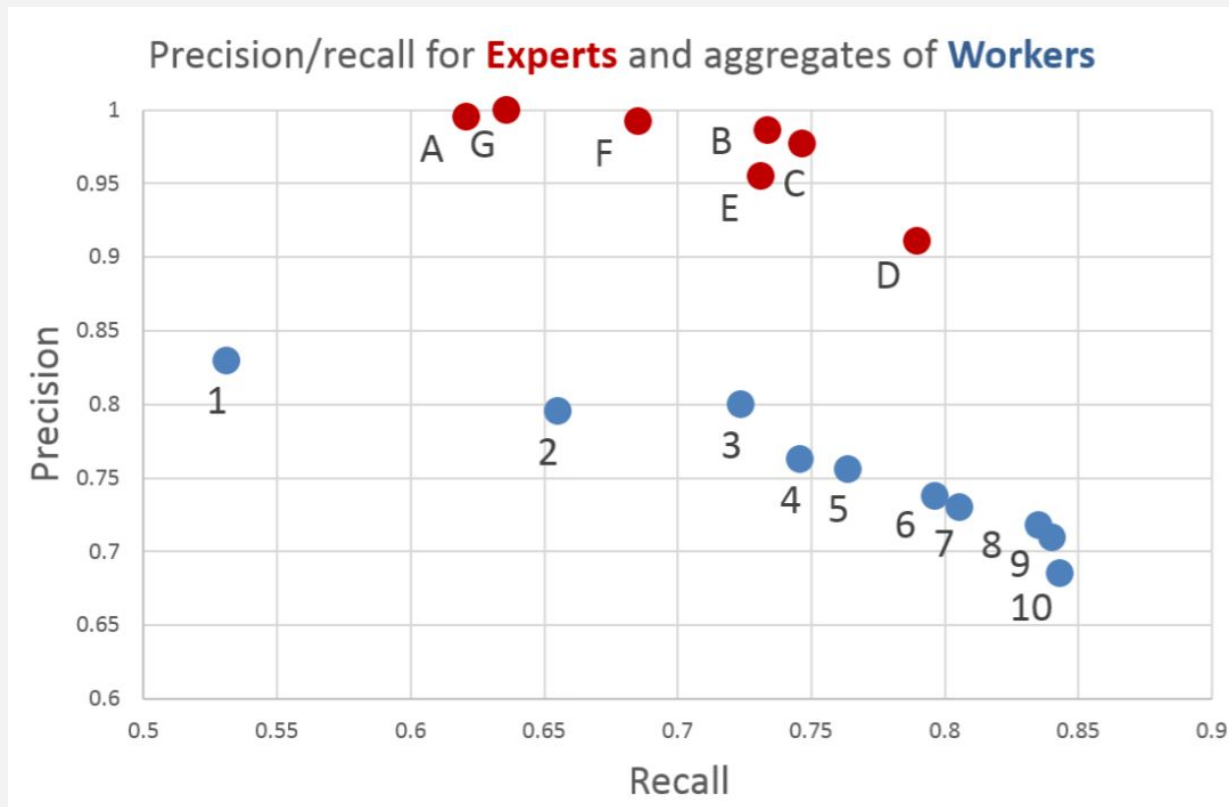
Crowdsourcing.

An expensive task. Only 1 worker per image.

Training stage enforced.



Data Verification



Tools

Data Source

Google/Bing Search, Flickr, Instagram, Google Map/Streetview, Satellite

Visual Annotation

[VGG Image Annotator](#), [Video Annotation Tool](#), [Scalabel](#)

Or Build your own (HTML+JS)

Crowdsourcing

Amazon MTurk

Today

- Overview
- Existing Dataset
- Build A Dataset
- **Data Loading and Processing in PyTorch**
 - Dataloader Class
 - Transforms
 - torchvision
- Amazon MTurk Tutorial

https://pytorch.org/tutorials/beginner/data_loading_tutorial.html

<https://stanford.edu/~shervine/blog/pytorch-how-to-generate-data-parallel>

PyTorch Dataloader

The basic class provides useful tools to load and prepare data.

```
import torch
from torch.utils import data

class Dataset(data.Dataset):
    'Characterizes a dataset for PyTorch'
    def __init__(self, data_files, labels):
        'Initialization'
        self.labels = labels # list of labels for each data sample
        self.data_files = data_files # list of file names storing the data

    def __len__(self):
        'Denotes the total number of samples'
        return len(self.data_files)

    def __getitem__(self, index):
        'Generates one sample of data'
        ...
```

PyTorch Dataloader

The basic class provides useful tools to load and prepare data.

```
import torch
from torch.utils import data

class Dataset(data.Dataset):

    ...

    def __getitem__(self, index):
        'Generates one sample of data'
        # Select sample
        datafile = self.data_files[index]

        # Load data and get label
        im = torch.load(datafile)
        label = self.labels[index]

        return im, label
```

PyTorch DataLoader

The basic class provides useful tools to load and prepare data.

```
import torch
from torch.utils import data

class Dataset(data.Dataset):

    ...

    def __getitem__(self, index):
        'Generates one sample of data'
        # Select sample
        datafile = self.data_files[index]

        # Load data and get label
        im = torch.load(datafile)
        label = self.labels[index]

        return im, label

training_set = Dataset(data_files, labels)
training_loader = data.DataLoader(training_set, **params)
```

PyTorch Dataloader Params

<https://pytorch.org/docs/stable/data.html#torch.utils.data.DataLoader>

Commonly used:

- ❑ dataset (*Dataset*) – dataset from which to load the data.
- ❑ batch_size (*int, optional*) – how many samples per batch to load (default: 1).
- ❑ shuffle (*bool, optional*) – set to `True` to have the data reshuffled at every epoch (default: `False`).
- ❑ num_workers (*int, optional*) – how many subprocesses to use for data loading. 0 means that the data will be loaded in the main process. (default: 0)
- ❑ drop_last (*bool, optional*) – set to `True` to drop the last incomplete batch, if the dataset size is not divisible by the batch size. If `False` and the size of dataset is not divisible by the batch size, then the last batch will be smaller. (default: `False`)

Use PyTorch Dataloader

Get data in the training loop:

```
training_set = Dataset(partition['train'], labels)
training_loader = data.DataLoader(training_set, **params)

for (i, batch_data) in enumerate(training_loader):
    # process the data
    output = net(batch_data)
```

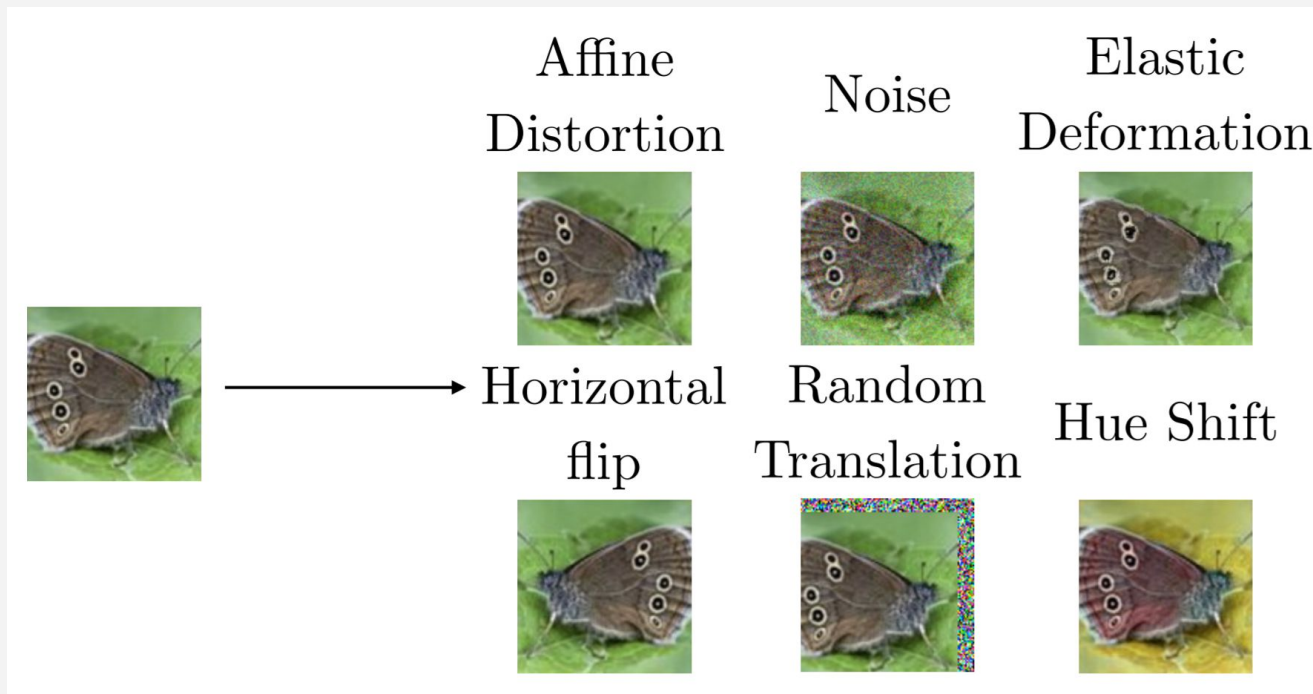
Get data in one batch (for debugging):

```
batch_data = next(iter(training_loader))

output = net(batch_data)
```


Data Transforms

Data Augmentation as regularization on deep neural networks



Data Transforms

```
import torch
from torch.utils import data

class Dataset(data.Dataset):

    ...

    def __getitem__(self, index):
        'Generates one sample of data'
        # Select sample
        datafile = self.data_files[index]

        # Load data and get label
        im = torch.load(datafile)
        label = self.labels[index]

        if self.transform:
            im = self.transform(im)

        return im, label
```

Data Transforms

```
import torch
from torch.utils import data

class Dataset(data.Dataset):

    ...

    def __getitem__(self, index):
        'Generates one sample of data'
        # Select sample
        datafile = self.data_files[index]

        # Load data and get label
        im = torch.load(datafile)
        label = self.labels[index]

        if self.transform:
            im = self.transform(im)

        return im, label
```

```
class RandomCrop(object):
    """Crop randomly the image in a sample.

    Args:
        output_size (tuple or int): Desired output size. If
        int, square crop
        is made.
    """
    def __init__(self, output_size):
        assert isinstance(output_size, (int, tuple))
        if isinstance(output_size, int):
            self.output_size = (output_size, output_size)
        else:
            assert len(output_size) == 2
            self.output_size = output_size

    def __call__(self, sample):
        image, landmarks = sample['image'], sample['landmarks']
        h, w = image.shape[:2]
        new_h, new_w = self.output_size
        top = np.random.randint(0, h - new_h)
        left = np.random.randint(0, w - new_w)

        image = image[top: top + new_h, left: left + new_w]
        landmarks = landmarks - [left, top]

        return {'image': image, 'landmarks': landmarks}
```

Data Transforms

```
import torch
from torch.utils import data

class Dataset(data.Dataset):
    def __init__(self, data_files, labels):
        'Initialization'
        self.transform = RandomCrop(out_size)

    def __getitem__(self, index):
        'Generates one sample of data'
        # Select sample
        datafile = self.data_files[index]

        # Load data and get label
        im = torch.load(datafile)
        label = self.labels[index]

        if self.transform:
            im = self.transform(im)

        return im, label
```

```
class RandomCrop(object):
    """Crop randomly the image in a sample.

    Args:
        output_size (tuple or int): Desired output size. If
        int, square crop
        is made.
    """
    def __init__(self, output_size):
        assert isinstance(output_size, (int, tuple))
        if isinstance(output_size, int):
            self.output_size = (output_size, output_size)
        else:
            assert len(output_size) == 2
            self.output_size = output_size

    def __call__(self, sample):
        image, landmarks = sample['image'], sample['landmarks']
        h, w = image.shape[:2]
        new_h, new_w = self.output_size
        top = np.random.randint(0, h - new_h)
        left = np.random.randint(0, w - new_w)

        image = image[top: top + new_h, left: left + new_w]
        landmarks = landmarks - [left, top]

        return {'image': image, 'landmarks': landmarks}
```

torchvision

A convenient package provides common dataset setting and transforms.

For example, `torchvision.datasets.ImageFolder` is a generic data loader where the images are arranged in this way:

`root/dog/xxx.png`

`root/cat/123.png`

`root/dog/xyy.png`

`root/cat/nsdf3.png`

`root/dog/xxz.png`

`root/cat/asd932_.png`

...

Also provides common vision datasets: MNIST, COCO, ImageNet, CIFAR...

torchvision

Common transforms <https://pytorch.org/docs/stable/torchvision/transforms.html#torchvision-transforms>

```
>>> transforms.Compose([  
  
>>>     transforms.CenterCrop(10),  
  
>>>     transforms.ToTensor(),  
  
>>> ])
```

Crop, ColorJitter, RandomAffine, ...

Today

- Overview
- Existing Dataset
- Build A Dataset
- Data Loading and Processing in PyTorch
- **Amazon MTurk Tutorial**

Amazon Mechanical Turk Tutorial



On Demand

Over 500K workers, 24x7

Speed

Work is done in parallel

Scalable

No minimum project size

Qualification

Set prerequisite to workers

MTurk Concepts

Requesters

Person creates tasks for Workers to work on.

Human Intelligence Tasks (HITs)

HIT is a single, self-contained task.

Assignment

Multiple Workers can be assigned to a single HIT.

A Worker can only accept a HIT once and submit one assignment per HIT.

Workers

Person completes assignments.

Approval and Payment

After assignment submission, if you approve the work, the HIT reward is draw from your MTurk account.

Qualification

Anyone can register as a worker. You can set qualification types such as approval rate to control the quality of submissions.

Common Use Cases

Image/Video Processing

MTurk is well-suited for processing images. While difficult for computers, it is a task that is extremely easy for people to do. In the past, companies have used MTurk to:



Tag objects found in an image to improve your search or advertising targeting



Review a set of images to select the best picture to represent a product



Audit user-uploaded images or videos to moderate content



Classify objects found in satellite imagery

Data Verification and Clean-up

Companies with large online directories or catalogs are using MTurk to identify duplicate entries and verify item details. Examples of this have included:



Removing duplicate content from business listings



Identifying incomplete or duplicate product listings in a catalog



Verifying restaurant details such as phone numbers or hours of operation



Converting unstructured data about locations into well-formed addresses

Common Use Cases

Information Gathering

The diversification and the scale of the MTurk workforce allows you to gather a breadth of information that would be almost impossible to do otherwise such as:



Allowing people to ask questions from a computer or mobile device about any topic and have Workers return the results

Example: Data Labeling Using MTurk

1. Setup

Python and Boto3 (AWS SDK).

2. Accounts

AWS and MTurk (Also need to link the two).

Purchasing Prepaid HITs.

3. Creating Tasks

Define a HIT and its reward.

4. Retrieving Results

Verify result, Add a bonus

No coding needed:

[Tutorial 1](#)

Command line approach:

[Tutorial 2](#)

Summary

- Overview
- Existing Dataset
- Build A Dataset
 - Data Collection
 - Annotation
 - Verification
 - Tools
- Data Loading and Processing in PyTorch
 - Dataloader Class
 - Transforms
 - torchvision
- Amazon MTurk Tutorial