

Automatic Music Tagging and Classification

Minz Won

minz.won@upf.edu



MTG
Music Technology
Group

About me

- PhD student at MTG
- Worked at:
 - Pandora (United States)
 - Naver Corp. (South Korea)
 - Kakao Corp. (South Korea)

Contents

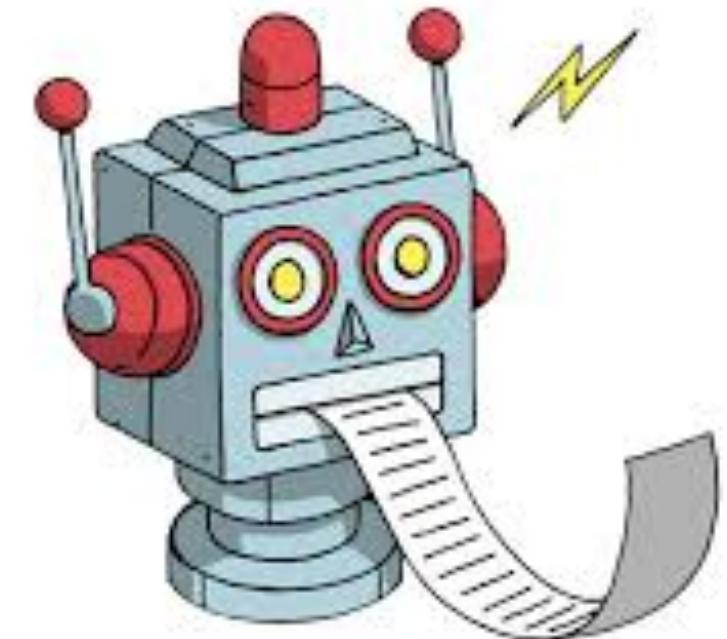
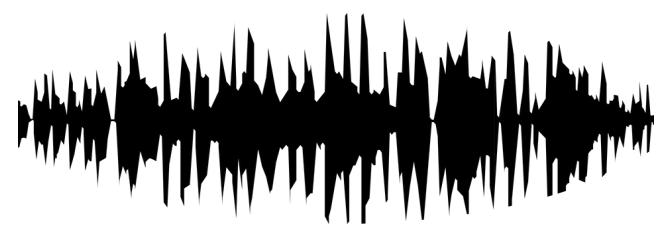
- Introduction
- Music tagging models
- Transfer learning
- Limitations
- Lab

Contents

● Introduction

- Music tagging models
- Transfer learning
- Limitations
- Lab

Music Classification



Blues

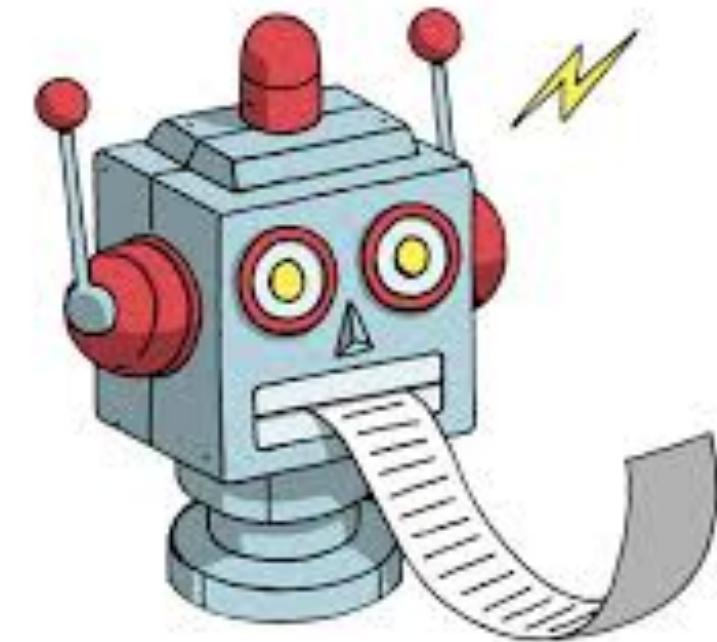
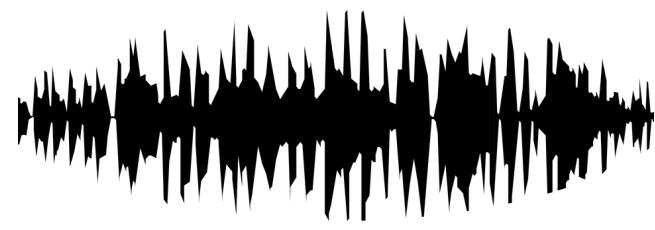
Classical

Hiphop

Jazz

Rock

Automatic Music Tagging



Rock

Guitar

Male vocal

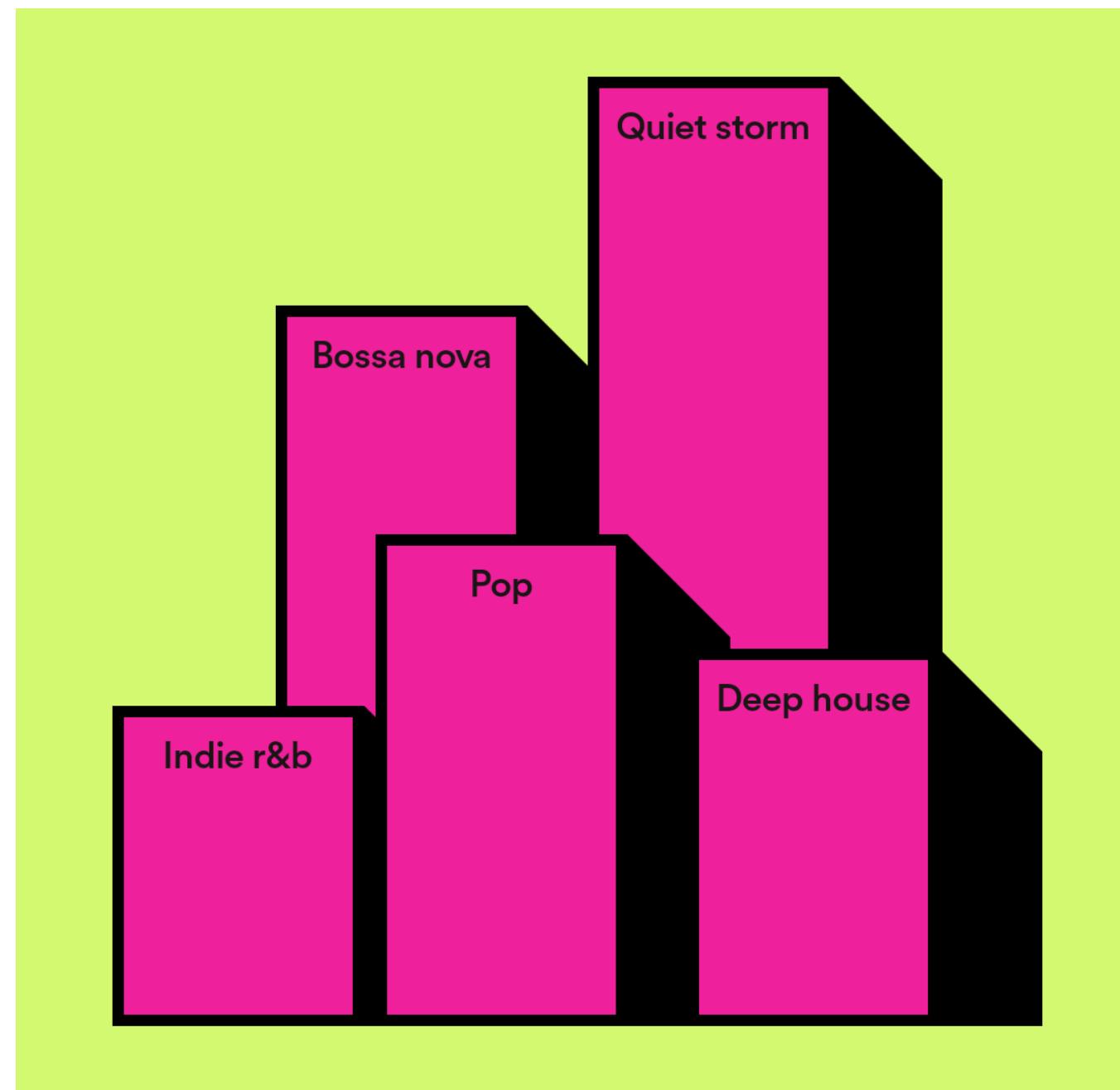
Female vocal

80s

90s

Why do we need this?

Why do we need this?



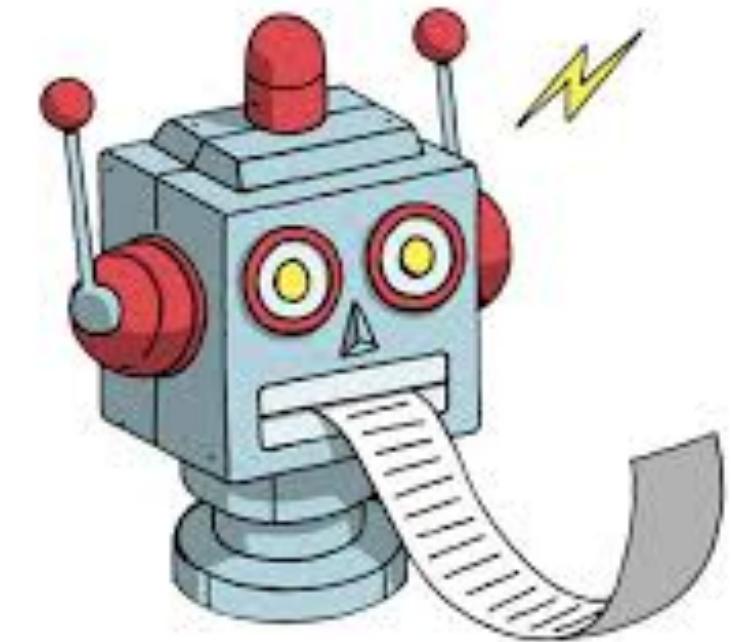
Recommendation

Why do we need this?

		Genre	Mood	Vocal	Save	Clear
	Title			Artist	Remixer	
▶	High U Gonna Feel (Original Mix)			Den Ishu		
▶	Replaced feat. Jaw (Norman Webers Back To The Ro...	Jaw, Enliven			Norman Weber	
▶	Trendy Jose (JT Donaldson Remix)		Ivaylo		JT Donaldson	
▶	Moments (Original Mix)		Saison			
▶	The Premise (Original Mix)		Kassian			
▶	In the Morning feat. Zakes Bantwini (Atjazz Remix)	Jazzanova, Zakes Bantwini	Atjazz			

Retrieval

Why do we need this?



Similarity-based search

Rock

Guitar

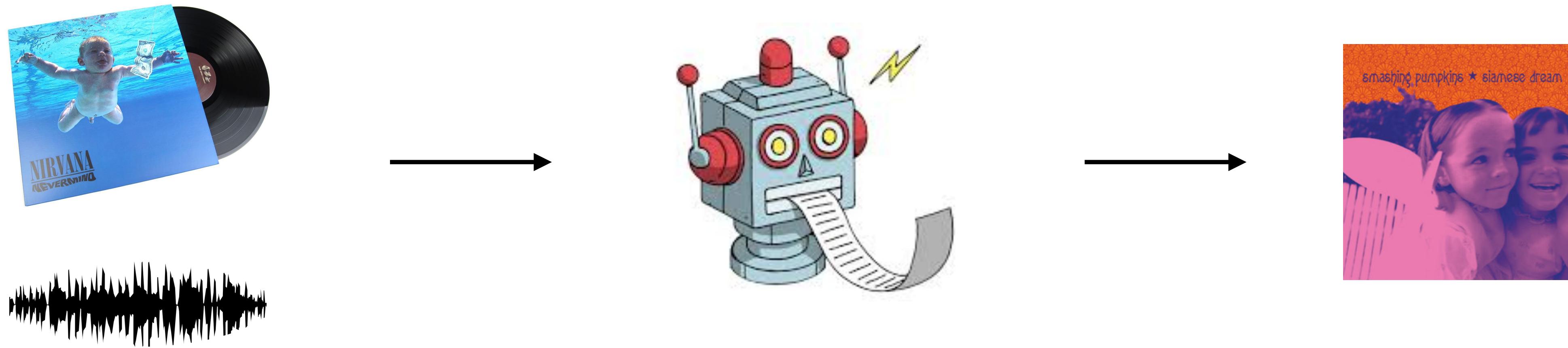
Male vocal

Female vocal

80s

90s

Why do we need this?



Similarity-based search

Why do we need this?

The screenshot shows a digital music player interface with a dark theme. At the top, there is a playback bar with a progress indicator at 00:19 / 01:00. Below the bar, Korean text provides keyboard shortcuts: 단축키: ← 이전 → 다음 ESC 재생/멈춤. It also says 미리듣기를 하고 싶은 곡의 앨범 이미지를 클릭해주세요.

The main area displays a list of five jazz songs:

- 01. How High The Moon - The Things You Did Last Summer New York Trio 4739476 Youtube 검색 [연관곡] [아티스트의 연관곡] [앨범의 연관곡]
- 02. Broadway - Lover Man The Jacky Terrasson Jazz Trio 6745138 Youtube 검색 [연관곡] [아티스트의 연관곡] [앨범의 연관곡]
- 03. Always - Secret Love Eddie Higgins Trio 4713506 Youtube 검색 [연관곡] [아티스트의 연관곡] [앨범의 연관곡]
- 04. Blues In The Closet - The Essen Jazz Festival Concert Bud Powell 2805632 Youtube 검색 [연관곡] [아티스트의 연관곡] [앨범의 연관곡]
- 05. I Will Wait For You - I Will Wait For You Steve Kuhn Trio 3721382 Youtube 검색 [연관곡] [아티스트의 연관곡] [앨범의 연관곡]

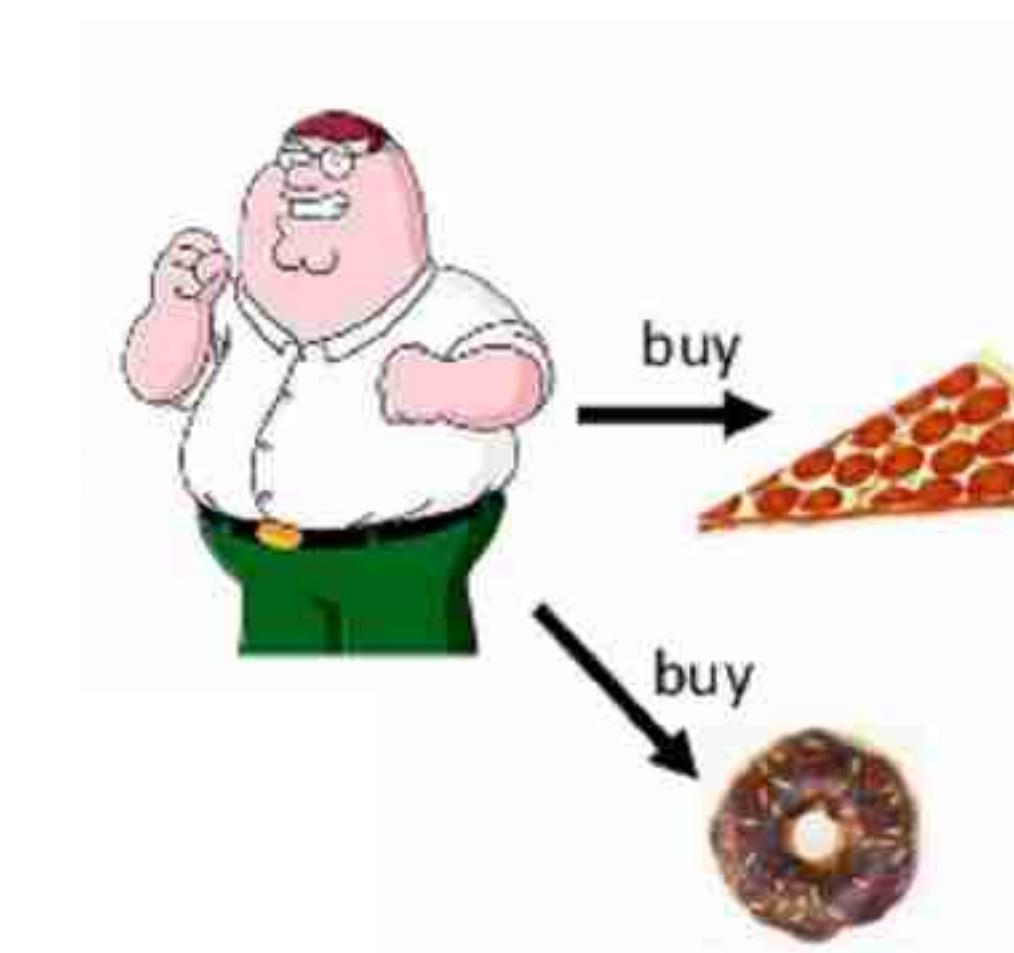
A large watermark "Input (Jazz)" is visible across the bottom of the list.

Similarity-based search

Why do we need this?

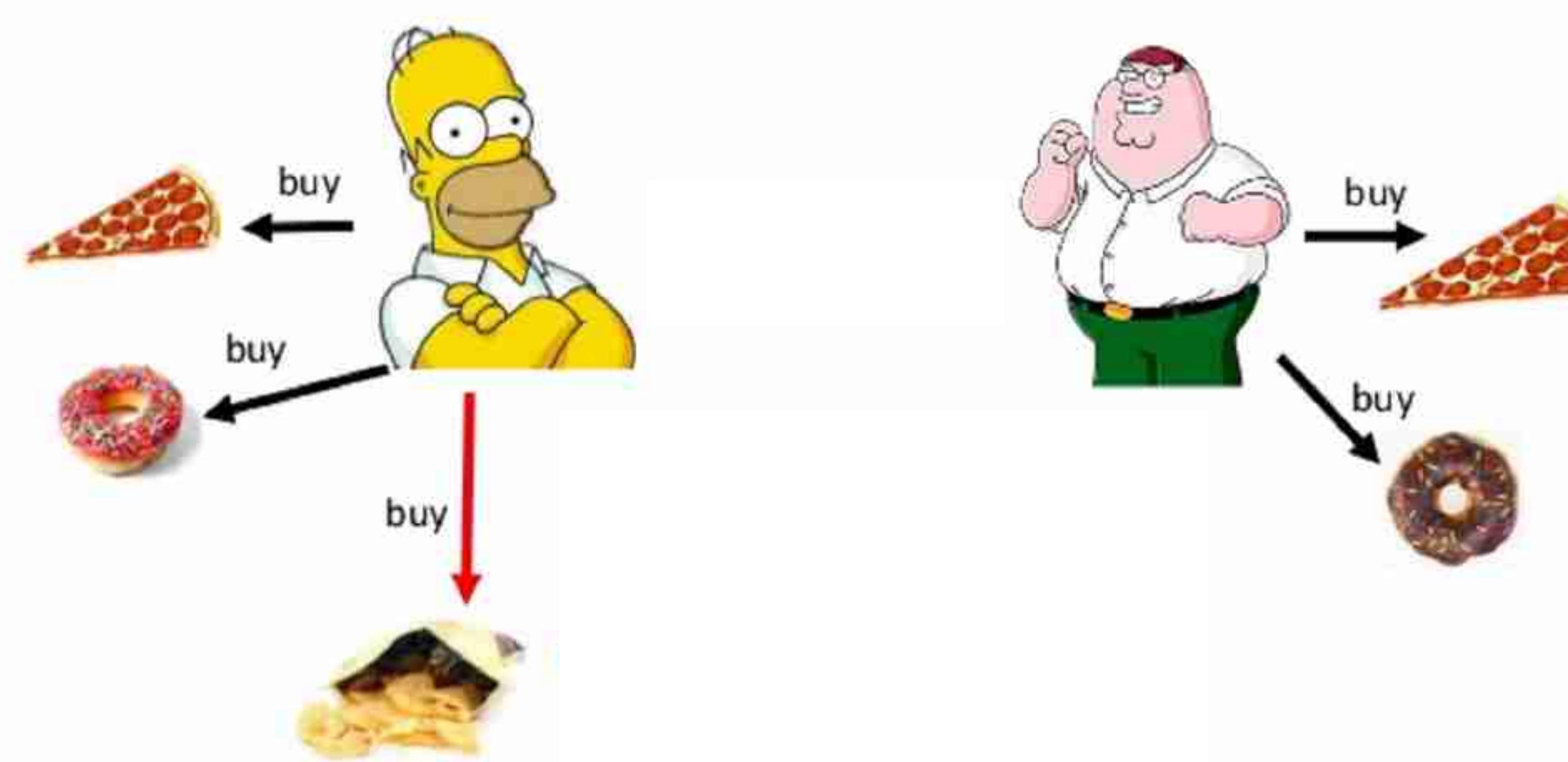
It's fully content-based!

Why do we need this?



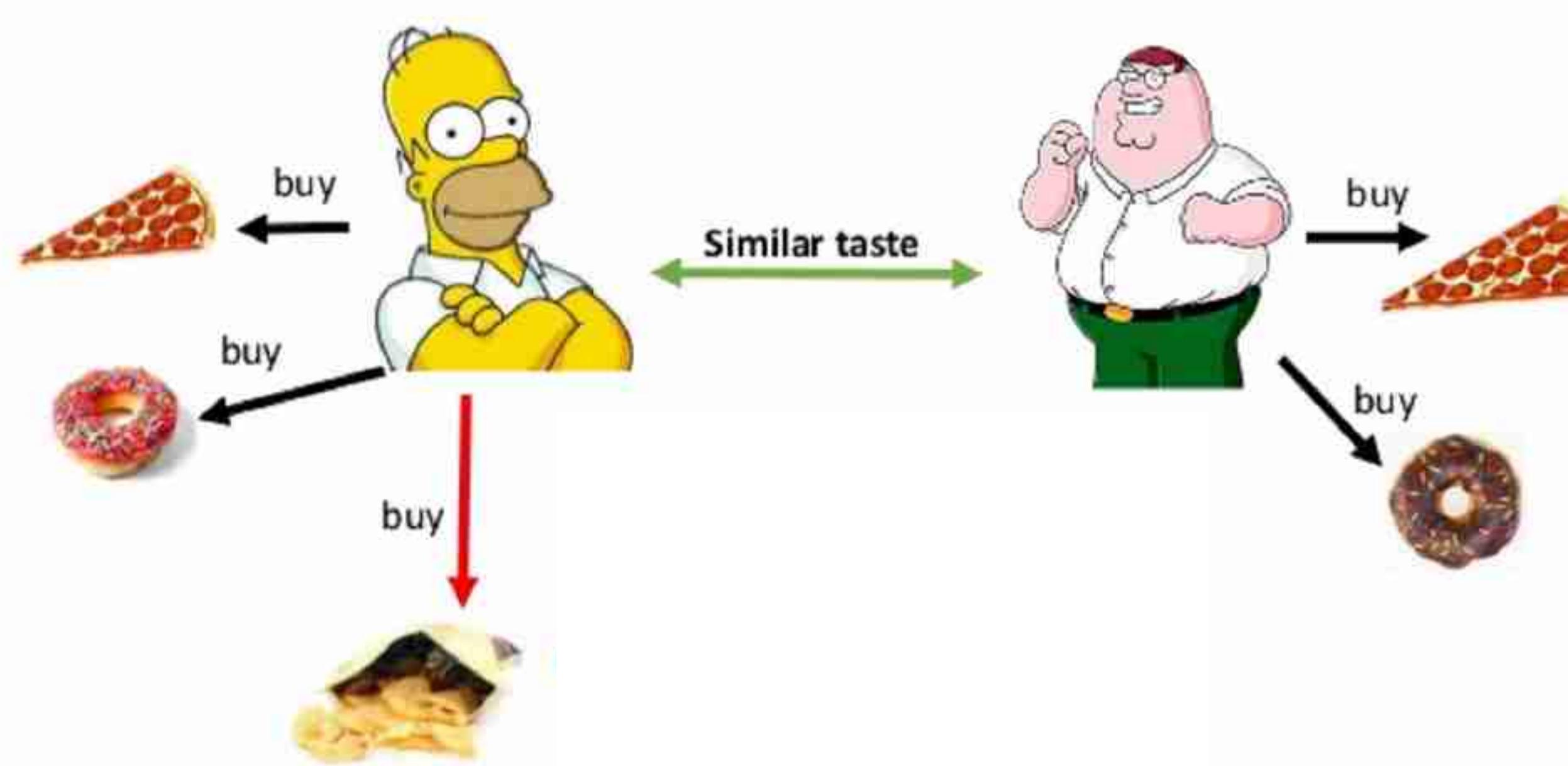
Collaborative filtering

Why do we need this?



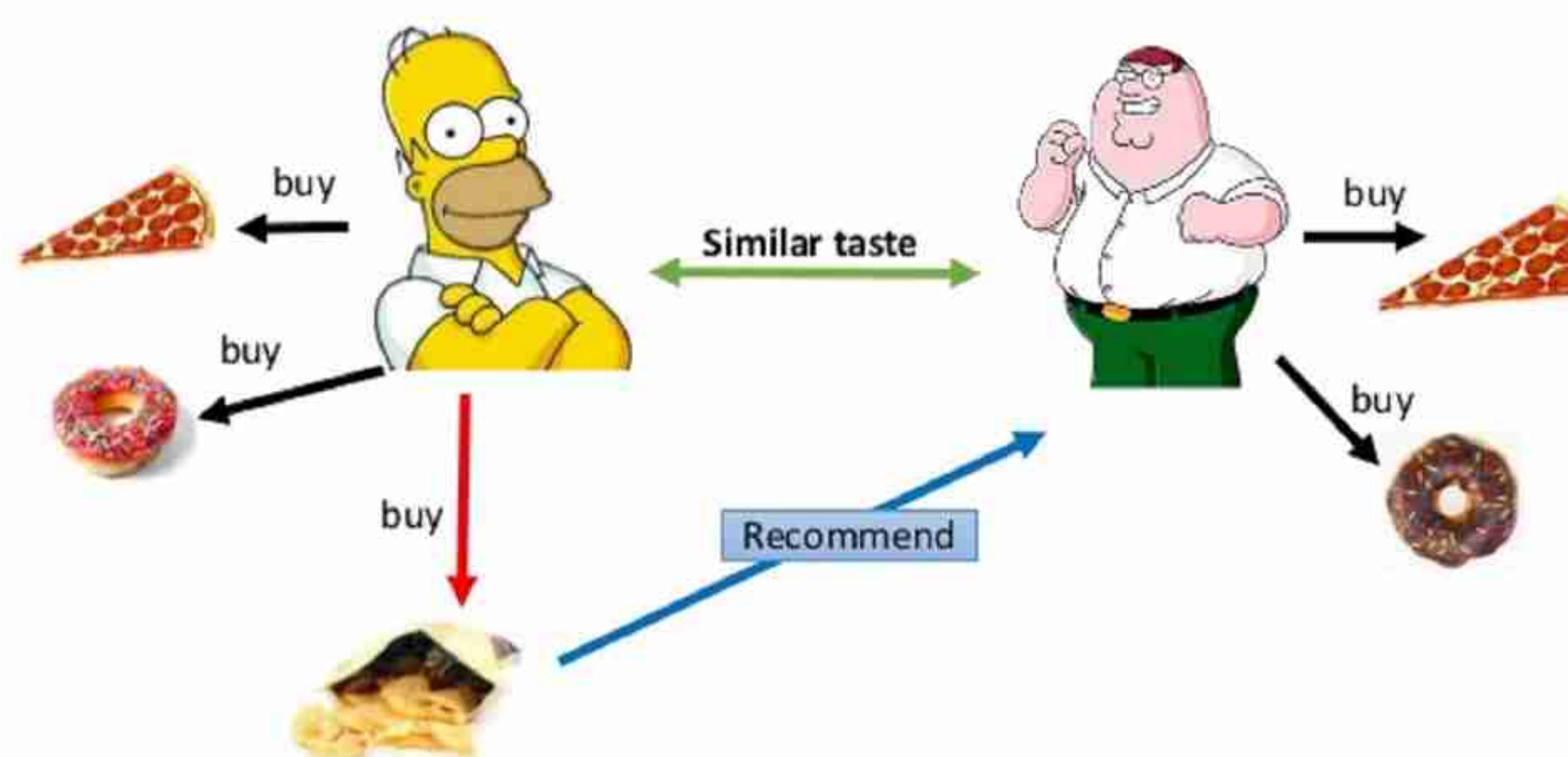
Collaborative filtering

Why do we need this?



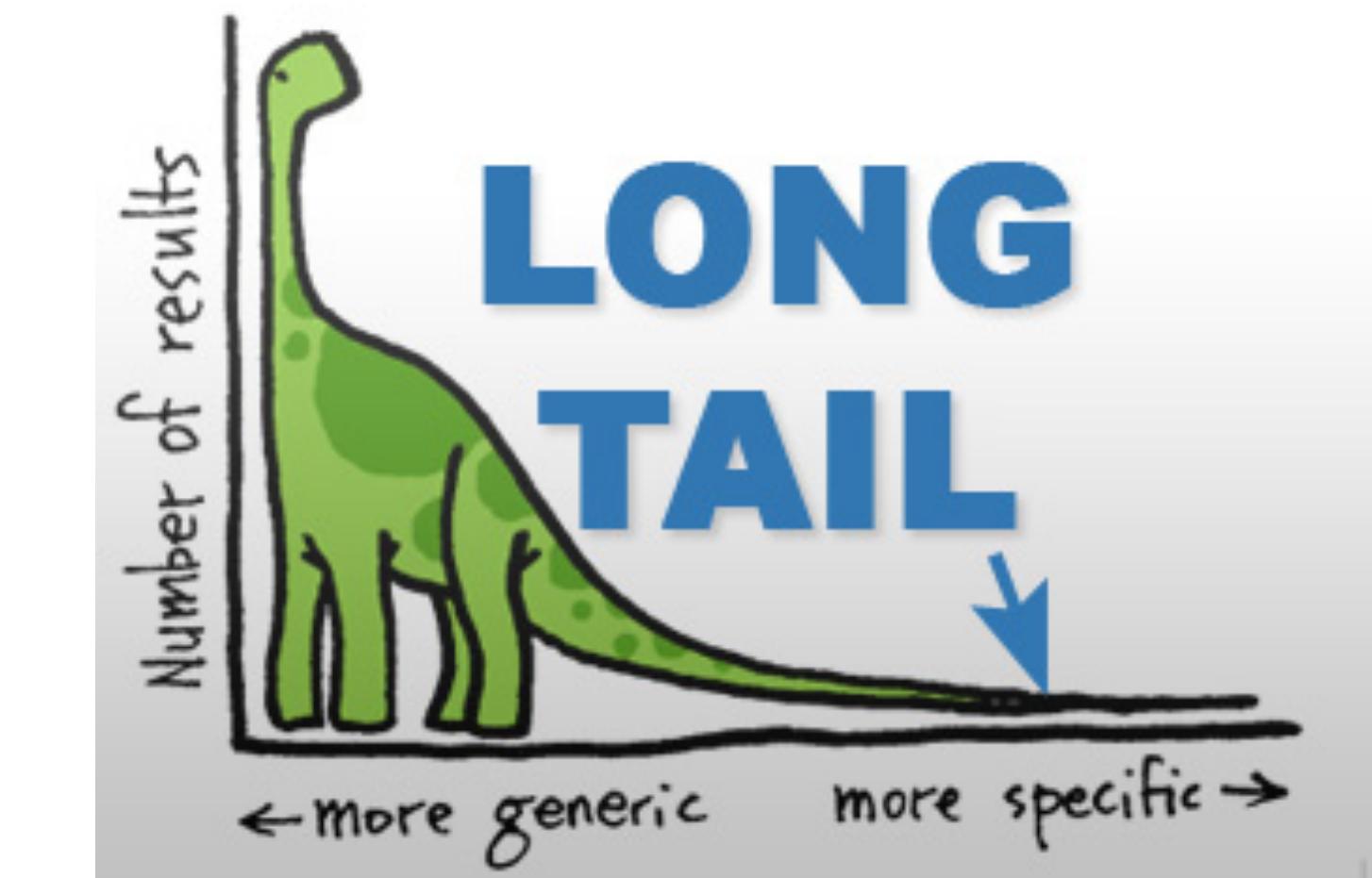
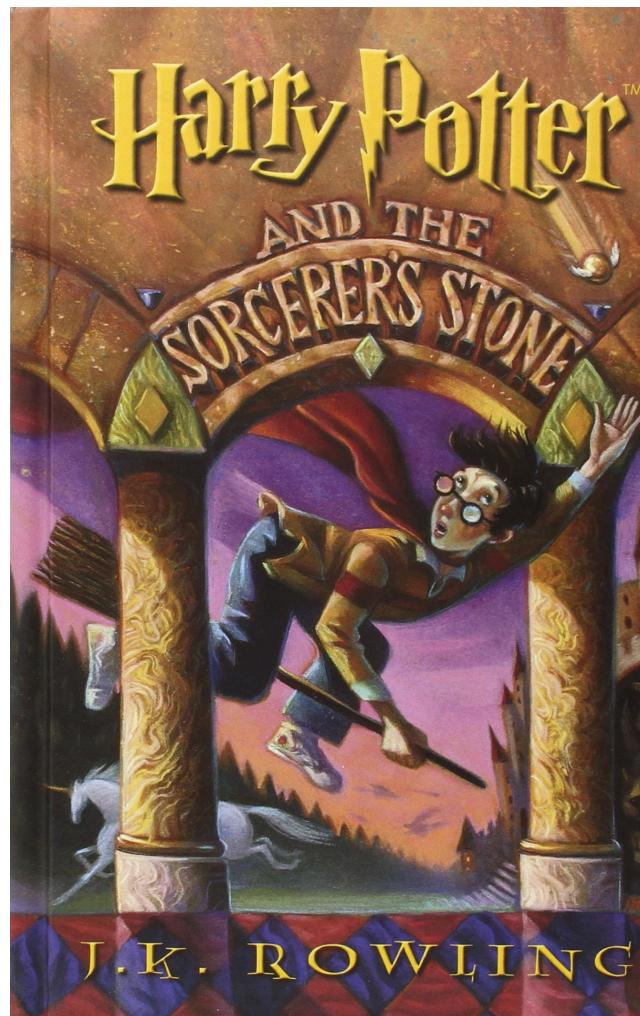
Collaborative filtering

Why do we need this?



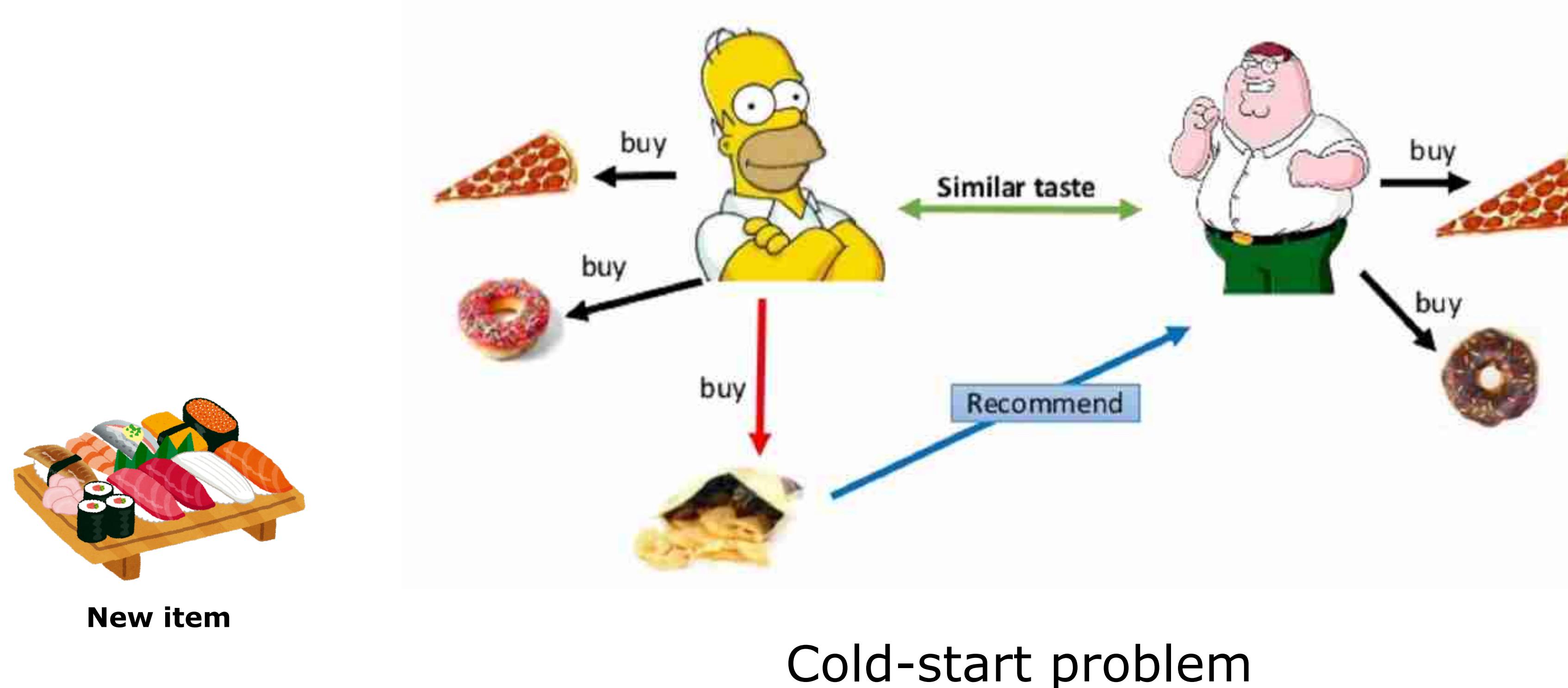
Collaborative filtering

Why do we need this?

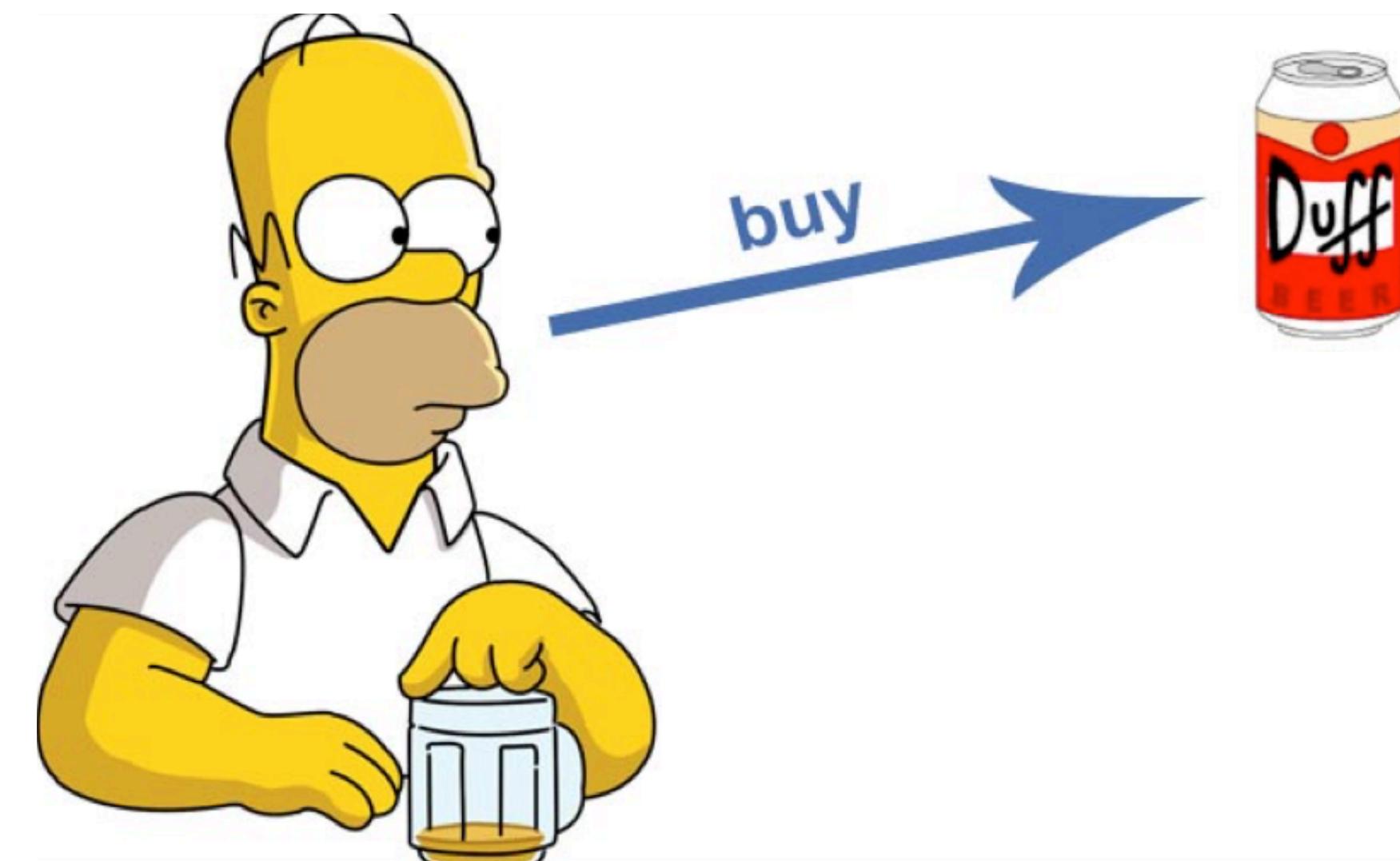


Popularity bias

Why do we need this?

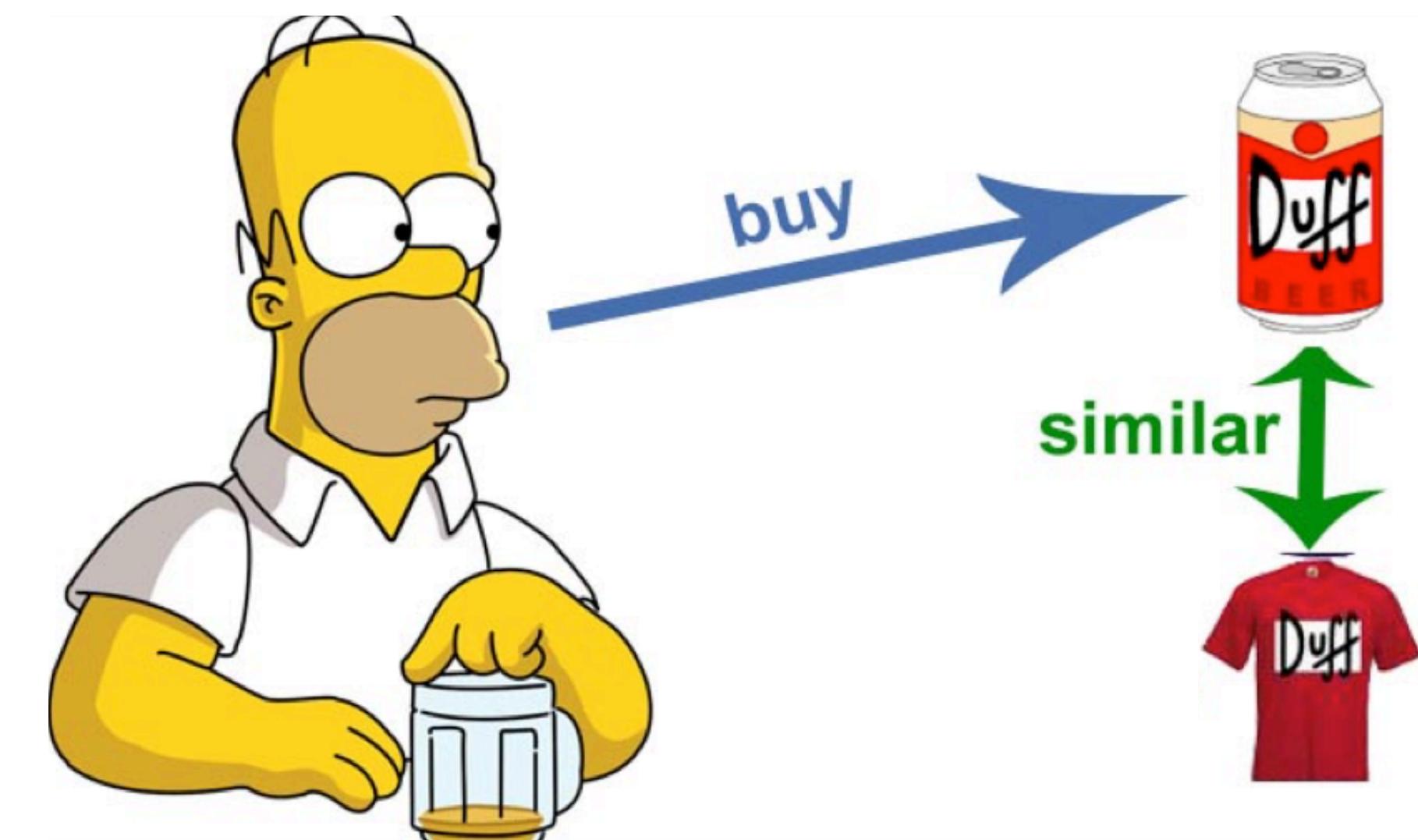


Why do we need this?



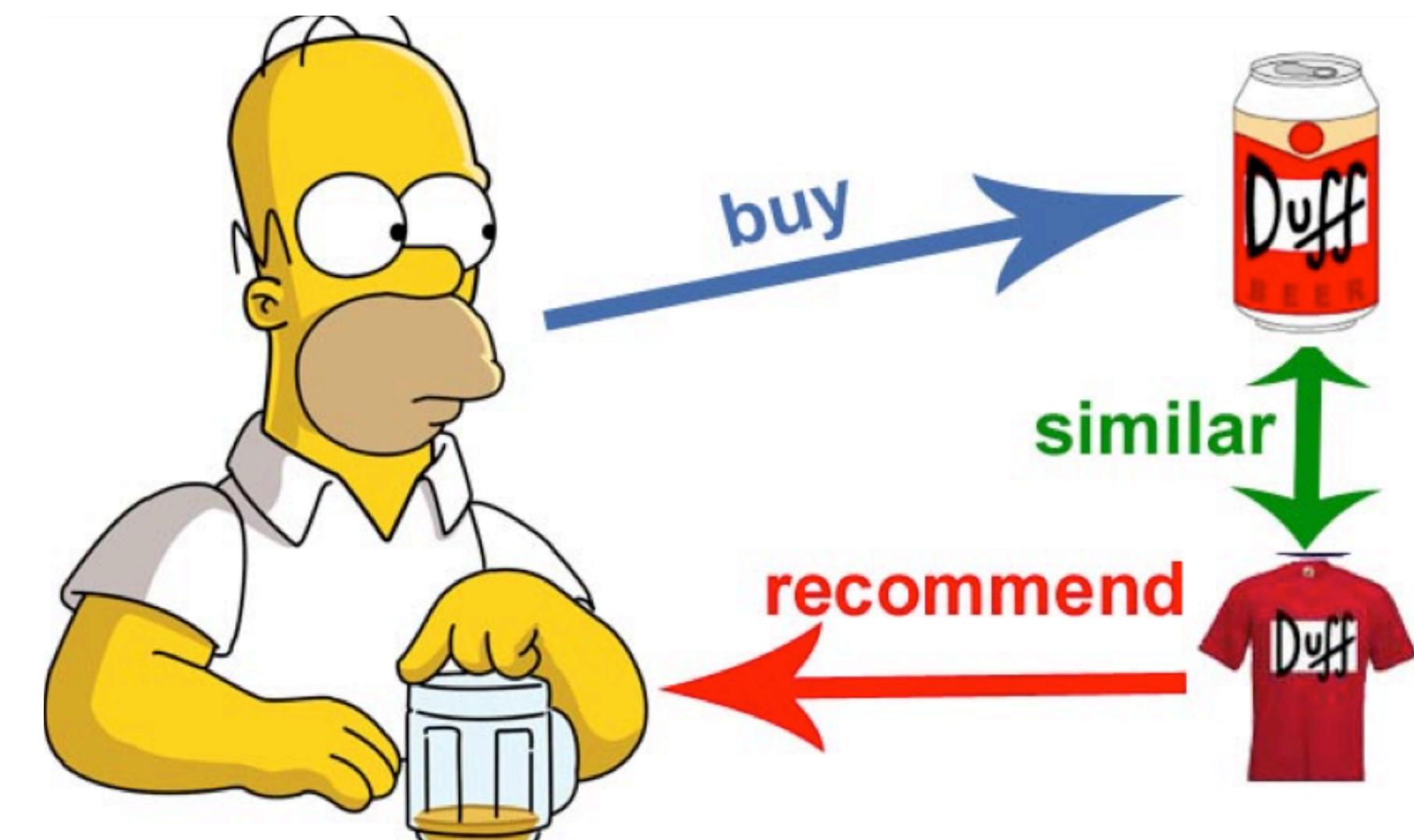
Content-based

Why do we need this?



Content-based

Why do we need this?



Content-based

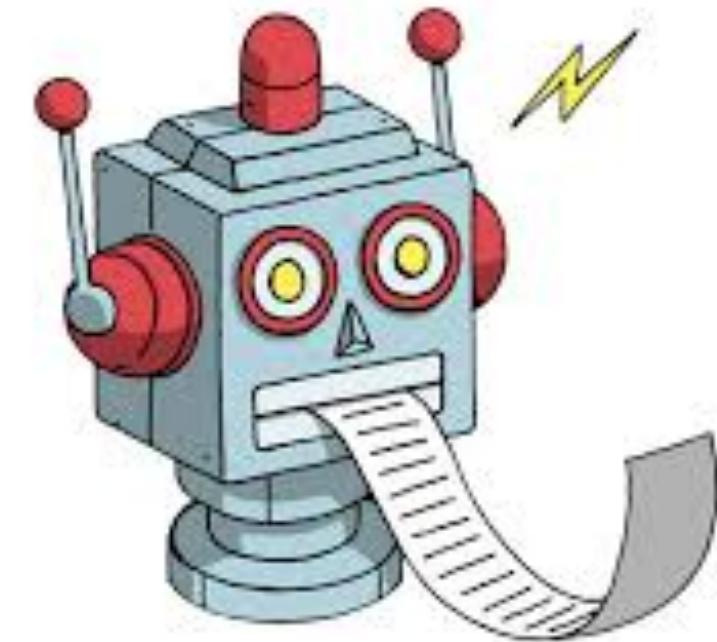
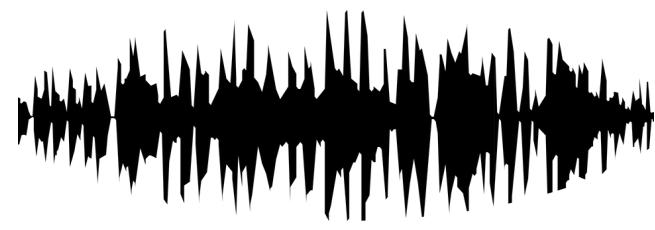
Contents

- Introduction

● **Music tagging models**

- Transfer learning
- Limitations
- Lab

Automatic Music Tagging



Rock

Guitar

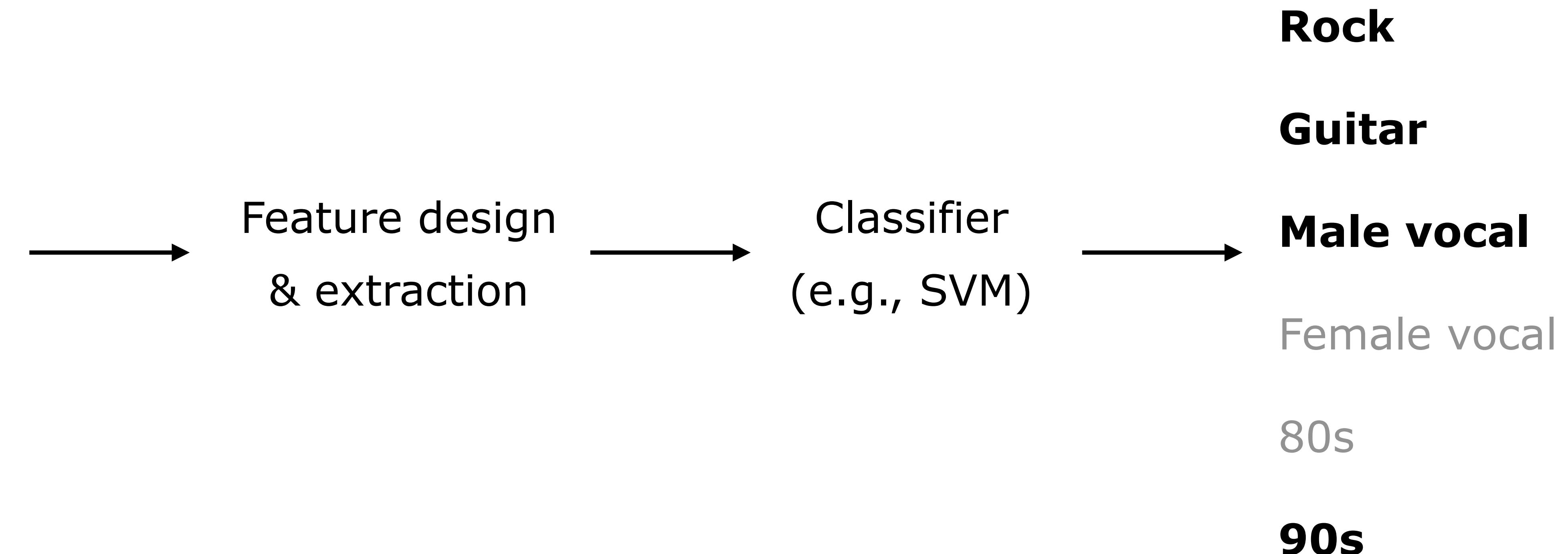
Male vocal

Female vocal

80s

90s

Traditional approaches



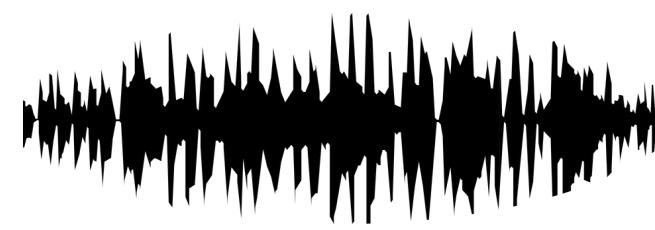
Recent approaches

People with no idea about AI
saying it will take over the world:

My Neural Network:



Recent approaches



Deep neural
networks



Rock

Guitar

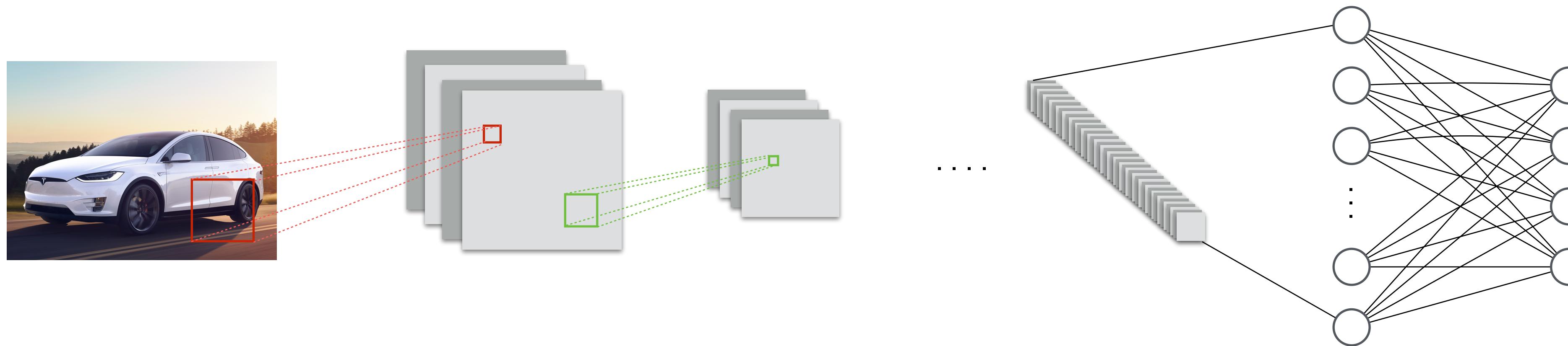
Male vocal

Female vocal

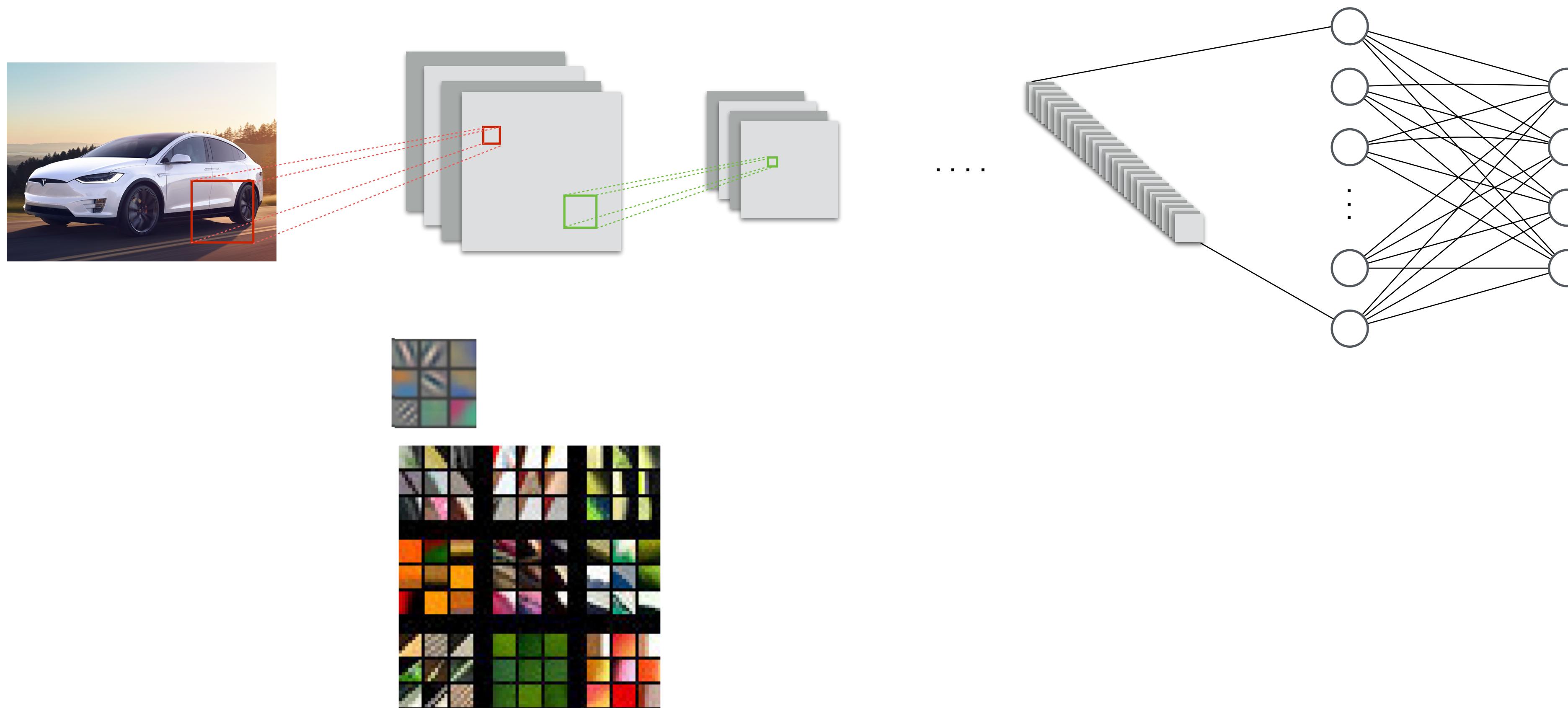
80s

90s

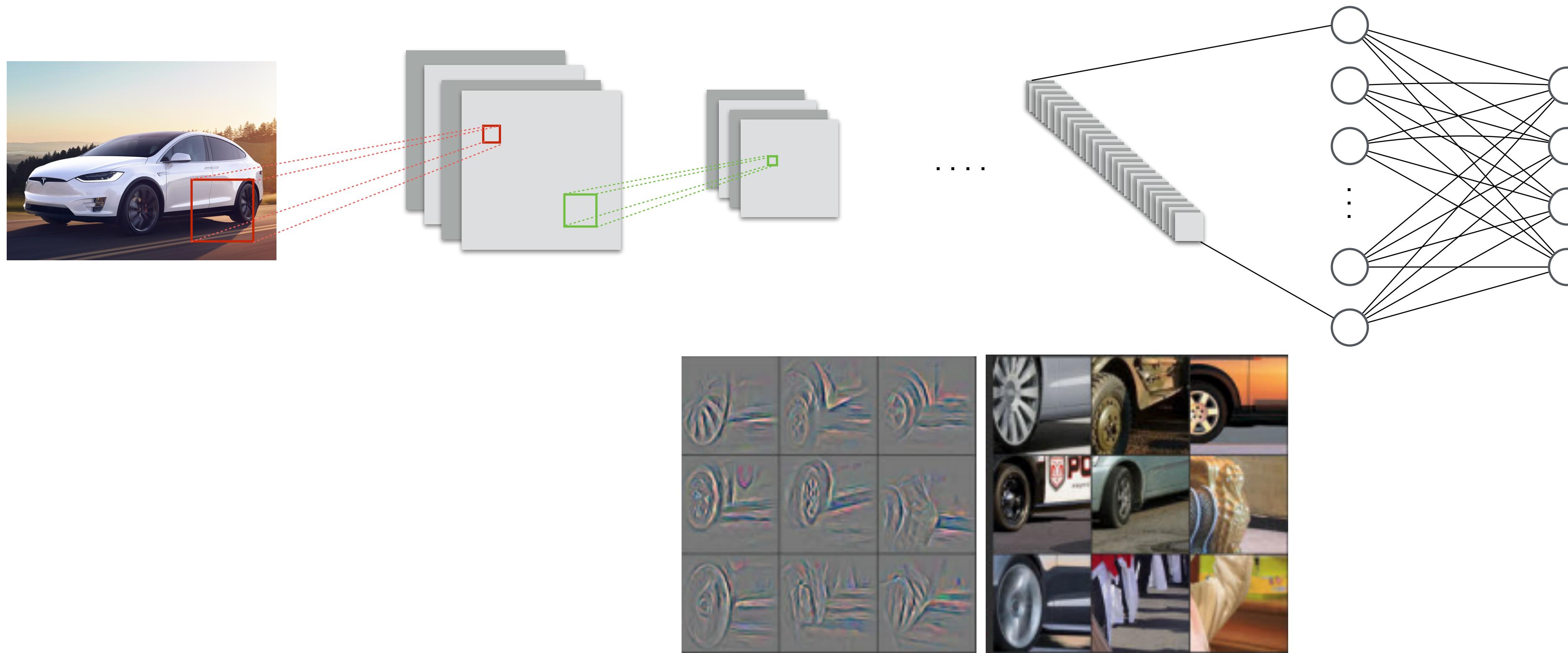
Convolutional Neural Networks



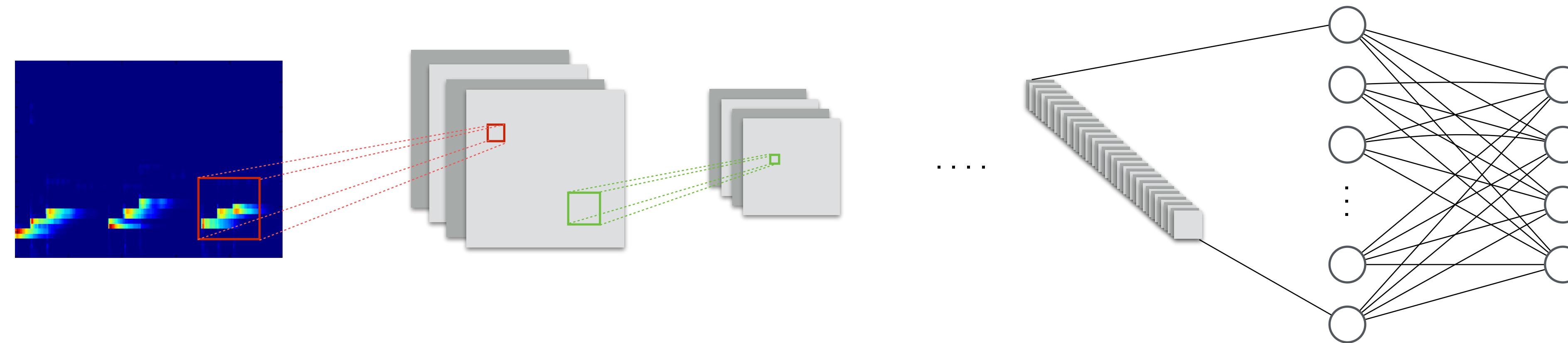
Convolutional Neural Networks



Convolutional Neural Networks

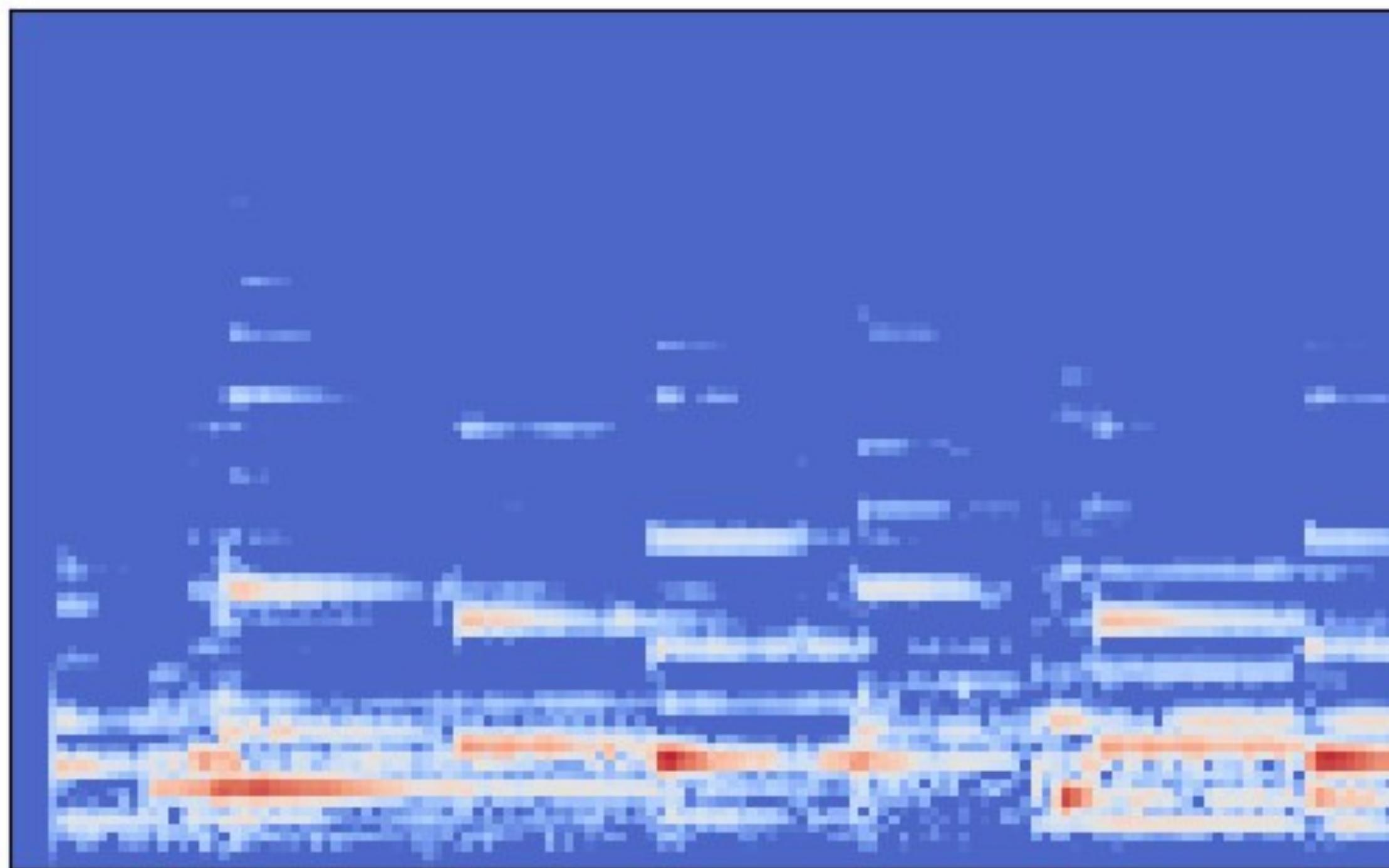


Convolutional Neural Networks

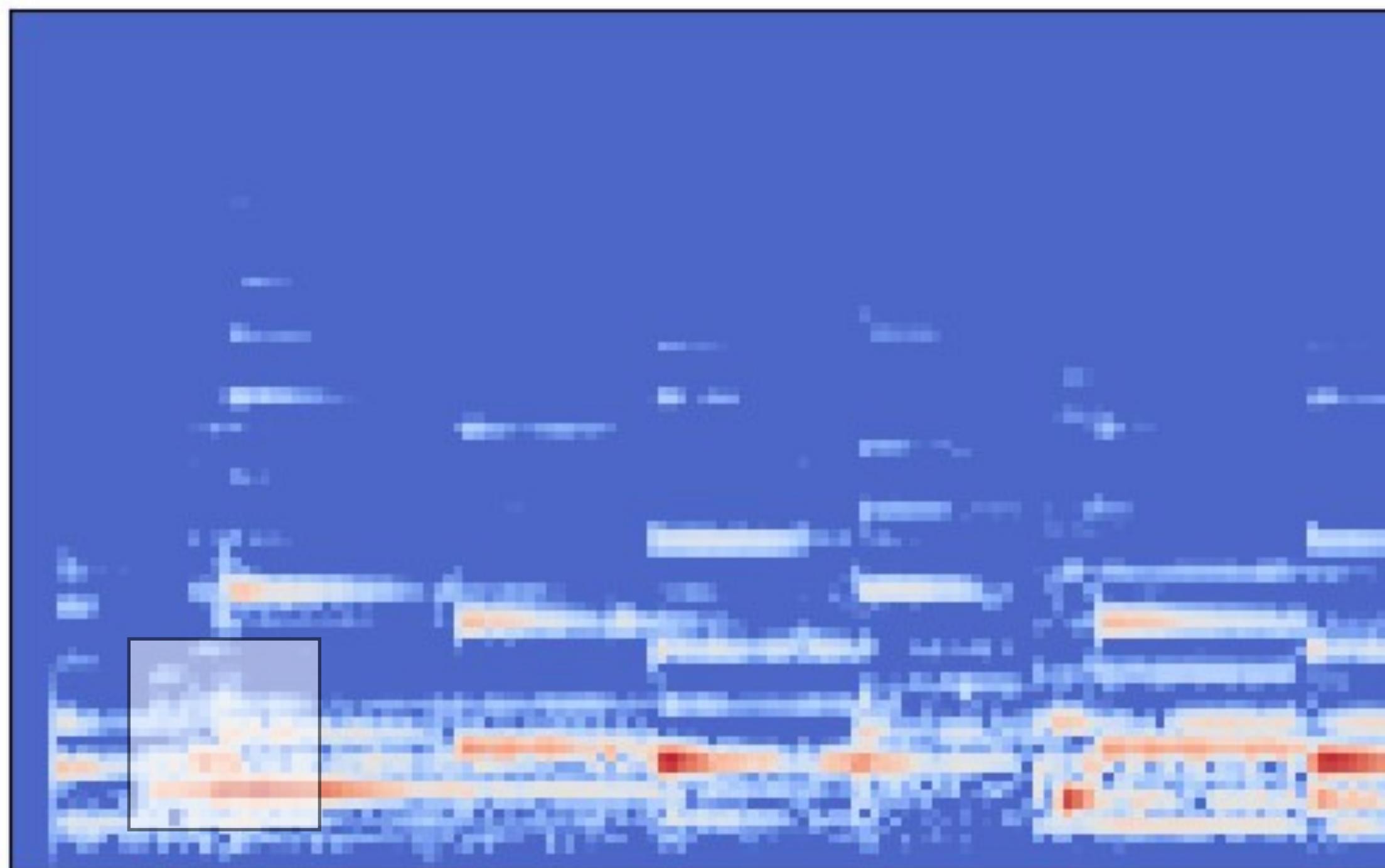


Fully Convolutional Networks (Choi et al., 2016)

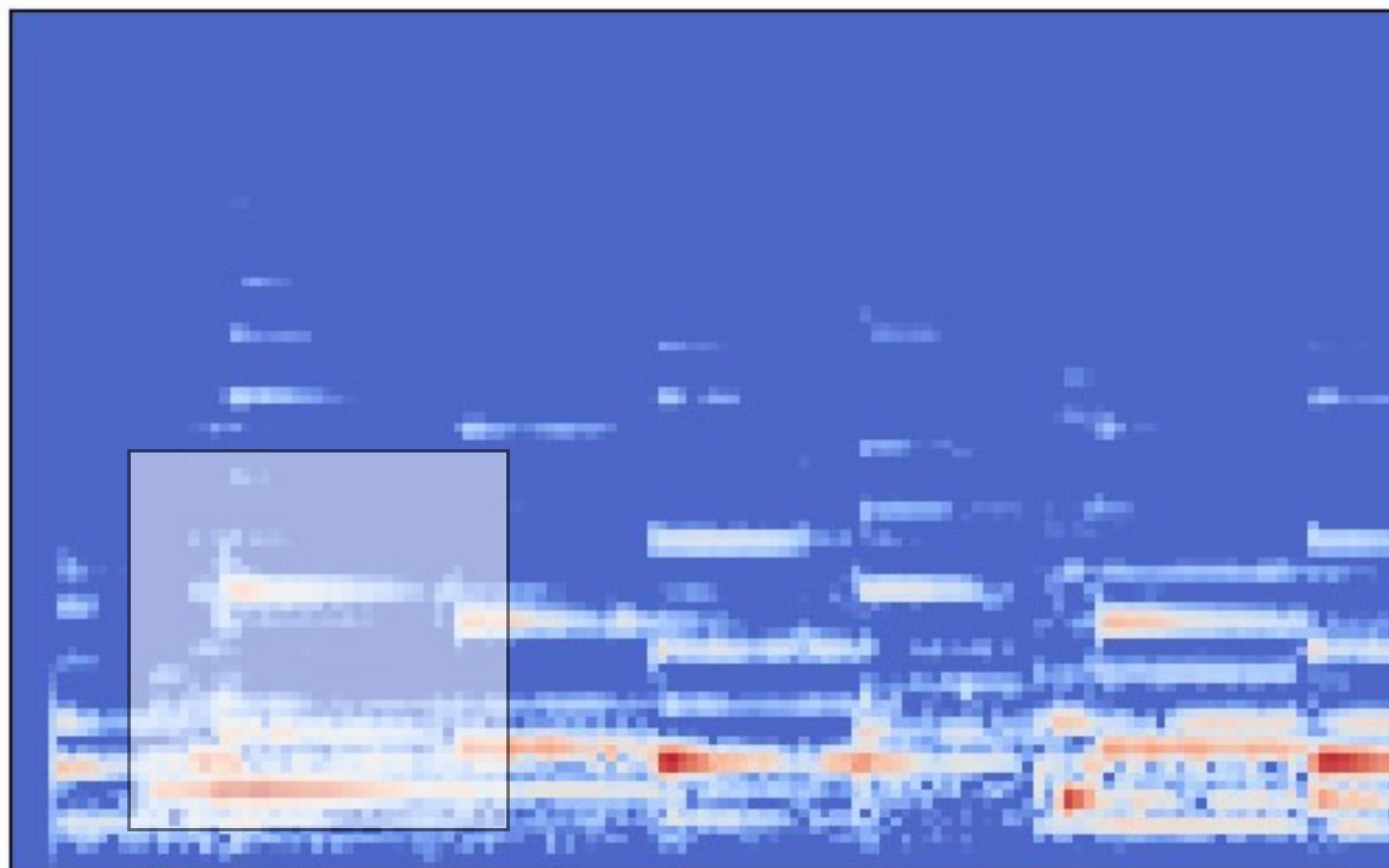
Convolutional Neural Networks



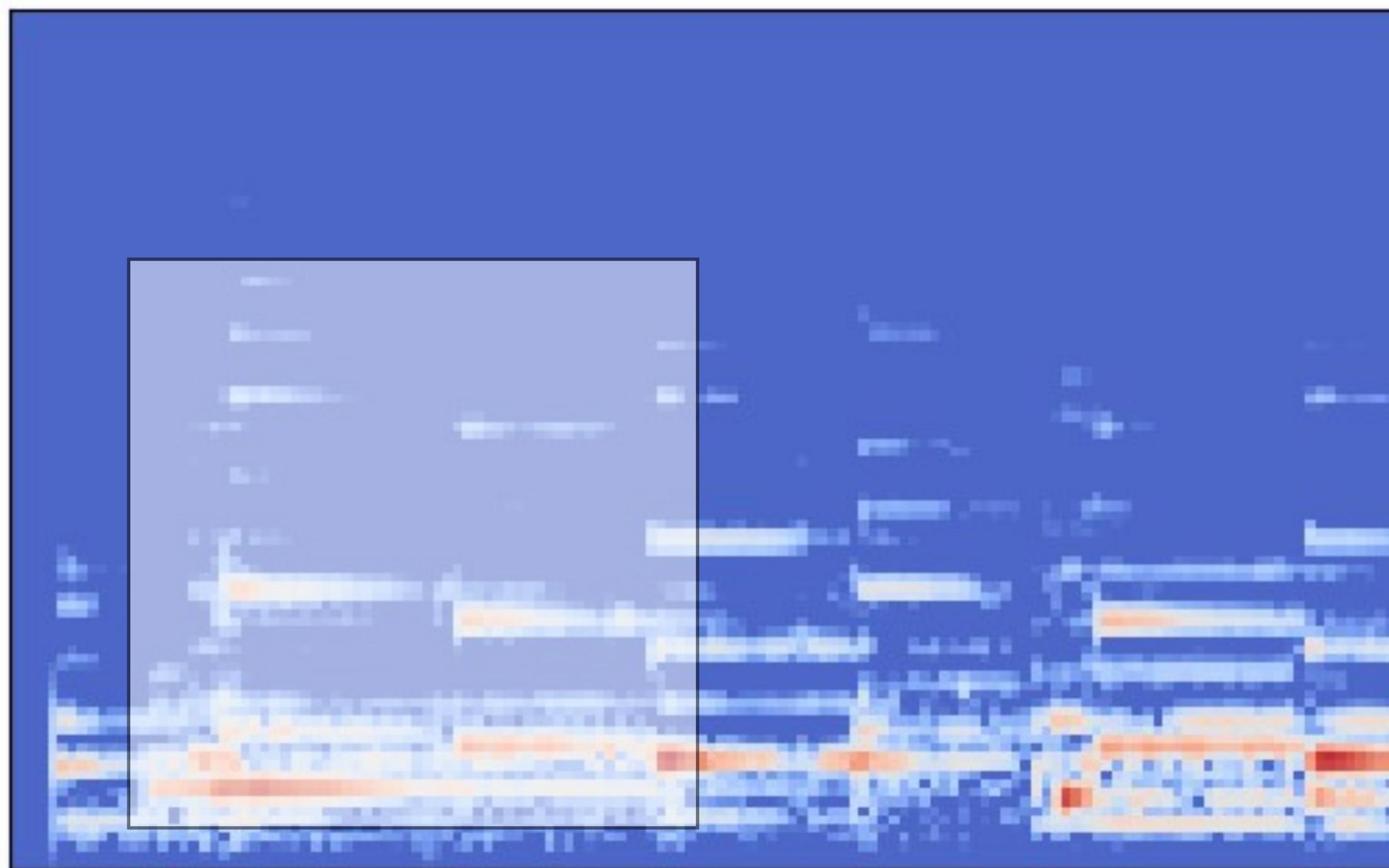
Convolutional Neural Networks



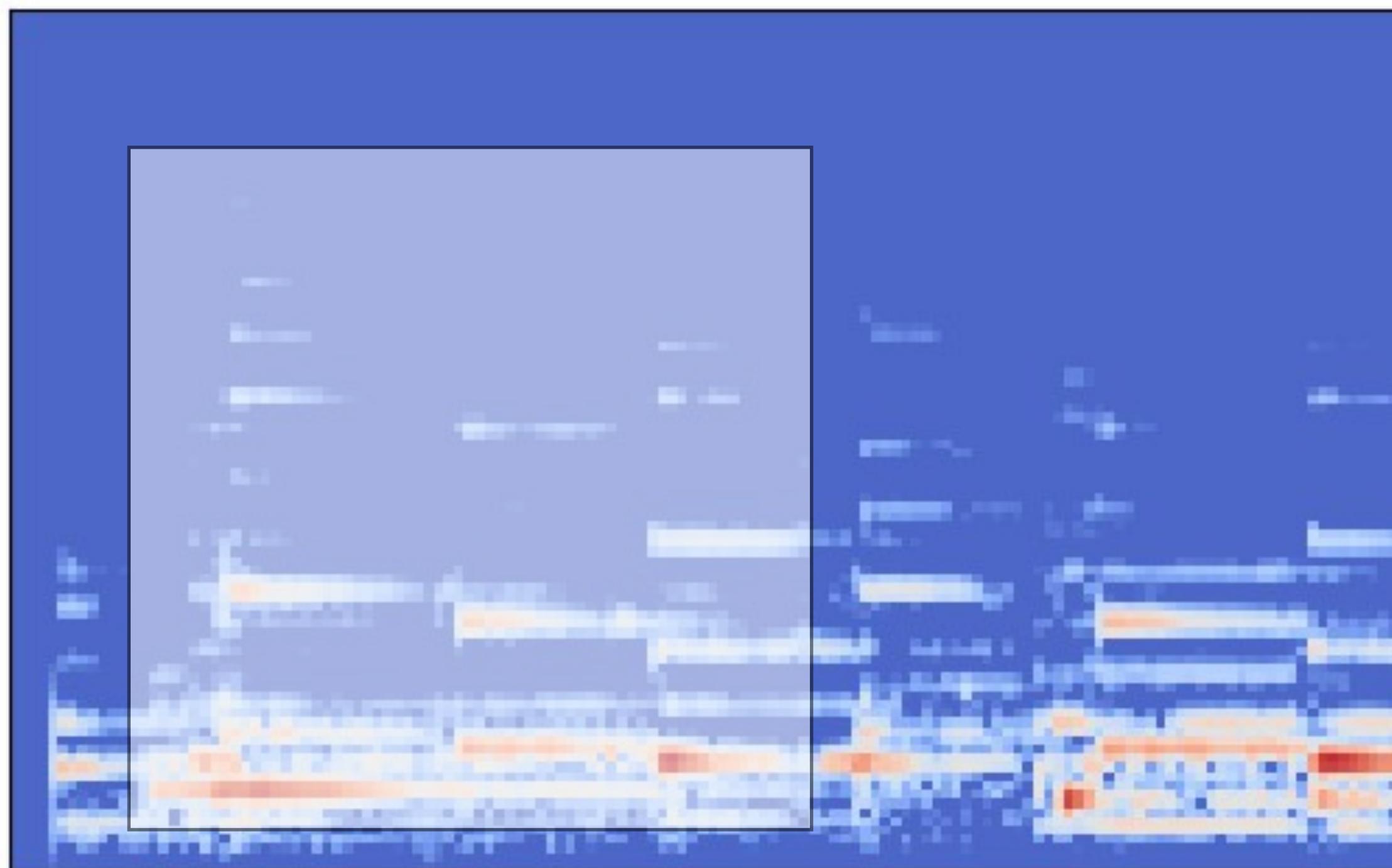
Convolutional Neural Networks



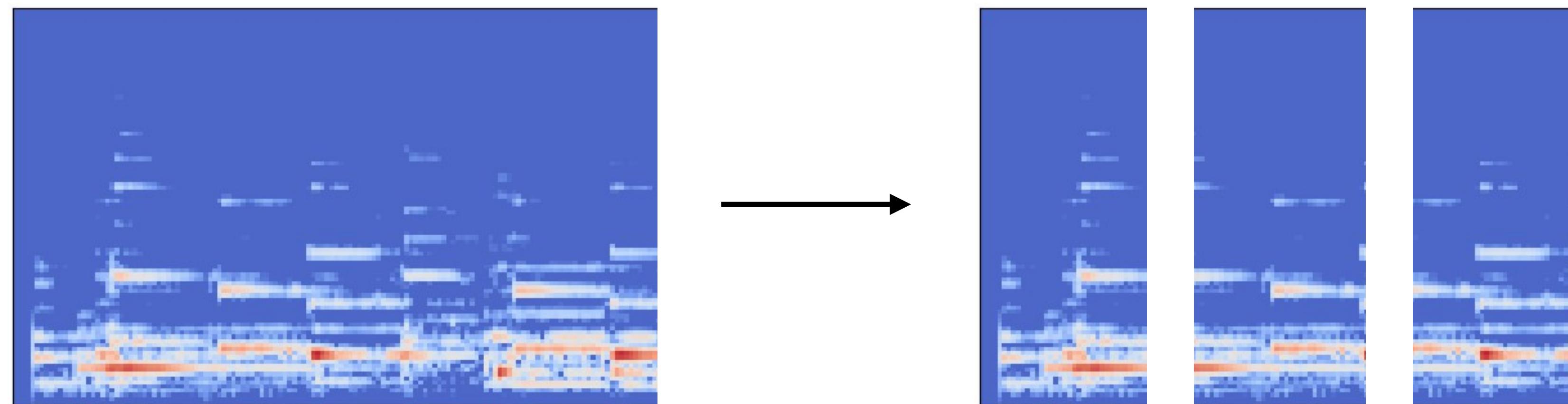
Convolutional Neural Networks



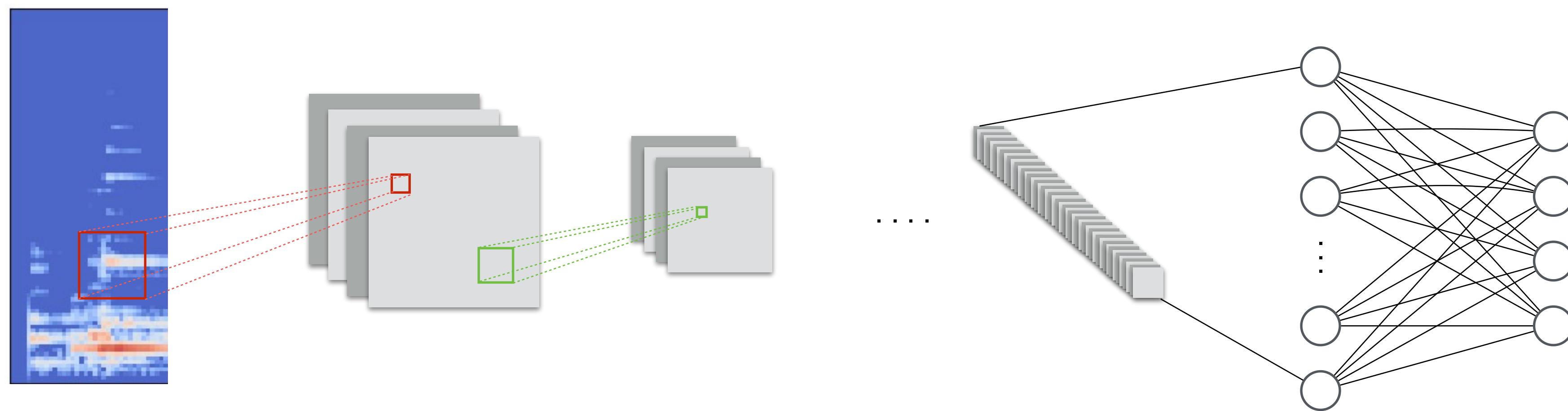
Convolutional Neural Networks



Bag-of-chunk



Bag-of-chunk



Bag-of-chunk



Rock



Guitar

Male vocal
Female vocal

80s
90s



Rock

Guitar
Male vocal

80s
90s

Guitar
Male vocal

80s
90s



Rock

Guitar

Male vocal

Female vocal

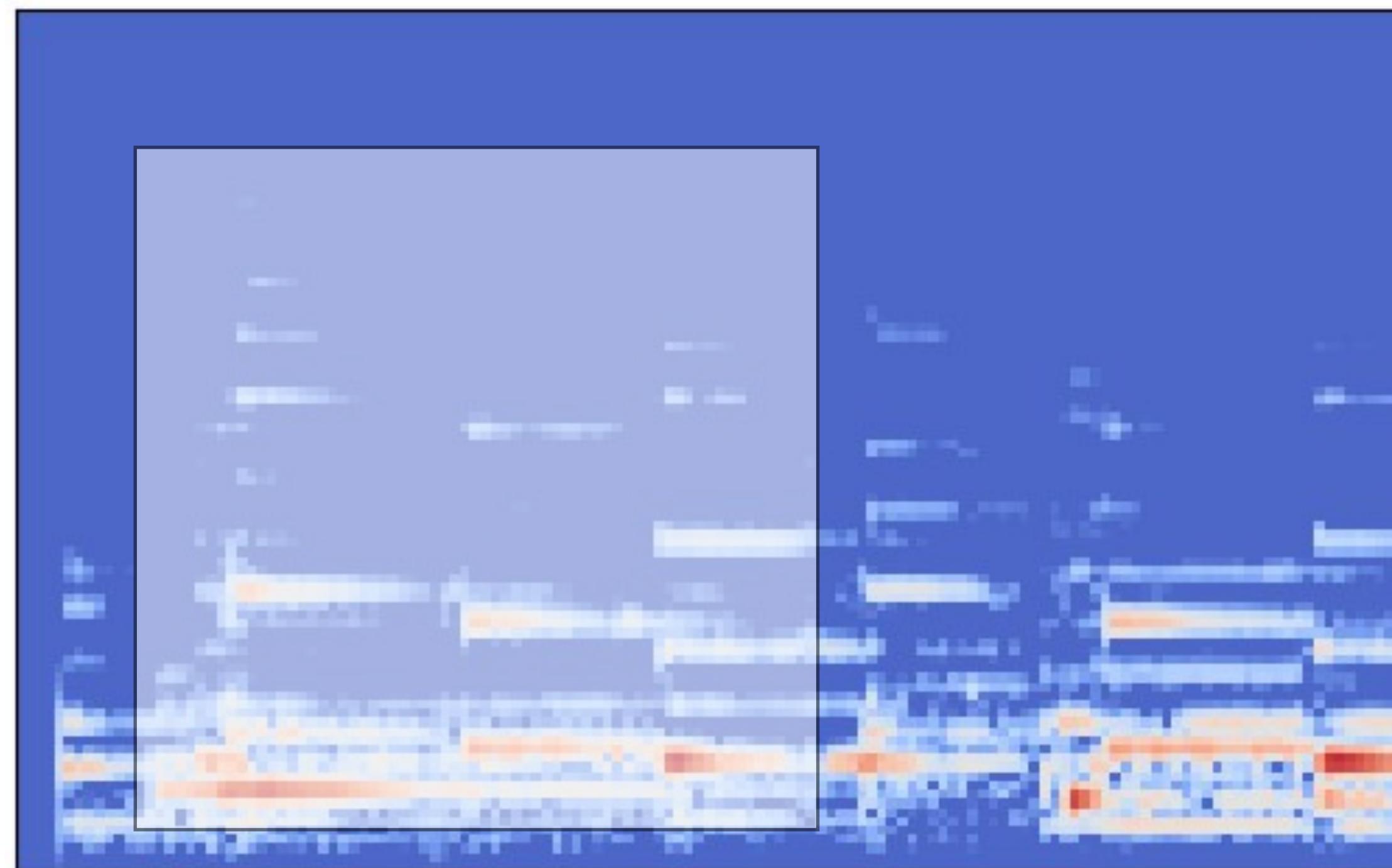
80s

90s

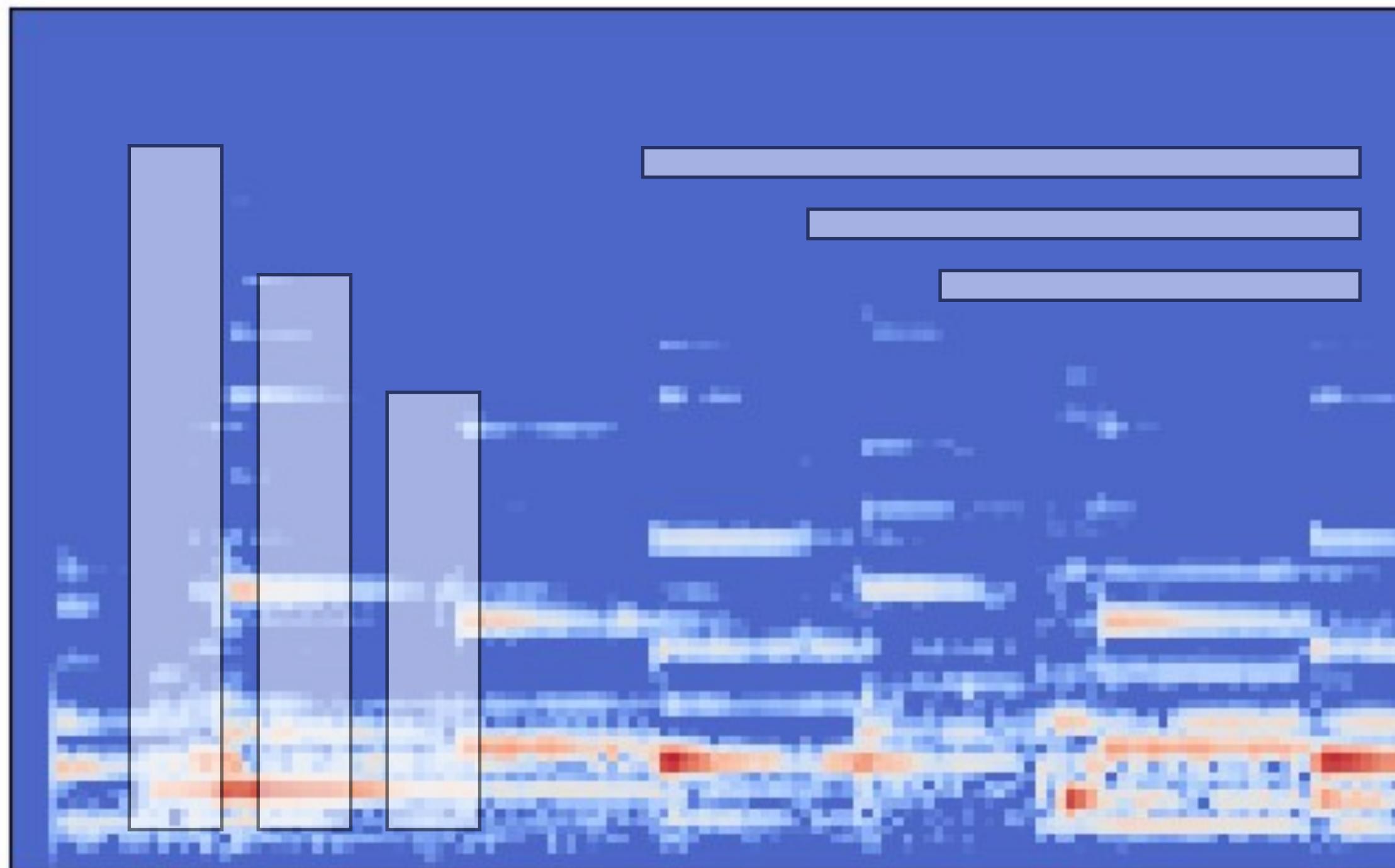
Bag-of-chunk

- Higher performance
- More robust
- A sort of data augmentation
- e.g., 1 million songs x 10 chunks = 10 million chunks

Let's put some domain knowledge



Let's put some domain knowledge

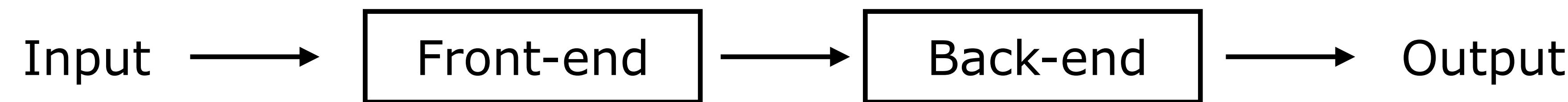


Musicnn (Pons et al., 2018)

Let's put some domain knowledge

- More efficient
- Good performance with a limited size of data
- However, with sufficient amount of data and computing power,
assumption-free models still work better

Music is sequential!



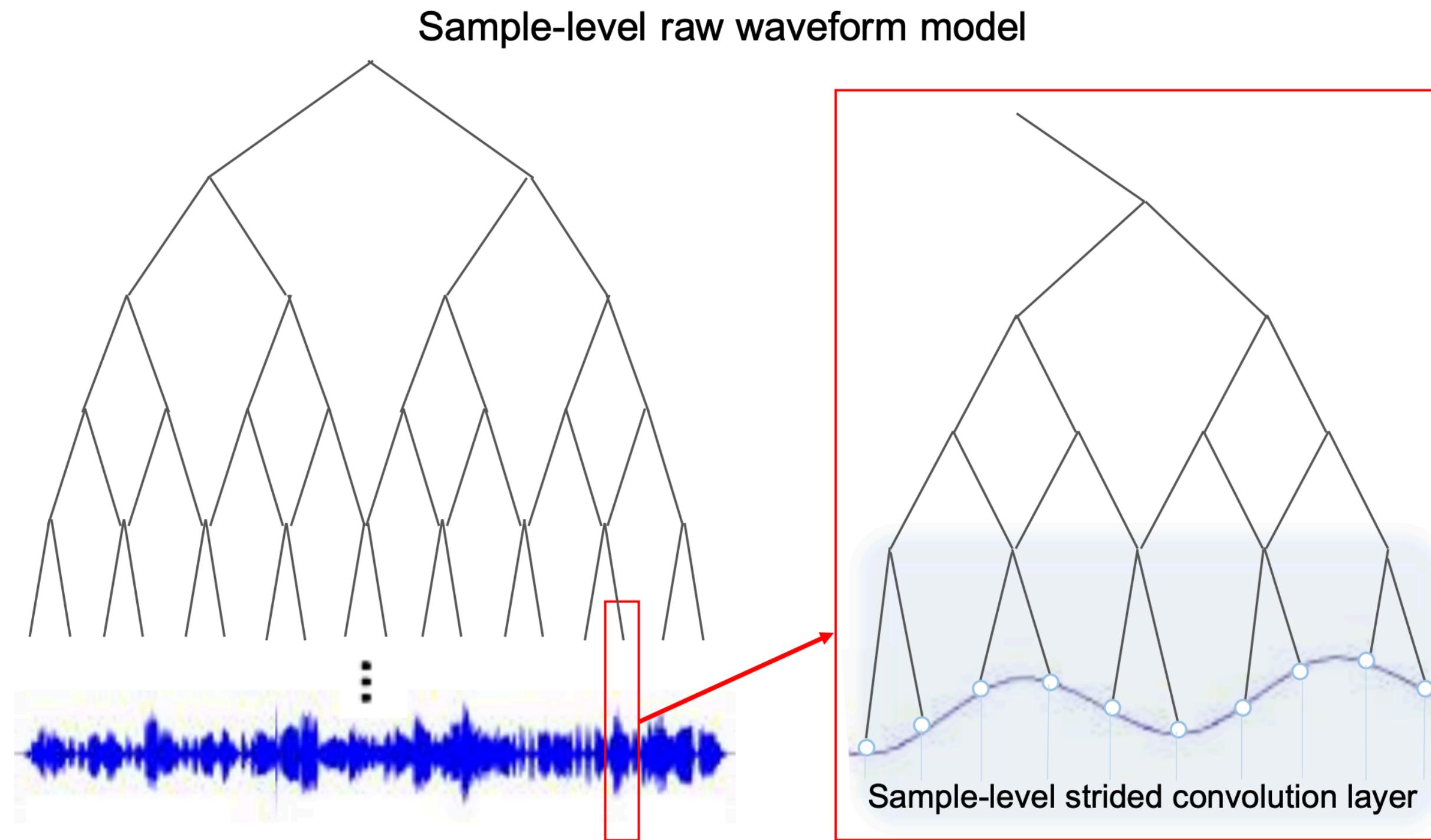
Music is sequential!

- Front-end: capture acoustic features with CNN
- Back-end: temporally summarize extracted features
- CNN front-end + RNN back-end (Choi et al. 2017)
- CNN front-end + self-attention back-end (Won et al. 2019)

Choi, Keunwoo, et al. "Convolutional recurrent neural networks for music classification." (2017)

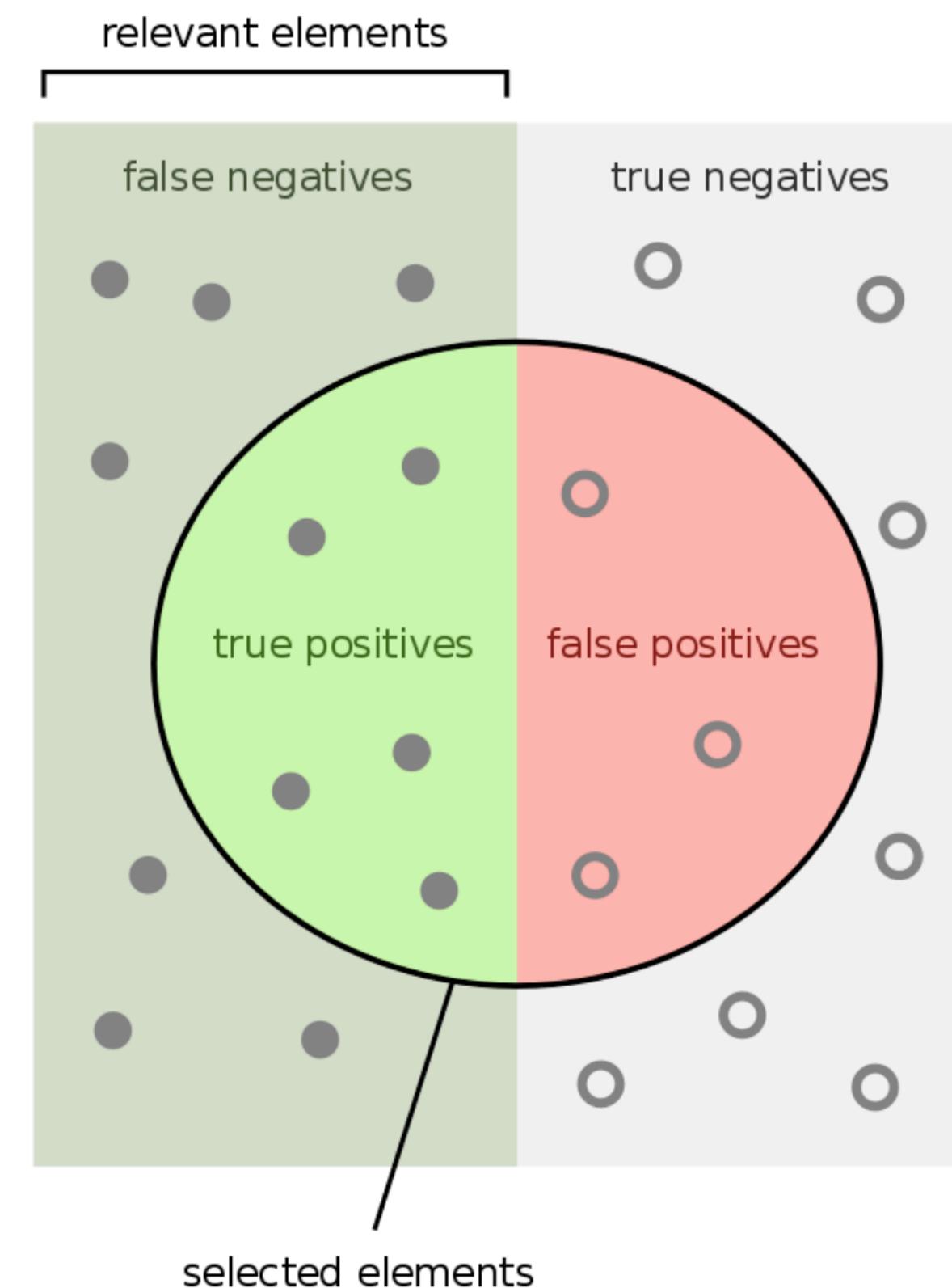
Won, Minz, Sanghyuk Chun, and Xavier Serra. "Toward interpretable music tagging with self-attention." (2019).

Let's make it end-to-end



Sample-level CNN (Lee et al., 2017)

Evaluation



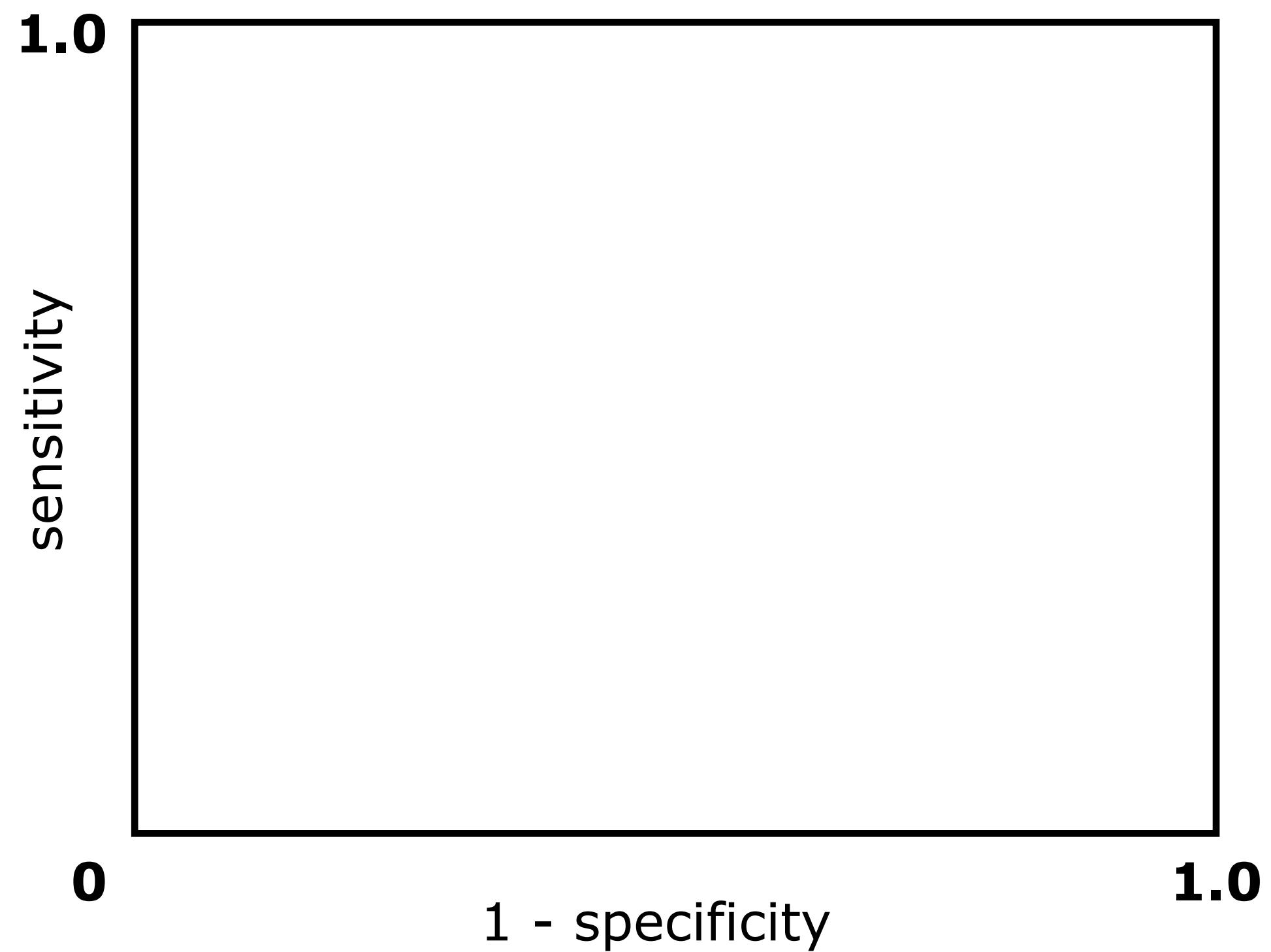
How many relevant items are selected?
e.g. How many sick people are correctly identified as having the condition.

Sensitivity = $\frac{\text{true positives}}{\text{relevant elements}}$

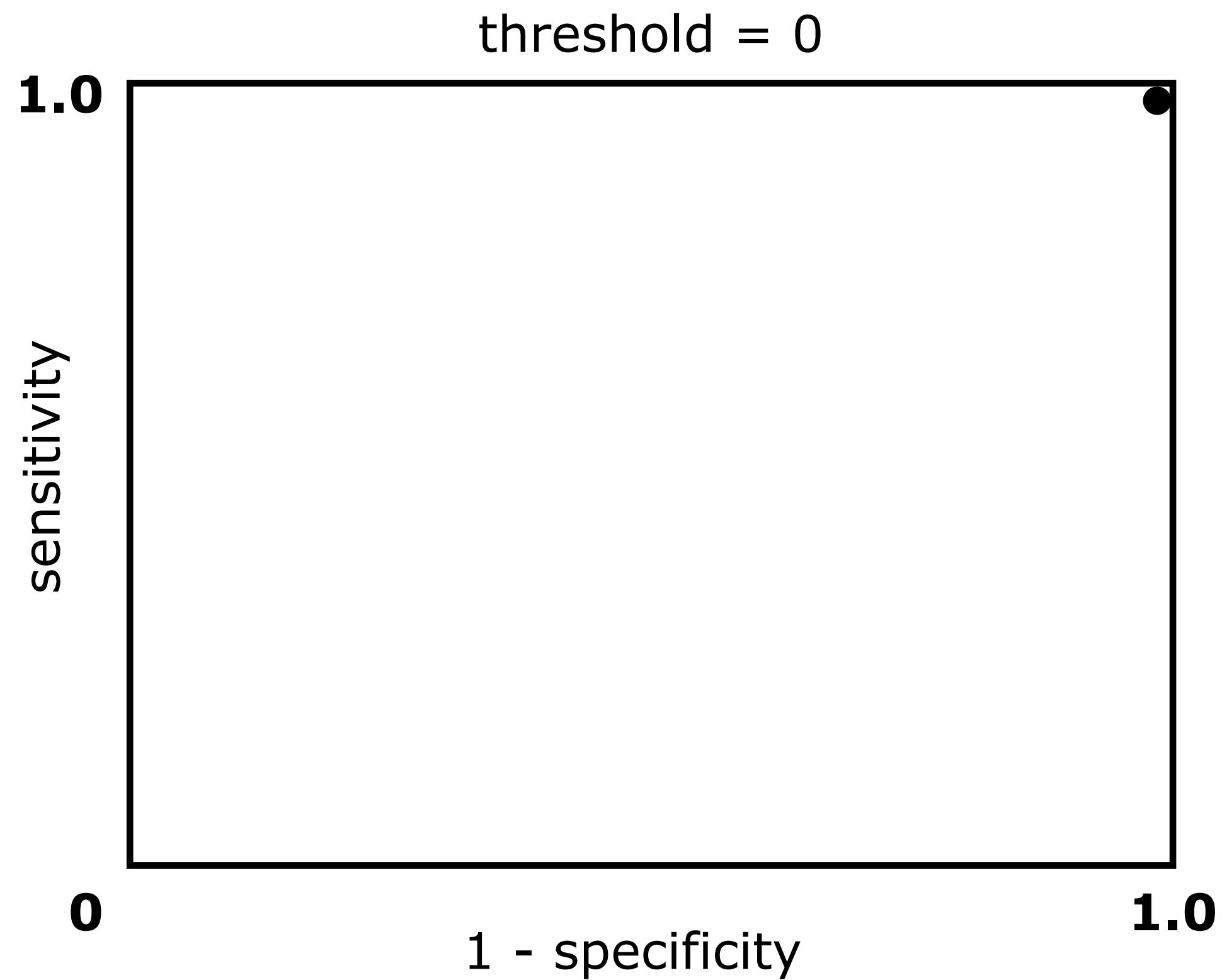
How many negative selected elements are truly negative?
e.g. How many healthy people are identified as not having the condition.

Specificity = $\frac{\text{true negatives}}{\text{non-relevant elements}}$

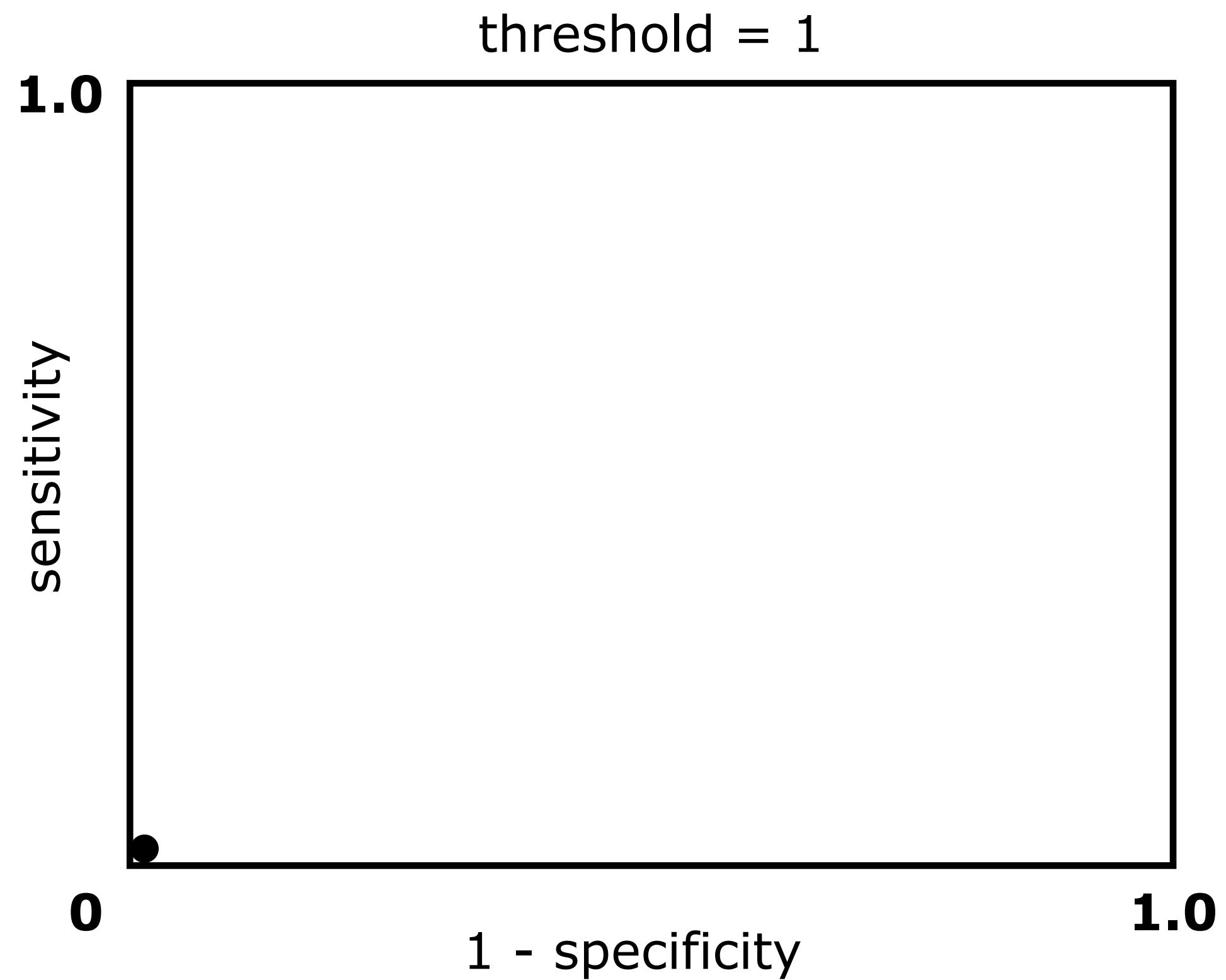
Evaluation



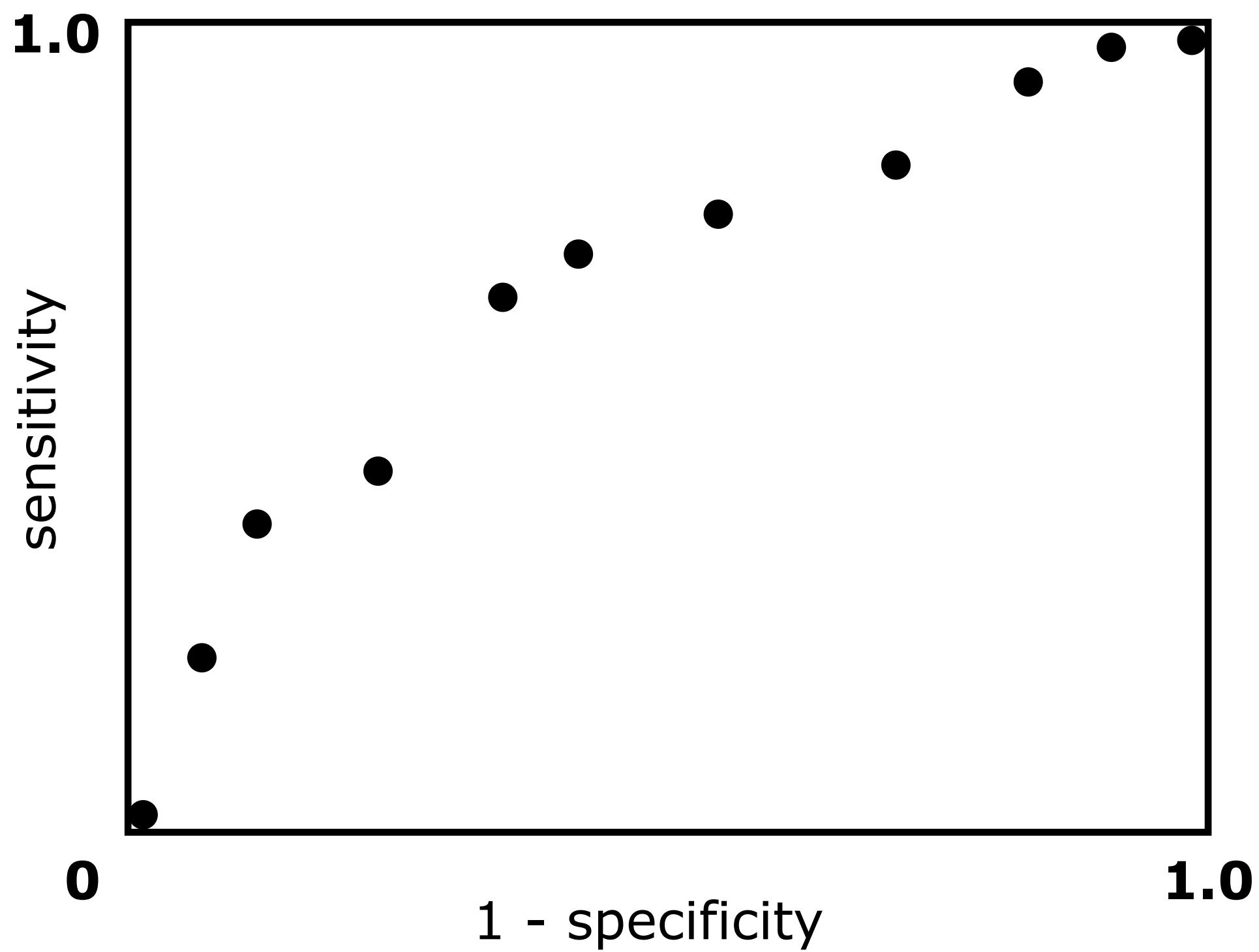
Evaluation



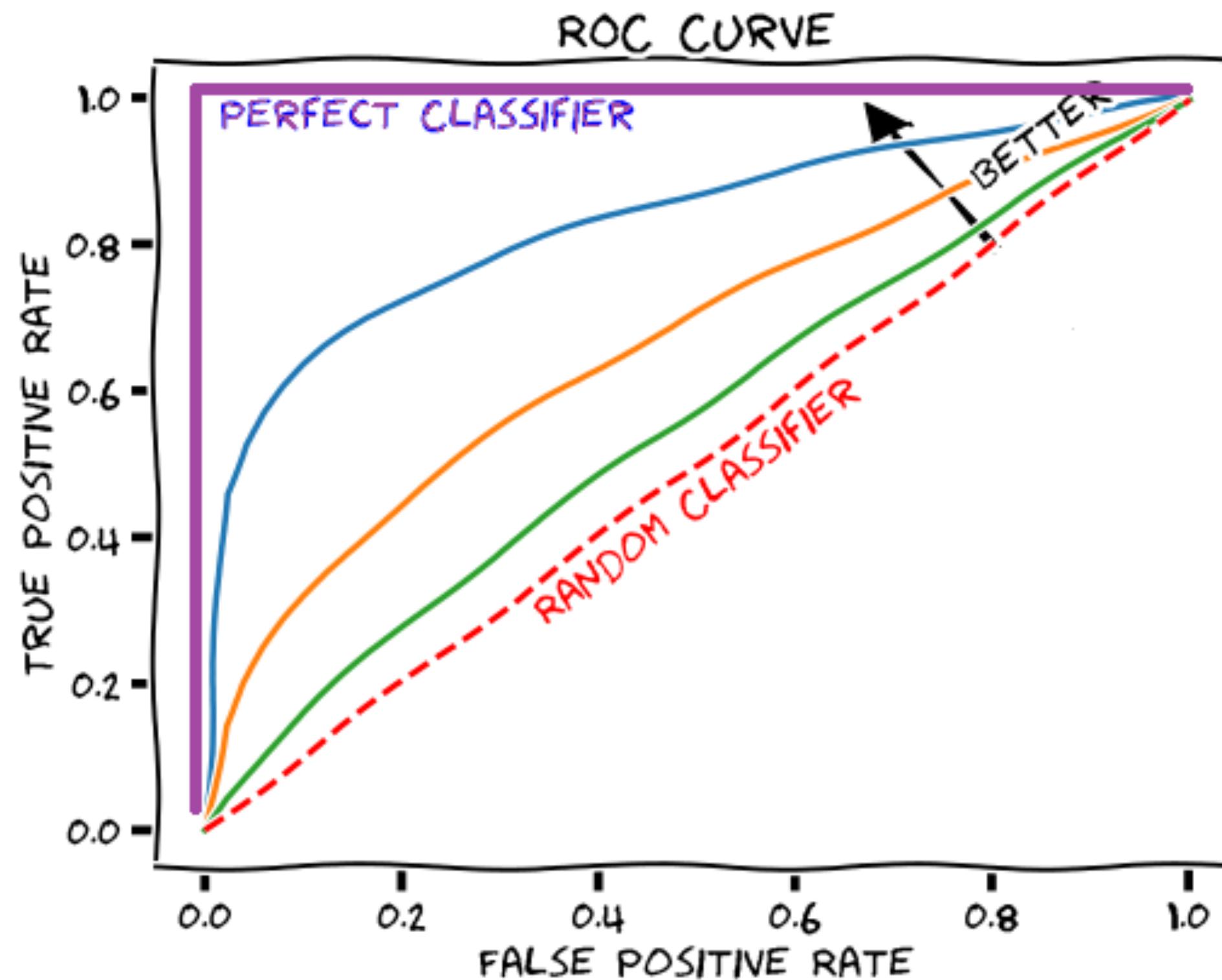
Evaluation



Evaluation

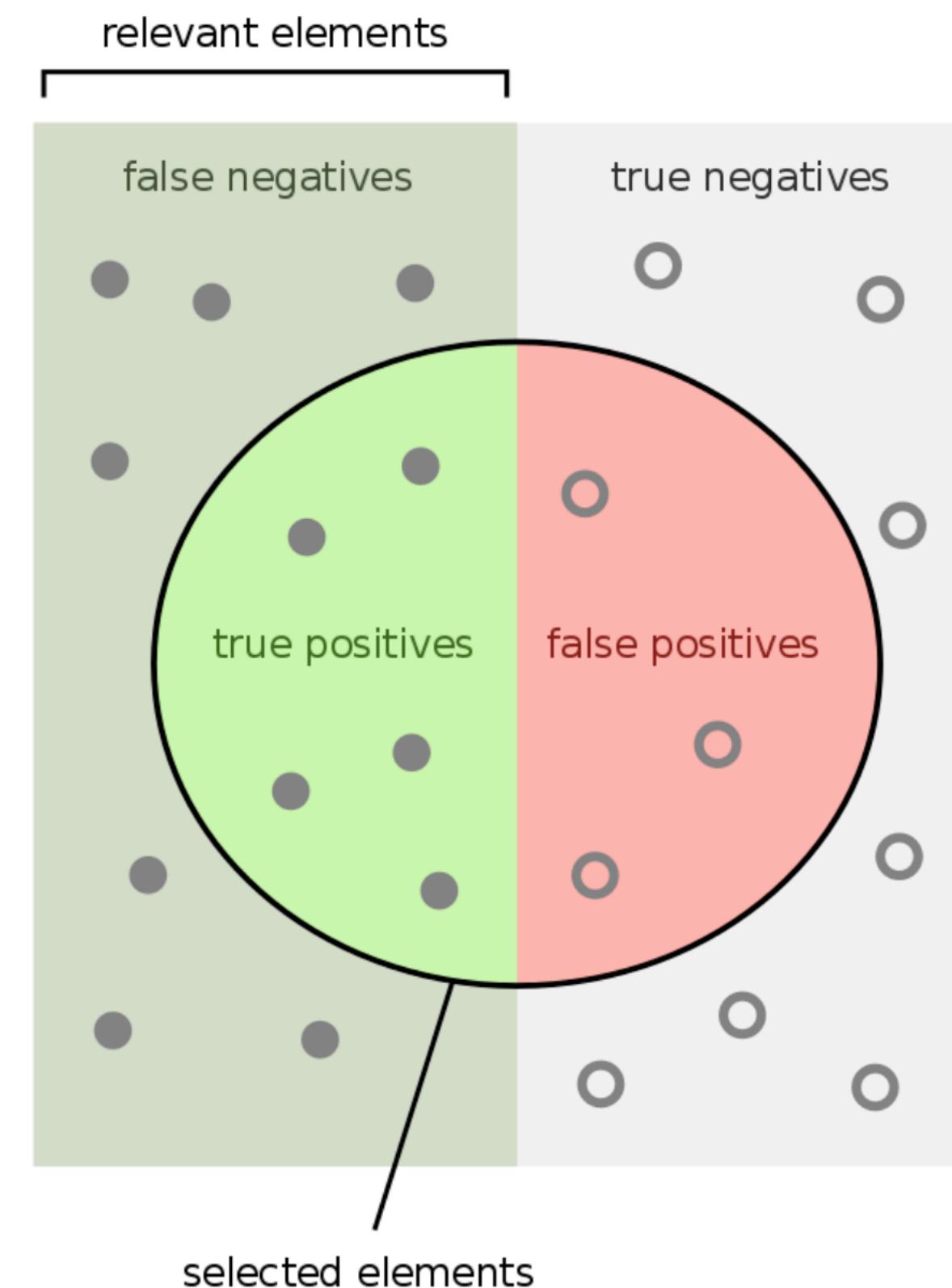


Evaluation



Receiver operating characteristic curve

Evaluation



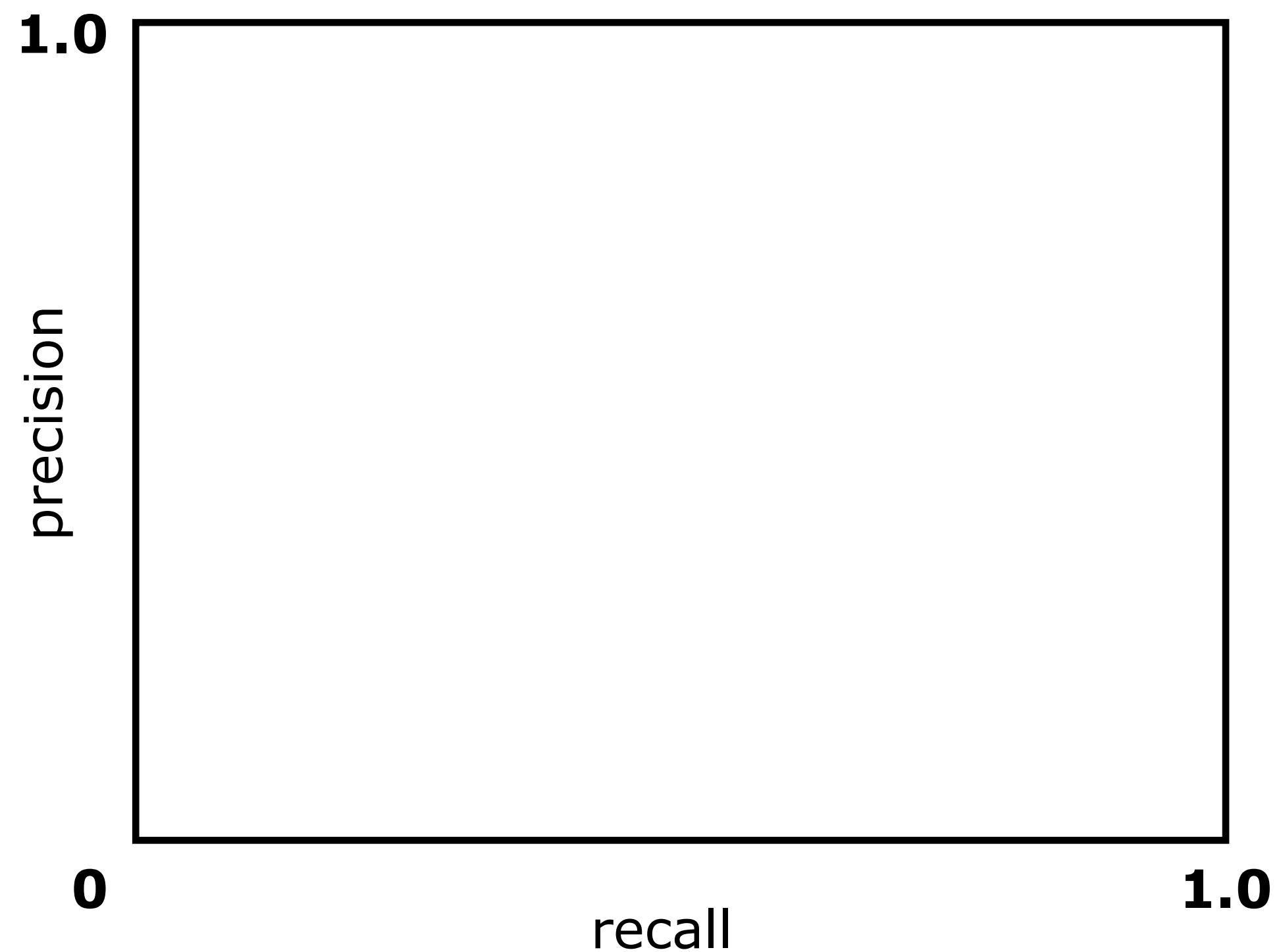
How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{selected elements}}$$

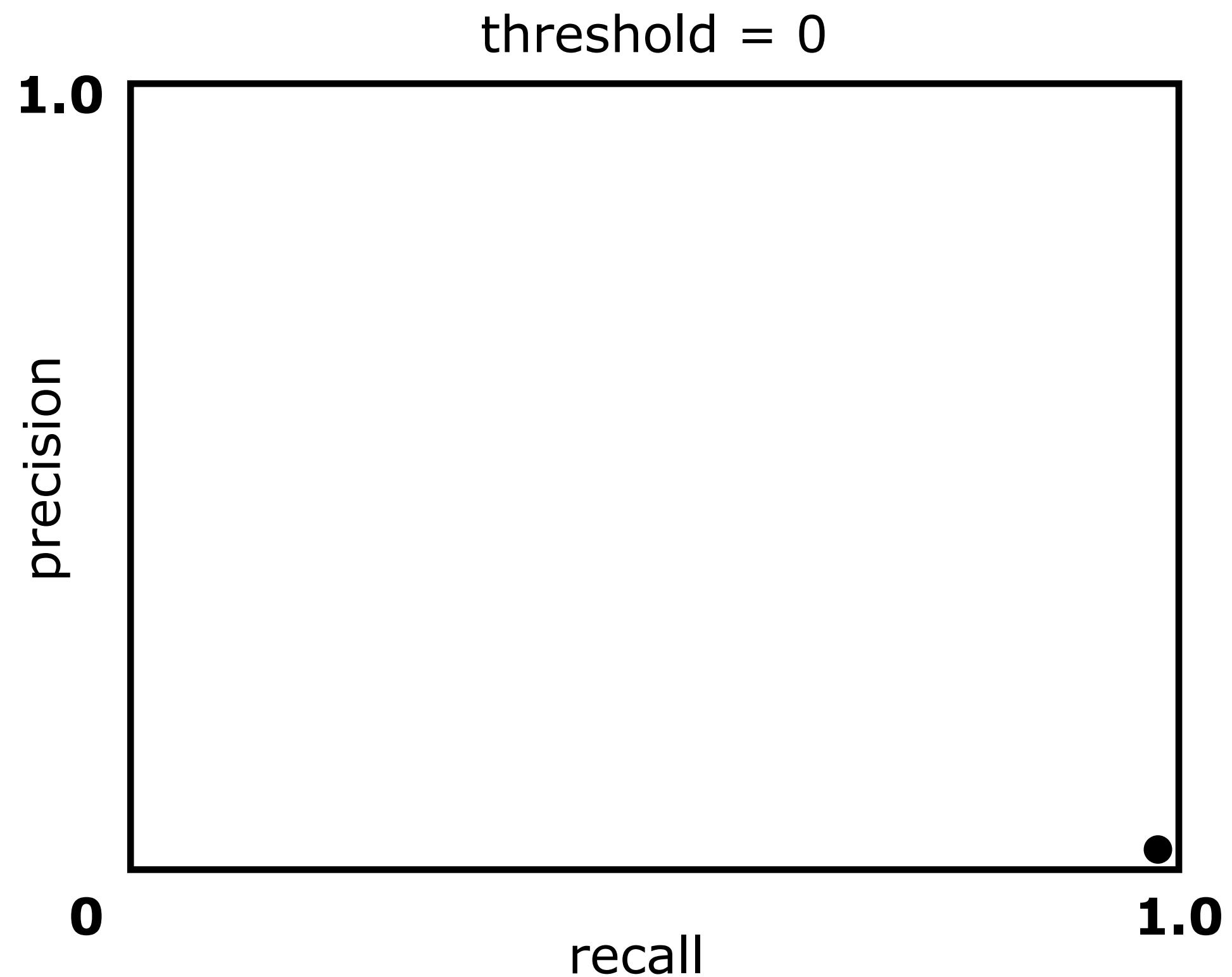
How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{relevant elements}}$$

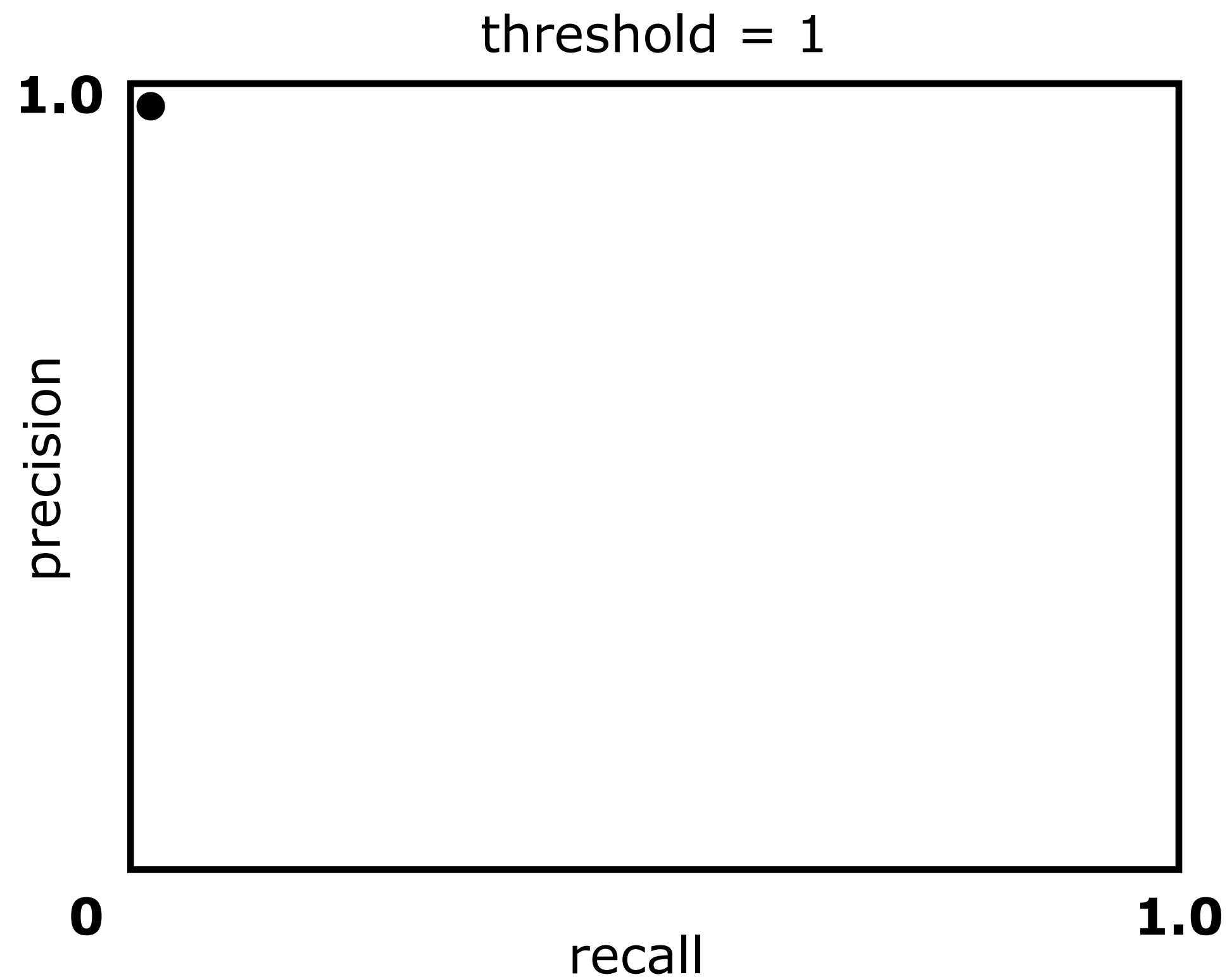
Evaluation



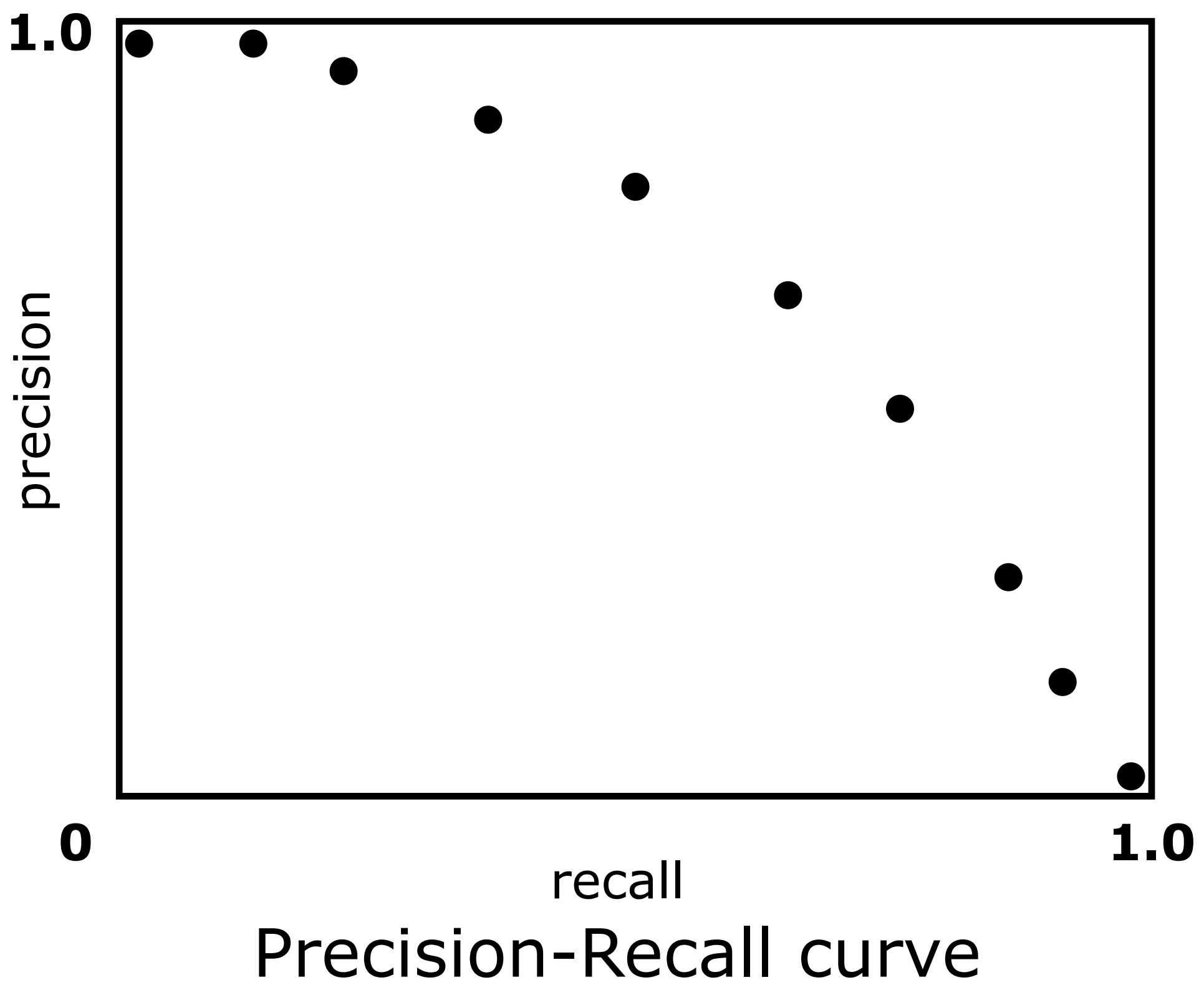
Evaluation



Evaluation



Evaluation



Benchmark

Methods	MTAT		MSD		MTG-Jamendo	
	ROC-AUC	PR-AUC	ROC-AUC	PR-AUC	ROC-AUC	PR-AUC
FCN [1]	0.9020	0.4367	0.8744	0.2970	0.8255	0.2801
Musicnn [2]	0.9107	0.4511	0.8803	0.2983	0.8226	0.2713
Sample-level [3]	0.9065	0.4428	0.8789	0.2959	0.8208	0.2742
Sample-level + SE [4]	0.9097	0.4540	0.8838	0.3109	0.8233	0.2784
CRNN [6]	0.8729	0.3599	0.8499	0.2469	0.7978	0.2358
Self-attention [7]	0.9096	0.4485	0.8810	0.3103	0.8261	0.2883
BoC-CNN + Res	0.9129	0.4614	0.8898	0.3280	0.8316	0.2951

Contents

- Introduction
- Music tagging models

● Transfer learning

- Limitations
- Lab

Transfer Learning



Rock

Guitar

Male vocal

Female vocal

80s

90s

Transfer Learning

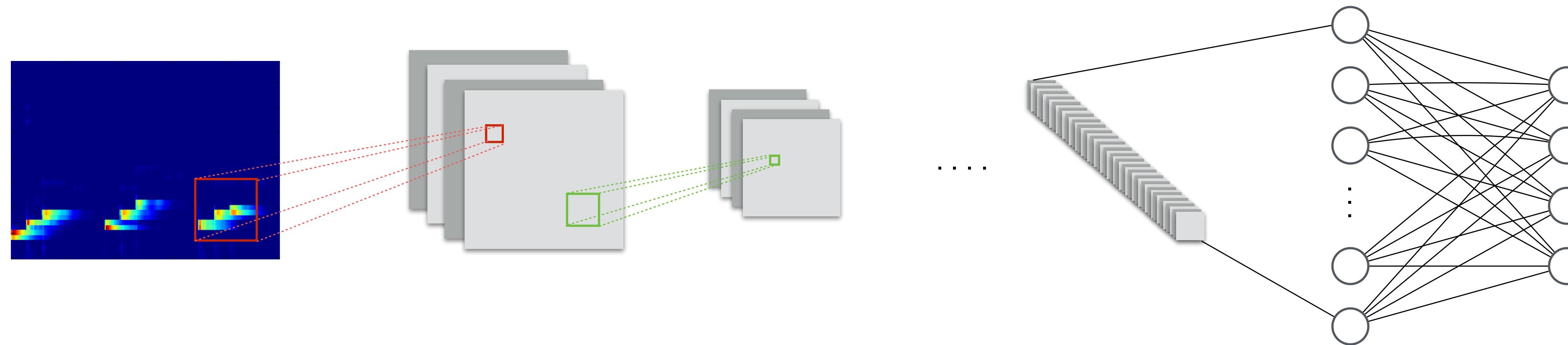


Transfer Learning

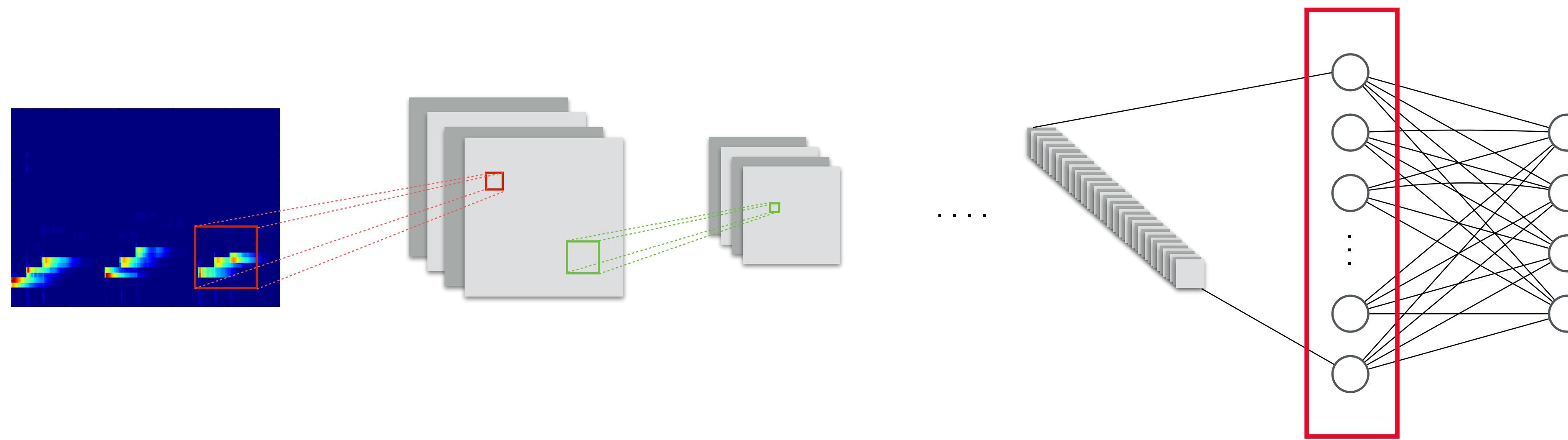


Genre classification → Music similarity

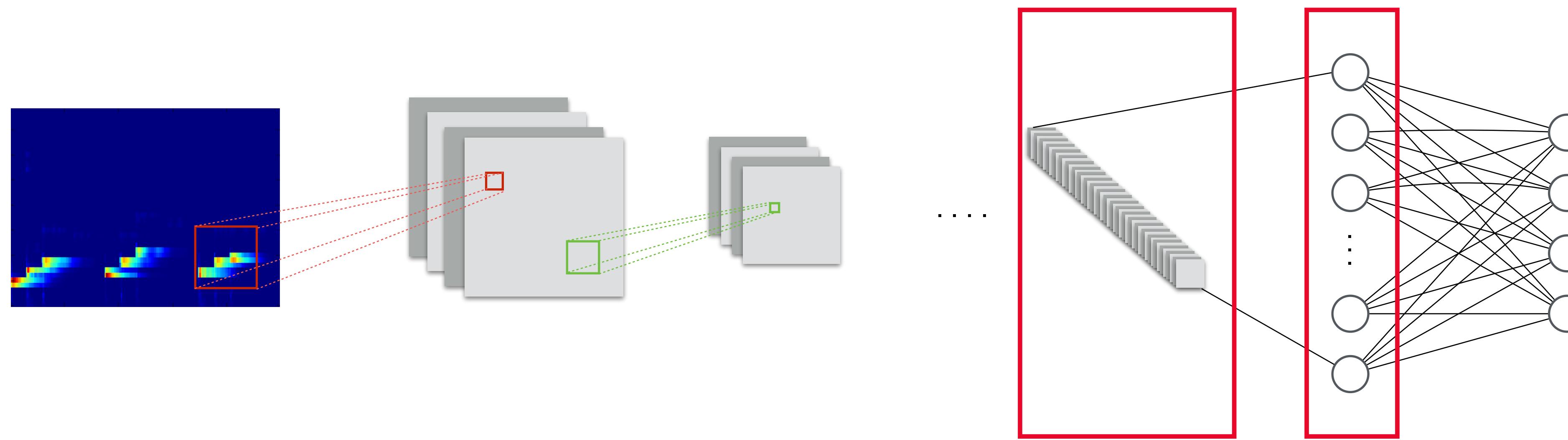
Transfer Learning



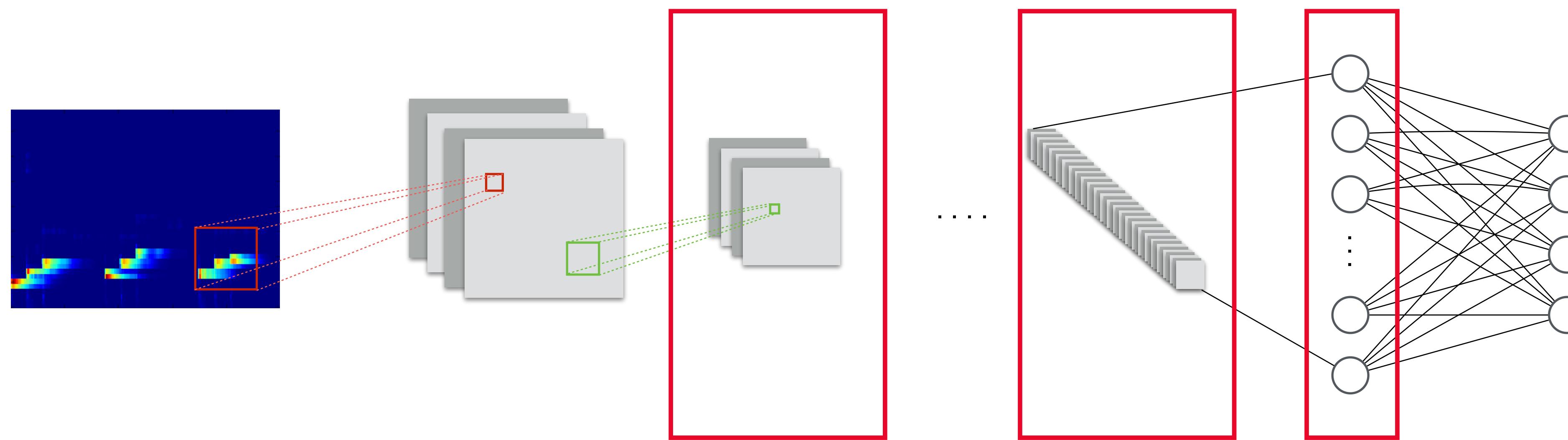
Transfer Learning



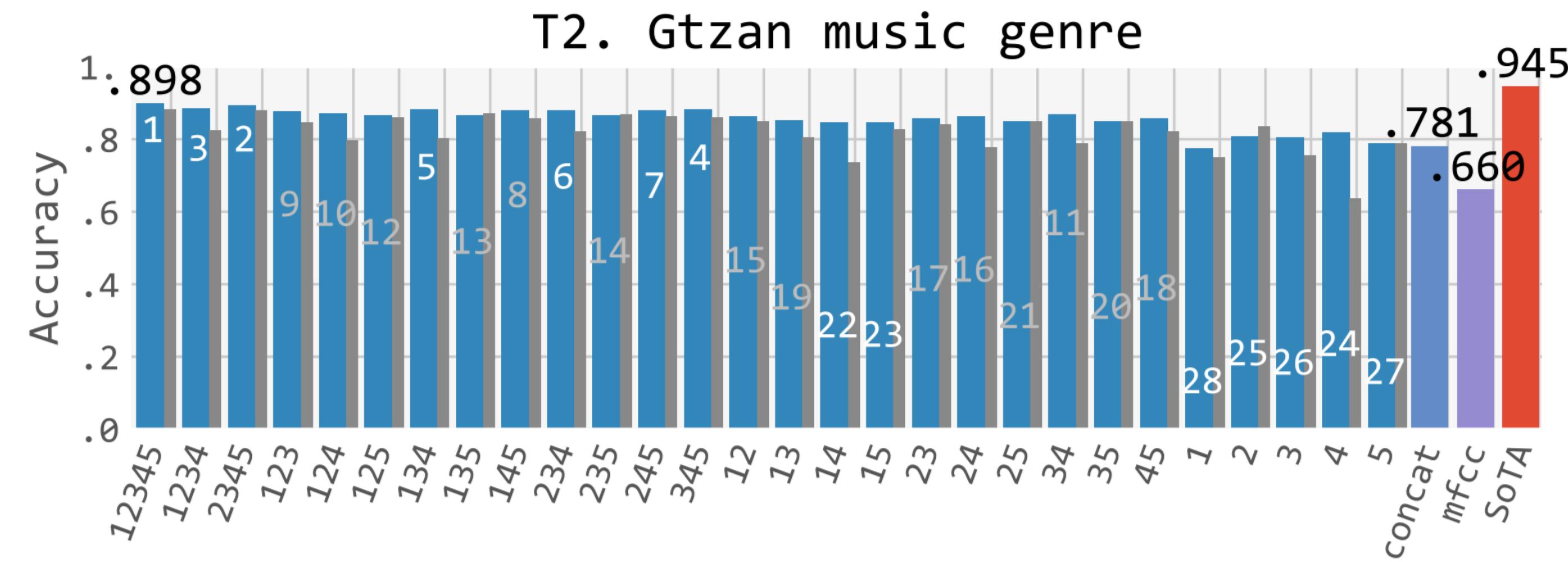
Transfer Learning



Transfer Learning



Transfer Learning



Contents

- Introduction
- Music tagging models
- Transfer learning

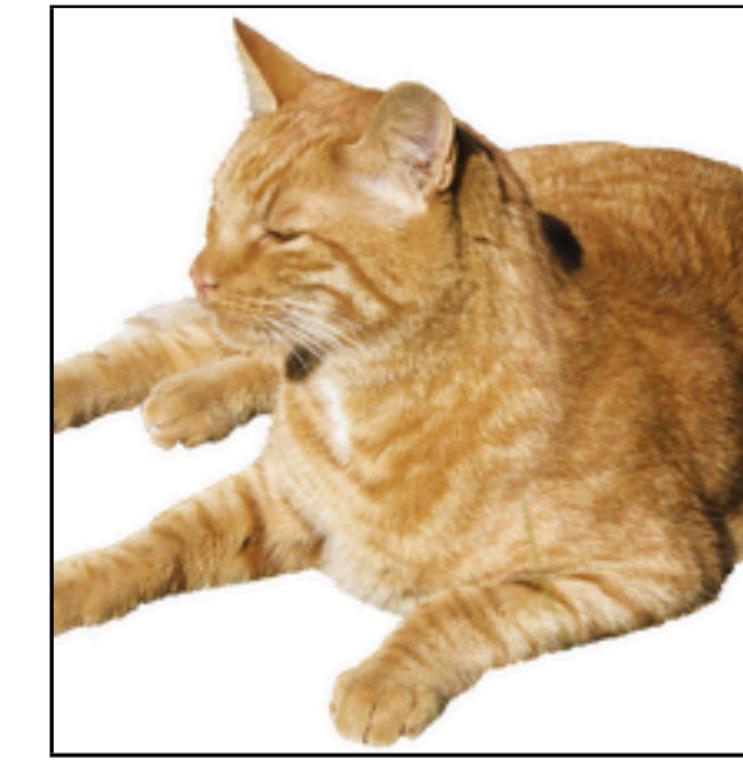
● **Limitations**

- Lab

Limitations



(a) Texture image
81.4% **Indian elephant**
10.3% indri
8.2% black swan



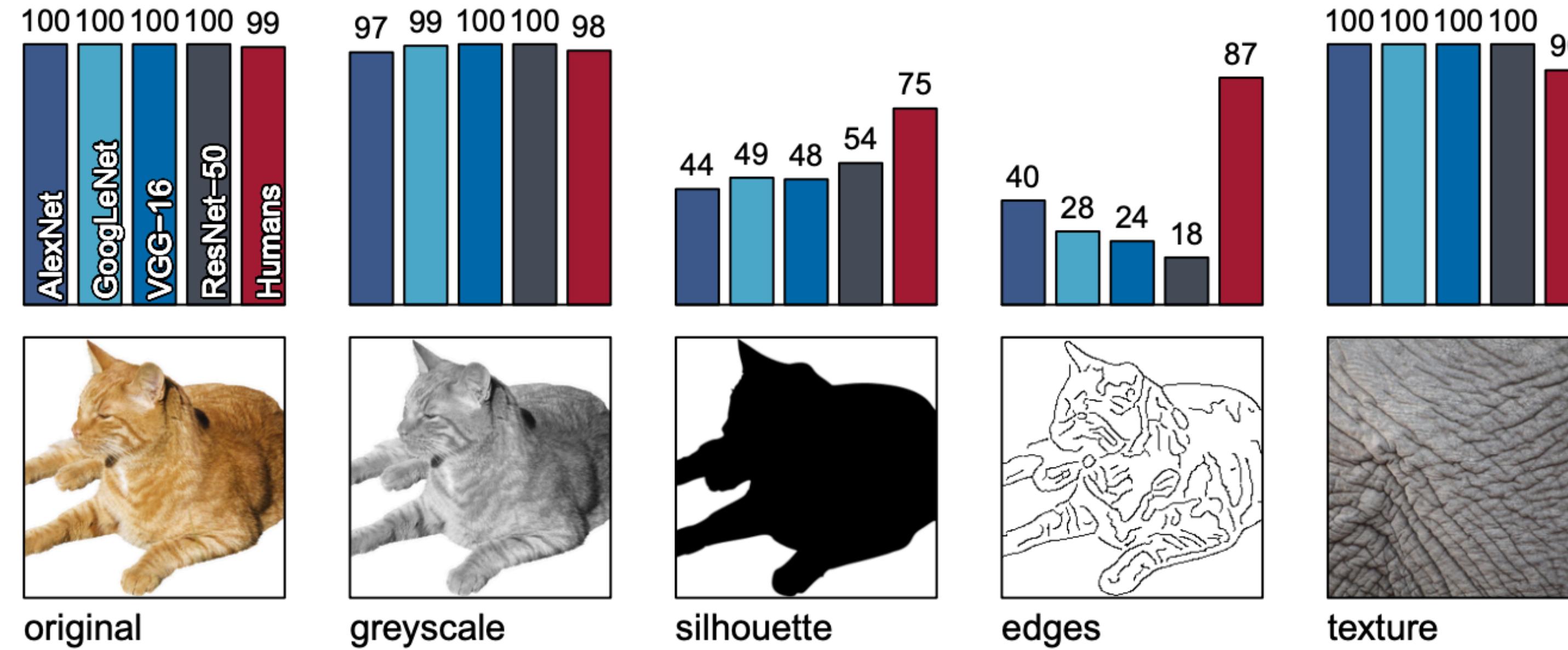
(b) Content image
71.1% **tabby cat**
17.3% grey fox
3.3% Siamese cat



(c) Texture-shape cue conflict
63.9% **Indian elephant**
26.4% indri
9.6% black swan

Inherent texture bias in CNN

Limitations



Inherent texture bias in CNN

Limitations

- Data augmentation can be useful
- e.g., Time stretch, Pitch shift, Dynamic range compression, Add noise
- But still it cannot overcome inherent texture bias. For example, tasks like cover song identification rely more on chroma features.

Limitations

common
images



turtle in grass
baseline: turtle
ours: turtle

uncommon
images



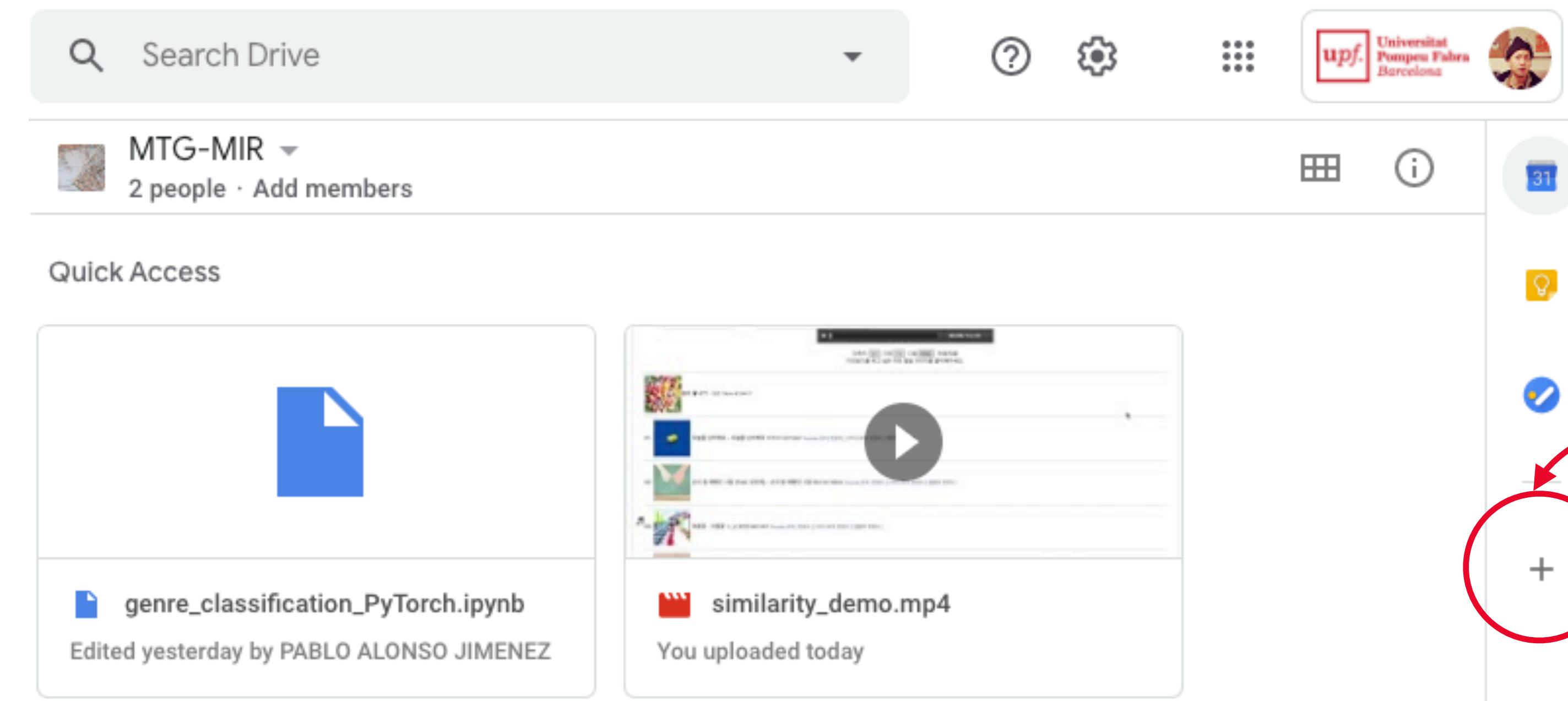
crab in grass
baseline: turtle
ours: crab

Bias in dataset

Contents

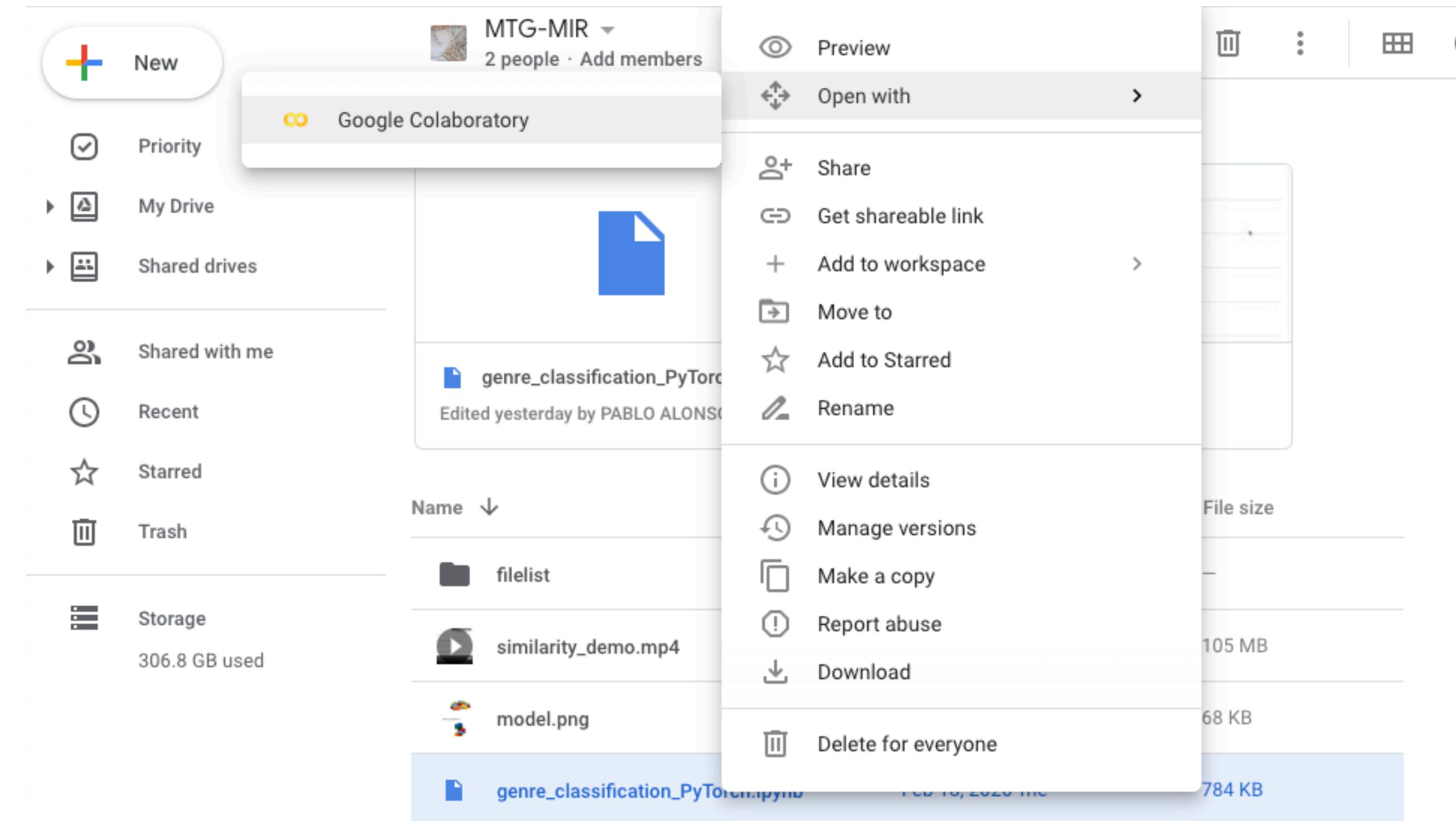
- Introduction
- Music tagging models
- Transfer learning
- Limitations
- **Lab**

Lab



Click!
Add Colaboratory

Lab



Lab

The screenshot shows a Jupyter Notebook interface with a single code cell containing the following Python code:

```
force_remount=True)
root_dir = '/content/gdrive/Shared drives/MTGxMIP/'
```

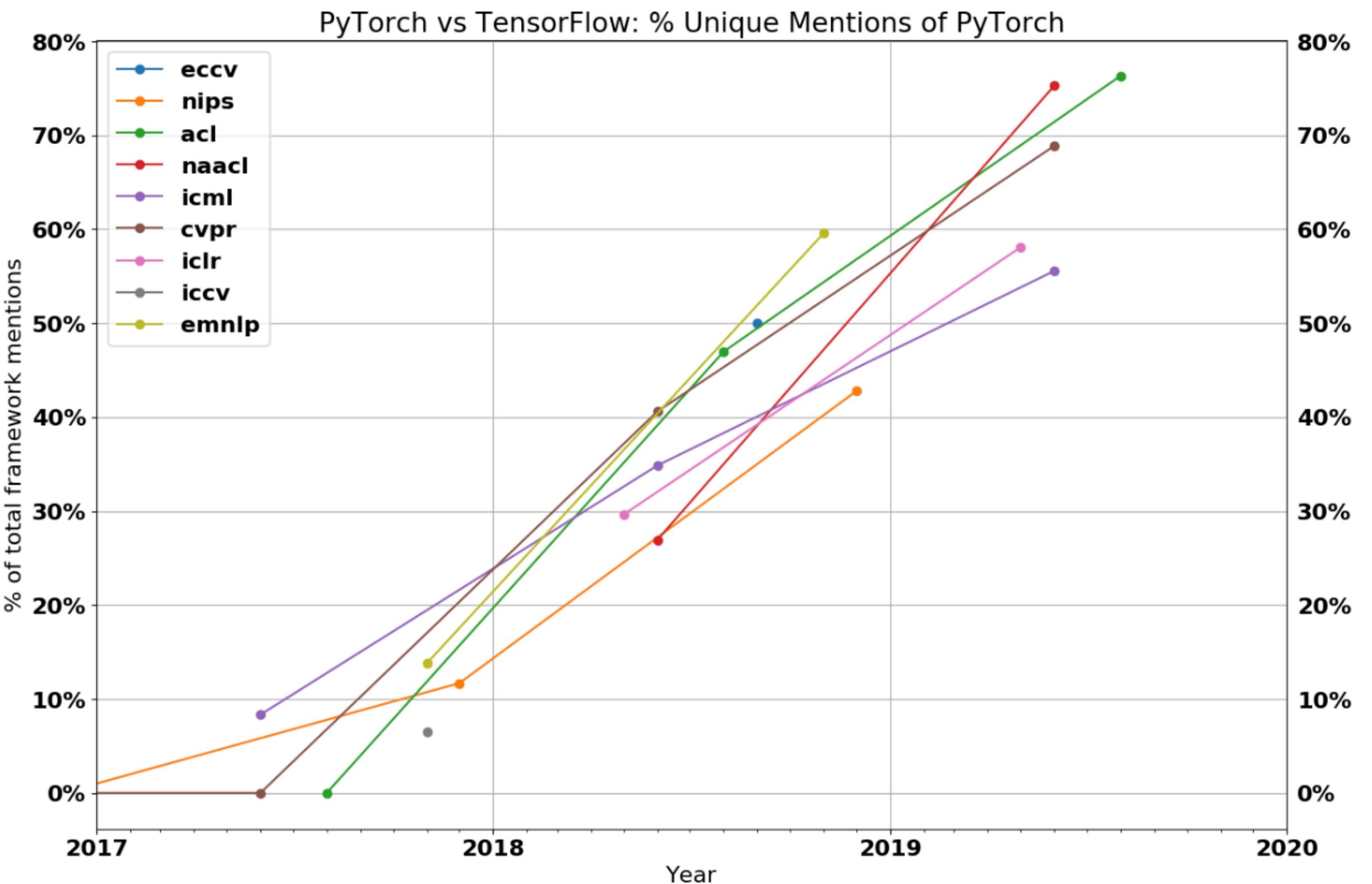
The output of the cell is displayed below it, showing the value `55.308)`. The notebook title is `genre_classification_PyTorch.ipynb`. The file menu is open, showing options like Save a copy in Drive..., Save, and Print.

PyTorch

- Easy to install
- Easy to prototype
- Easy to debug
- More pythonic
- Dynamic graph

PyTorch

- Easy to install
- Easy to prototype
- Easy to debug
- More pythonic
- Dynamic graph



PyTorch

```
class MyModel(nn.Module):
```

```
    def __init__(self):
```



```
        def forward(self, x):
```

