DLI Accelerated Data Science Teaching Kit

# Lecture 2.4 - Data Annotation and Data Quality

# Data Annotation

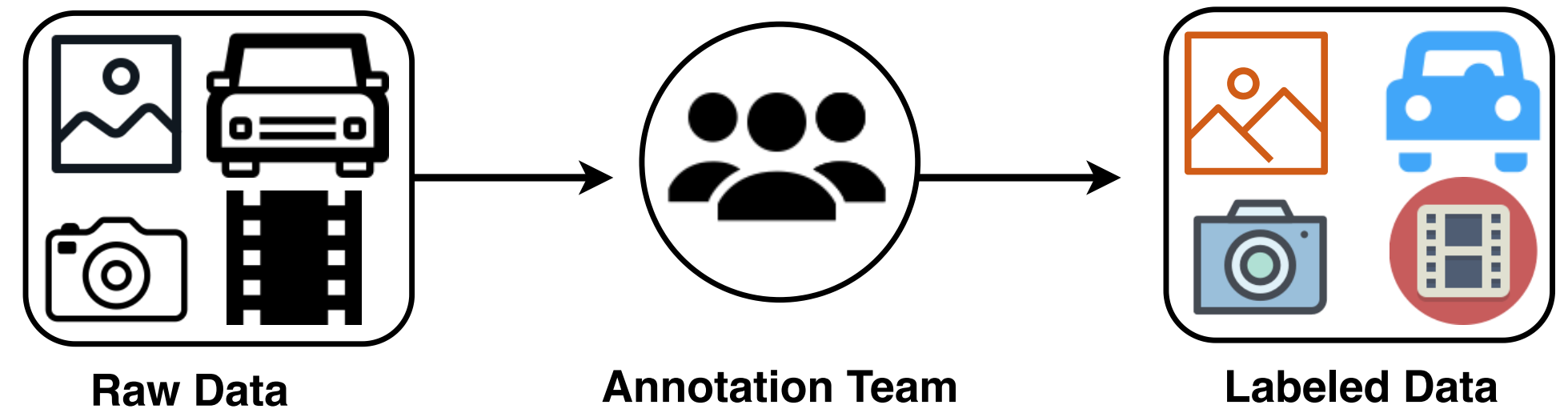Data annotation is the labelling of collected data for various applications.

- Image Processing

- Natural Language Processing

- Smart Grid

- … …



**Raw Data**          **Annotation Team**          **Labeled Data**

# Image Processing

Examples of data annotation

**Image Classification**



CAT

**Object Detection**



DOG, DOG, CAT

**Semantic Segmentation**
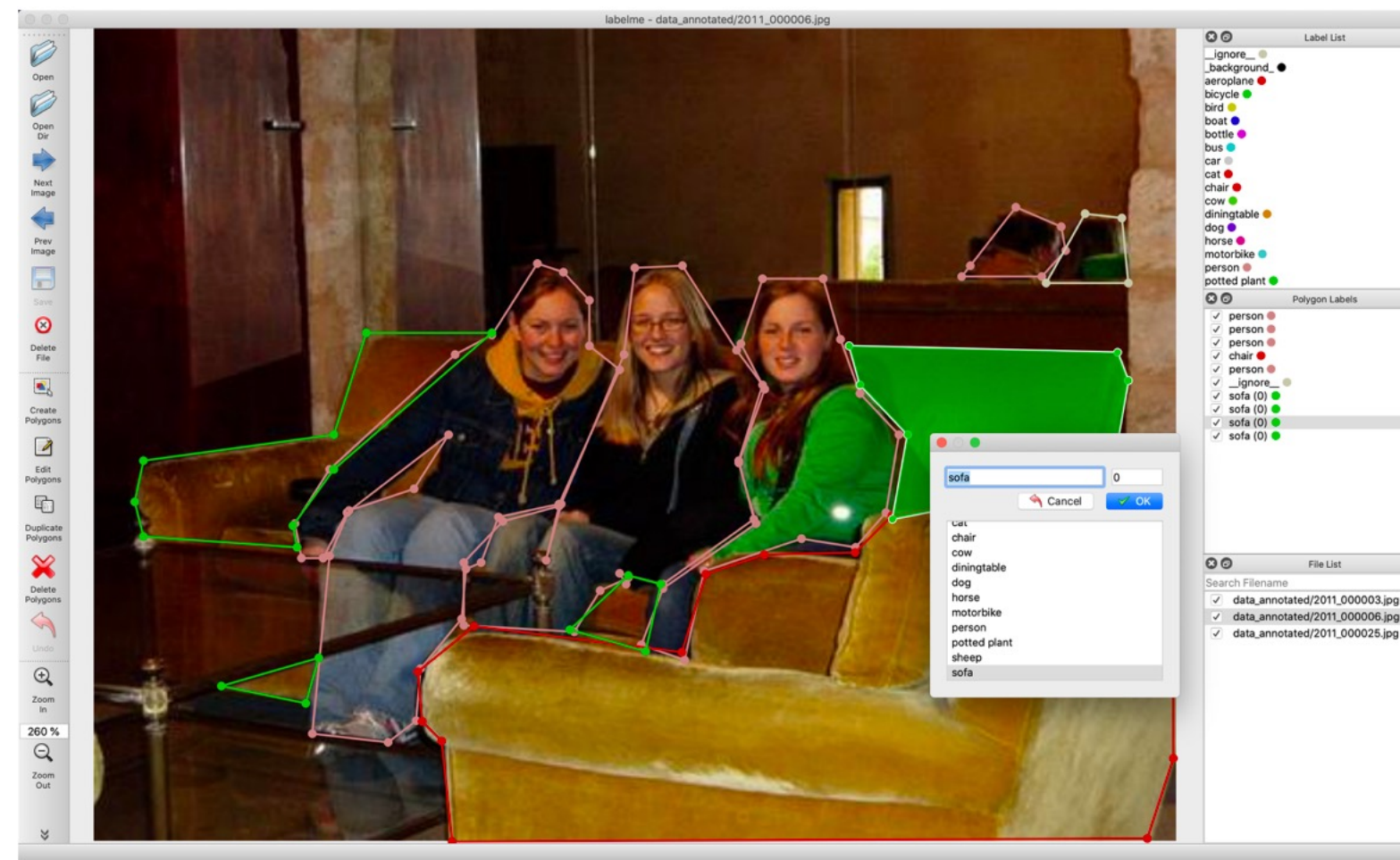




Sky

Trees

Cat

Grass

# LabelMe

A graphical image annotation tool inspired by http://labelme.csail.mit.edu.

It is written in Python and uses Qt for its graphical interface.
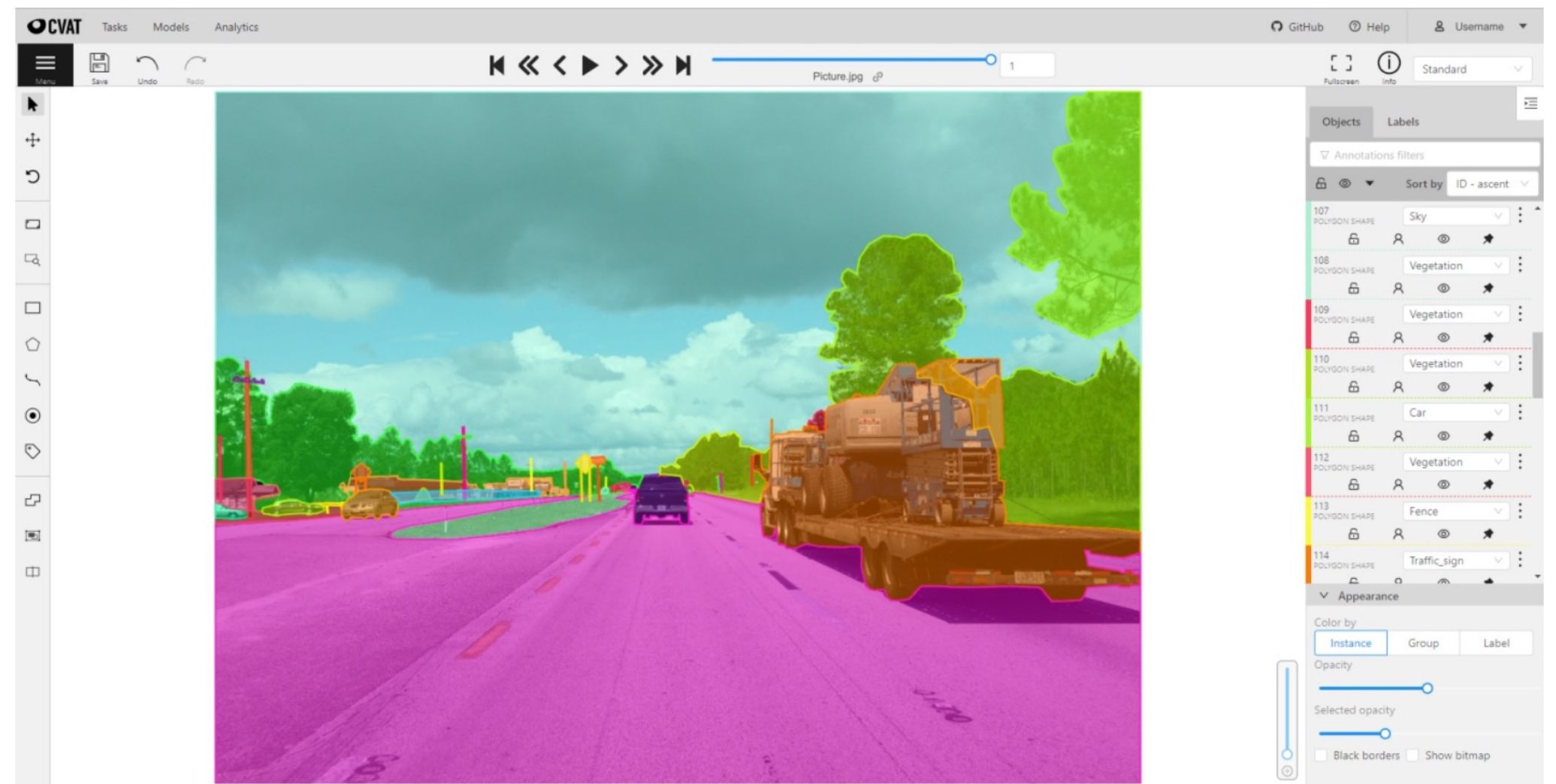
**Source: https://github.com/wkentaro/labelme**

# Computer Vision Annotation Tool (CVAT)

Open source

Auto annotation using trained DL models

Collaborative

Easy to deploy and maintain

**Source: https://github.com/openvinotoolkit/cvat**

# Einstein Vision Object Detection

Open Source

Image Size Change and Move

Label Change and Remove

Easy to Use

**Source: https://github.com/misu007/einstein-vision-object-detection-annotation-tool-sample**

# Natural Language Processing (NLP)

Sentiment Classification on Texts



Named Entity Recognition

# Quality Control for Data Annotation

Calculating the accuracy of an annotator compared to ground-truth data.

Calculating the overall agreement and reliability of a dataset.

Measuring inter-annotator agreement on a per-task basis to generate a confidence score for each training data label.

Designing architectures that incorporate subject matter experts into the annotation workflow.

Breaking up a task into simpler subtasks to improve accuracy, efficiency and quality control.

**Source: https://livebook.manning.com/book/human-in-the-loop-machine-learning/chapter-8/v-8/1**

# Data Quality (DQ)

**Data quality** refers to the state of qualitative or quantitative pieces of information.

Data is of high quality if it is "Fit for Use" in their intended operational, decision-making and other roles.

Quality Properties

- Relevance – does data meet basic needs?

- Accuracy – are key data elements correct?

- Timeliness – is data current?

# Data Quality (DQ)

**Data quality** refers to the state of qualitative or quantitative pieces of information.

Data is of high quality if it is "Fit for Use" in their intended operational, decision-making and other roles.

Quality Properties

- Comparability
  - Can several databases be combined into a data warehouse?

- Completeness
  - No missing records, no missing data elements

# Why Care About Data Quality?

High-quality data can be a major business asset, a unique source of competitive advantage.

Poor data quality can lower customer satisfaction.

Poor data quality can lower employee job satisfaction too, leading to excessive turnover and the resulting loss of key process knowledge.

Poor data quality can also breed organizational mistrust and make it hard to mount efforts that lead to needed improvements.

# An Example

Insurance companies need accurate and complete databases in order to:

- Accomplish Mandated Accounting Tasks

- Compute Insurance-in-force

- Compute Claim Rates

- Calculate Premiums

- Construct Models to Predict Future Results

**Source: https://redrockslocksmith.com/do-insurance-companies-cover-car-lockouts/**

# How Do You Obtain High-Quality Data?

Prevention: Keep Bad Data Out of the Database/List

Detection: Proactively Look for Bad Data Already Entered

Repair: Let the Bad Data Find You and Then Fix Things

Allocating Resources: How Much for Prevention, Detection, and Repair

DLI Accelerated Data Science Teaching Kit

# Thank You