# Machine Listening for Music and Sound Analysis

# Lecture 3 – Music Information Retrieval I

Dr.-Ing. Jakob Abeßer

Fraunhofer IDMT

Jakob.abesser@idmt.fraunhofer.de

https://machinelistening.github.io

Fraunhofer

IDMT

# Overview

- Music Information Retrieval

- Music Tagging

- Music Similarity

- Tempo Estimation

# Music Information Retrieval
## Examples

■ Examples:

   ■ Musical Instrument

   AUD-1     AUD-2

   ■ Musical Genre / Tempo

   AUD-3     AUD-4

Fraunhofer
IDMT

# Music Information Retrieval
## Motivation

- Large music collections

- Mobile device apps / instruments

- Music industry shifts almost completely to online products & services
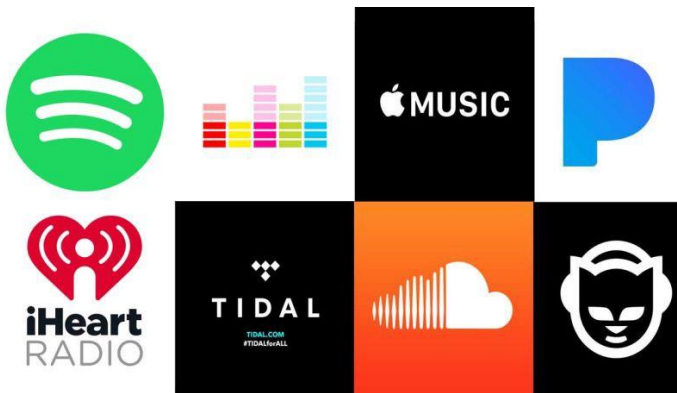
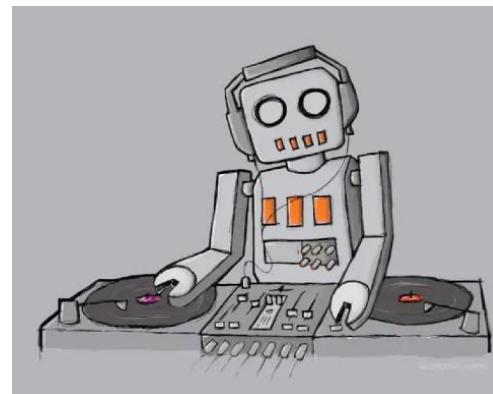- Growing market of music streaming services



Fig. 1



Fig. 2

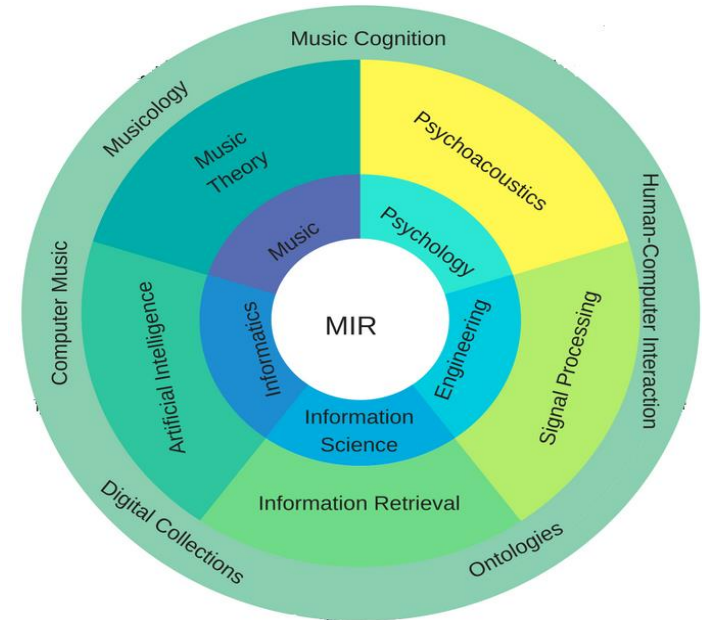# Music Information Retrieval
## Typical Research Tasks

- What's that song again? Who's singing that?

  - Audio identification

- I want to learn that song on my instrument!

  - Automatic music transcription

- What songs are similar? How to generate a playlist?

  - Audio similarity search

- How to organize my music? Which genre / style?

  - Audio classification

# Music Information Retrieval
## Research Landscape

- Interdisciplinary research community since 2000

- Conferences

  - ISMIR (International Society for Music Information Retrieval)

  - IEEE ICASSP, DAFx, AES, ICMC, SMC

- MIREX competition (Music Information Retrieval Evaluation eXchange)
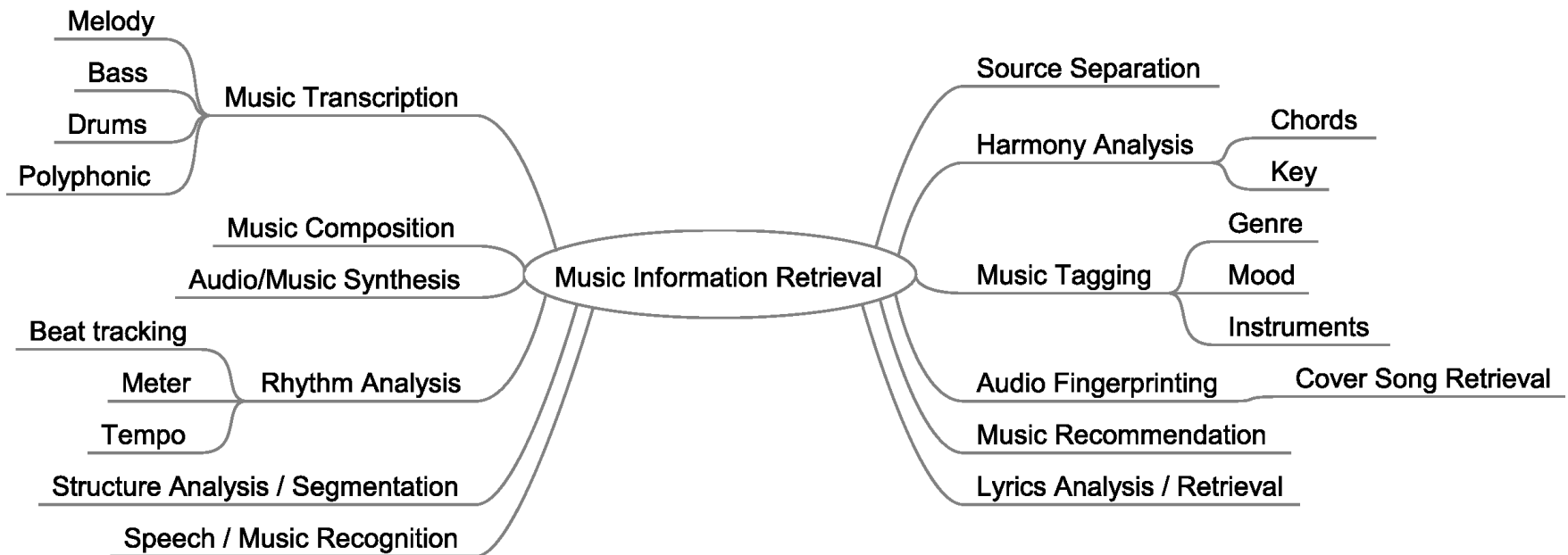
# Music Information Retrieval
## Research Landscape

- MIR @ Fraunhofer IDMT

    - Semantic music technologies (SMT) group

        - Staff + PhD / master / bachelor students + interns

- National / international research groups

    - International Audio Laboratories Erlangen, Germany

    - Centre for Digital Music, Queen Mary University, London, UK

    - Universitat Pompeu Fabra, Barcelona, Spain

    - Institute for music/acoustic research and coordination (IRCAM), Paris, France

    - USA, China, Taiwan, Japan, Korea, etc.

Fraunhofer
IDMT

# Music Information Retrieval
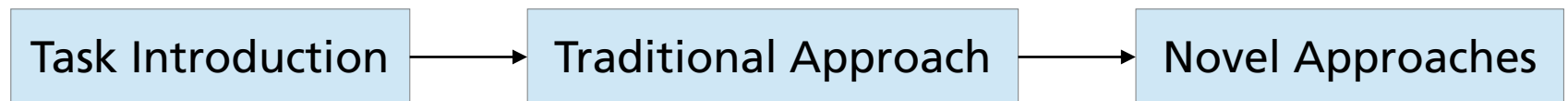## Research Task Taxonomy

Fraunhofer

IDMT

# Music Information Retrieval
## Case Studies

- MIR 1 lecture
    - Music tagging / music similarity → general tasks
    - Tempo estimation → rhythm
- MIR 2 lecture
    - Pitch detection → pitch / tonality
    - Source separation & instrument recognition → timbre

- Teaching Concept

| Task Introduction | → | Traditional Approach | → | Novel Approaches |
|---|---|---|---|---|

Fraunhofer
IDMT

# Music Tagging
## Introduction

- Tags

    - Textual (objective / subjective) annotations of songs

    - Examples

        - Instruments        (drums, bass, guitar, vocals ...)

        - Genre        (classical, electro, hip hop)

        - Mood        (mellow, romantic, angry, happy)
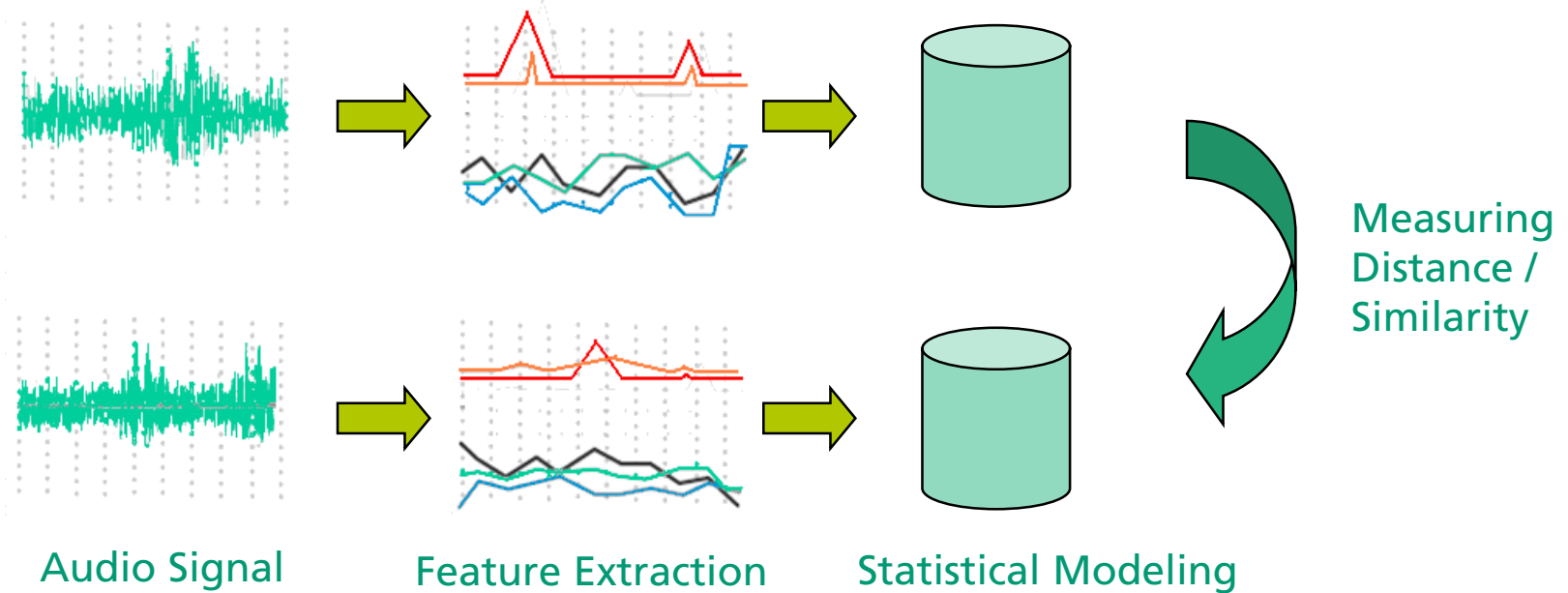
        - Miscellaneous    (noise, loud, ambient)

- Challenge

    - Music pieces change their characteristics over time

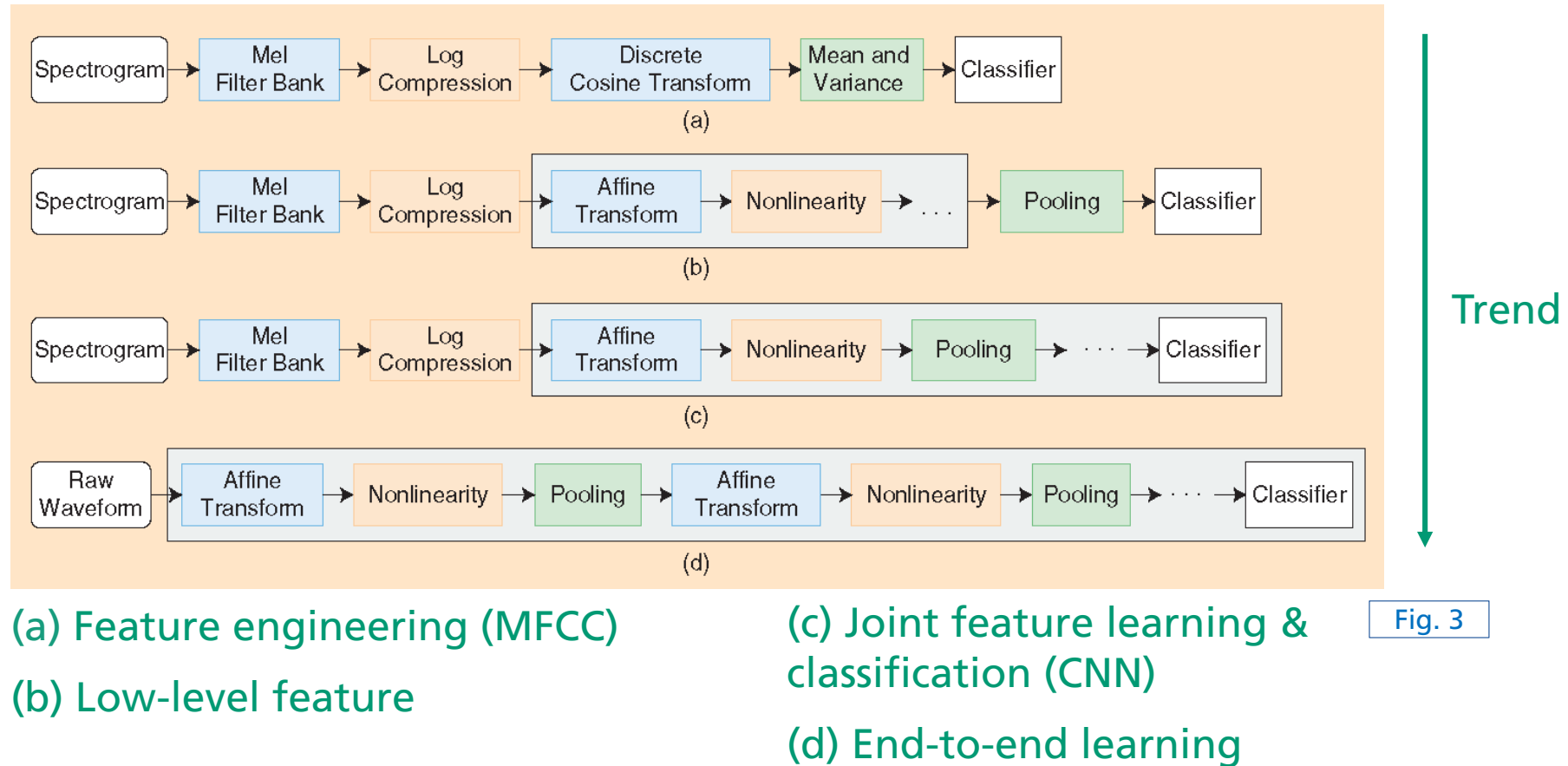        - E.g.: trumpet plays only in the chorus (jazz)

# Music Tagging
## Traditional Approach

- Audio feature engineering & music domain knowledge
- Standard classification methods (GMM, SVM, kNN)

Measuring Distance / Similarity

Audio Signal          Feature Extraction          Statistical Modeling

Fraunhofer
IDMT

# Music Tagging
## Novel Approaches



Fig. 3

(a) Feature engineering (MFCC)

(b) Low-level feature

(c) Joint feature learning & classification (CNN)

(d) End-to-end learning

# Music Tagging
## Novel Approaches

- Joint representation learning & classification using CNNs
    - Input: spectrograms (2D) or audio samples (1D end-to-end)

- Integrate musical knowledge in network design (e.g., filter shapes)
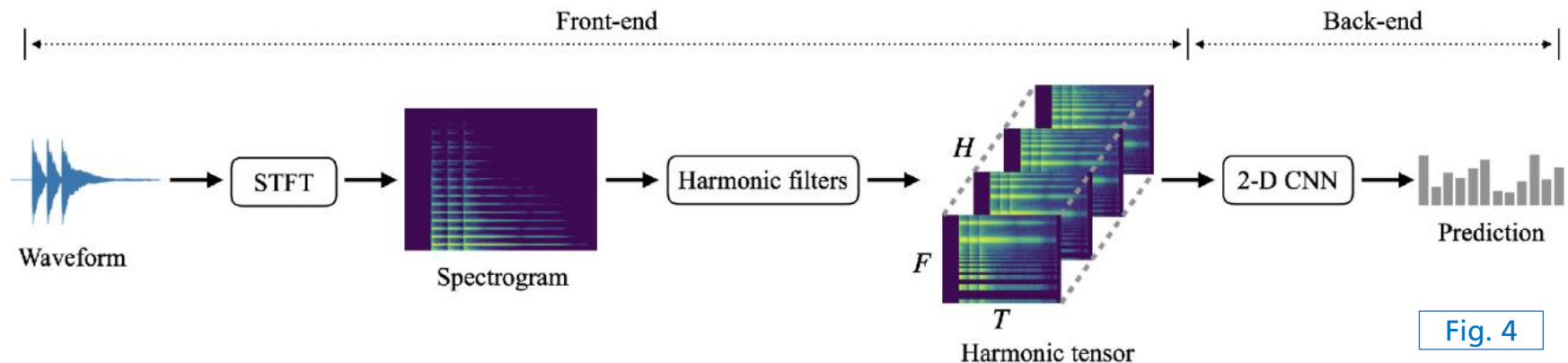


Fig. 4

# Music Tagging
## Novel Approaches

- End-to-end Learning

    - Model input is low-level representation (audio waveform)

    - No pre-processing / assumptions required

    - Not restricted to spectral magnitudes $\rightarrow$ can model phase!

    - Requires large amounts of training data
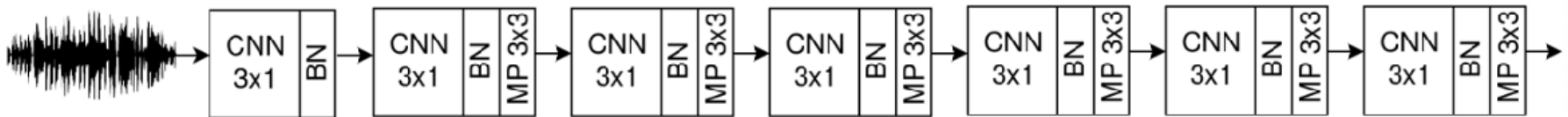


Fig. 6

Fraunhofer
**IDMT**

# Music Tagging
## Novel Approaches

- **Transfer Learning**

  - **Pre-train model on source task (lot of data available)**

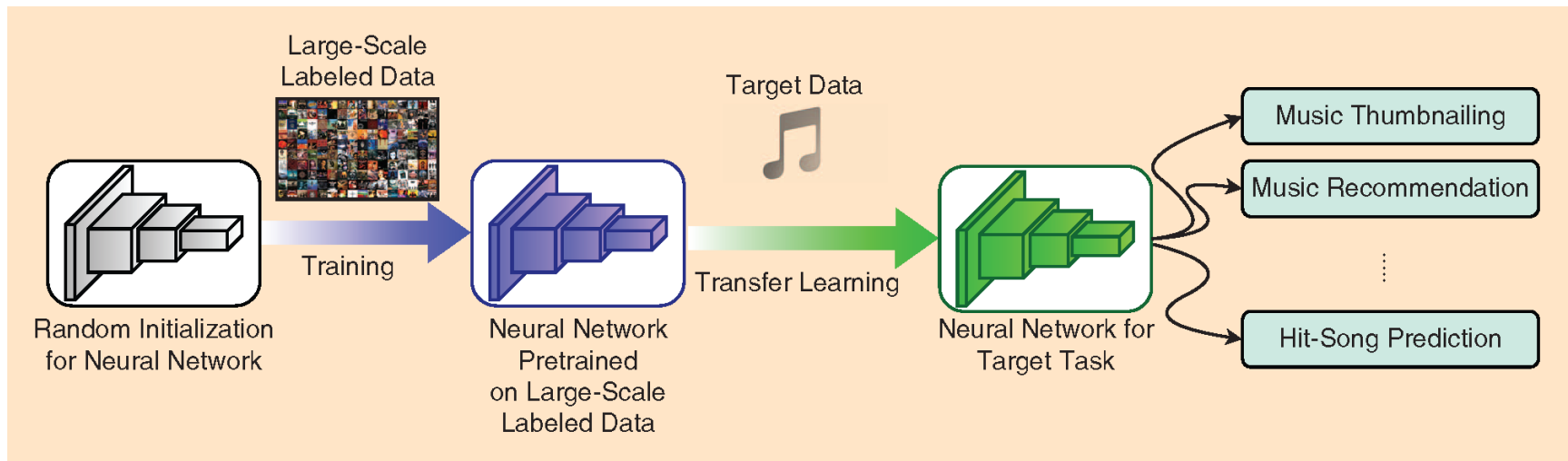  - **Fine-tune model on target task (only little data available)**



Large-Scale Labeled Data

Target Data

Music Thumbnailing

Music Recommendation

Hit-Song Prediction

Training

Transfer Learning

Random Initialization for Neural Network

Neural Network Pretrained on Large-Scale Labeled Data

Neural Network for Target Task

Fig. 5

- **Source model (CNN) → Target model (embeddings + shallow classifier)**

© Fraunhofer IDMT

Fraunhofer
IDMT

# Music Similarity
## Introduction

- Music → inherently multi-dimensional

- Challenge

  - Large music databases
  - Incomplete / missing metadata



Fig. 7

- Query by example → general retrieval approach

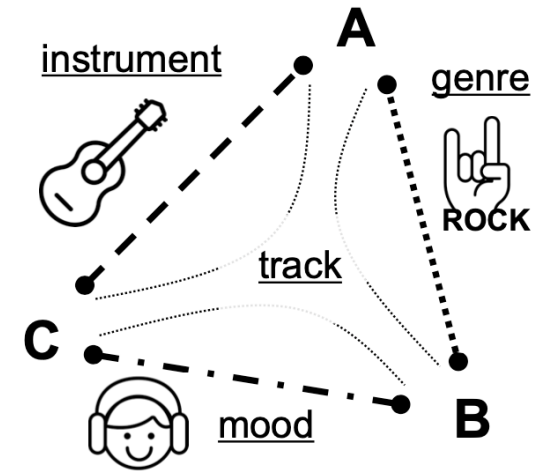  - Retrieval most similar song *S* given a query song *Q*

# Music Similarity
## Introduction

- Retrieval tasks

    - Music fingerprinting (retrieve title, artist, e.g., Shazam app)

    - Cover song identification (similar text, chord progressions …)

    - Music replacement (similar style, instrumentation)

- Specificity of different tasks



High    **Specificity**    Low

Cover Song
Identification

Music
Finger-printing

Music
Replacement

Instrument-based
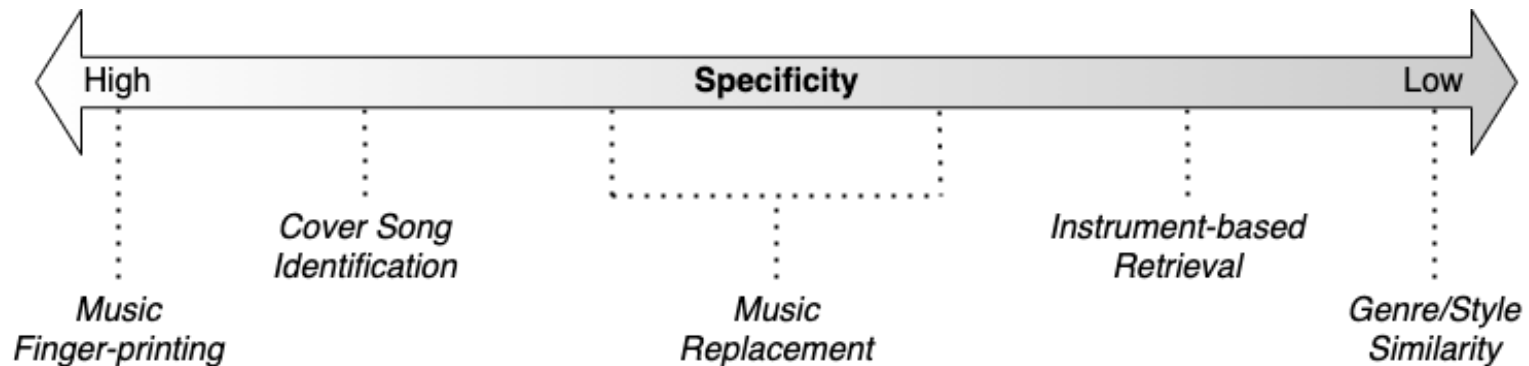Retrieval

Genre/Style
Similarity

Fig. 8

Fraunhofer
IDMT

# Music Similarity
## Traditional Approaches

- Different dimensions of music similarity

  - Melodic similarity (pitch contours)

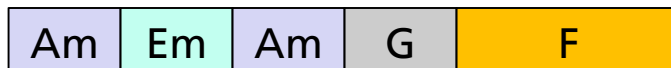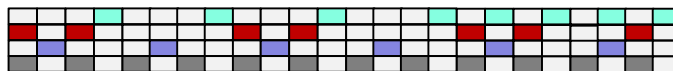  - Timbral similarity (instrumentation)

    Piano ▬  Guitar ▬  Vocals ▬

  - Structural / harmonic similarity (segments, chords)

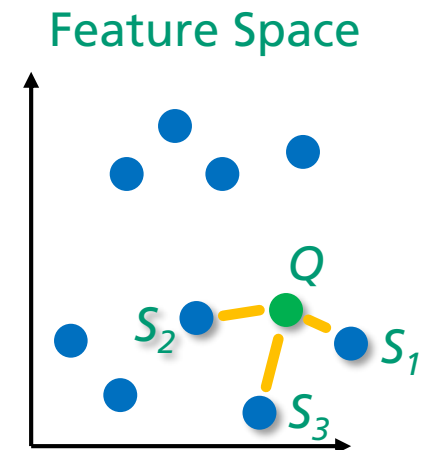    | Am | Em | Am | G | F |
    |----|----|----|---|---|

  - Rhythmic similarity (patterns)

# Music Similarity
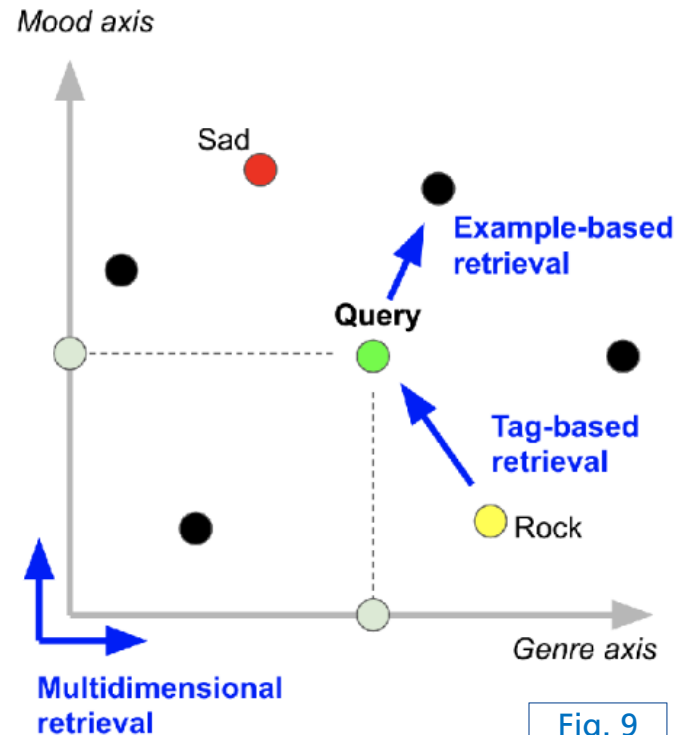## Novel Approaches

- Metric learning
    - Model (abstract) notion of similarity between data instances
        - Pair-wise distance between feature representations

- Training → Preserve similarity in the feature space
    - Proximity between similar instances
    - Distance between dissimilar instances

- Distance measures (Euclidean, cosine)

- Query $Q$ → Ranked list of most similar items ($S_i$)

Feature Space

$Q$

$S_2$

$S_1$

$S_3$

Fraunhofer

**IDMT**

# Music Similarity
## Novel Approaches

- Disentanglement learning

    - Similarity → multiple semantic concepts (e.g., genre, instrument, mood)

        - learnt jointly

        - remain separable in the embedding space

- Improves tagging (classification) and recommendation (similarity)
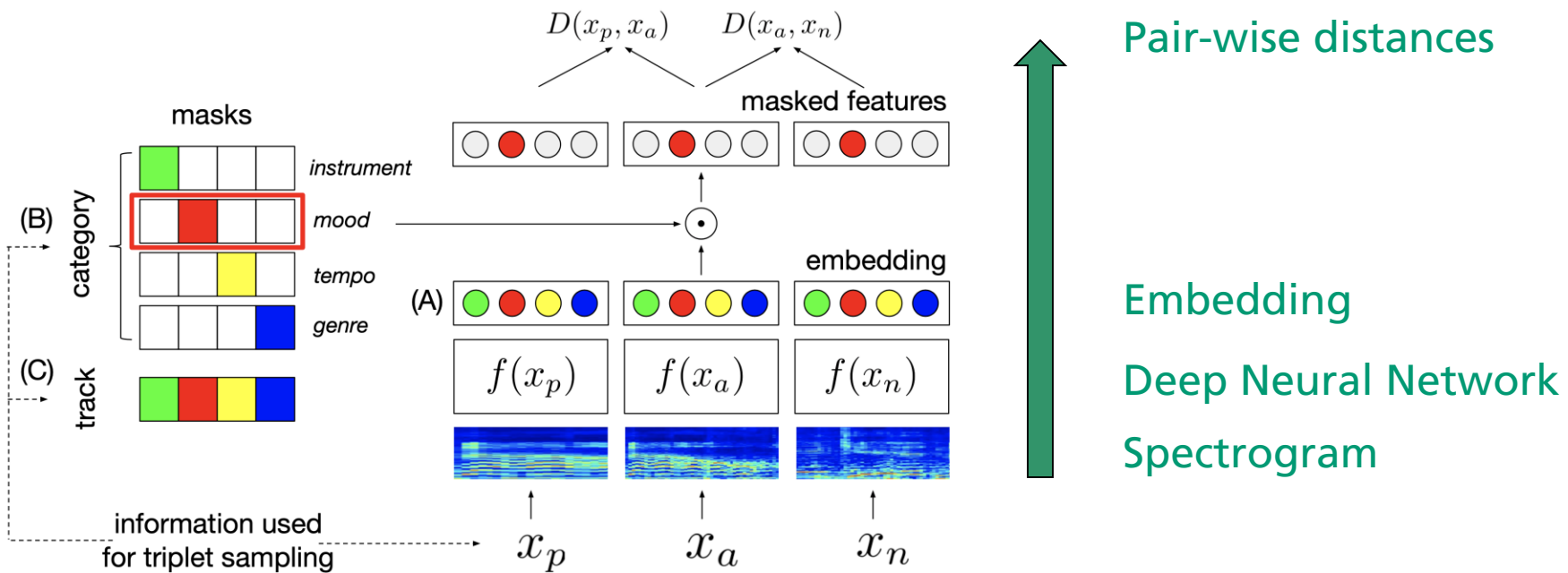


Fig. 9

# Music Similarity
## Novel Approaches

- Triplet-based Training

  - Conditional Similarity Networks (CSN) [Lee, 2020]



Pair-wise distances

Embedding

Deep Neural Network

Spectrogram

Fig. 10

# Tempo Detection
## Introduction

- Tempo [beats / minute]

  - Frequency with which humans tap along the beat



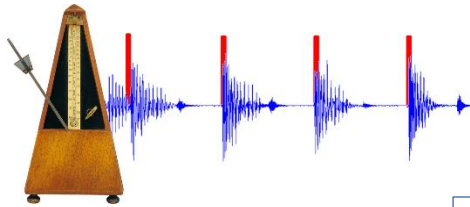Fig. 11

- Beat tracking

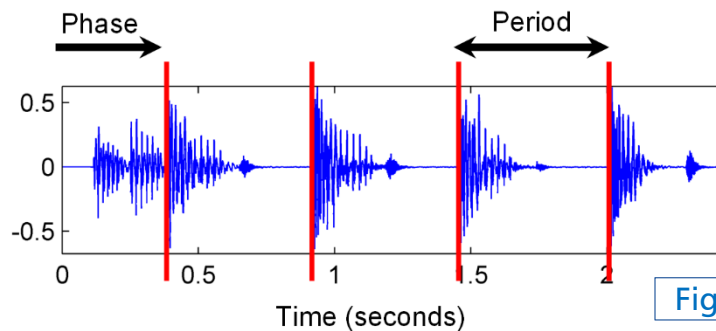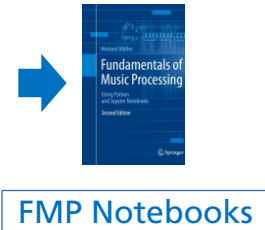  - Estimating precise beat positions



Phase          Period

Time (seconds)

Fig. 12



FMP Notebooks

Fraunhofer
IDMT

# Tempo Detection
## Introduction

- Note onsets → note beginning times

    - Clearly defined for plucked string and percussion instruments

    - Ambiguous for wind & brass instruments

- Onset detection
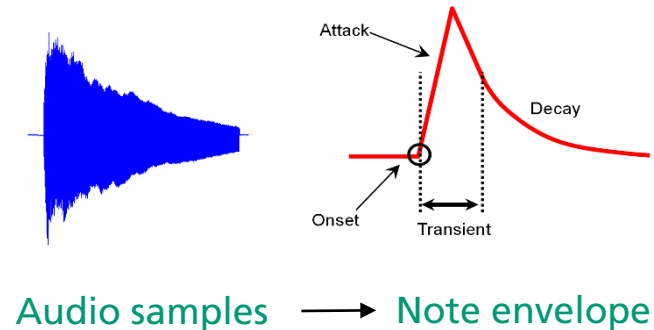
    - Onset detection function

    - Peak picking



Audio samples ⟶ Note envelope

Fig. 13

Audio samples

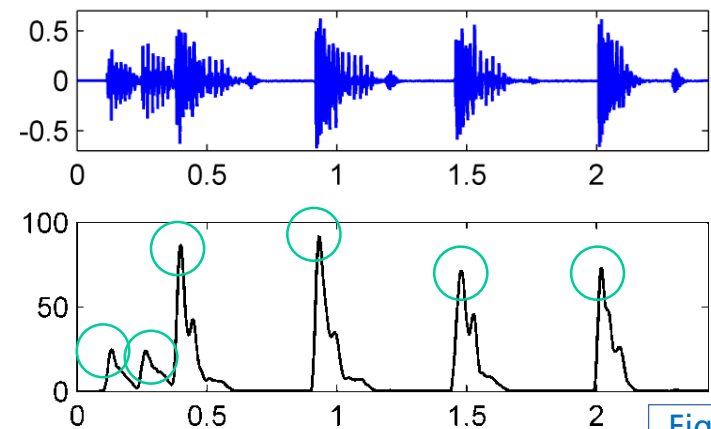Onset detection function & peaks



Fig. 14

Fraunhofer
**IDMT**

# Tempo Detection
## Traditional Methods

■ Predominant local pulse (PLP)

    ■ Correlation with local (windowed) periodic patterns

■ Tempogram [Grosche & Müller, 2009]

    ■ Local likelihood of different tempo candidates

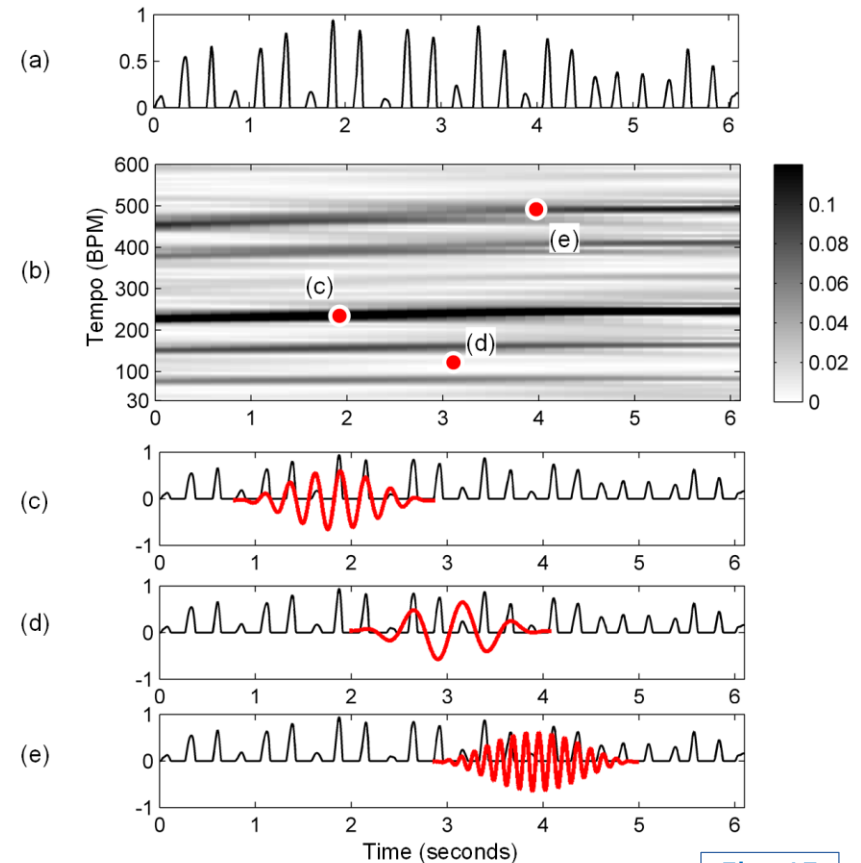    ■ Allows to follow tempo changes (classical music)

FMP Notebooks



Fig. 15

# Tempo Detection
## Novel Methods
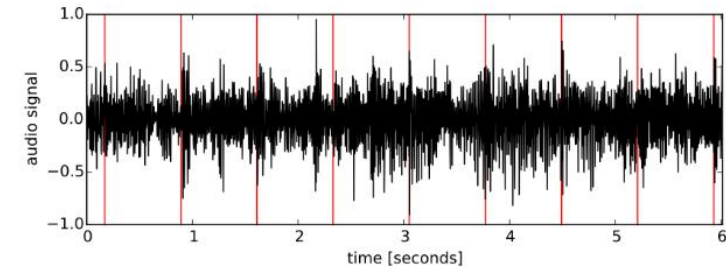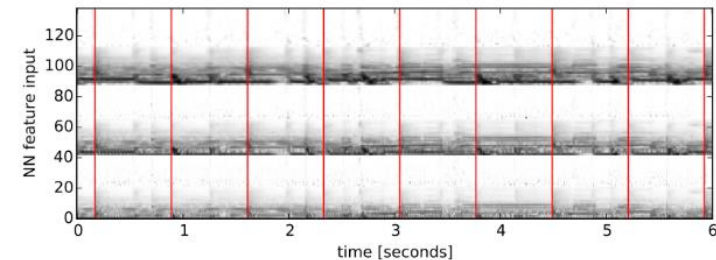
- Approach [Böck et al., 2015]



Fig. 16



(a) Input audio signal

- Signal representation

  - Stacking of 3 STFT magnitude spectrograms (N=1024, 2048, 4096)

  - Log-amplitude & log-frequency
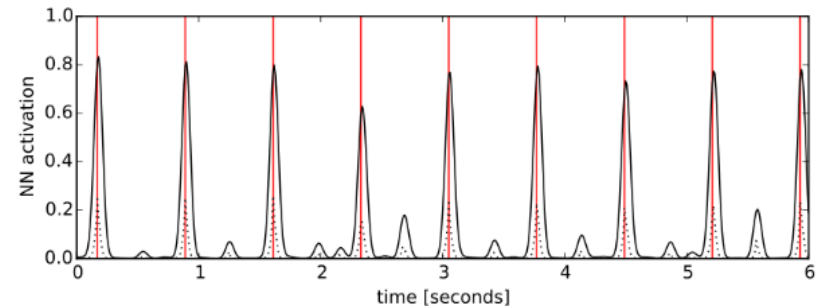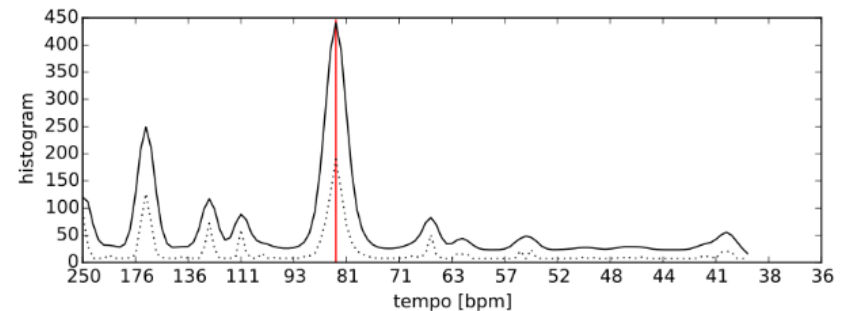


(b) Input to the neural network

Fig. 17

# Tempo Detection
## Novel Methods

- Neural Network
    - Recurrent (bi-directional LSTM) layer
    - Outputs beat activation function
- Comb filter bank
    - Multiple comb filters → detect periodicities
- Estimate tempo from histogram maximum

(c) Neural network output (beat activation function)

(f) Weighted histogram with summed maxima

Fig. 18

# Tempo Detection
## Novel Methods

- **Approach** [Schreiber & Müller, 2018]

  - Sample rate ~ 11 kHz, 40-band mel spectrogram

- **Tempo estimation → classification (256 classes: 30 – 285 bpm)**

- **Neural Network**

  - 3 layers (short filters) → onsets

  - 4 multi-filter modules (parallel conv layers) → compress along frequency & find periodicities

  - Dense layers → tempo classification
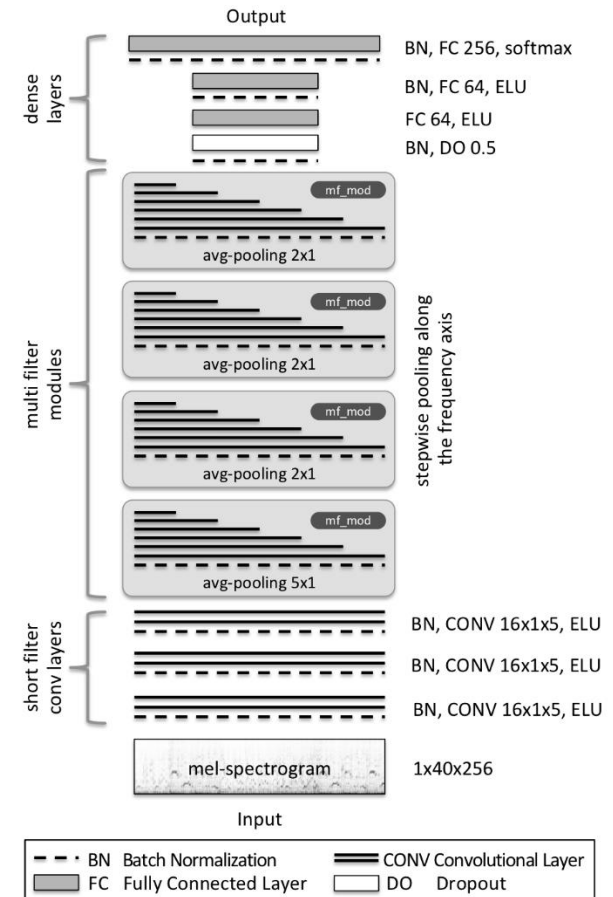


Fig. 19

# Summary

- Music Information Retrieval

- Music Tagging

- Music Similarity

- Tempo Estimation

- Main trends
    - Adapt (data-driven) deep learning methods to music domain
    - Incorporate music domain knowledge

# References

Böck, S., Krebs, F., & Widmer, G. (2015). Accurate tempo estimation based on recurrent neural networks and resonating comb filters. *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, 625–631.

Grosche, P., & Müller, M. (2009). A mid-level representation for capturing dominant tempo and pulse information in music recordings. *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, 189–194.

Lee, J., Bryan, N. J., Salamon, J., Jin, Z., & Nam, J. (2020). Disentangled Multidimensional Metric Learning for Music Similarity. *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 6–10. Barcelona, Spain.

Lee, J., Bryan, N. J., Salamon, J., Jin, Z., & Nam, J. (2020). Metric learning vs classification for disentangled music representation learning. *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, 439–445. Montréal, Canada.

Müller, M. (2021). *Fundamentals of Music Processing - Using Python and Jupyter Notebooks* (2nd ed.). Springer.

Nam, J., Choi, K., Lee, J., Chou, S. Y., & Yang, Y. H. (2019). Deep Learning for Audio-Based Music Classification and Tagging: Teaching Computers to Distinguish Rock from Bach. *IEEE Signal Processing Magazine*, *36*(1), 41–51.

Pons, J., Nieto, O., Prockup, M., Schmidt, E., Ehrmann, A., & Serra, X. (2018). End-to-End Learning for Music Audio Tagging at Scale. *Proceedings of the International Society for Music Information Retrieval (ISMIR)2*, 637–644. Paris, France.

**Fraunhofer**

**IDMT**

# References

Ribecky, S. (2021). *Disentanglement Representation Learning for Music Annotation and Music Similarity*. Technische Universität Ilmenau.

Schreiber, H., & Müller, M. (2018). A Single-Step Approach to Musical Tempo Estimation using a Convolutional Neural Network. *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, 98–105. Paris, France.

Won, M., Chun, S., Nieto, O., & Serra, X. (2020). Data-Driven Harmonic Filters for Audio Representation Learning. *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 536–540. Barcelona, Spain.

# Images

Fig. 1: https://www.synchtank.com/wp-content/uploads/2018/06/1476277072027.jpg

Fig. 2: https://miro.medium.com/max/800/1*cC1KOdyzzt1nazak42cBdg.jpeg

Fig. 3: [Nam, 2019], p. 42, Fig. 1

Fig. 4: [Won, 2020], p. 537, Fig. 1a

Fig. 5: [Nam, 2019], p. 48, Fig. 4

Fig. 6: [Pons, 2018], p. 639, Fig. 2 (top left)

Fig. 7: [Lee, 2020, ICASSP], p. 1, Fig. 1

Fig. 8: [Ribecky, 2021], p. 26, Fig. 2.11

Fig. 9: [Lee, 2020, ISMIR], p. 1, Fig. 1

Fig. 10: [Lee, 2020, ICASSP], p. 2, Fig. 2

Fig. 11: [Müller, 2021], p. 309, chapter 6 (cover image)

Fig. 12: [Müller, 2021], p, 310, Fig. 6.1(b)

Fig. 13: [Müller, 2021], p. 311, Fig. 6.2

Fig. 14: [Müller, 2021], p. 313, Fig. 6.3(a)&(b)

# Images

Fig. 15: [Grosche & Müller, 2009], p. 2, Fig. 1(e-g) & p. 3, Fig. 2 (a)

Fig. 16: [Böck et al., 2015], p. 2, Fig. 1

Fig. 17: [Böck et al., 2015], p. 3, Fig. 2 (a) & (b)

Fig. 18: [Böck et al., 2015], p. 3, Fig. 2 (c) & (f)

Fig. 19: [Schreiber & Müller, 2018], p. 3, Fig. 2

Fraunhofer

**IDMT**

# Sounds

**AUD-1:** Mr Smith – Black Top (2021), https://freemusicarchive.org/music/mr-smith/studio-city/black-top

**AUD-2:** Crowander – Humbug (2021), https://freemusicarchive.org/music/crowander/from-the-piano-solo-piano/humbug

**AUD-3:** Bumy Goldson: Keep Walking (2021), https://freemusicarchive.org/music/bumy-goldson/parlor/keep-walking

**AUD-4:** Cloudjumper: Mocking the god (2016), https://freemusicarchive.org/music/Cloudjumper/Memories_of_Snow/05_Cloudjumper_-_Mocking_the_gods

Fraunhofer
IDMT

# Thank you!

- Any questions?

Dr.-Ing. Jakob Abeßer

Fraunhofer IDMT

Jakob.abesser@idmt.fraunhofer.de

https://machinelistening.github.io

Fraunhofer

IDMT