

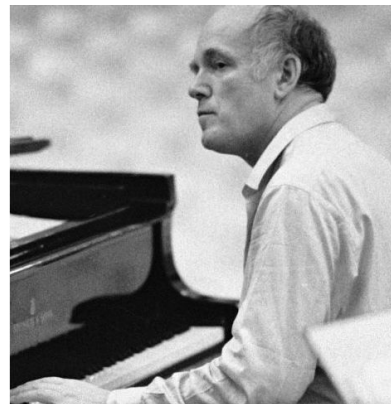


# TOWARDS EXPLAINING EXPRESSIVE QUALITIES IN PIANO RECORDINGS: Transfer of Explanatory Features via Acoustic Domain Adaptation

Shreyan Chowdhury, Gerhard Widmer

# *Prologue*





*“The search for audio features that capture the expressive perceptual qualities of performed music”*

# The Story So Far

## Con Espressione Game

*On the Characterization of Expressive Performance in Classical Music: First Results of the Con Espressione Game* (ISMIR 2020)

[C. Cancino-Chacón, S. Peter, S. Chowdhury, A. Aljanaki, G. Widmer]

## Mid-level features

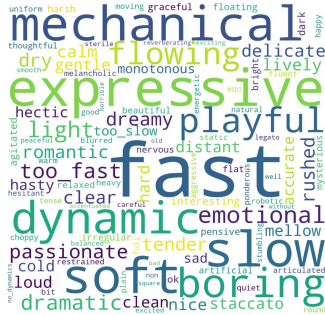
*Towards Explainable Music Emotion Recognition: The Route via Mid-level Features* (ISMIR 2019)

[S. Chowdhury, A. Vall, V. Haunschmid, G. Widmer]

**JKU**  
JOHANNES KEPLER  
UNIVERSITY LINZ



# The Con Espressione Dataset



“gentle”

“agitated”

“ponderous”

“hasty”

“graceful”

“staccato”

“stumbling”

“tender”

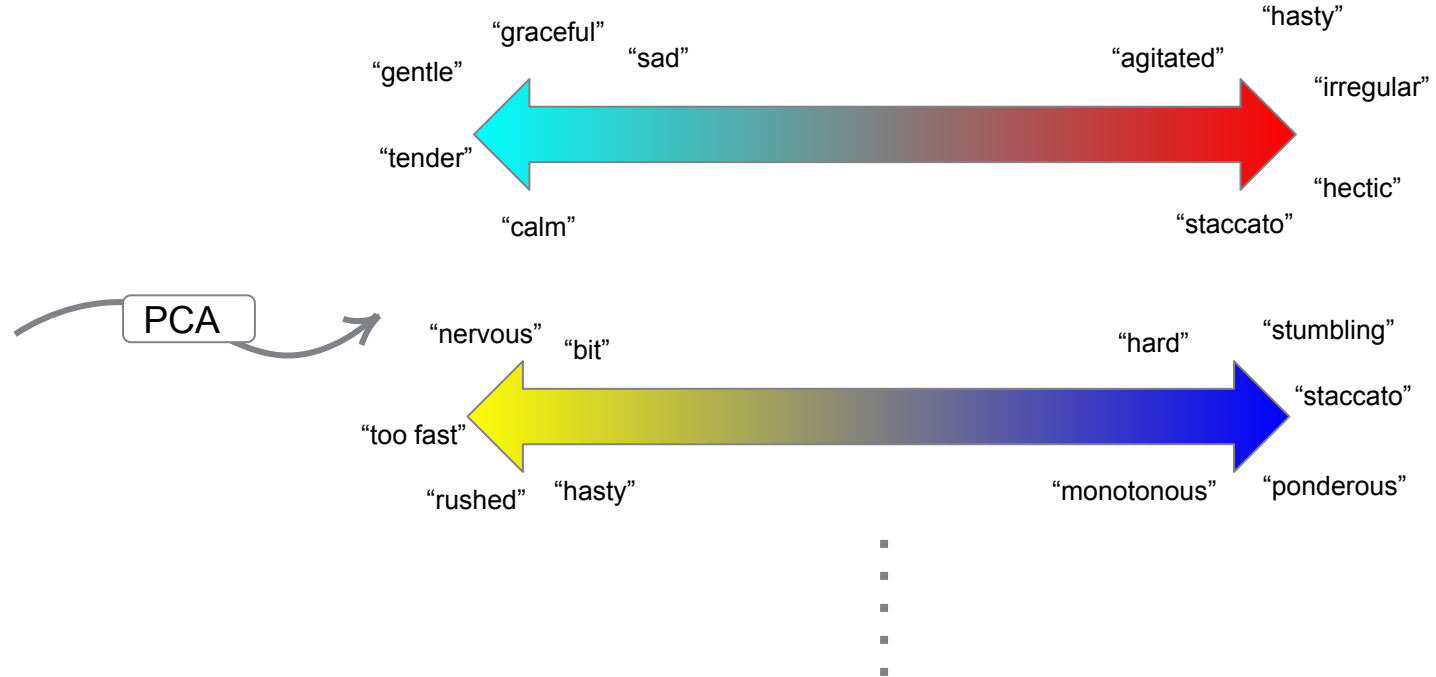
“nervous”

“too fast”

“rushed”

“monotonous”

# The Con Espressione Dataset





# The Con Espressione Dataset

*Can the embedding dimensions obtained  
from free-text descriptions of expressive piano performances  
be modeled using audio features?*

# Mid-level Features<sup>1</sup>

Low-level features,  
such as pitch

Building blocks of musical  
signals

Melodiousness  
Articulation  
Rhythm complexity  
Rhythm stability  
Dissonance  
Tonal stability  
Minorness

Perceptual and subjective, but  
make intuitive musical sense

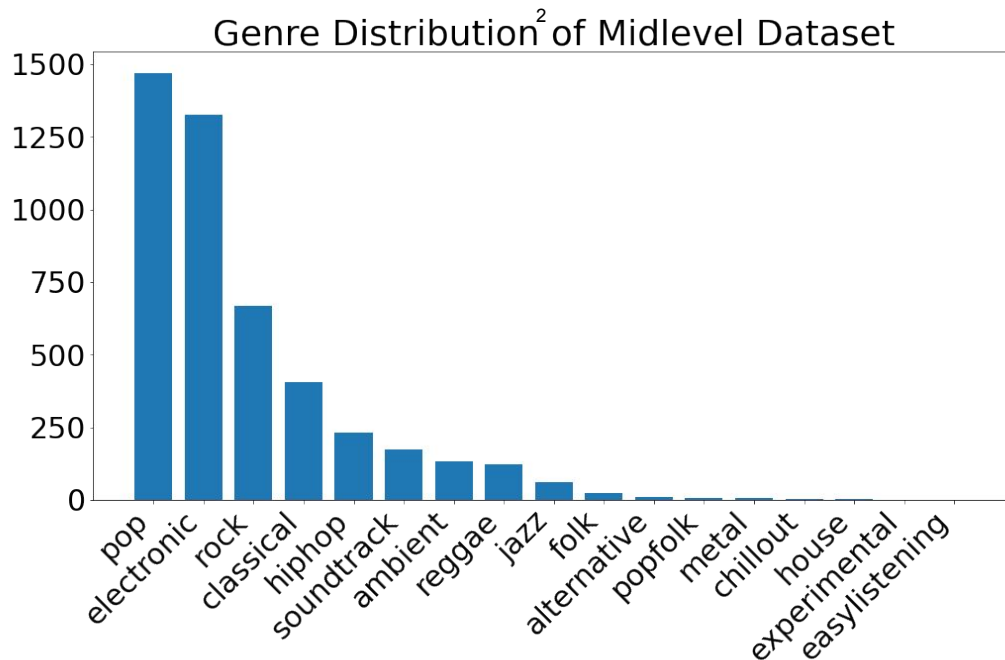
High-level  
features, such as  
emotion

Subjective, abstract  
descriptions

# Learning and Transferring Mid-level Features

## Mid-level Dataset

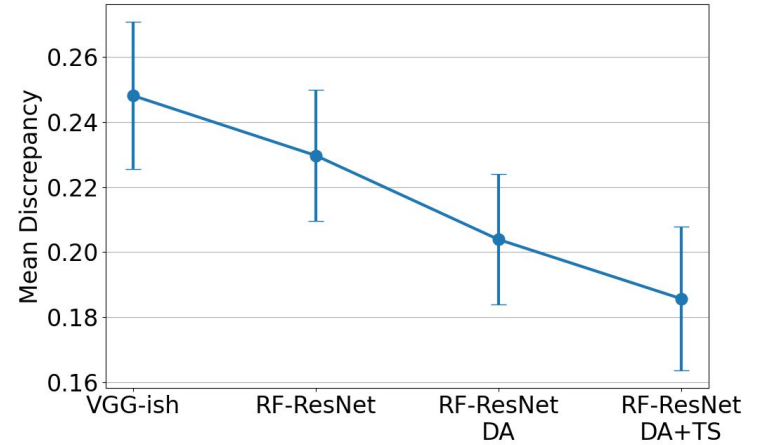
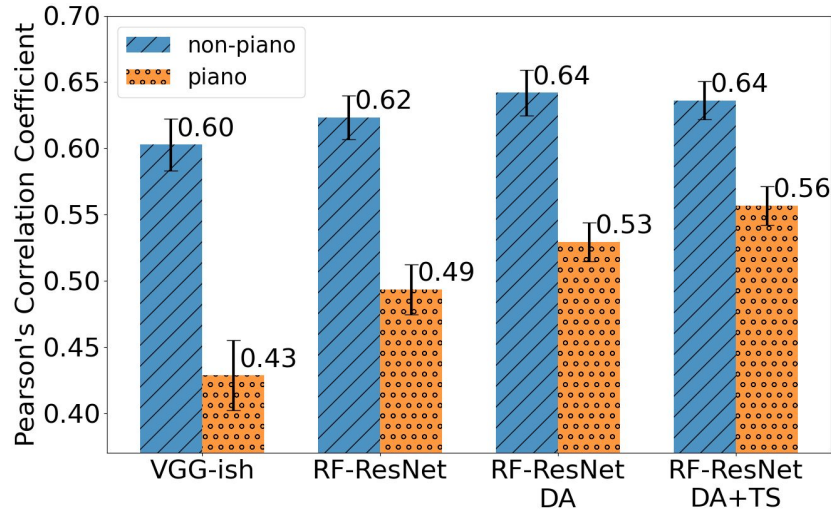
- 5000 snippets
- 15-second clips
- Crowdsourced annotation



# Transferring to our Domain of Interest (Piano)

- Setup
  - Test-set: manually created from piano recordings in the Mid-level Dataset
  - RF-ResNet architecture
- Learn Mid-level Features with Domain Adaptation
  - Unsupervised Domain Adaptation (UDA) by Backpropagation<sup>3</sup>
- Teacher-Student Refinement
  - Multiple DA teachers
  - Student learns from Mid-level combined with pseudo-labeled piano dataset

# Results on the Test-set



# Results – Transferring to Con Espressione Recordings

	Dim 1	Dim 2	Dim 3	Dim 4
VGG-ish	0.35	0.10	0.22	0.32
RF-ResNet	0.36	0.07	0.28	0.33
RF-ResNet DA	<b>0.40</b>	0.09	<b>0.29</b>	0.32
RF-ResNet DA+TS	0.35	<b>0.15</b>	<b>0.29</b>	<b>0.34</b>

Coefficient of determination (R2-score)

RF-ResNet		RF-ResNet DA+TS	
Feature	<i>r</i>	Feature	<i>r</i>
articulation	0.47	melodiousness	– 0.39
rhythmic complexity	0.41	articulation	0.46
		rhythmic complexity	0.41
		dissonance	0.40

Pearson's correlation (*r*) for mid-level features with description embedding dimension 1.

Features with  $p < 0.05$  and  $|r| > 0.20$  are selected.

This dimension has positive loadings for words like “hectic”, “irregular”, and negative loadings for words like “sad”, “gentle”, “tender”.

Looking forward to see you at the poster session **AUD-20!**