, Rome GPE was the 11th ORDINAL most visited city in the world with 10.1 million CAR

s, the third ORDINAL most visited in the European Union ORG , and the most popular tourist destinatio

s historic centre is listed by UNESCO ORG as a World Heritage ORG Site.[14] Host city for the 1960

mer Olympics EVENT , Rome GPE is also the seat of several specialised agencies of the United Natio

# DATA ANNOTATION

## VAYIANOS PERTSAS

ORG , and national and international banks such as Unicredit ORG and BNL ORG . Rome GPE

business district is the home of many companies involved in the oil industry, the pharmaceutical industry, an

es. The presence of renowned international brands in the city have made Rome GPE an important centre

# INTRODUCTION

# WHAT IS ANNOTATION?

➡ Annotation is the process of producing extra information and associating it with a particular point in a document or other piece of information

➡ In Machine Learning, annotation is the process of labelling individual elements of data

# WHY DO WE USE ANNOTATION?

➡ To enhance our data with more information regarding particular data elements

➡ In Machine Learning, annotation is used in order to train ML algorithms by showing them the outcome we want them to predict

# TYPES OF ANNOTATION:

➡ Categorization / Classification

➡ Semantic Segmentation / Entity Annotation

➡ Semantic Association / Entity Linking

# TYPES OF DATA ANNOTATION:

➡ Image Annotation
➡ Video Annotation
➡ Audio Annotation
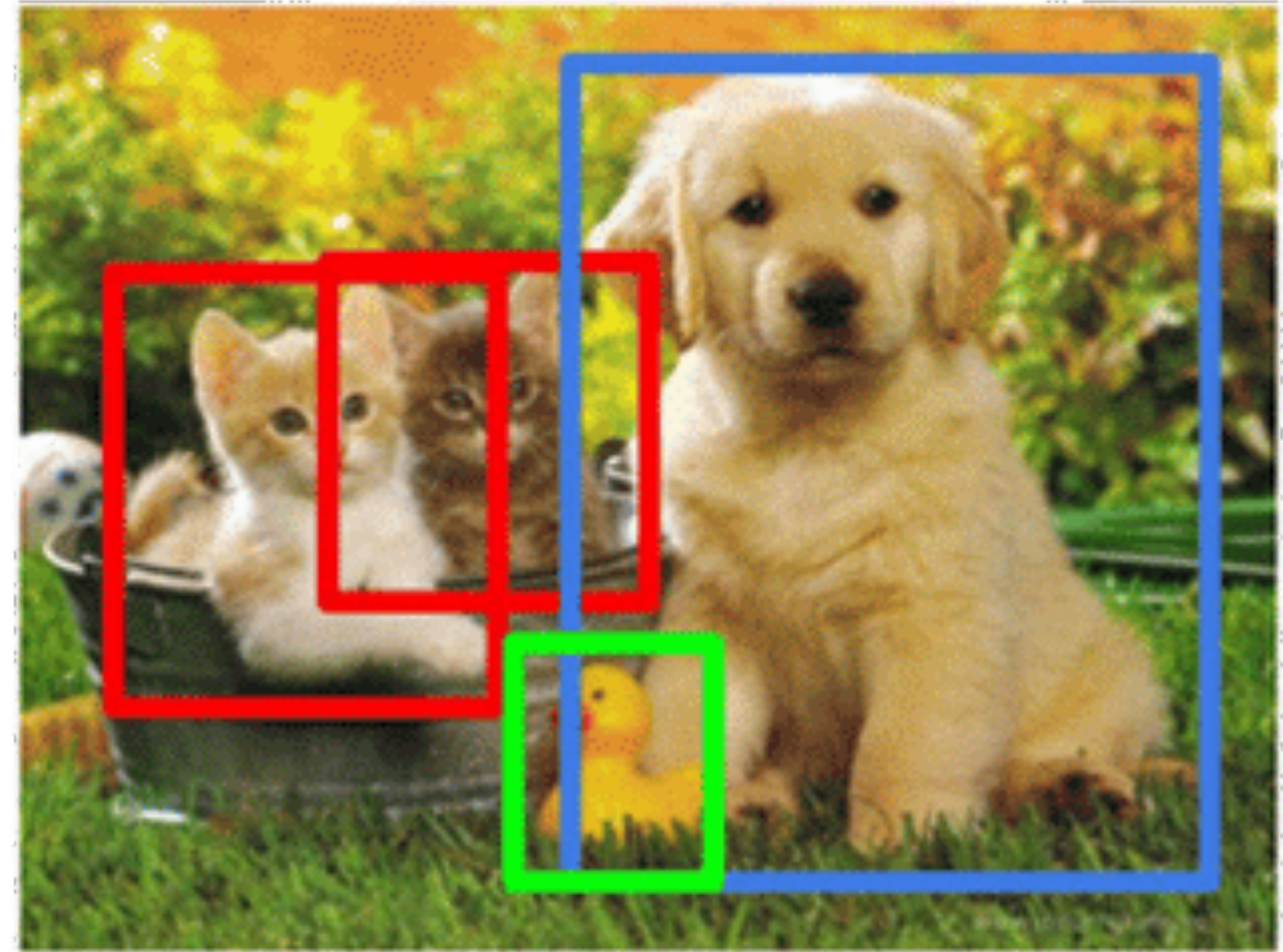➡ Text Annotation

# USE CASES OF DATA ANNOTATION:

➡ Image Annotation
   ➡ **Image Classification**
   ➡ Object Detection
   ➡ Image Captioning
   ➡ Optical Character Recognition

➡ Video Annotation

➡ Audio Annotation

➡ Text Annotation



CAT

# USE CASES OF DATA ANNOTATION:

➡ Image Annotation
  ➡ Image Classification
  ➡ **Object Detection**
  ➡ Image Captioning
  ➡ Optical Character Recognition

➡ Video Annotation
➡ Audio Annotation
➡ Text Annotation



CAT, DOG, DUCK

# USE CASES OF DATA ANNOTATION:

➡ Image Annotation
  ➡ Image Classification
  ➡ Object Detection
  ➡ **Image Captioning**
  ➡ Optical Character Recognition

➡ Video Annotation
➡ Audio Annotation
➡ Text Annotation



A couple of people standing next to an elephant.

# USE CASES OF DATA ANNOTATION:

➡️ Image Annotation
  ➡️ Image Classification
  ➡️ Object Detection
  ➡️ Image Captioning
  ➡️ **Optical Character Recognition**

➡️ Video Annotation

➡️ Audio Annotation

➡️ Text Annotation

# USE CASES OF DATA ANNOTATION:

➡ Image Annotation
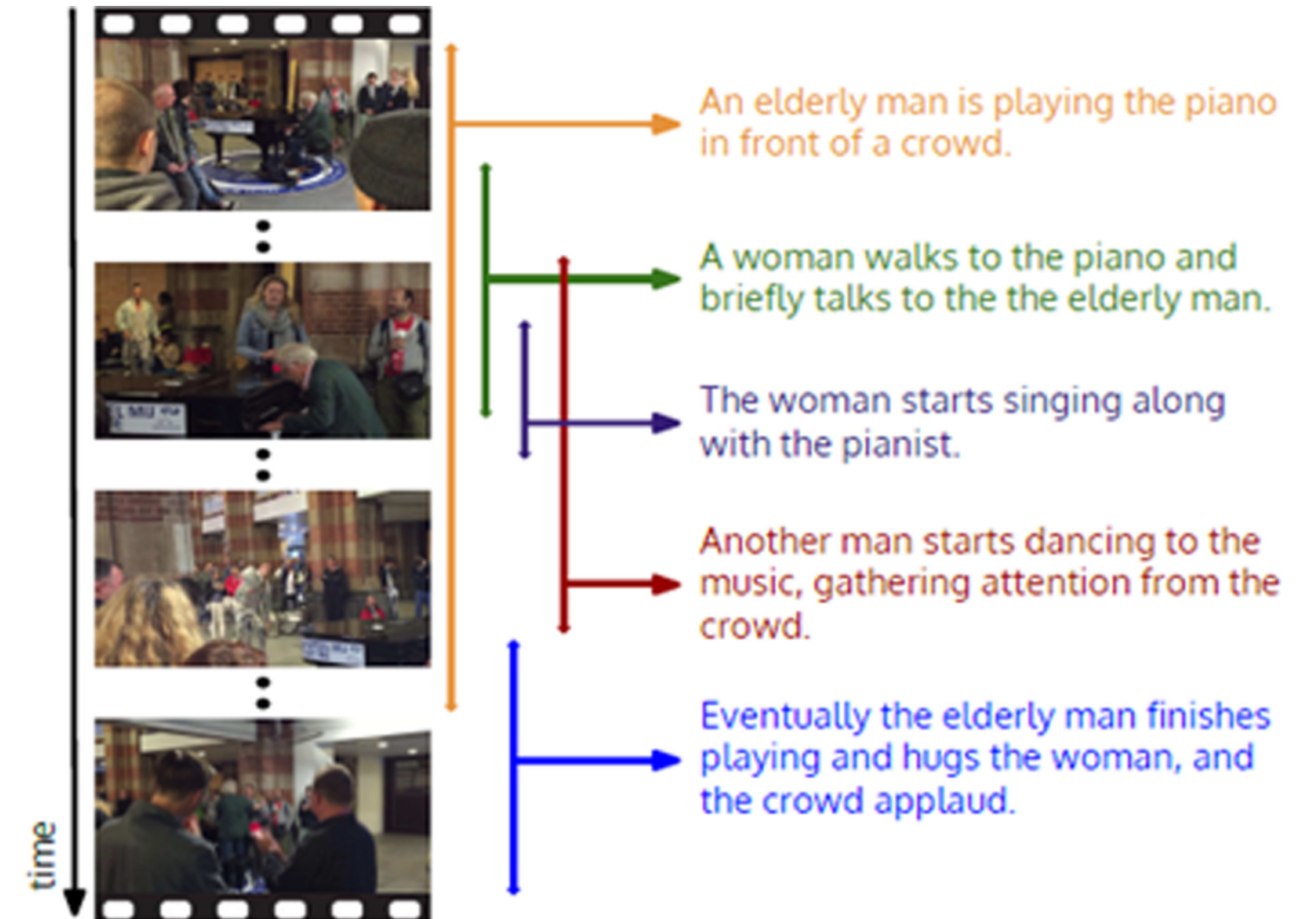➡ Video Annotation
➡ Audio Annotation
➡ Text Annotation

# USE CASES OF DATA ANNOTATION:

➡ Image Annotation
➡ Video Annotation
  ➡ **Video classification**
    ➡ Video captioning
    ➡ Video object detection and tracking

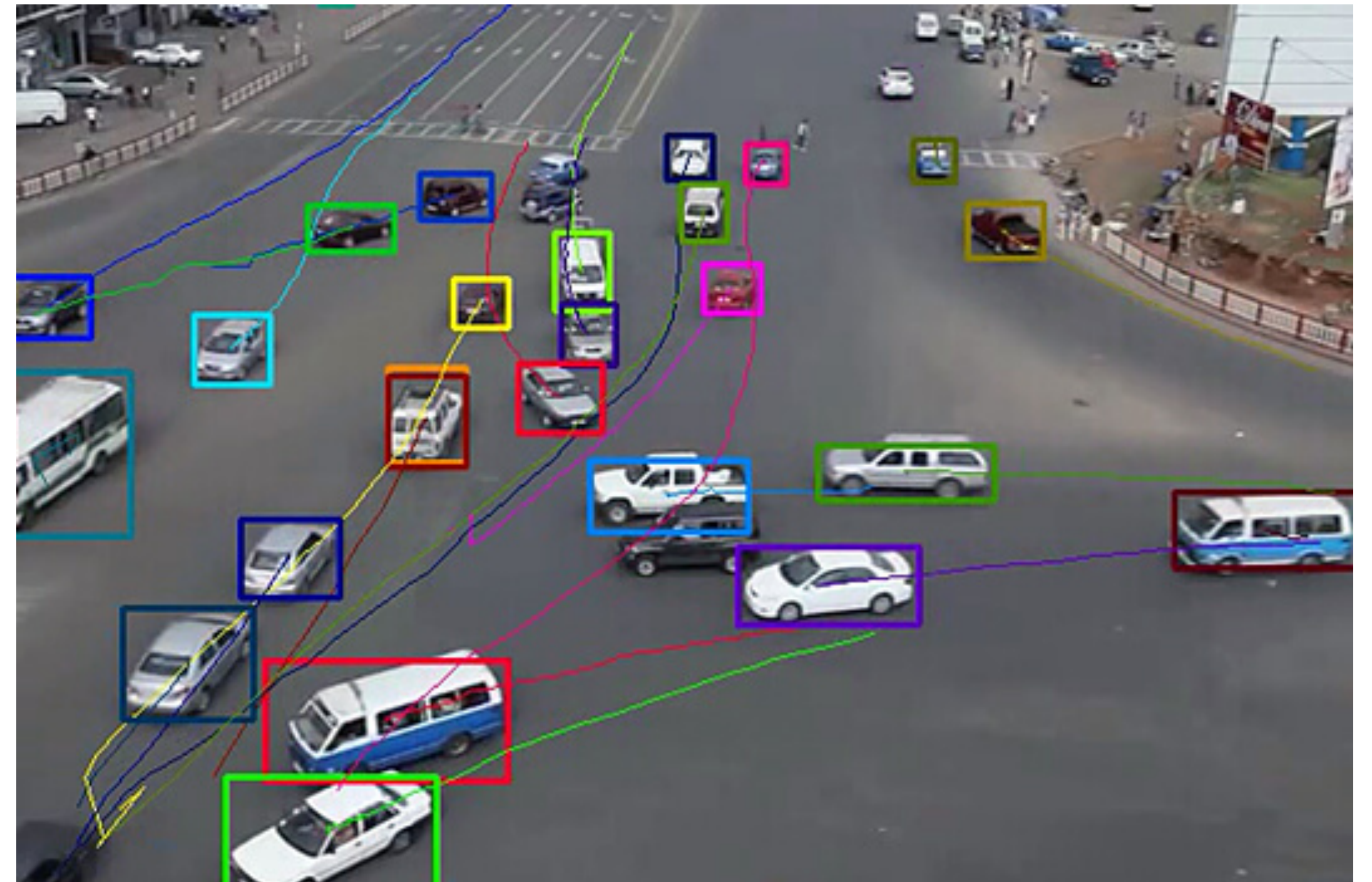➡ Audio Annotation
➡ Text Annotation

# USE CASES OF DATA ANNOTATION:

➡ Image Annotation

➡ Video Annotation

   ➡ Video classification

   ➡ **Video captioning**

   ➡ Video object detection and tracking

➡ Audio Annotation

➡ Text Annotation



An elderly man is playing the piano in front of a crowd.

A woman walks to the piano and briefly talks to the the elderly man.

The woman starts singing along with the pianist.

Another man starts dancing to the music, gathering attention from the crowd.

Eventually the elderly man finishes playing and hugs the woman, and the crowd applaud.

# USE CASES OF DATA ANNOTATION:

➡️ Image Annotation
➡️ Video Annotation
    ➡️ Video classification
    ➡️ Video captioning
    ➡️ **Video object detection and tracking**

➡️ Audio Annotation
➡️ Text Annotation

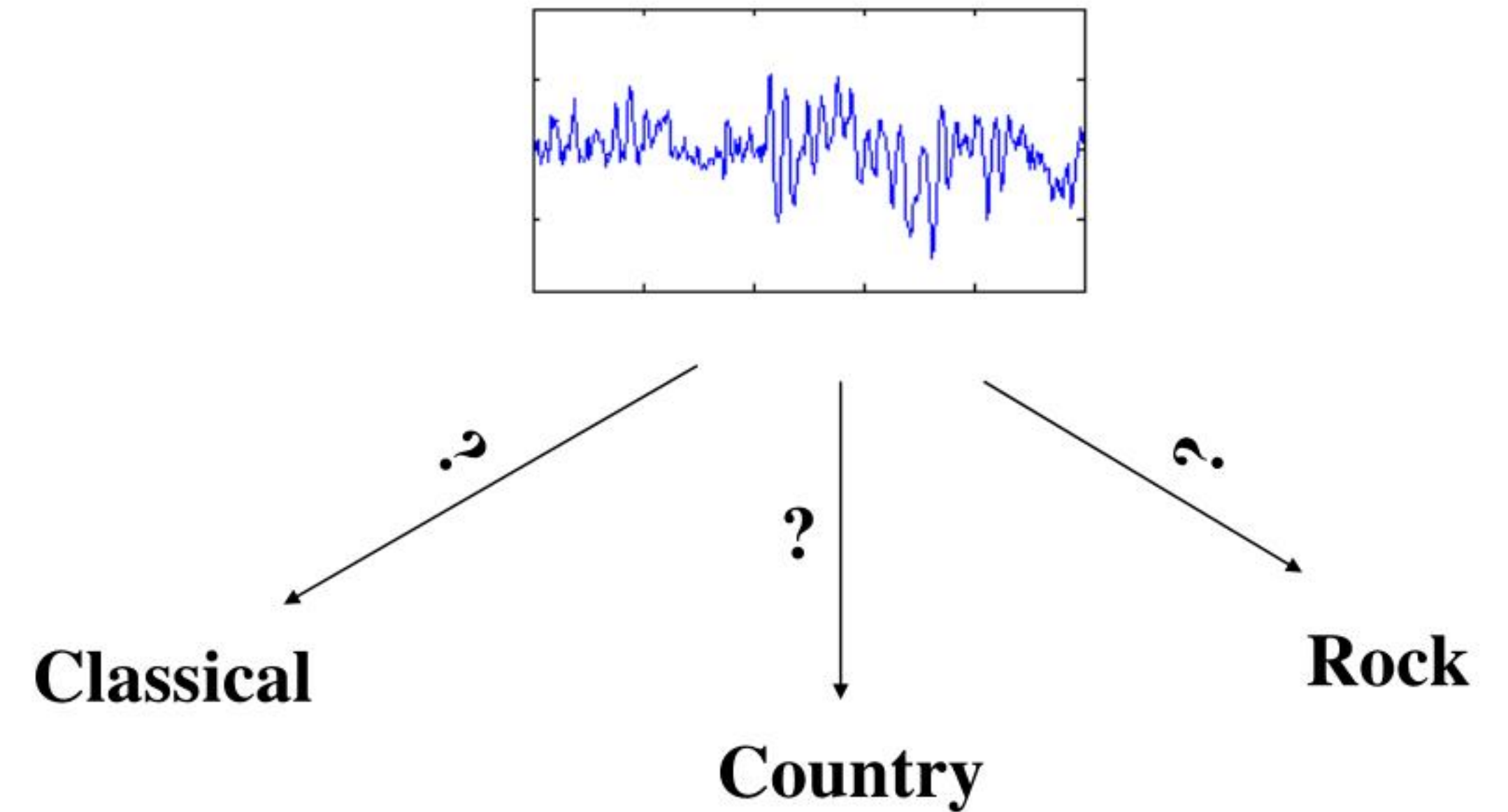# USE CASES OF DATA ANNOTATION:

➡ Image Annotation
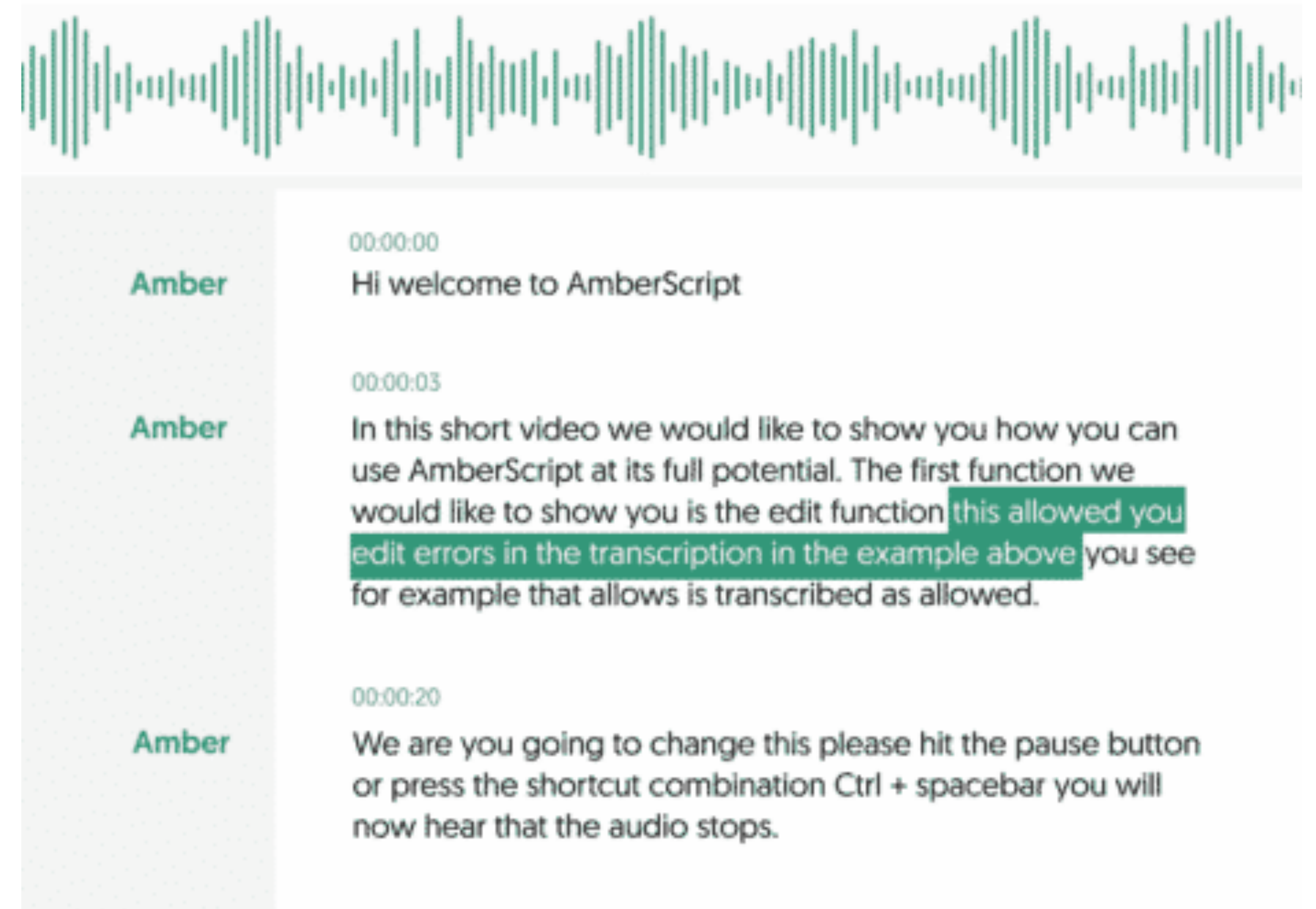➡ Video Annotation
➡ Audio Annotation
➡ Text Annotation

# USE CASES OF DATA ANNOTATION:

➡️ Image Annotation

➡️ Video Annotation

➡️ Audio Annotation

    ➡️ **Audio Classification**

    ➡️ Audio Transcription
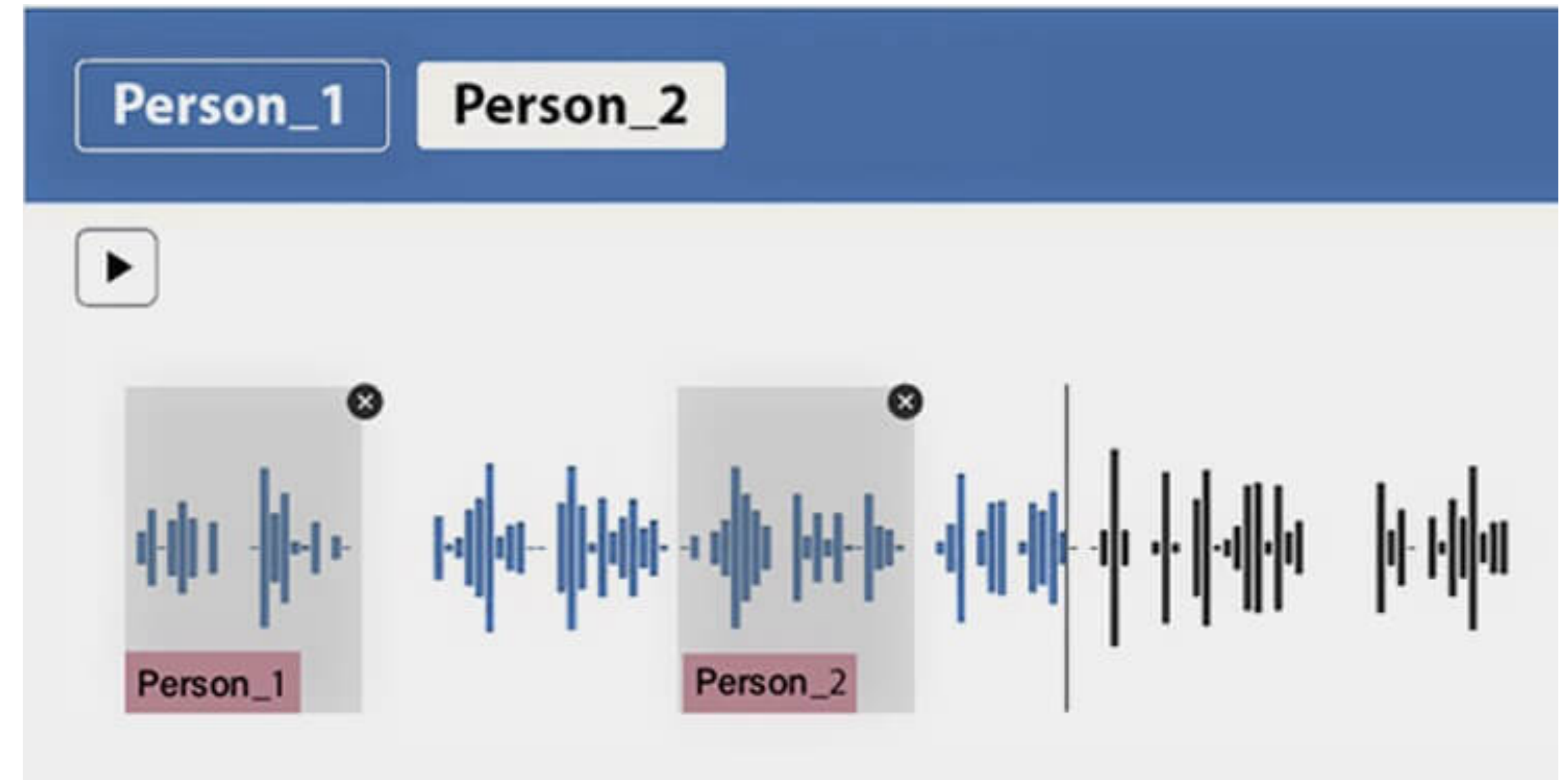
    ➡️ Speaker Detection

➡️ Text Annotation

# USE CASES OF DATA ANNOTATION:

➡️ Image Annotation

➡️ Video Annotation

➡️ Audio Annotation
  ➡️ Audio Classification
  ➡️ **Audio Transcription**
  ➡️ Speaker Detection

➡️ Text Annotation

# USE CASES OF DATA ANNOTATION:

➡️ Image Annotation

➡️ Video Annotation

➡️ Audio Annotation

   ➡️ Audio Classification

   ➡️ Audio Transcription

   ➡️ **Speaker Detection**

➡️ Text Annotation

# USE CASES OF DATA ANNOTATION:

➡ Image Annotation
➡ Video Annotation
➡ Audio Annotation
➡ Text Annotation

# USE CASES OF DATA ANNOTATION:

➡️ Image Annotation
➡️ Video Annotation
➡️ Audio Annotation
➡️ Text Annotation
➡️ **Text classification**
➡️ Language translation
➡️ Entity recognition
➡️ Entity linking

NLP in Health: A comprehensive look

| ☑ | health | 1 |
| ☑ | natural-language-processing | 2 |
| ☐ | computer-vision | 3 |
| ☐ | other | 4 |

# USE CASES OF DATA ANNOTATION:

➡️ Image Annotation

➡️ Video Annotation

➡️ Audio Annotation

➡️ Text Annotation

   ➡️ Text classification

   ➡️ **Language translation**

➡️ Entity recognition

➡️ Entity linking

---

**Source**

Dies ist ein deutscher Satz, der abre Fehler enthältt.

**Target**

This is a German sentence, which contains mistakes.

**Comment**

Korrektur: Da Fehler im deutschen Text!

# USE CASES OF DATA ANNOTATION:

➡ Image Annotation

➡ Video Annotation

➡ Audio Annotation

➡ Text Annotation

    ➡ Text classification

    ➡ Language translation

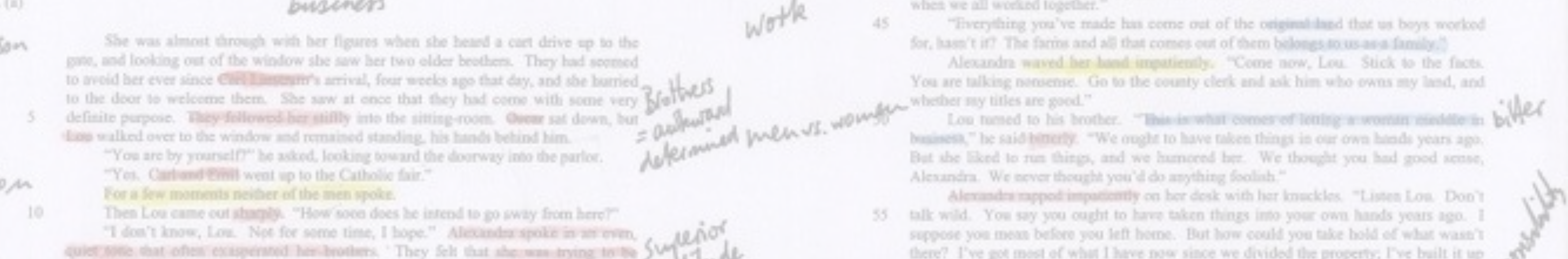    ➡ **Entity recognition**

    ➡ Entity linking

| Person | p | Loc | l | Org | o | Event | e | Date | d | Other | z |
|---|---|---|---|---|---|---|---|---|---|---|---|

Barack Hussein Obama II ✕ (born August 4, 1961 ✕ ) is an American ✕ attorney and politician who served as the 44th President of the United States ✕ from January 20, 2009 ✕ , to January 20, 2017 ✕ . A member of the Democratic Party ✕ , he was the first African American ✕ to serve as president. He was previously a United States Senator ✕ from Illinois ✕ and a member of the Illinois State Senate ✕ .

# USE CASES OF DATA ANNOTATION:

➡ Image Annotation

➡ Video Annotation

➡ Audio Annotation

➡ Text Annotation

  ➡ Text classification

  ➡ Language translation

  ➡ Entity recognition

  ➡ **Entity linking**

# ANNOTATION FORMATS

```
<NE id="i0" type ="building">
 The Massachussetts State
House</NE> in <NE id="i1"
type="city">Boston, MA</NE>
houses the offices of many
important state figures,
including <NE in="i2" type="
title">Governor</NE><NE id="
i3" type="Person">Deval
Patrick</NE>and those of the
<NE id="i4" type="
organization">Massachussetts
General Court</NE>.
```

# INLINE ANNOTATION

➡ Annotations (XML Tags) physically surround the extend that the tag refers to

## CONS:

➡ Changes the formatting of the original text

➡ Difficult to read by humans

➡ Difficult to merge with other annotating tag sets (e.g. POS tags)

➡ Difficult for multi tagging & group tagging

## PROS:

➡ Used by many programs

➡ No need for position tracking of the annotation

# STAND-OFF ANNOTATION BY TOKENS

| TOKEN | SENT_ID | TOKEN_ID |
|-------|---------|----------|
| The | 1 | 1 |
| Massachusetts | 1 | 2 |
| State | 1 | 3 |
| House | 1 | 4 |
| in | 1 | 5 |
| Boston | 1 | 6 |
| , | 1 | 7 |
| MA | 1 | 8 |
| houses | 1 | 9 |
| … | | |

| TAG | START_SENT_ID | START_TOKEN_ID | END_SENT_ID | END_TOKEN_ID |
|-----|---------------|----------------|-------------|--------------|
| NE_building | 1 | 1 | 1 | 4 |
| NE_city | 1 | 6 | 1 | 8 |

➡ Text needs to be tokenized

➡ Text is identified by assigning an ID to each token.

➡ Other IDs (paragraph section, etc.) can be assigned too.

➡ Annotation data is stored separately in a tab-separated file

➡ It is necessary to keep the associations between the IDs and the tokens.

## PROS:

➡ Different annotations on the same data can be easily merged (due to separation from the actual data)

## CONS:

➡ Doesn't allow for annotating parts of the word

➡ Relatively difficult to retrieve the original text

# STAND-OFF ANNOTATION BY CHARACTERS

```
The Massachusetts State
House in Boston, MA houses
the offices of many
important state figures,
including Govenor Deval
Patrick and those of the
Massachusetts General Court.

<NE id="N0" start="5" end ="
31" text="Massachusetts
State House" type="building"
/>
<NE id="N1" start="35" end="
45" text="Boston, MA" type="
city" />
<NE is="N2" start="118" end="
131" text="Deval Patrick"
type="person" />
```

➡ Start and end offsets declare the position of each annotation in the text

➡ Character encoding is crucial and must be maintained throughout the annotation process

➡ Technically only the offsets and the tag attributes suffice to retrieve the annotation but the actual annotated text is also kept for redundancy

➡ Original text can be very easily retrieved

# LINKED EXTENT ANNOTATIONS

The Massachusetts State House in Boston, MA houses the offices of many important state figures, including Govenor Deval Patrick and those of the Massachusetts General Court.

```
<NE id="N0" start="5" end ="31" text="Massachusetts State House" type="building" />
<NE id="N1" start="35" end="45" text="Boston, MA" type="city" />
<NE is="N2" start="118" end="131" text="Deval Patrick" type="person" />

<L-LINK id="L0" fromID="N2" toID="N0" relationship="worksIN"/>
```
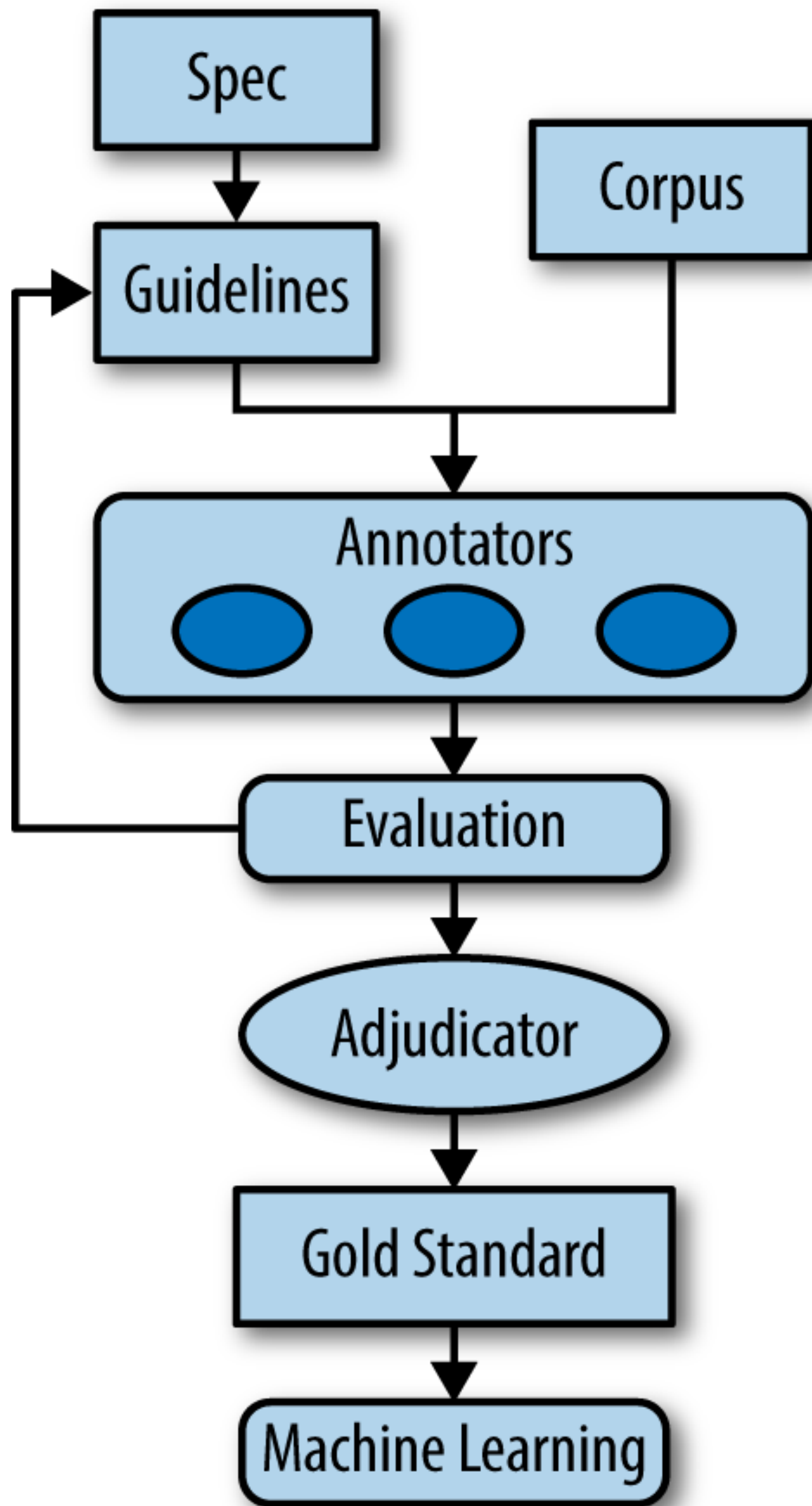
➡ Use the ID of the tags as anchors to represent the relationships between them

➡ Represent directionality by using fromID / toID attributes of the annotation

➡ Can work with both token-based and character-based stand-off annotations

➡ The NEs have to be annotated first in order to create the anchor IDs
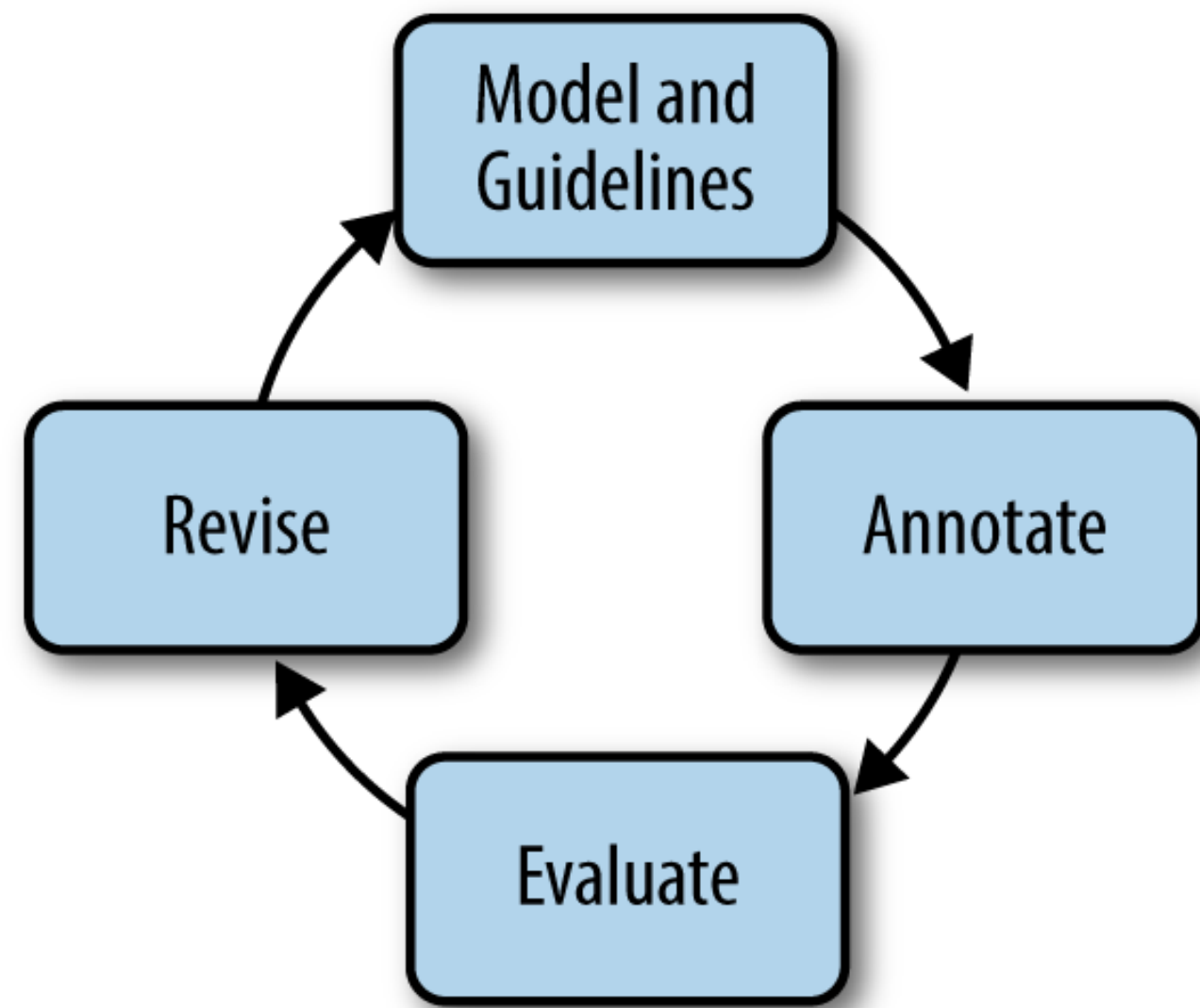
ANNOTATION WORKFLOW

Source: Natural Language Annotation for Machine Learning

# GUIDELINES & SPECIFICATIONS

➡ Guidelines show how Specifications (schema) is applied to the data

➡ Provide instructions to annotators with examples and use cases

➡ Are designed specifically for the particular specification and dataset

➡ Are designed specifically for the particular ML task

# THE M.A.M.A. CYCLE



Source: Natural Language Annotation for Machine Learning

➡ Supervisor creates guidelines based on the model and the annotation task

➡ Annotators use the guidelines and create annotations for the same batch

➡ After each batch, annotators gather and discuss their differences

➡ Revisions are based on the (dis)agreement of annotators

➡ Each revision leads to refinements for the guidelines and /or specification

➡ Continuous revisions should lead to higher inter annotator agreement (IAA)

➡ Once IAA reaches a sufficient score, each annotator uses different batch

# EVALUATING ANNOTATORS

➡ conducted before creation of gold standard
➡ based on the measurement of inter-annotator agreement (IAA)
➡ Good IAA doesn't necessary mean that the dataset will produce good results when used in ML
➡ Good IAA indicates that the annotation task can be easily reproduced by many people and lead to bigger dataset
➡ IAA must take into account random chance agreements
➡ Cohen's Kappa measures the IAA among a pair of annotators
➡ Fleiss's Kappa measures the IAA among more than two annotators

# INTERPRETING IAA SCORES

| κ | Agreement level |
|---|---|
| < 0 | poor |
| 0.01–0.20 | slight |
| 0.21–0.40 | fair |
| 0.41–0.60 | moderate |
| 0.61–0.80 | substantial |
| 0.81–1.00 | perfect |

➡ Depends on the complexity and objectivity of the task
➡ Should be taken in context with other scores in relevant tasks
➡ Annotation Charts can provide fruitful information regarding annotators' behaviour
➡ Poor initial results are normal especially in difficult tasks
➡ Sparse entities should be taken into account
➡ Use small batches and conduct as many as needed in order to increase the IAA

# ANNOTATION TOOLS

# CHOOSING AN ANNOTATION ENVIRONMENT

➡ Supported types of annotations

➡ Architecture

➡ Supported formats

➡ Support for multi-session / groups

➡ Support for workflows / automations

➡ Metrics

# USEFUL RESOURCES:

**Textbooks:**

- J. Pustejovsky. Natural Language Annotation for Machine Learning. O' Reilly 2013
- Alex M. PattersonThe Art of Data Annotation: Transforming Raw Data into Machine Learning Gold. Kindle Editions 2023
- Anthony Sarkis. Training Data for Machine Learning. O' Reilly 2023

**Web Resources:**

- Annotation Studio: a suite of collaborative web-based annotation tools currently under development at MIT
- Diigo for highlighting and bookmarking web pages
- Hypothes.is: web browser extension of annotating online documents (web pages, pdfs and docs)
- Perusall: social reading tool
- BRAT: web based annotation tool for texts
- Prodigy: Annotation tool for Machine Learning with support for multiple types of annotations

# Thank you!

vpertsas@aueb.gr