

PAT 498/598 (Winter 2025)

Music & AI

Lecture 11: Music Classification

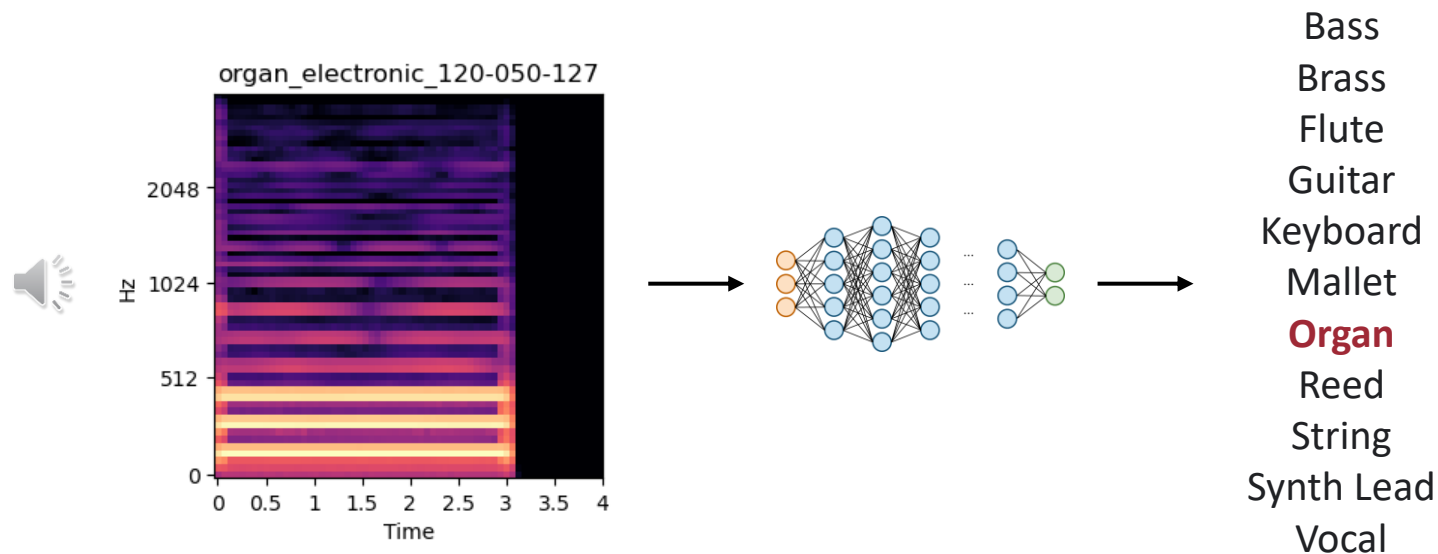
Instructor: Hao-Wen Dong



SCHOOL OF MUSIC, THEATRE & DANCE
PERFORMING ARTS TECHNOLOGY
UNIVERSITY OF MICHIGAN

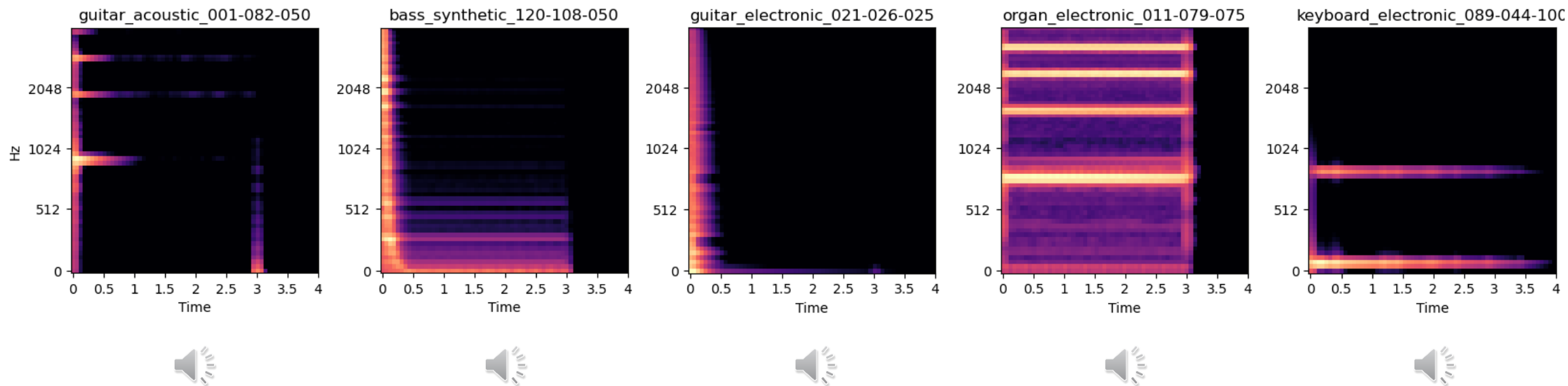
Homework 3: Musical Note Classification using CNNs

- Train a CNN that can classify audio files into their **instrument families**
 - **Input:** 64x64 mel spectrogram
 - **Output:** 11 instrument classes
 - Using the **NSynth** dataset (Engel et al., 2017)



NSynth Dataset

- A collection of 305,979 **single-shot musical notes** (Engel et al., 2017)
 - Produced from 1,006 **commercial sample libraries**
 - With different **MIDI pitches** (21–108) and **velocities** (25, 50, 75, 100, 127)



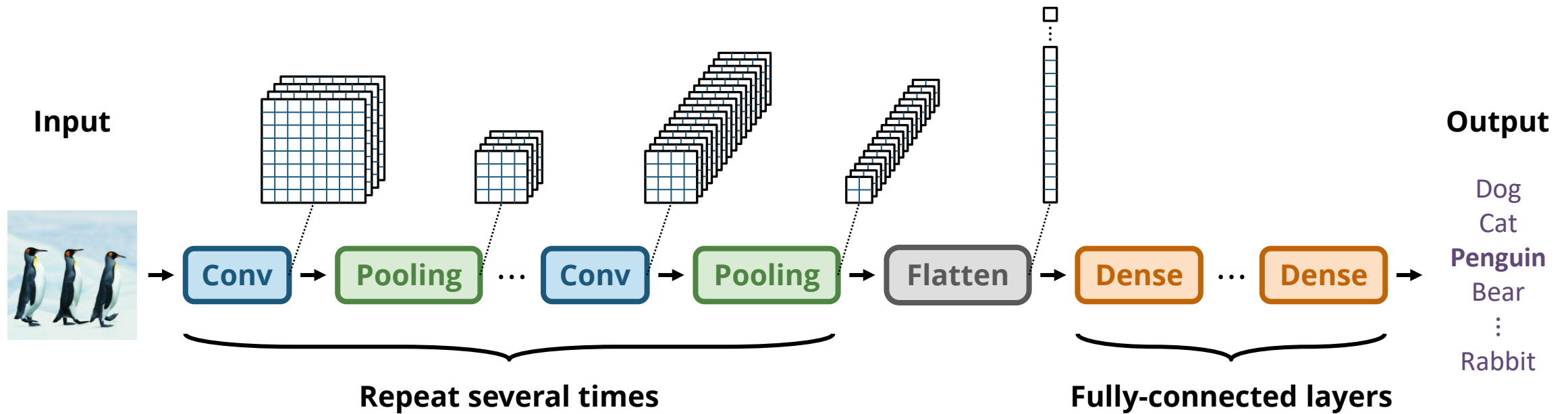
Homework 3: Musical Note Classification using CNNs

- Instructions will be released on Gradescope
- Due at **11:59pm ET** on **February 17**
- Late submissions: **1 point deducted per day**

(Recap) Reusable Pattern Detectors

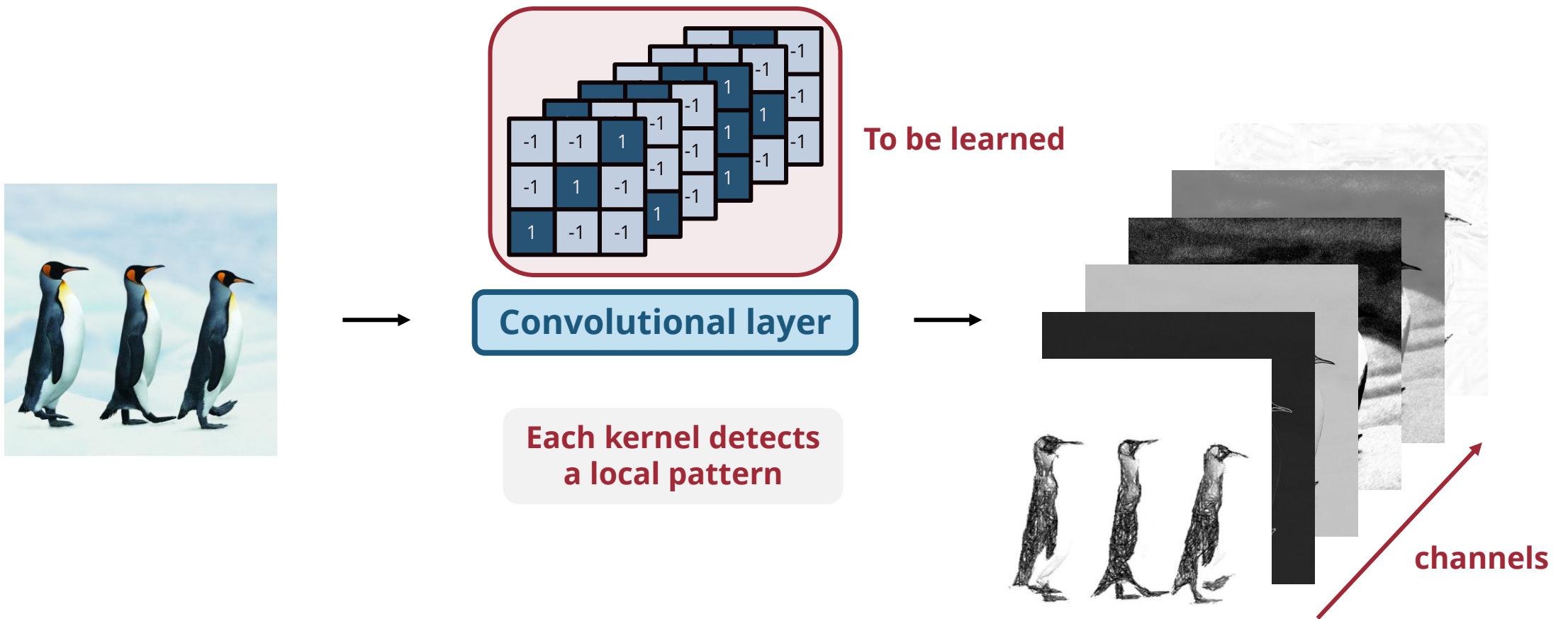


(Recap) Convolutional Neural Network (CNNs)



(Recap) Convolutional Layer

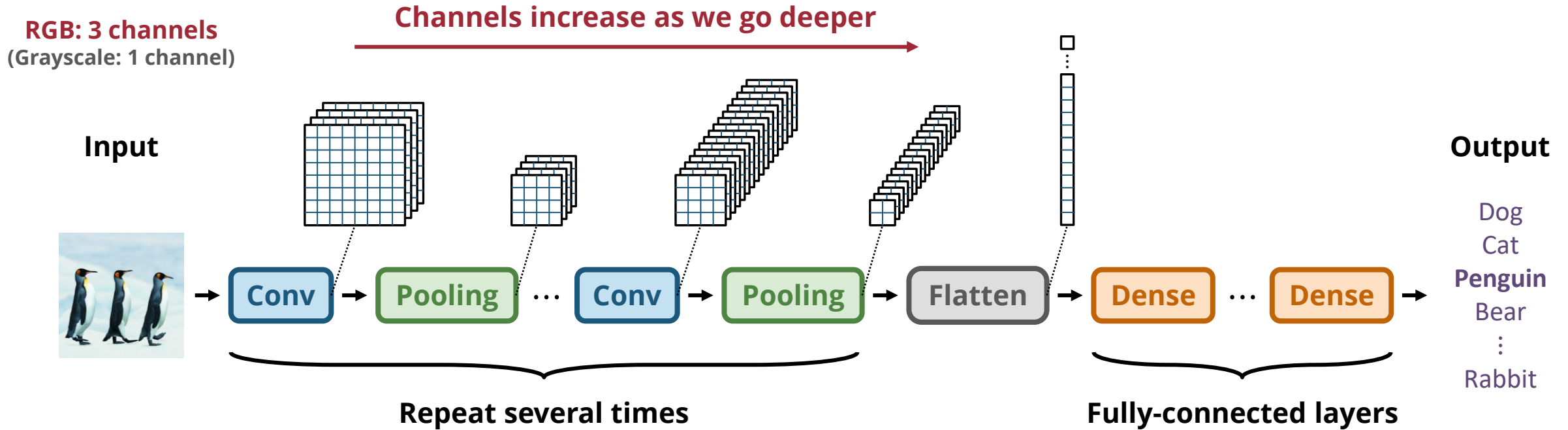
- A convolutional layer consists of many **learnable kernels** (channels)



(Recap) Benefits of CNNs

- Learn **local patterns**
- **Invariant to shifts**
 - Also called *translational invariance*
- **Reuse the learned filters** across
 - Different parts of the image
 - Across different images
- **Higher parameter-efficiency** against fully-connected neural network

(Recap) Convolutional Neural Network (CNNs)



(Recap) A Real Example

```
class CNN(nn.Module):  
    """A basic convolutional neural network."""  
    def __init__(self):  
        super().__init__()  
        self.conv1 = nn.Conv2d(1, 16, 3, padding="same")  
        self.conv2 = nn.Conv2d(16, 32, 3, padding="same")  
        self.conv3 = nn.Conv2d(32, 64, 3, padding="same")  
        self.conv4 = nn.Conv2d(64, 128, 3, padding="same")  
        self.pool = nn.MaxPool2d(2, 2)  
        self.fc = nn.Linear(128 * 4 * 4, n_classes)
```

```
def forward(self, x):  
    x = self.pool(F.relu(self.conv1(x)))  
    x = self.pool(F.relu(self.conv2(x)))  
    x = self.pool(F.relu(self.conv3(x)))  
    x = self.pool(F.relu(self.conv4(x)))  
    x = torch.flatten(x, 1)  
    x = self.fc(x)  
    return x
```

Input channels

Output channels

Kernel size

Input: 1 x 64 x 64

16 x 32 x 32 = 16384

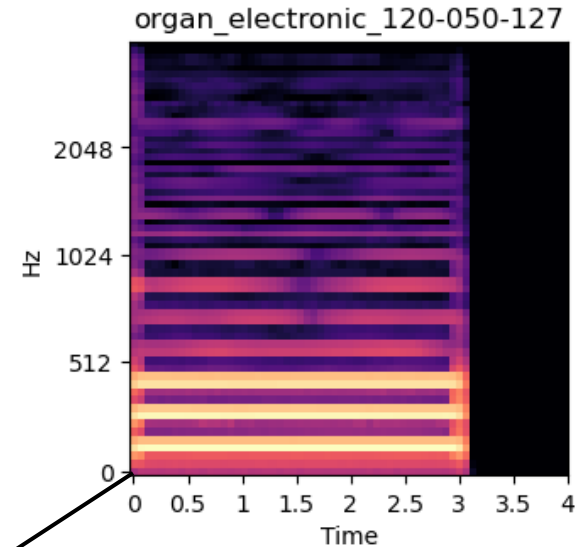
32 x 16 x 16 = 8192

64 x 8 x 8 = 4096

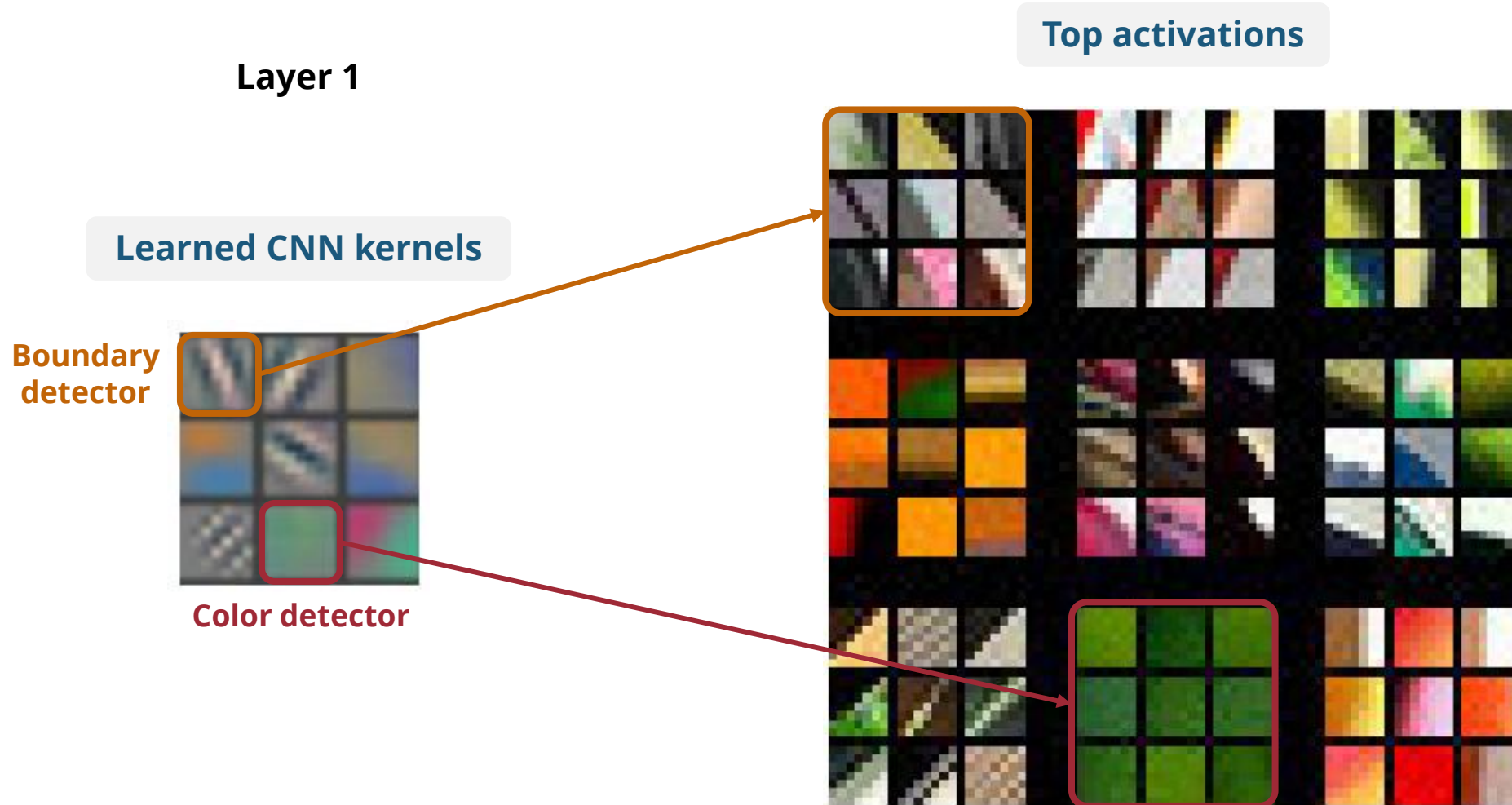
128 x 4 x 4 = 2048

**More channels,
lower resolution
as we go deeper**

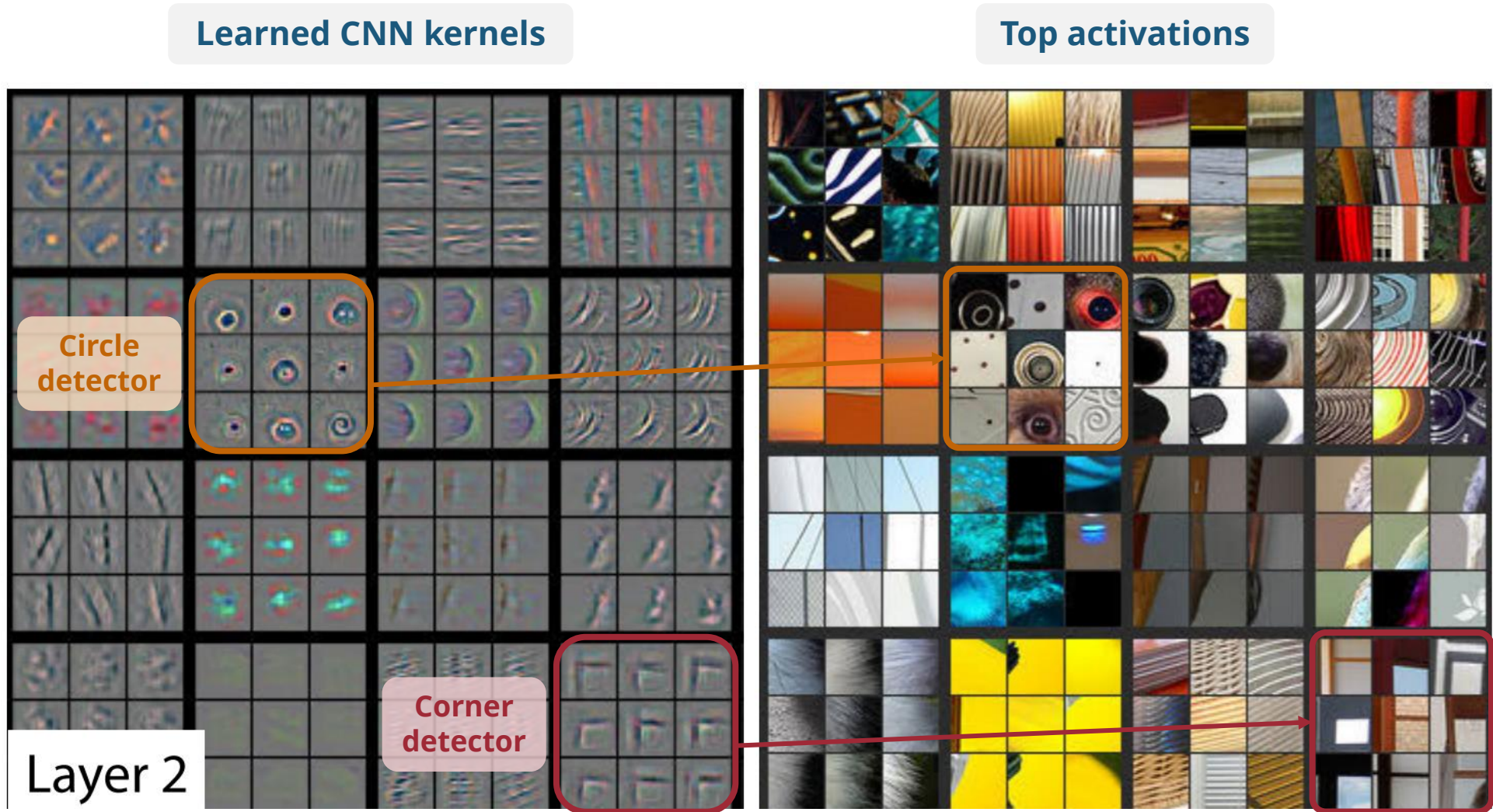
**Total number of
features decrease
as we go deeper**



(Recap) Learned CNN Kernels in a Trained AlexNet



(Recap) Learned CNN Kernels in a Trained AlexNet

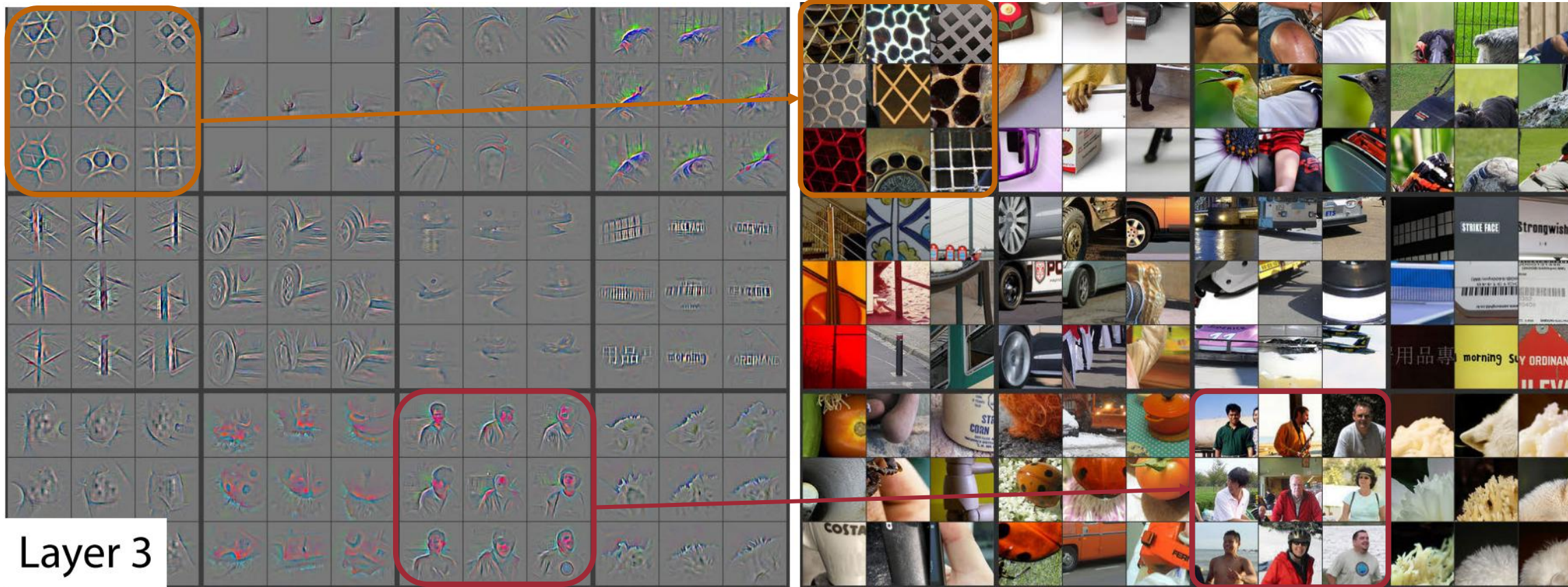


(Recap) Learned CNN Kernels in a Trained AlexNet

Learned CNN kernels

Top activations

Grid detector

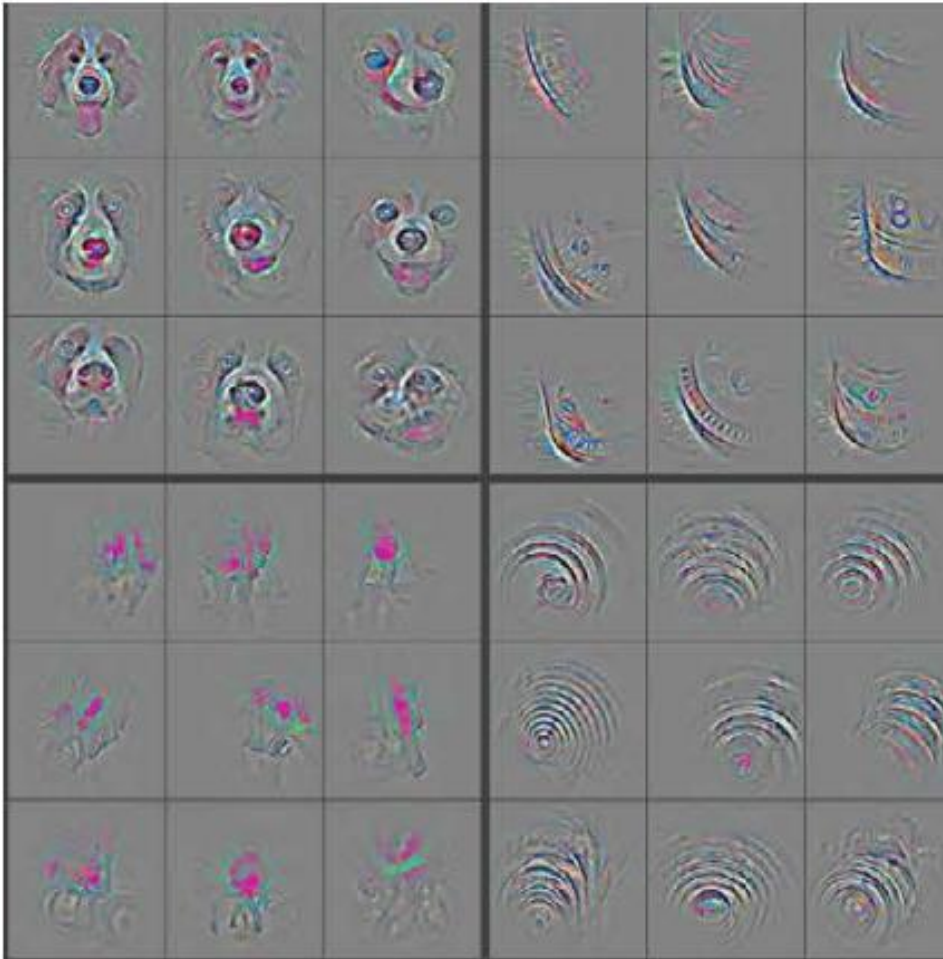


Layer 3

Human detector

(Recap) Learned CNN Kernels in a Trained AlexNet

Learned CNN kernels

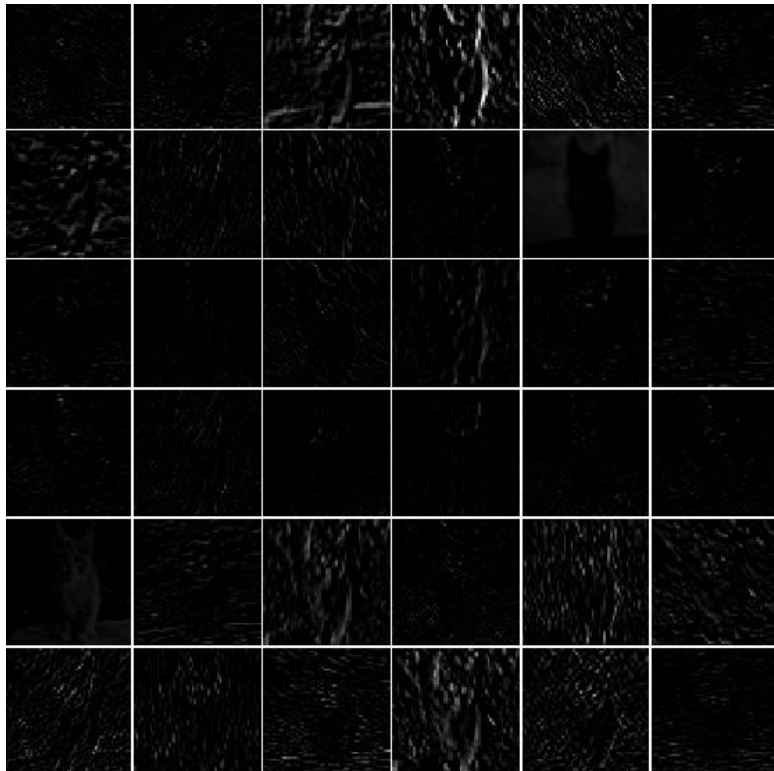


Top activations

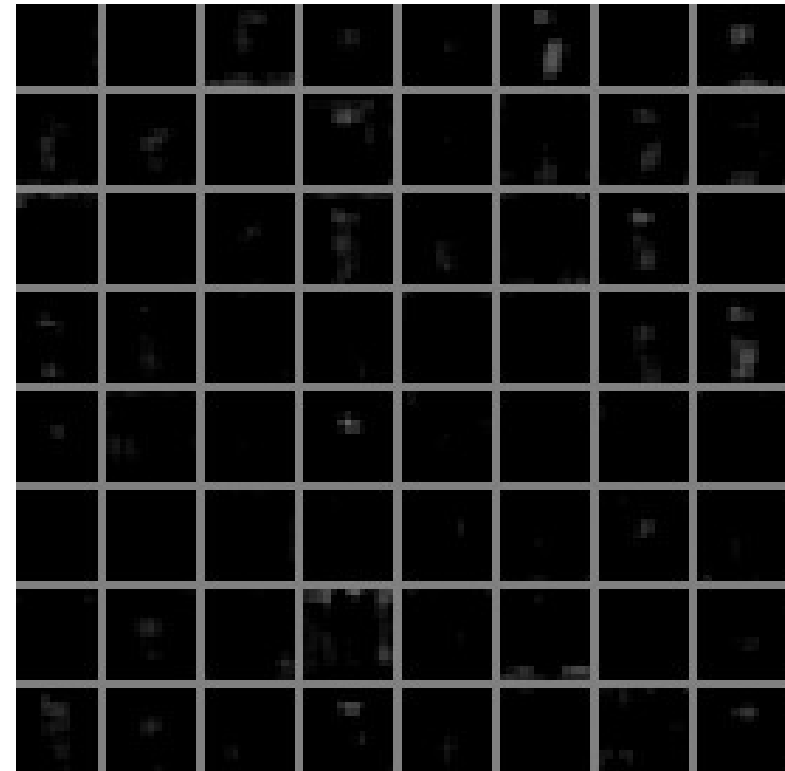


(Recap) Activations in a Trained AlexNet

1st convolutional layer



5th convolutional layer

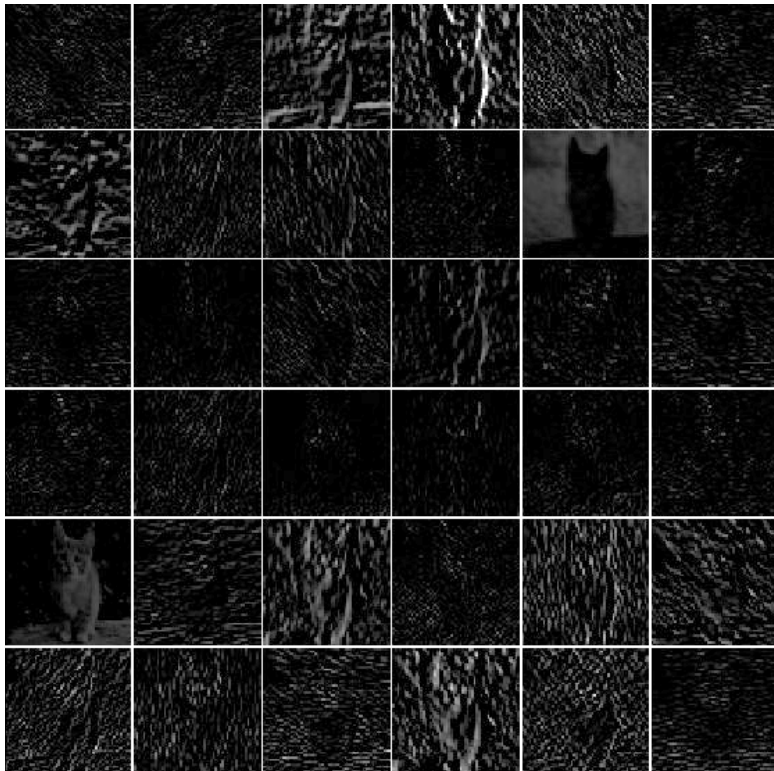


(Recap) Activations in a Trained AlexNet

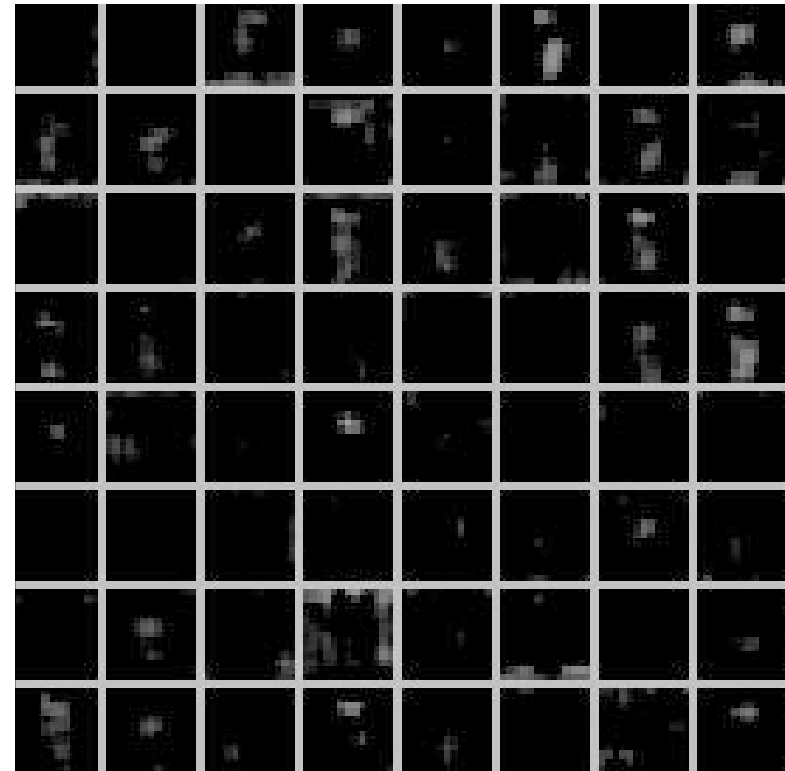
1st convolutional layer

5th convolutional layer

Brightness+
Contrast+

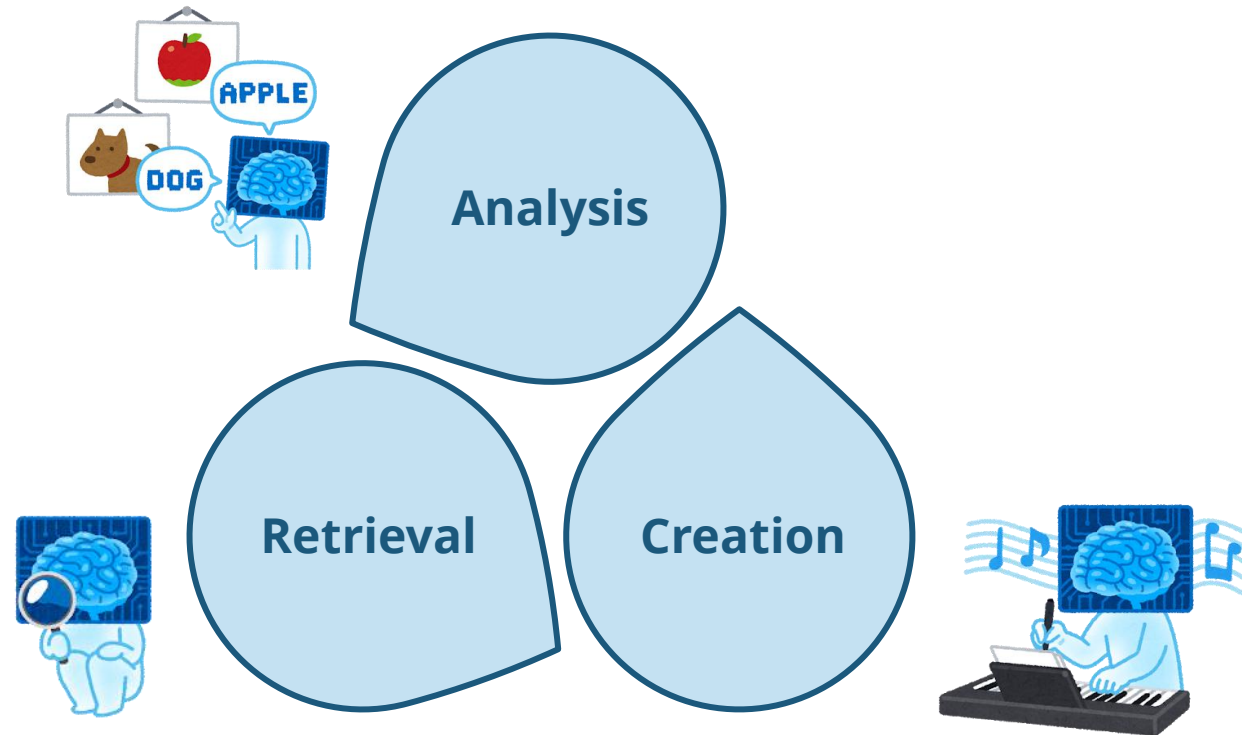


Brightness+
Contrast+



(Recap) Music Information Research (MIR)

- “Intelligent ways to analyze, retrieve and create music” (Yang 2018)



Music Classification

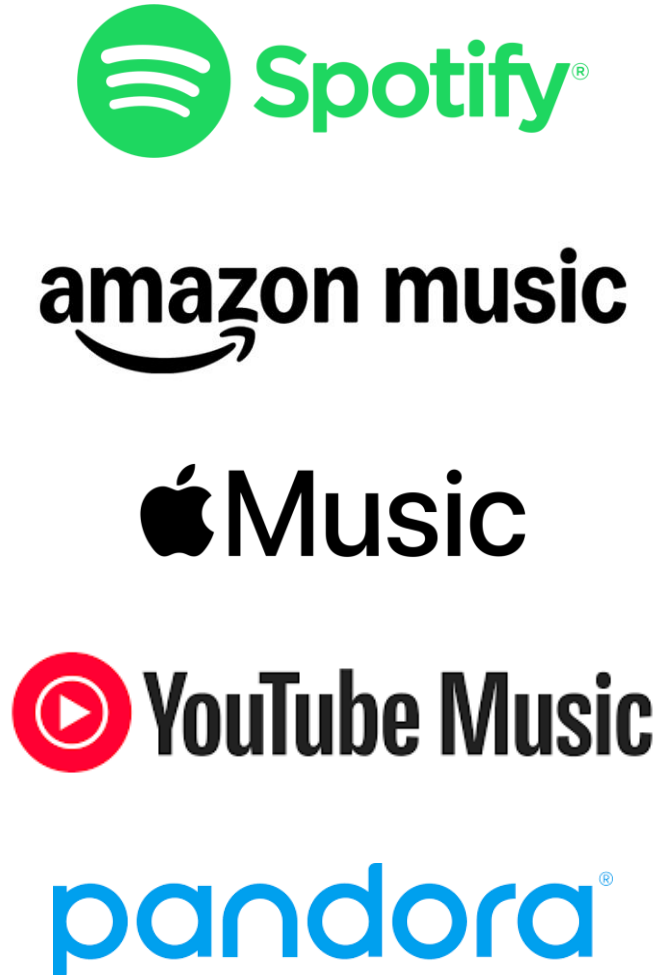
Music Classification Tasks

- **Genre** classification (pop, rock, r&b, jazz, hip-hop, classical, etc.)
- **Mood** classification (happy, sad, calm, aggressive, cheerful, etc.)
- **Instrument** recognition
- **Composer** identification
- **Key** detection
- **Chord** estimation
- **Music tagging** → **Can cover everything above!**

Applications of Music Classification Models

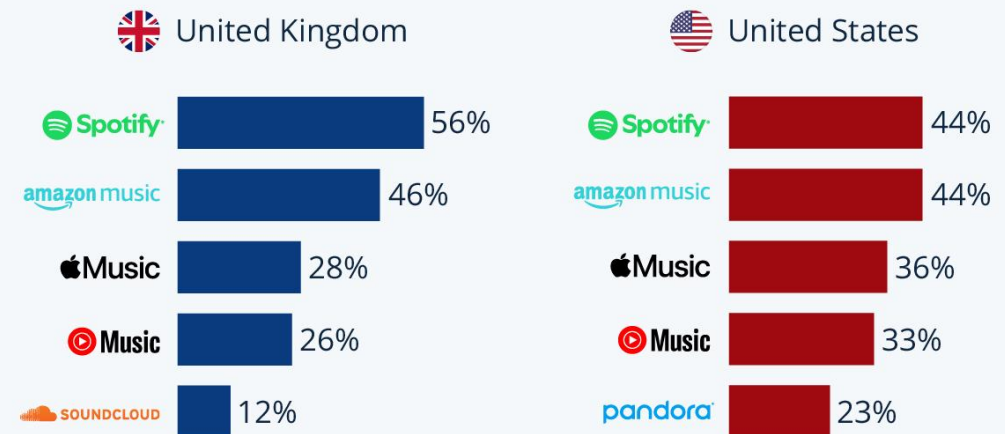
- Recommendation
- Curation
- Playlist generation
- Listening behavior analysis
- Musicology research

Music Classification for Recommendation



The Most Loved Digital Audio Streaming Platforms

Share of respondents who have paid for audio downloads or streaming services from the following platforms*



* in the 12 months prior to the survey
2,362 (UK)/4,944 (USA) respondents (18-64 y/o) surveyed Jul. 2023-Jun. 2024
Source: Statista Consumer Insights



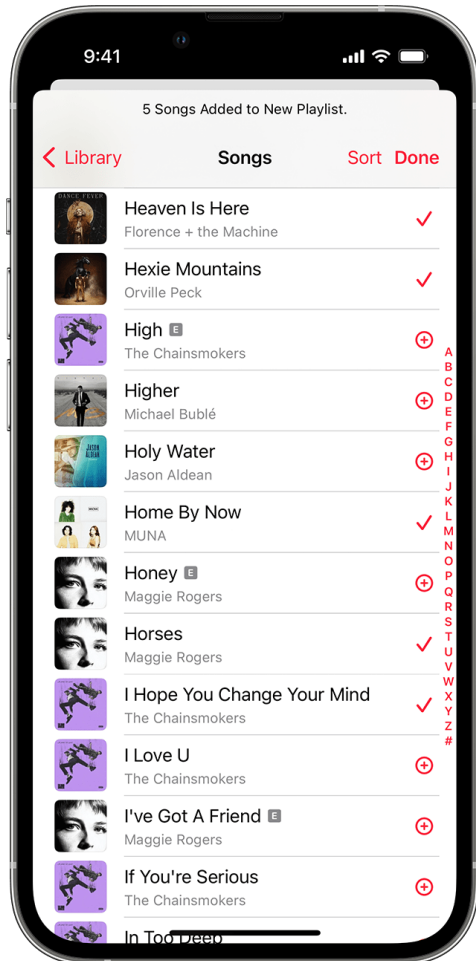
statista

Music Classification for Recommendation

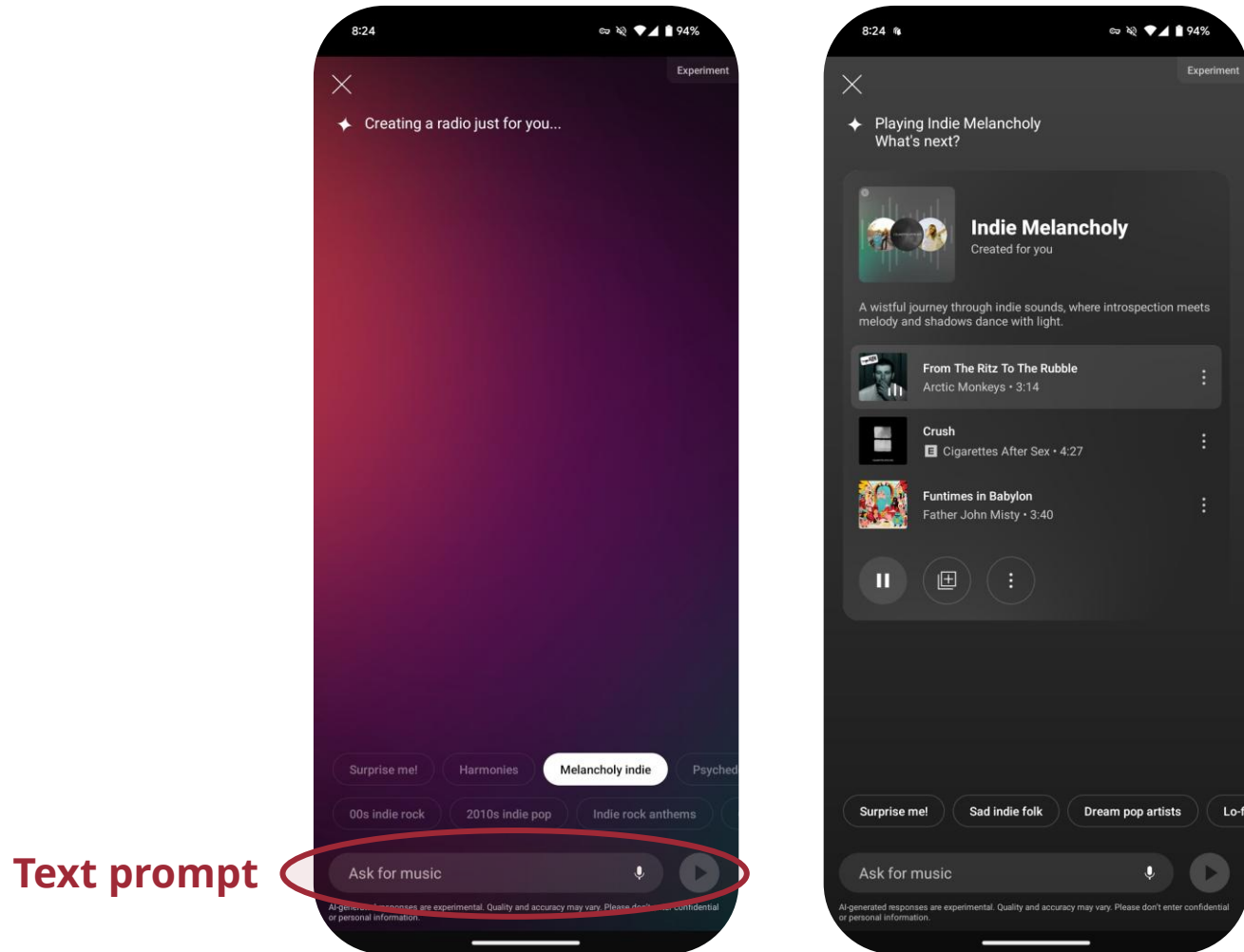
What to play next?



Music Classification for Playlist Generation

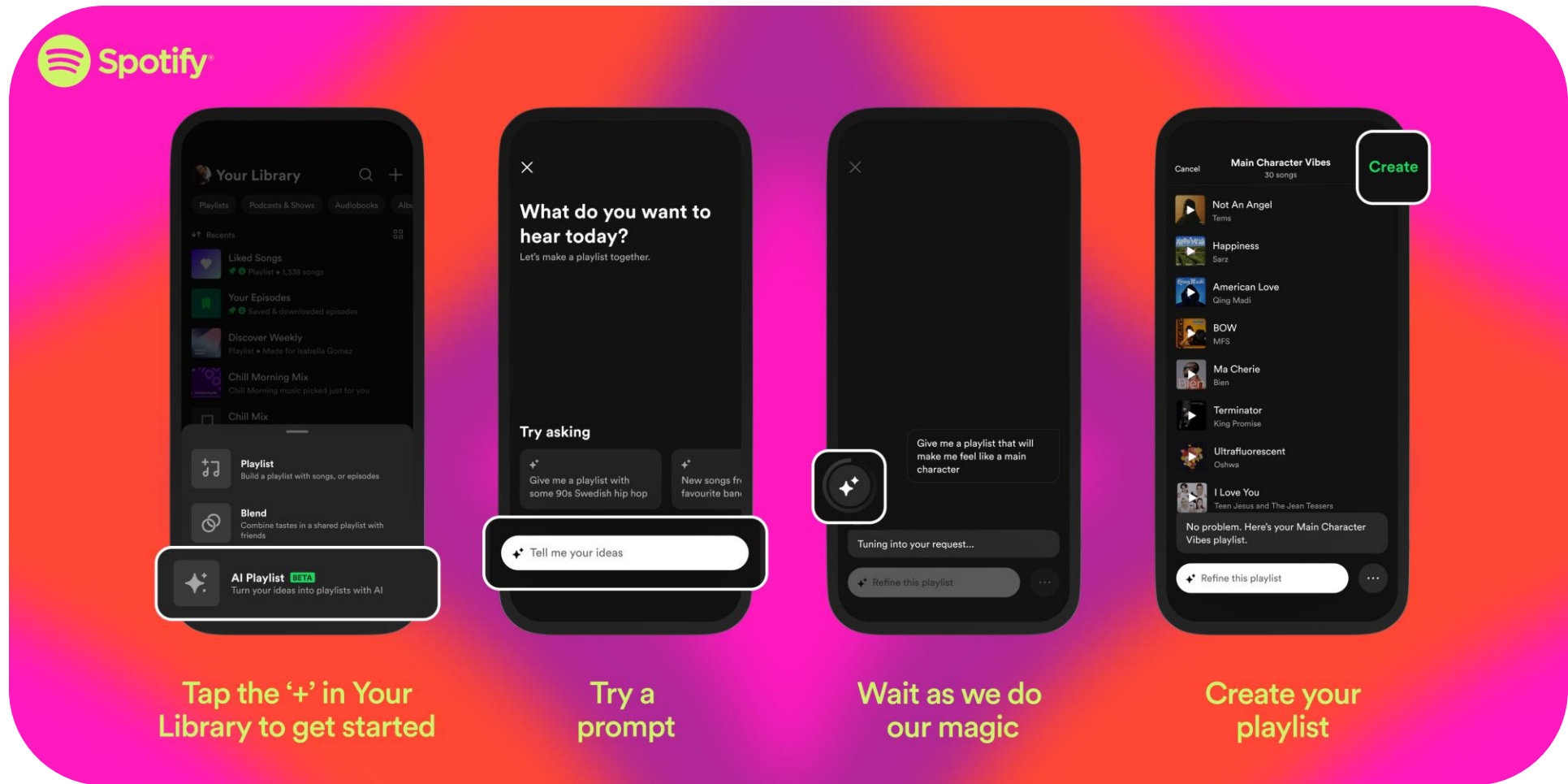


Ask Music (YouTube Music)



(Source: Android Police)

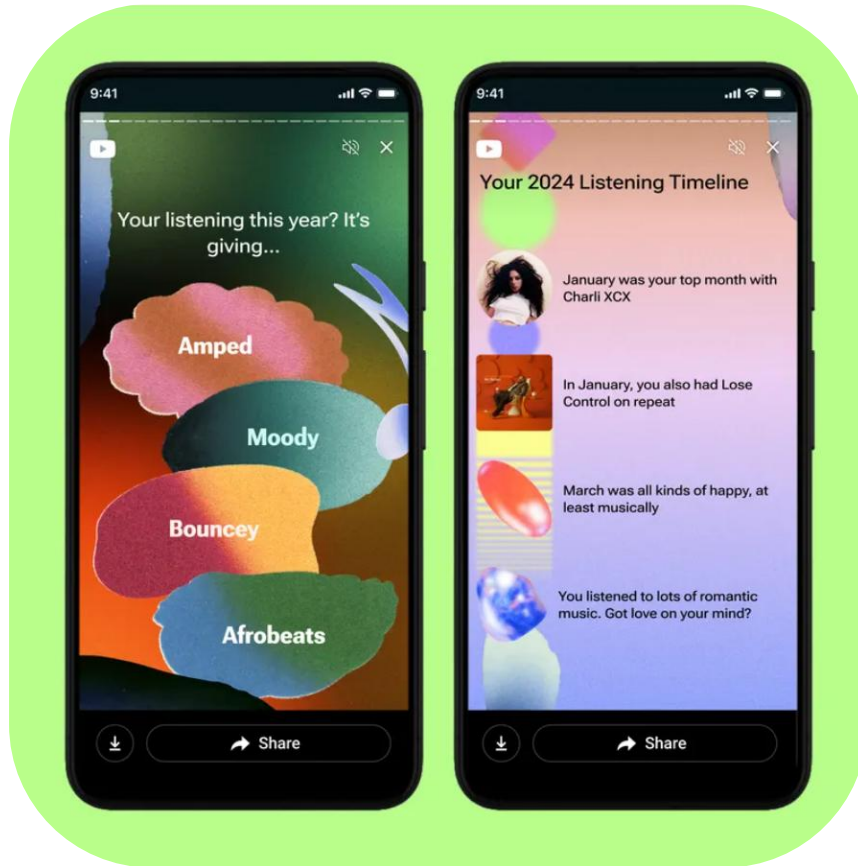
AI Playlist (Spotify)



(Source: Spotify)

Music Classification for Listening Behavior Analysis

YouTube's Music Recap



(Source: YouTube)

Spotify's Listening Personality



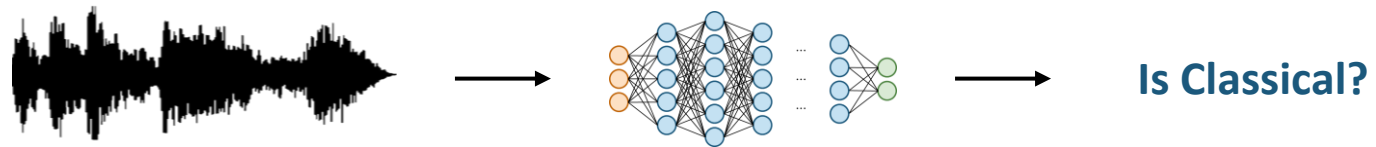
(Source: Spotify)

Types of Classification Tasks

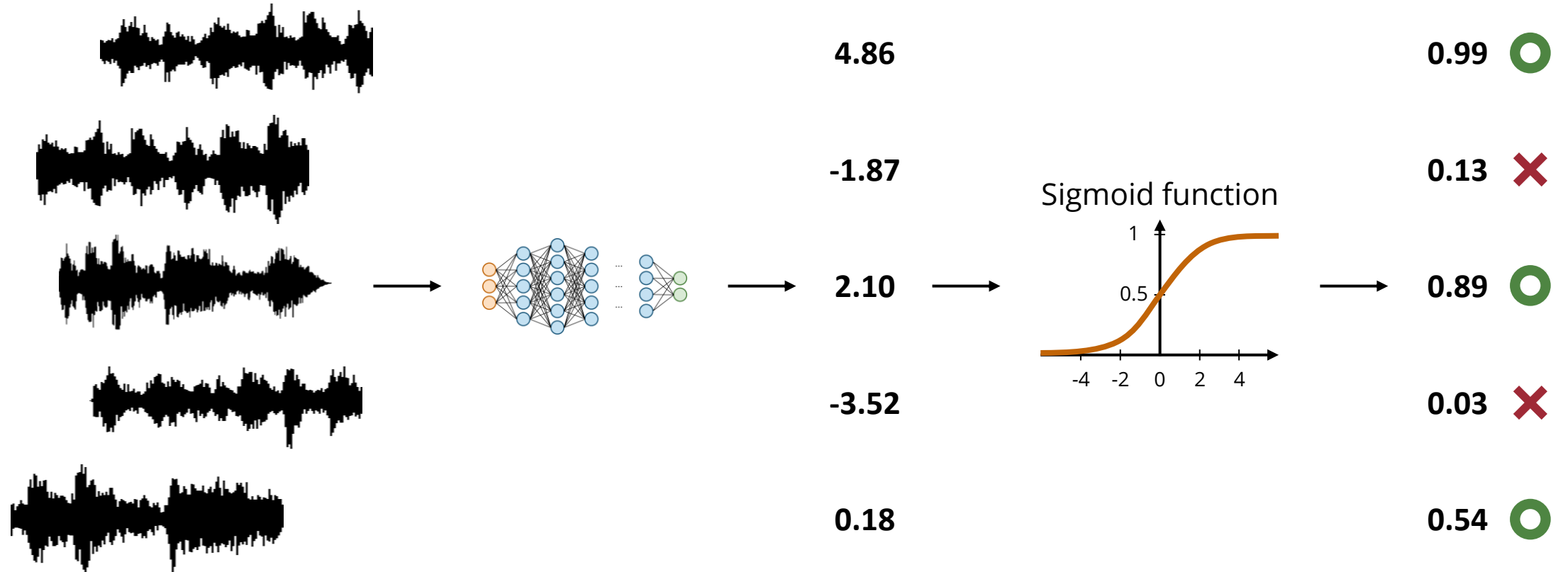
| Types of Classification Tasks

- **Binary** classification
- **Multiclass** classification
- **Multi-label** classification

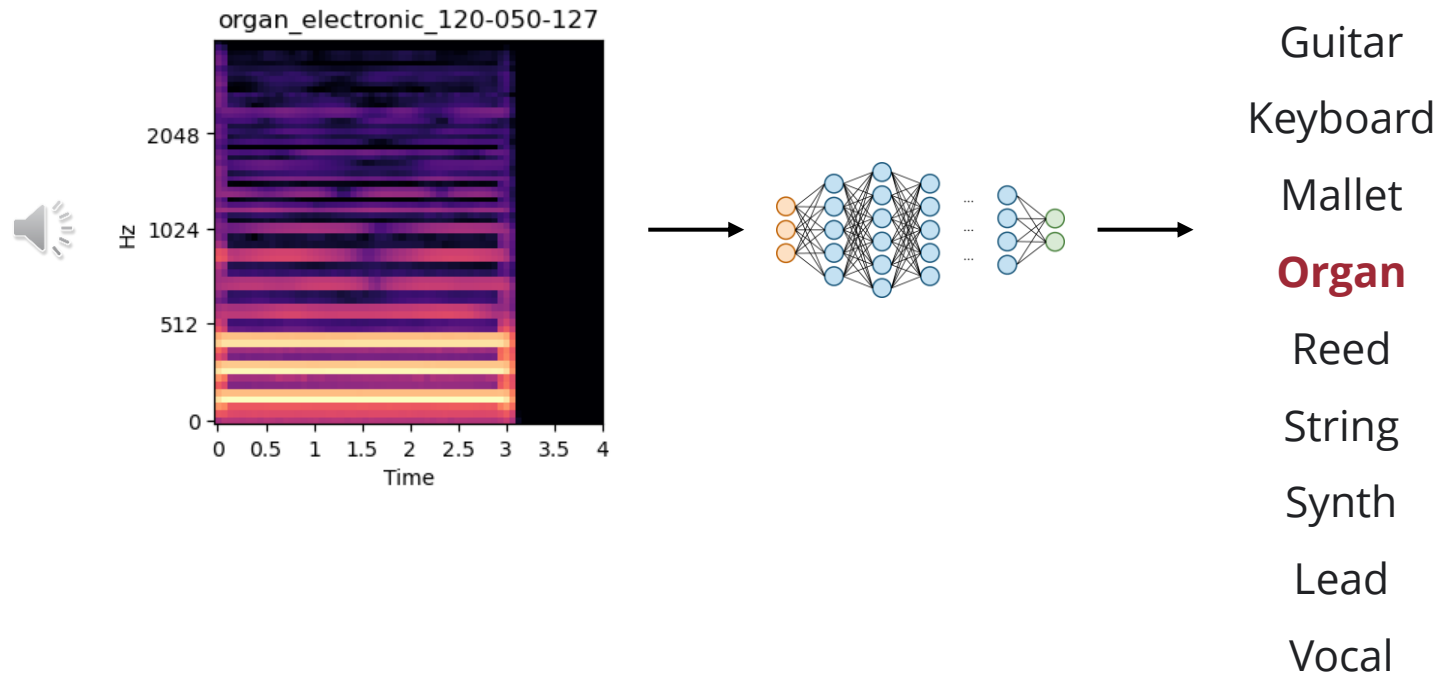
Binary Classification



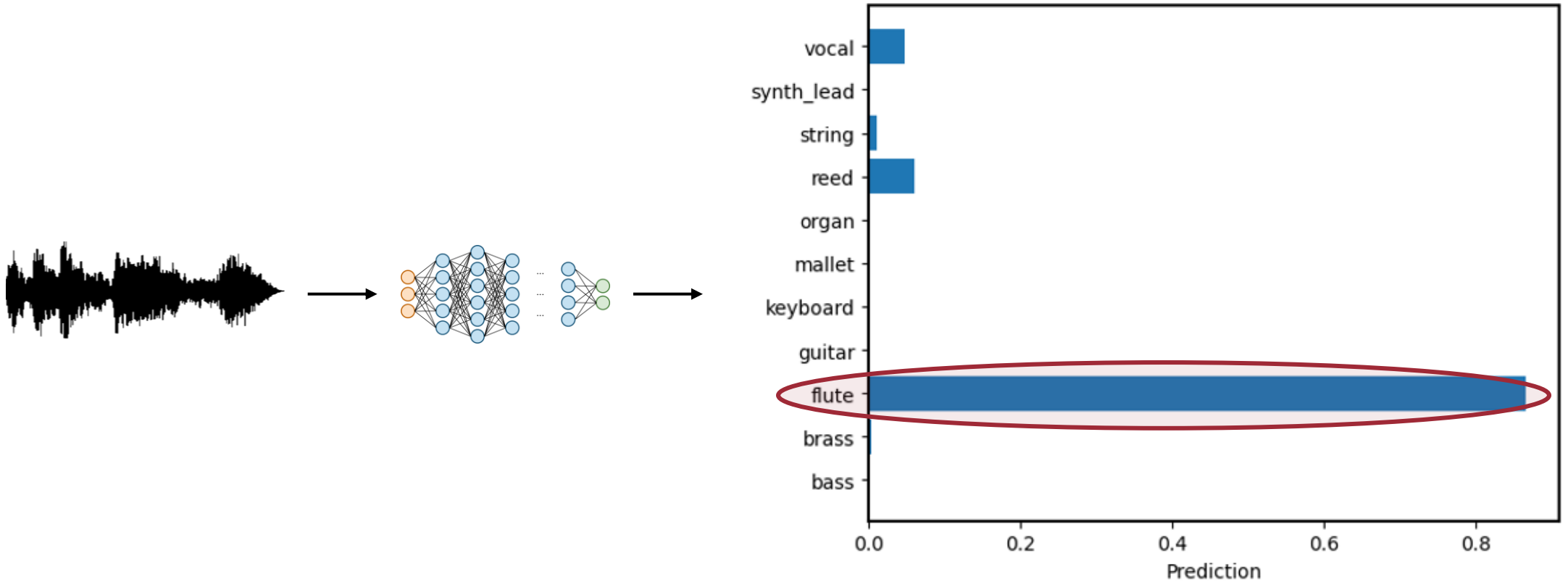
Binary Classification



Multiclass Classification



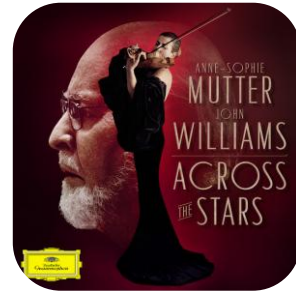
Multiclass Classification



Multi-label Classification



Soul, R&B, pop



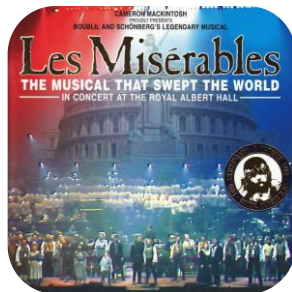
Soundtrack, classical



Rock



Classical



Musical



Alternative, rock

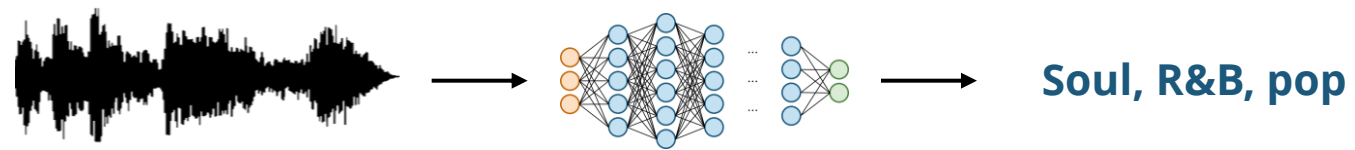


Pop

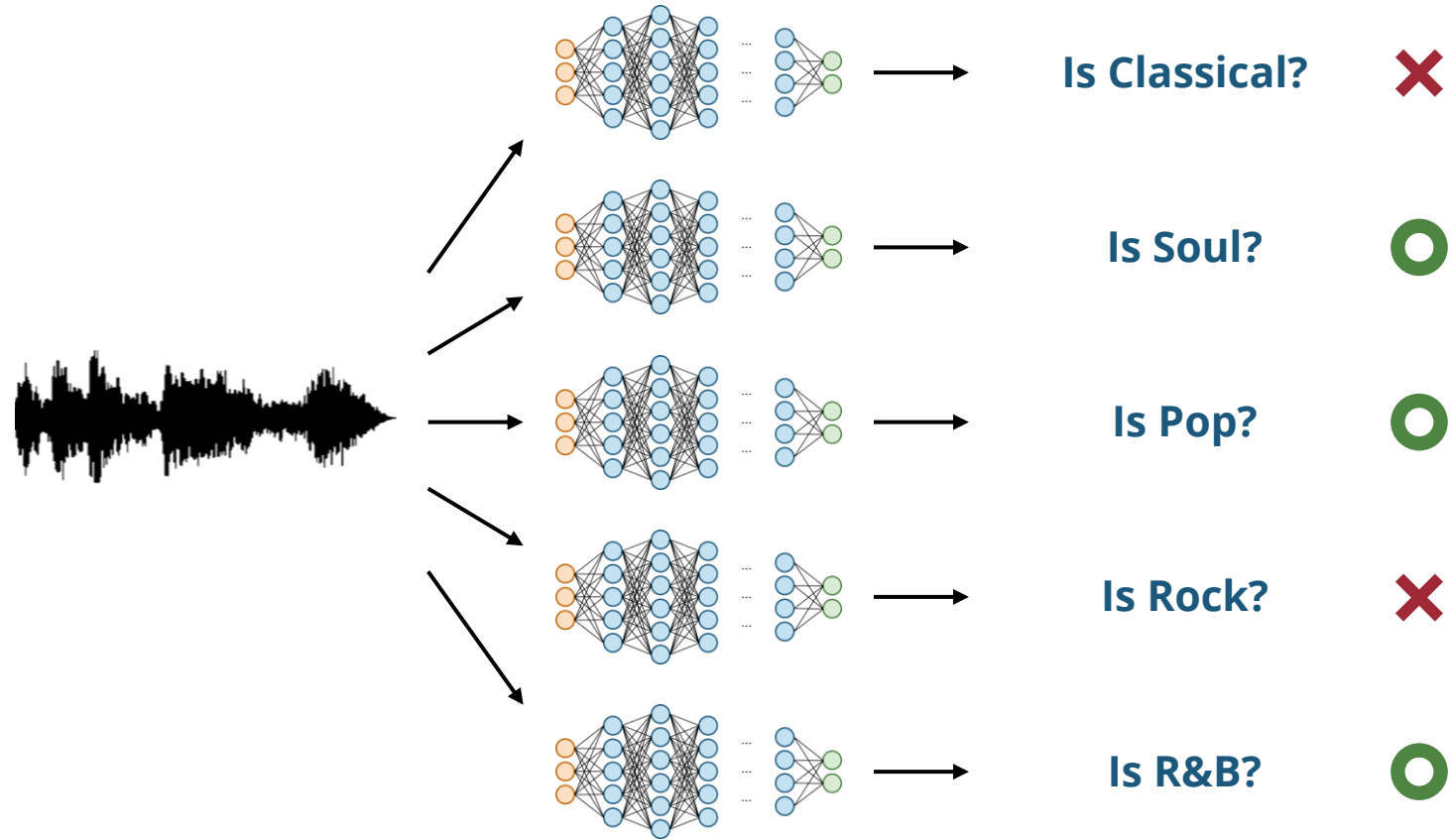


Pop, soul, R&B

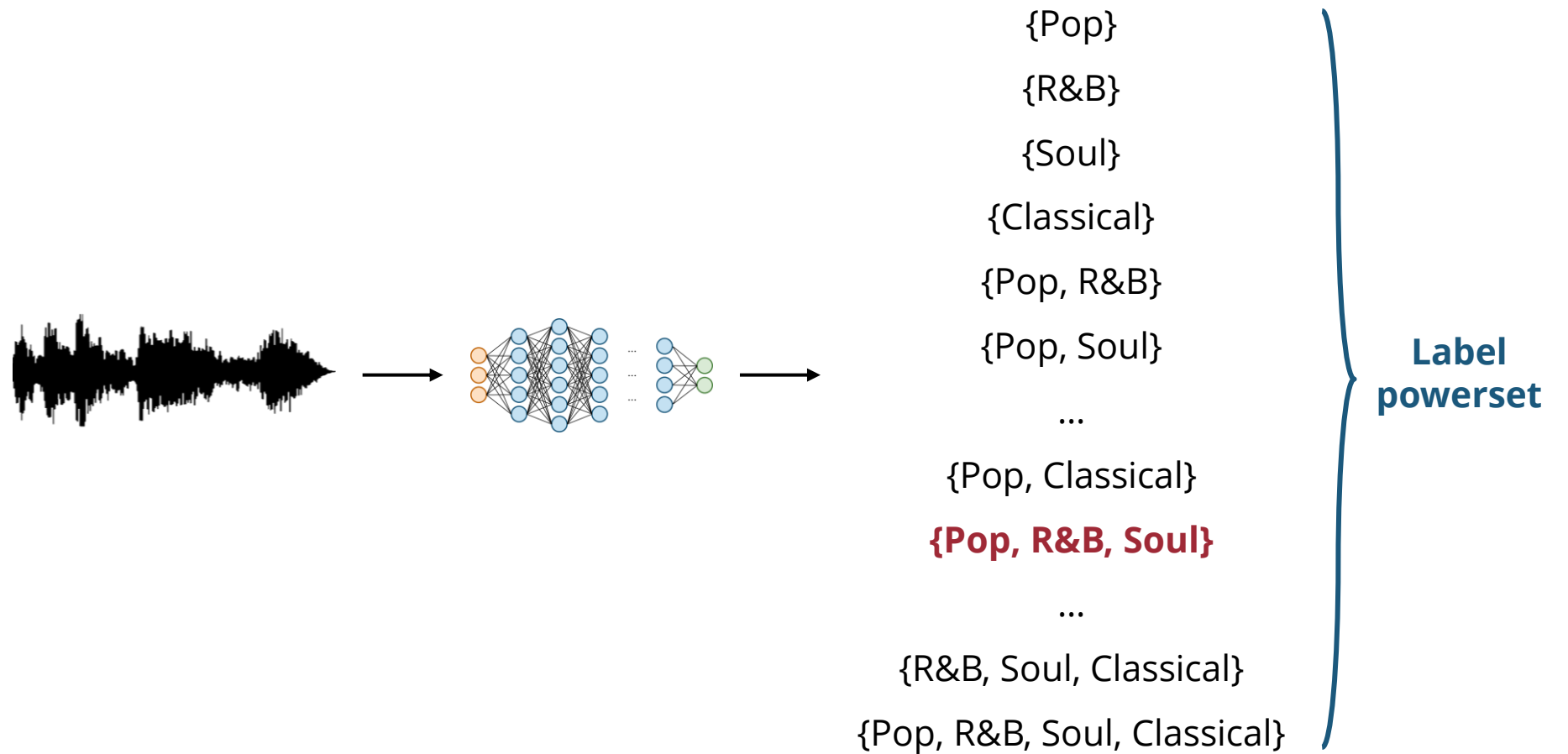
Multi-label Classification



Multi-label Classification as Binary Classification

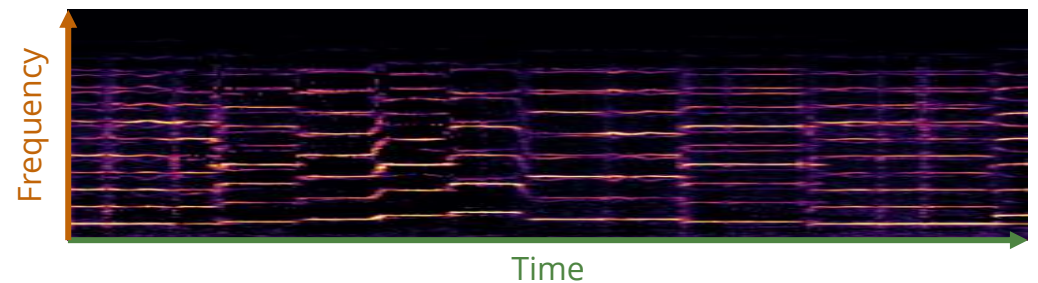


Multi-label Classification as Multi-class Classification



Input Features

- **Waveform**
- Time-frequency representation (**spectrograms**)
- **Hand-crafted features** or **features provided in metadata**
 - Acoustic: loudness, pitch, timbre
 - Rhythmic: beat, tempo, time signature
 - Tonal: key, scale, chords
 - Instrumentation, expressions, structures, etc.



Common Datasets

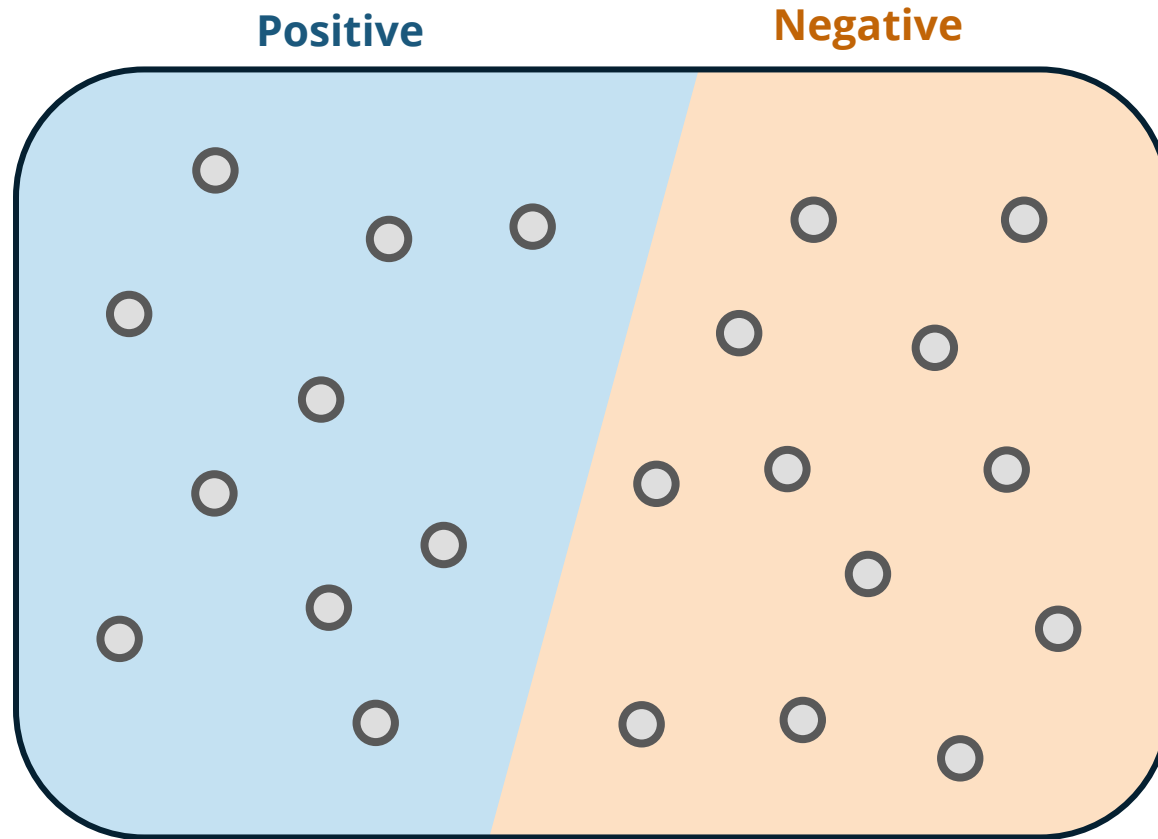
- **GTZAN**: 1,000 30-sec songs, 10 genres
- **MagnaTagATune**: 5,405 29-sec songs, 188 tags, 230 artists
- **Million Song Dataset (MSD)**: **1M** 30-sec songs, >500K tags, **tricky to access**
- **Free Music Archive (FMA)**: >10K **full** songs, 163 genres
- **MTG-Jamendo**: 55K **full** songs, 195 tags
- **AudioSet**: **1M** songs, YouTube URLs, **low-quality audio**
- **NSynth**: ~306K 4-sec instrument sounds

Evaluation Metrics

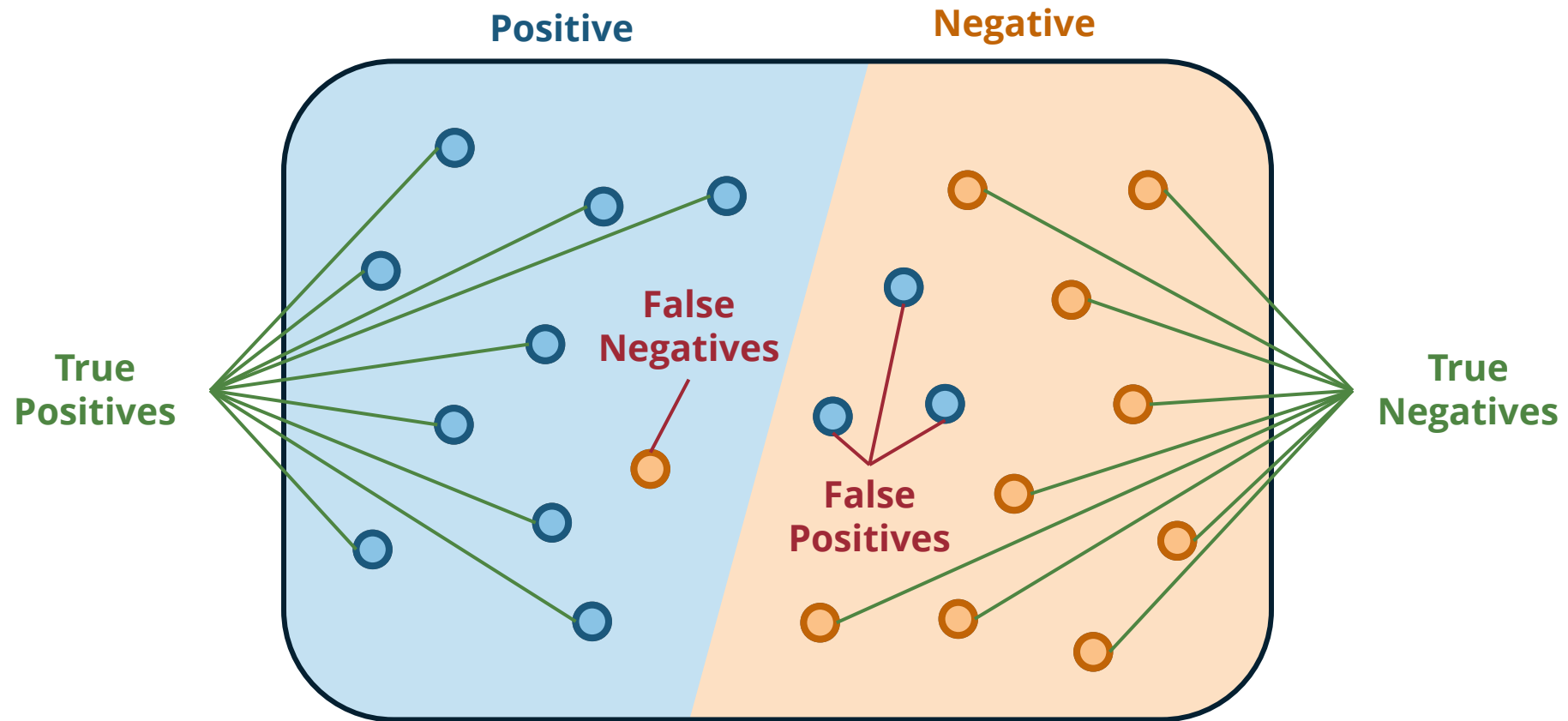
Evaluation Metrics

- **Key:** Capture what you care the most!
- The best evaluation metric depends on the actual use case
- Best to use several evaluation metrics to obtain a holistic view of your model's performance

Toy Example: Binary Classification

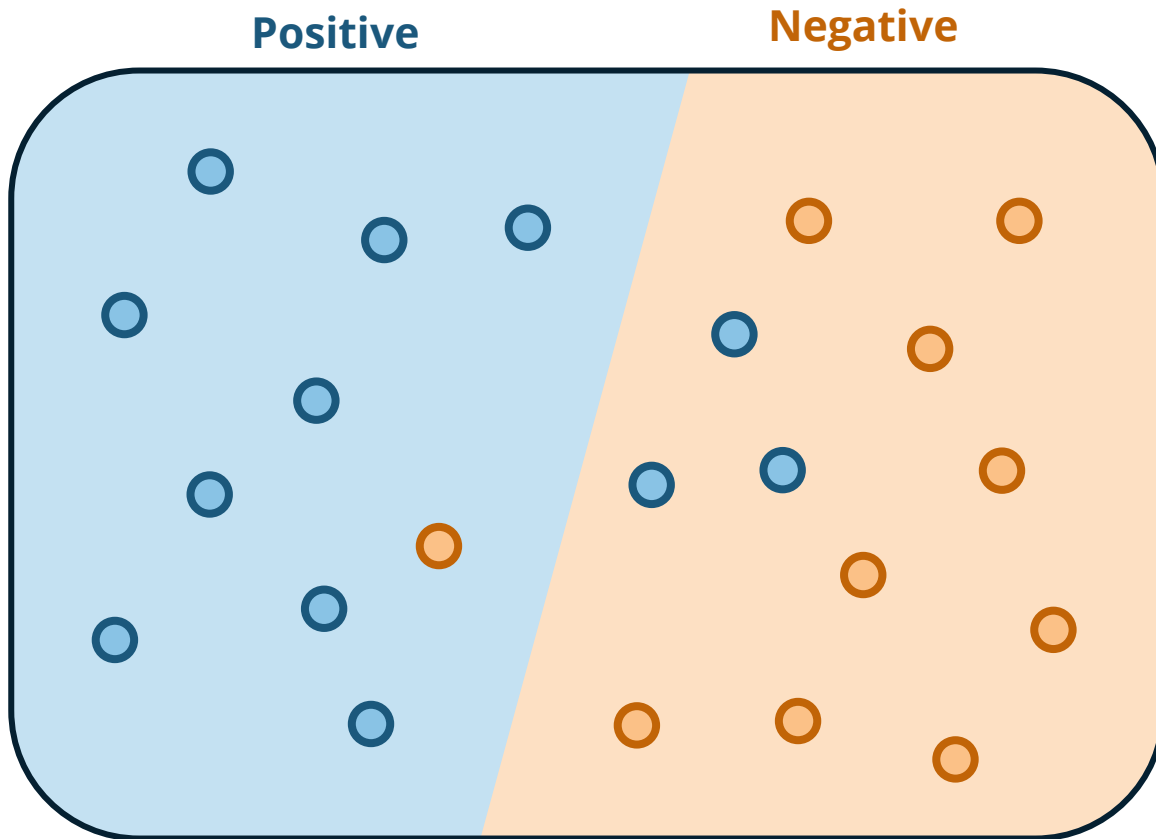


Toy Example: Binary Classification



Accuracy

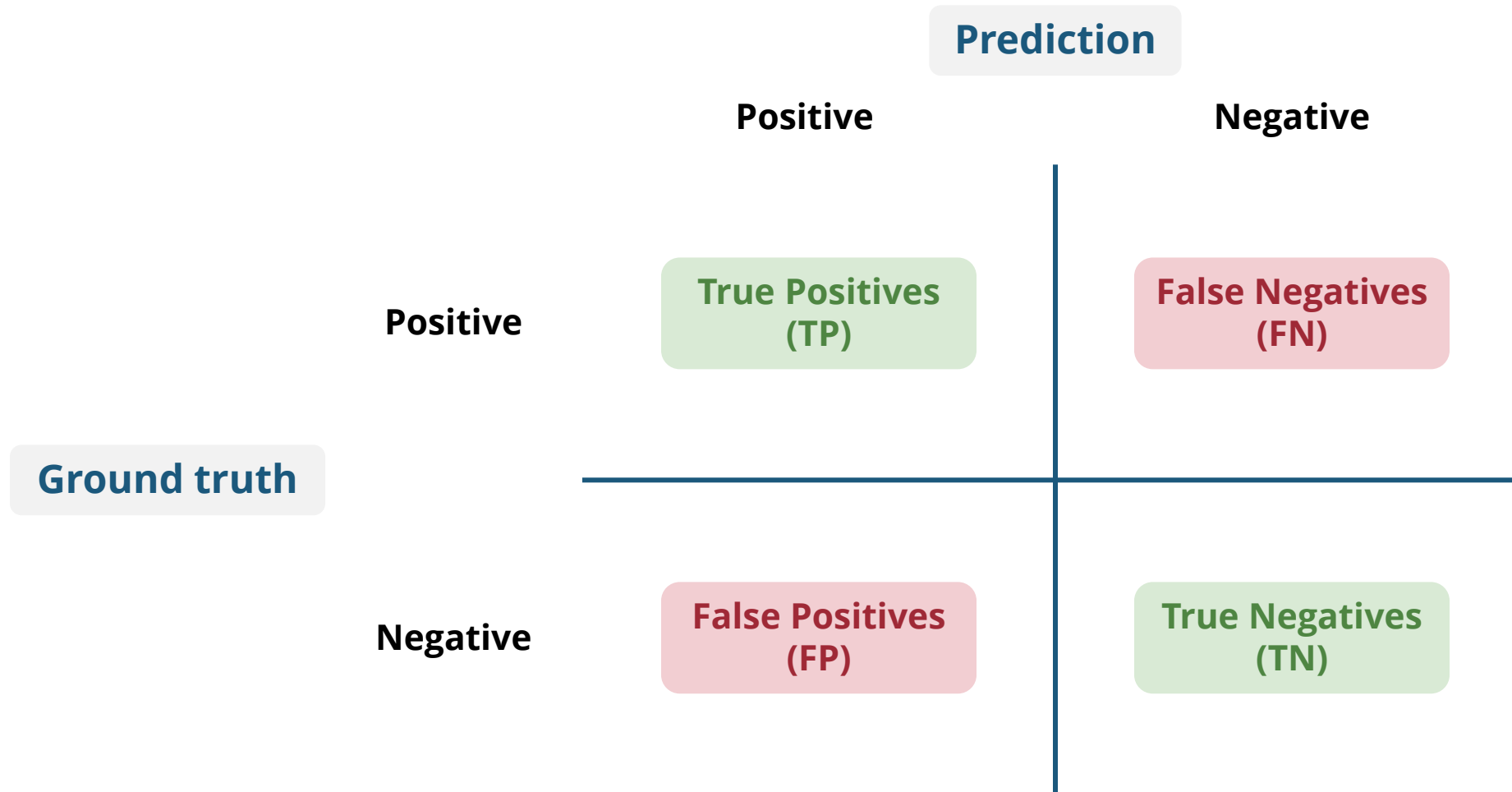
- **Definition:** Percentage of correct predictions across all classes



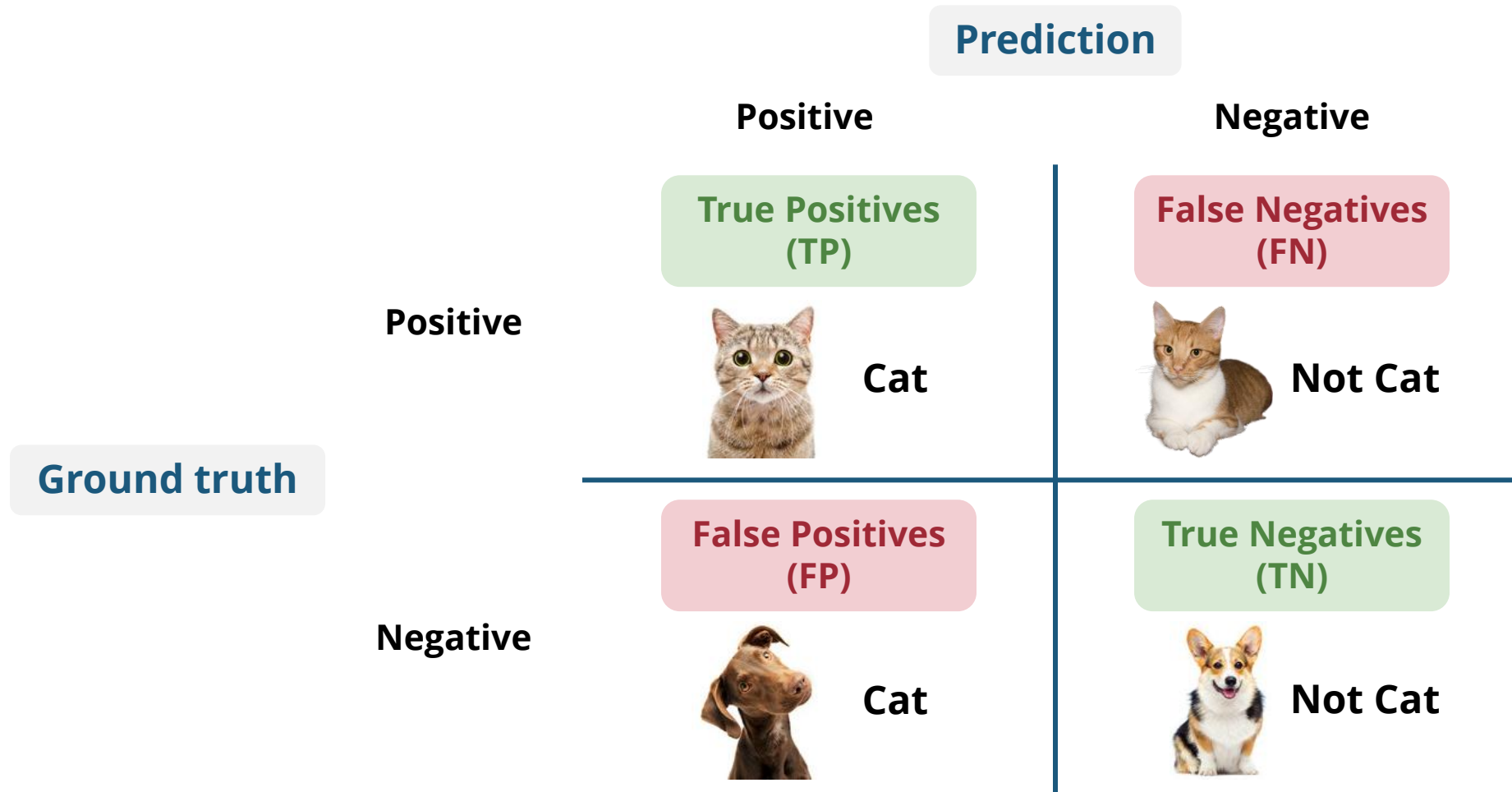
$$\text{Accuracy} = \frac{\text{Correct Predictions}}{\text{Total Predictions}} = \frac{14 + 13}{14 + 1 + 4 + 13} = \frac{27}{32} = 0.84375 \approx 0.82$$

The diagram illustrates the calculation of accuracy. The numerator shows 14 blue circles in a light blue box and 13 orange circles in a light orange box. The denominator shows 14 blue circles in a light blue box, 1 orange circle in a light blue box, 4 blue circles in a light orange box, and 13 orange circles in a light orange box.

Confusion Matrix for Binary Classification



Confusion Matrix for Binary Classification



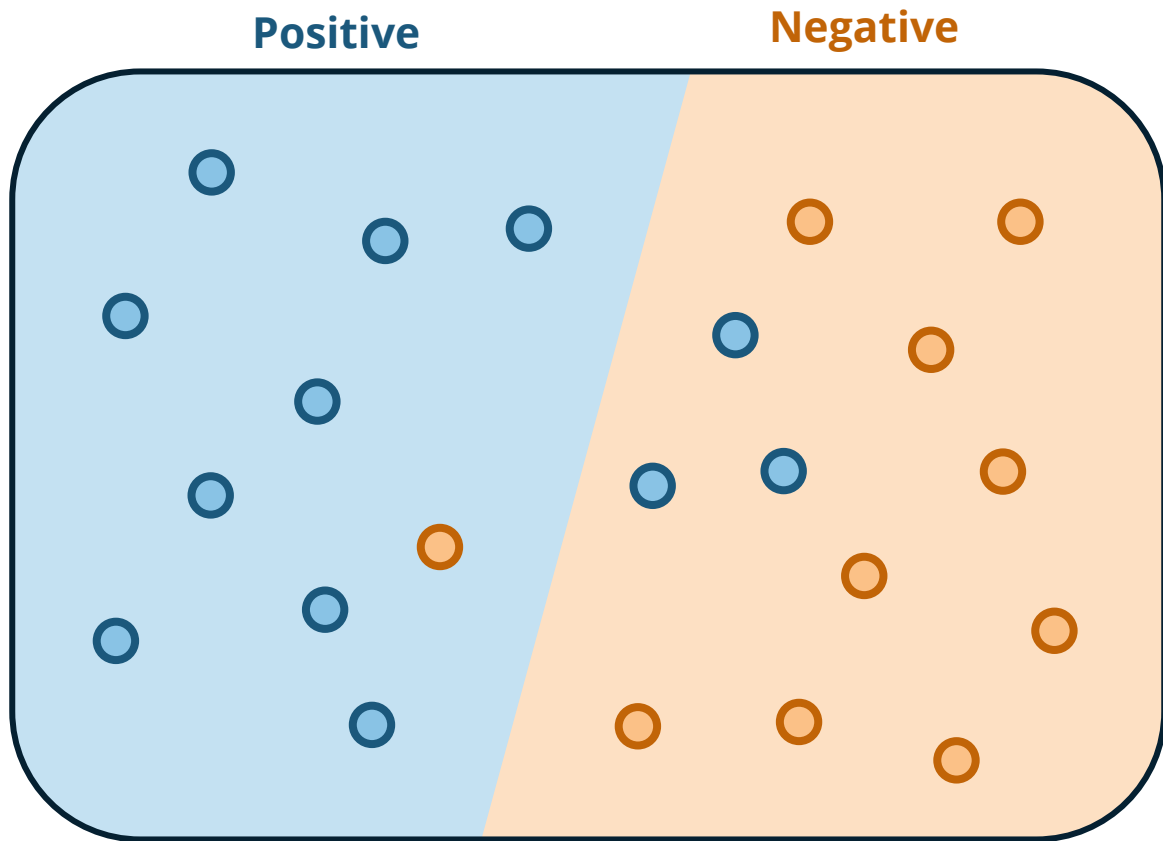
Accuracy on Imbalanced Datasets

- Accuracy does *not* work well on **imbalanced dataset**
- Take **a disease with a 1% prevalence** for example:
 - What if we simply say **negative to all diagnoses**?

		Prediction	
		Positive	Negative
Ground truth	Positive	True Positives (TP) 0	False Negatives (FN) 1
	Negative	False Positives (FP) 0	True Negatives (TN) 99

Accuracy = 0.99 🤔

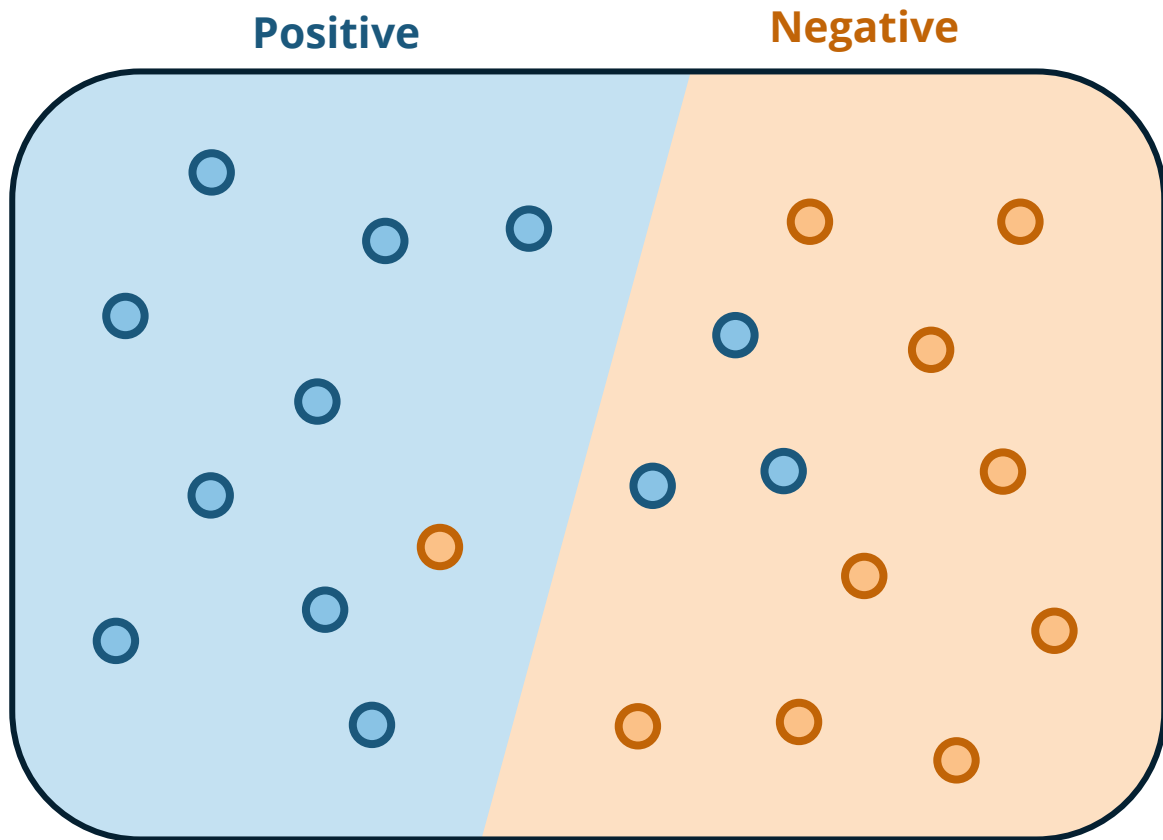
Precision



$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{\text{7}}{\text{7} + \text{2}} = 0.75$$

How often predictions for the positive are correct

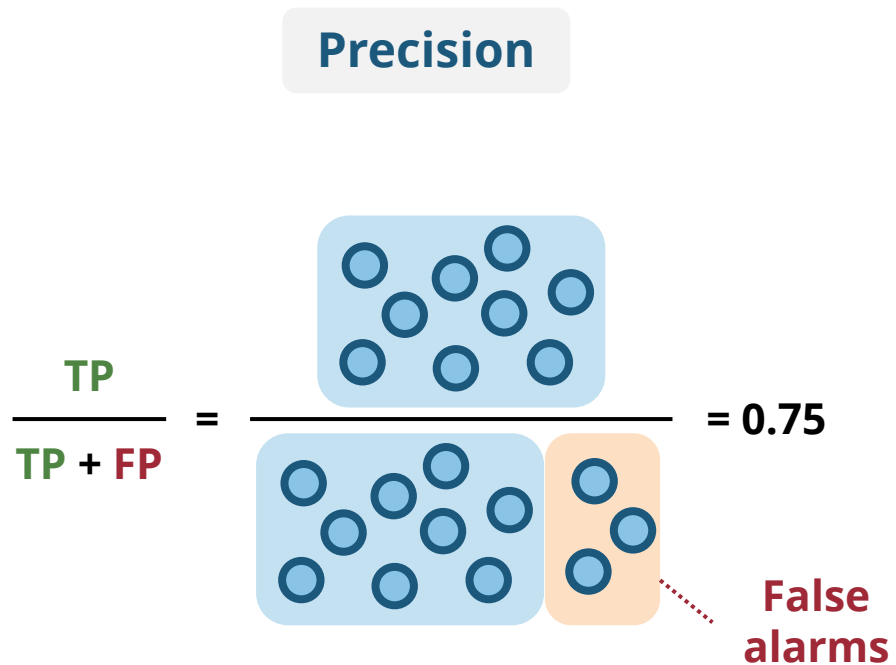
Recall



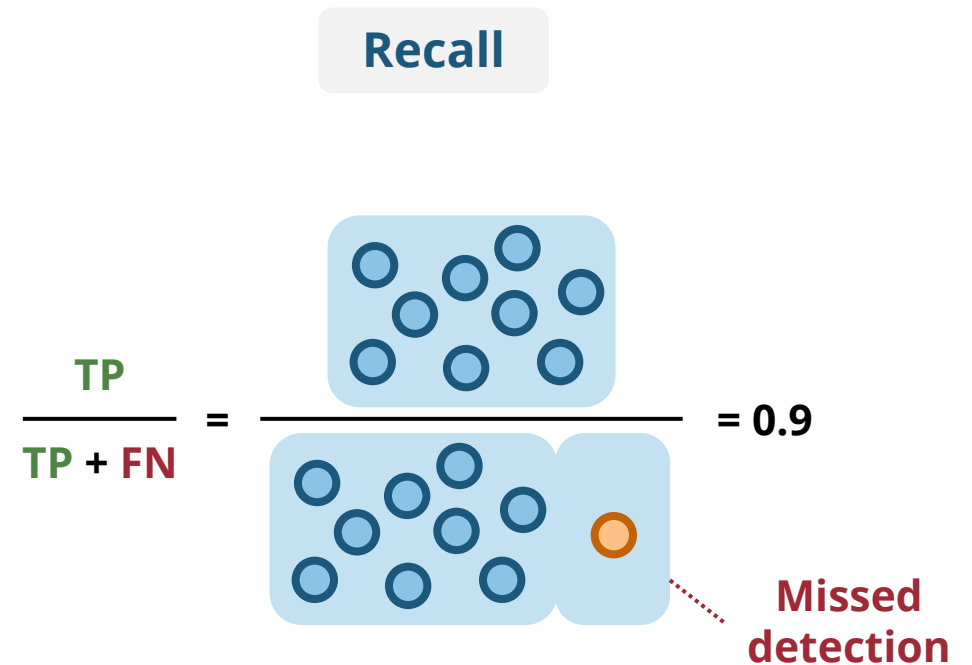
$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{\text{9 blue circles}}{\text{9 blue circles} + \text{1 orange circle}} = 0.9$$

How well the model finds all positive instances in the dataset

Precision vs Recall



How often predictions for the positive are correct



How well the model finds all positive instances in the dataset

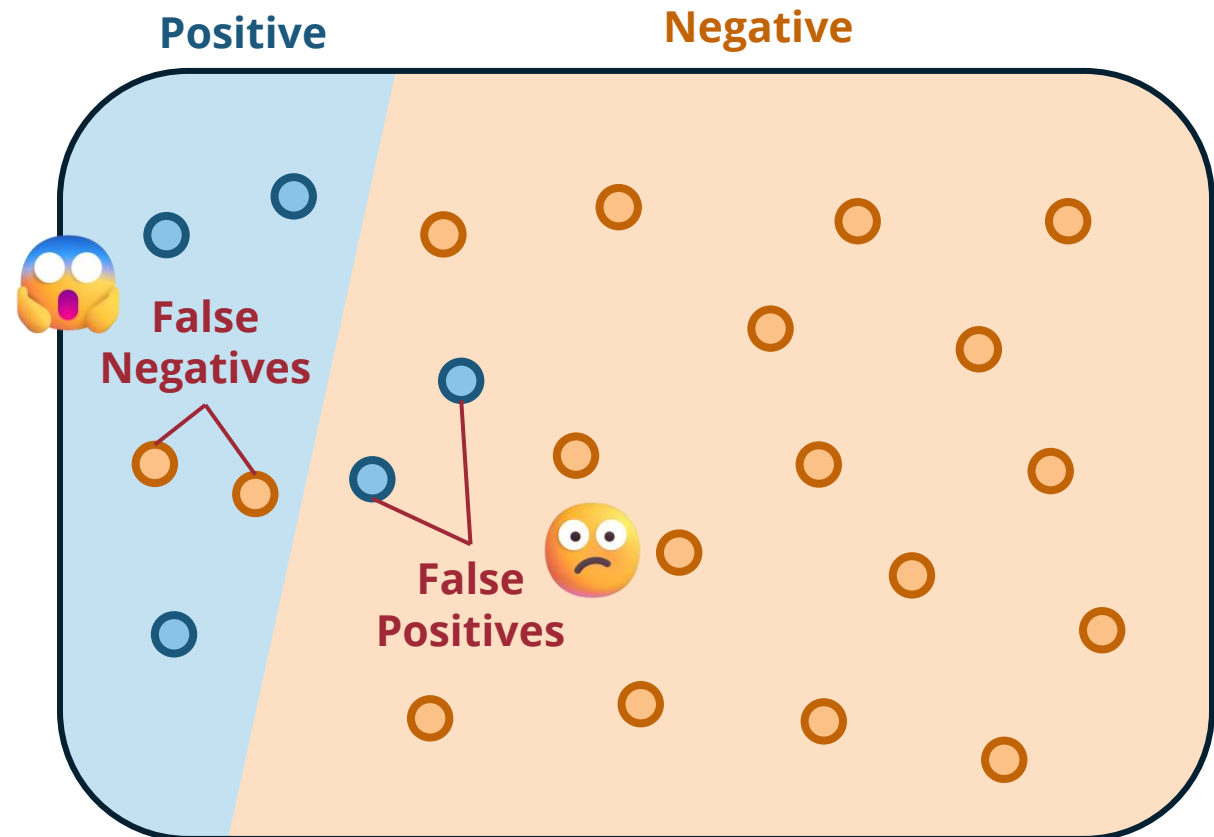
When should we care about Precision & Recall?

Rare cancer detection



Aim for high precision or high recall?

High recall ensures most cancer cases are identified.



False alarms vs Missed detections

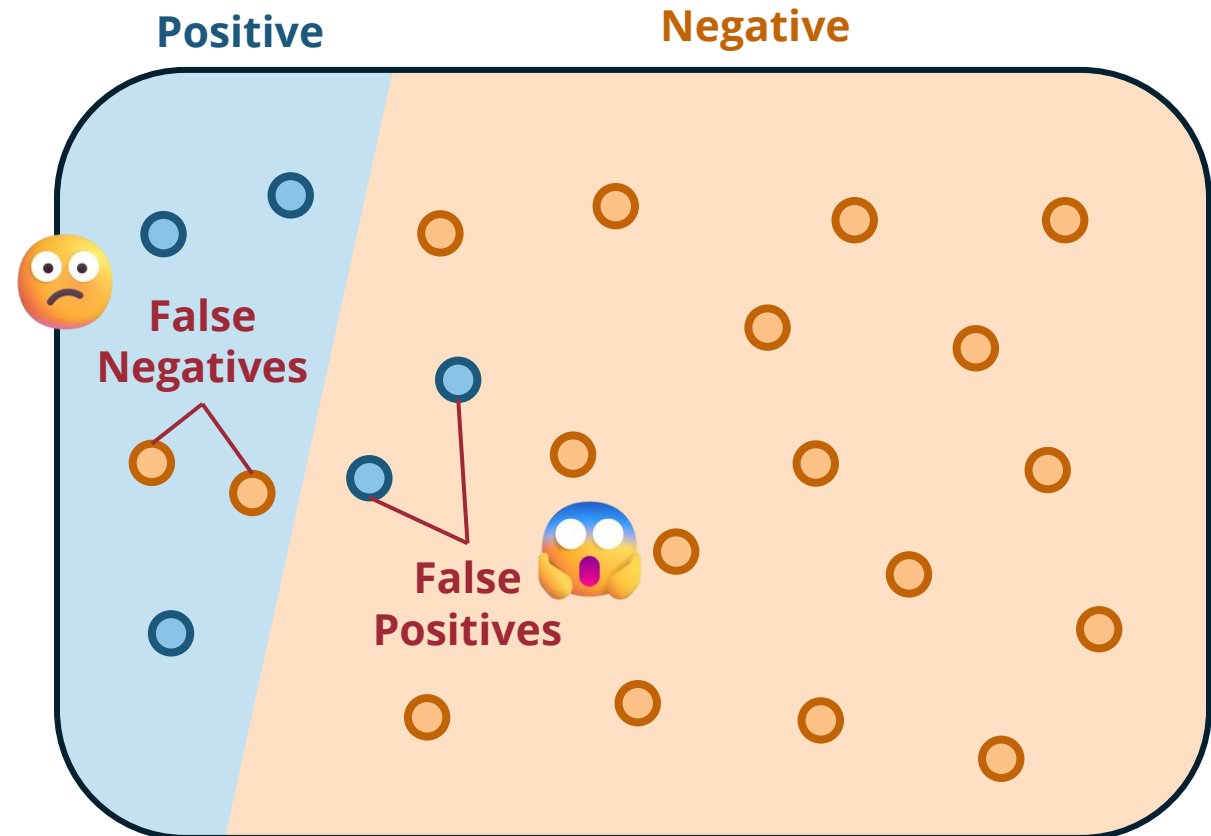
When should we care about Precision & Recall?

Music recommendation



Aim for high precision or high recall?

High precision ensures that the model won't recommend irrelevant items.



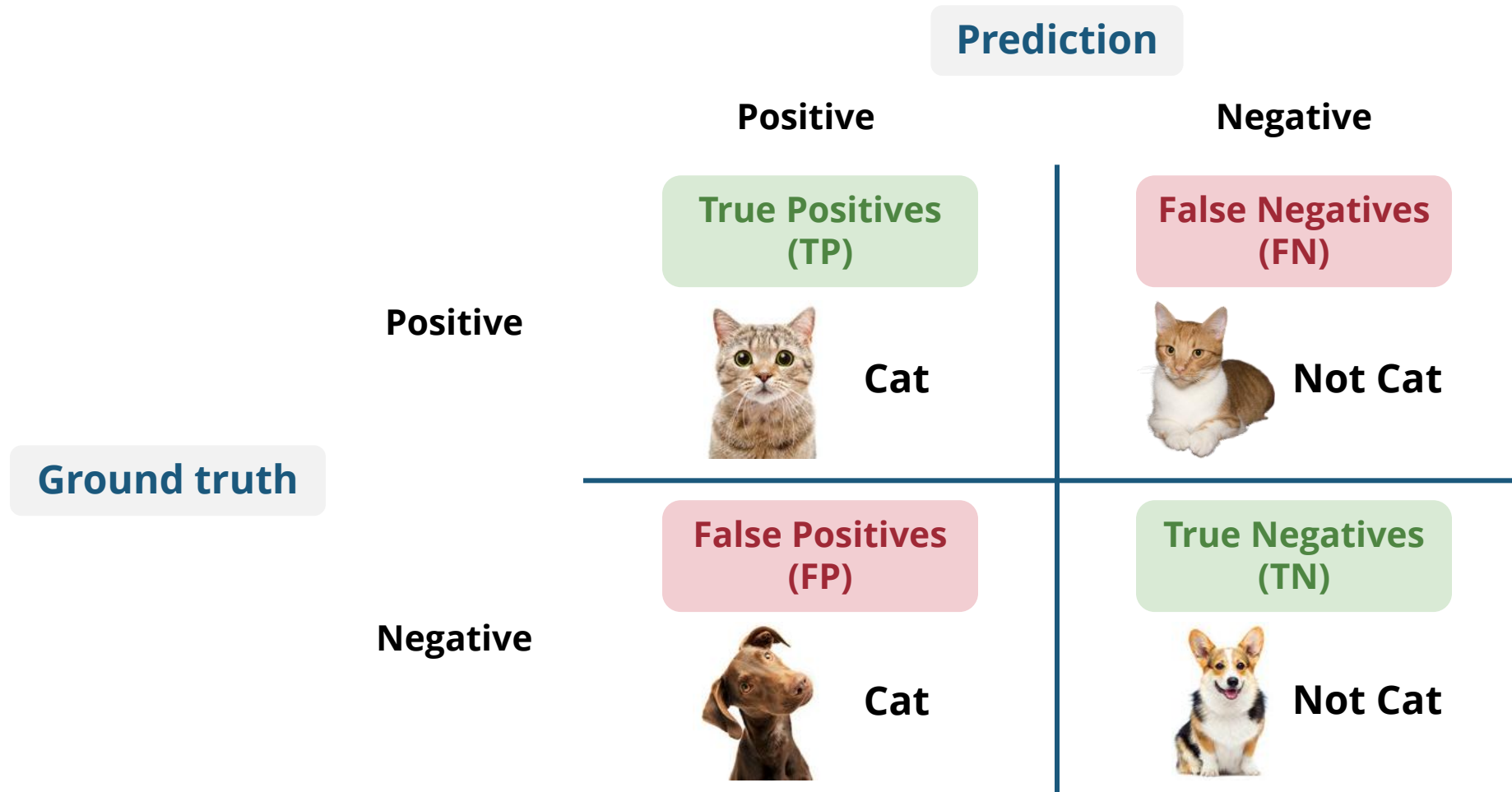
False alarms vs Missed recommendations

F1 Score: Considering both Precision & Recall

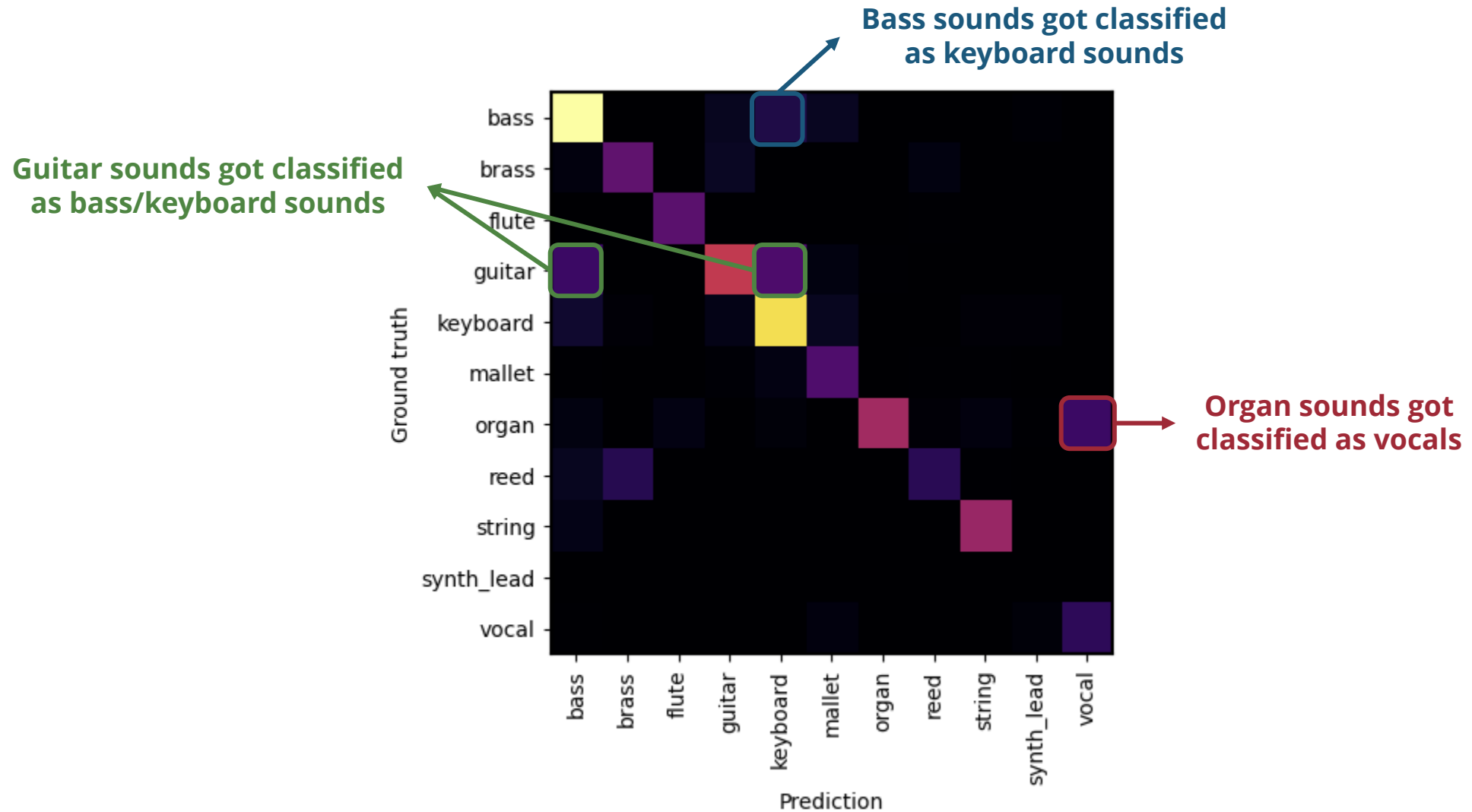
- Particularly useful for **imbalanced datasets**
 - Work better than accuracy when the dataset is imbalanced
 - For example, music search, retrieval, and recommendation

$$F_1 = \frac{2}{\frac{1}{Precision} + \frac{1}{Recall}}$$
$$= \frac{2 \cdot Precision \cdot Recall}{Precision + Recall}$$

Confusion Matrix for Binary Classification



Confusion Matrix for Multiclass Classification



Optional Reading

- Minz Won, Janne Spijkervet, and Keunwoo Choi, "[Music Classification: Beyond Supervised Learning, Towards Real-world Applications](#)," *Tutorials of ISMIR*, 2021.

Open Source Music Classification Models

- github.com/minzwon/sota-music-tagging-models
- github.com/jordipons/musicnn