

A Computational Approach to Analysis and Detection of Singing Techniques

January 29th, 2024 Ph.D. Defense
Yuya Yamamoto

Informatics Program, University of Tsukuba

Supervisors:
Nobutaka Suzuki
Hiroko Terasawa
Hiroyoshi Ito

Committee:
Nobutaka Suzuki
Hiroko Terasawa
Atsushi Toshimori
Shuichi Moritsugu
Juhan Nam

♪Can you listen to the music??

- **Introduction**
 - Background
 - Research Aims and Problem Statements
- **Research works**
 - Ch. 3: Exploration of Singing Techniques
 - Ch. 4: Singing Technique Analysis on Actual Vocal Performances
 - Ch. 5: Characteristics-aware Modeling for Singing Technique Classification
 - Ch. 6: Singing Technique Detection from Real-world Vocal Tracks
- **Conclusion**

Introduction

Combined Chapter 1 and Chapter 2 in the presentation
Chapter 1: Motivation, problems, and aims of the thesis
Chapter 2: Related works

Research area: Singing voices in music

4

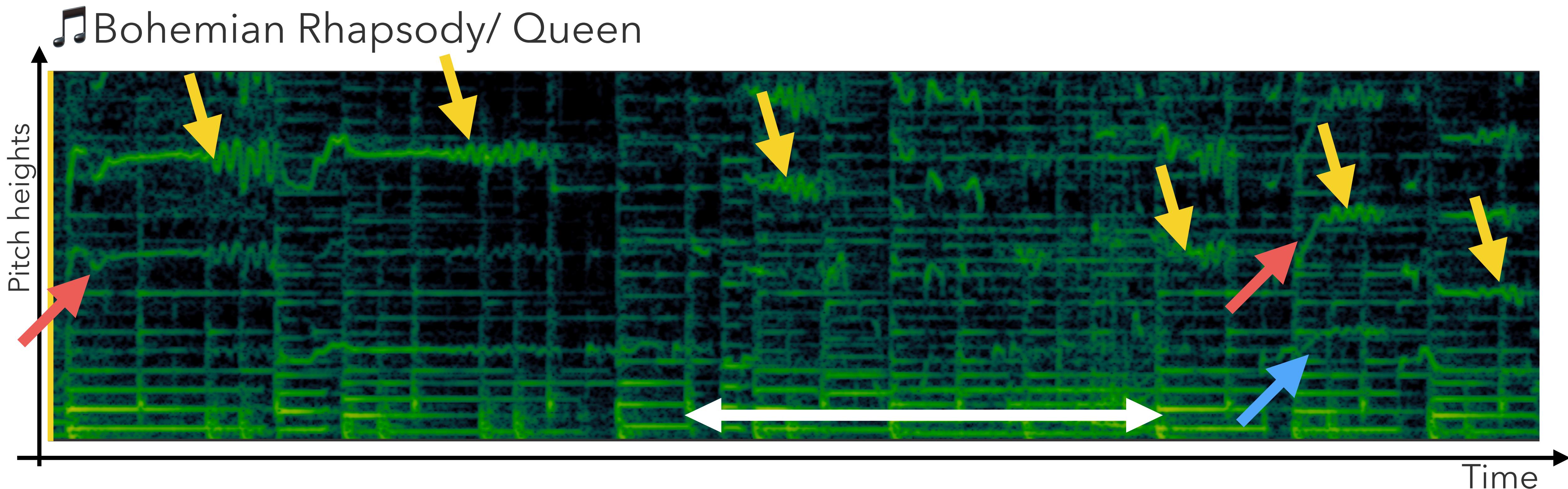
- Singing voice is one of the most important parts in music
- Singing, Listening, Creating etc. are fundamental cultural and artistic activities for human beings
 - Many research works have been done to clarify or to engineer about singing voice



Singing techniques

5

Singing technique is one of the ways to embody the expression



Techniques: Vibrato (yellow), Portamento (pink), Falsetto (blue), Raspy voice (White)

Sore Demo Shitai / Ken Hirai

♪ Playing...

Lyrics

A Na Ta Ha Ke-Shi Te

Tsu Ka Wa Na | No Yo

A Ta Shi No Su Ki Na

Bo Di So - Pu Wo

- **Vocal fry**: producing pulsive sound
- **Bend**: short going-around pitch bending
- **Vibrato**: periodic pitch modulation
- **Scooping**: short ascending pitch bending

Shin Ji Dai / Ado

♪ Playing...

Lyrics

Shin Ji Da I **Ha**

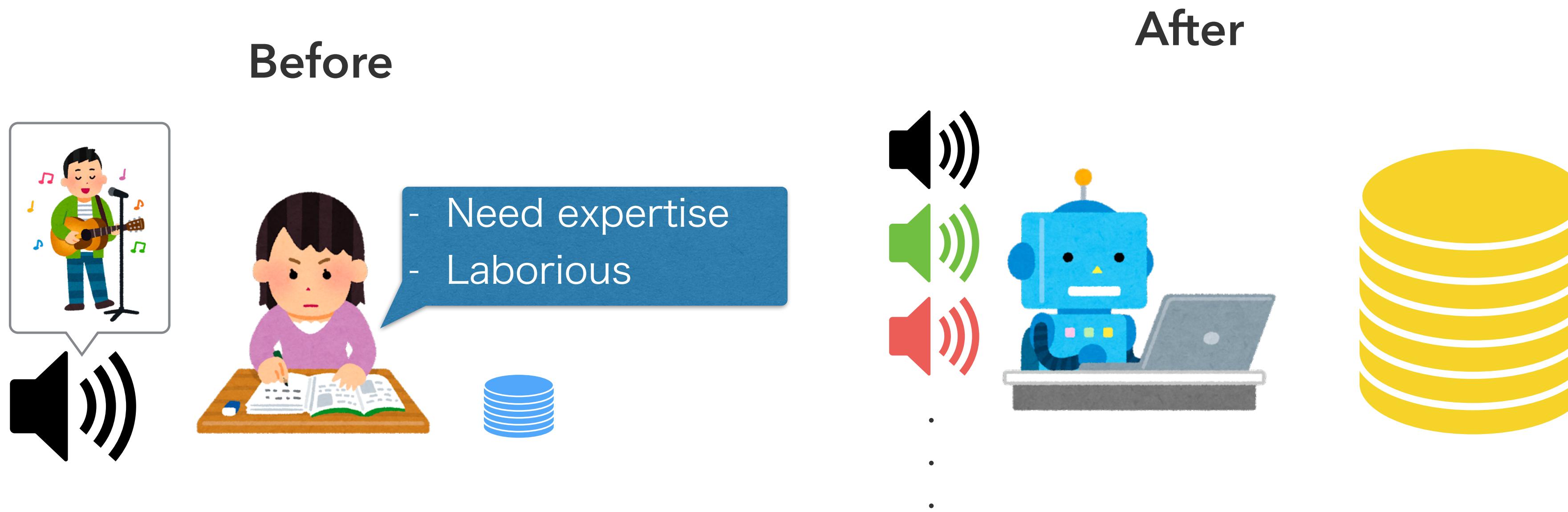
Ko No Mi **Ra** I Da

S(e) Kai Ju **U** Ze M **Bu**

Ka E Te Shi **Ma** E **Ba**...

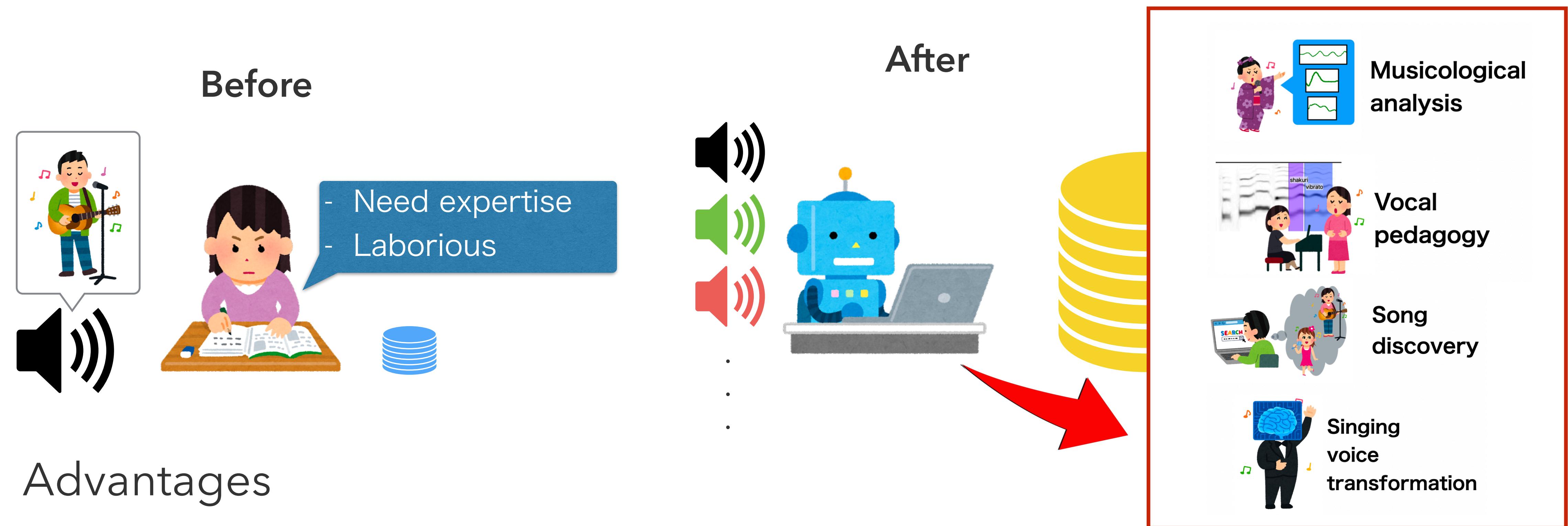
- **Falsetto:** different tone on higher pitch
- **Vibrato:** periodic pitch modulation
- **Scooping:** short ascending pitch bending
- **Hiccup:** short falsetto & pitch jump
- _, (): Enhancing/Omitting phonemes

Establish “Computational” framework for singing technique analysis



- Advantages
 - 1. Accelerates the singing technique analysis by automation

Establish “Computational” framework for singing technique analysis

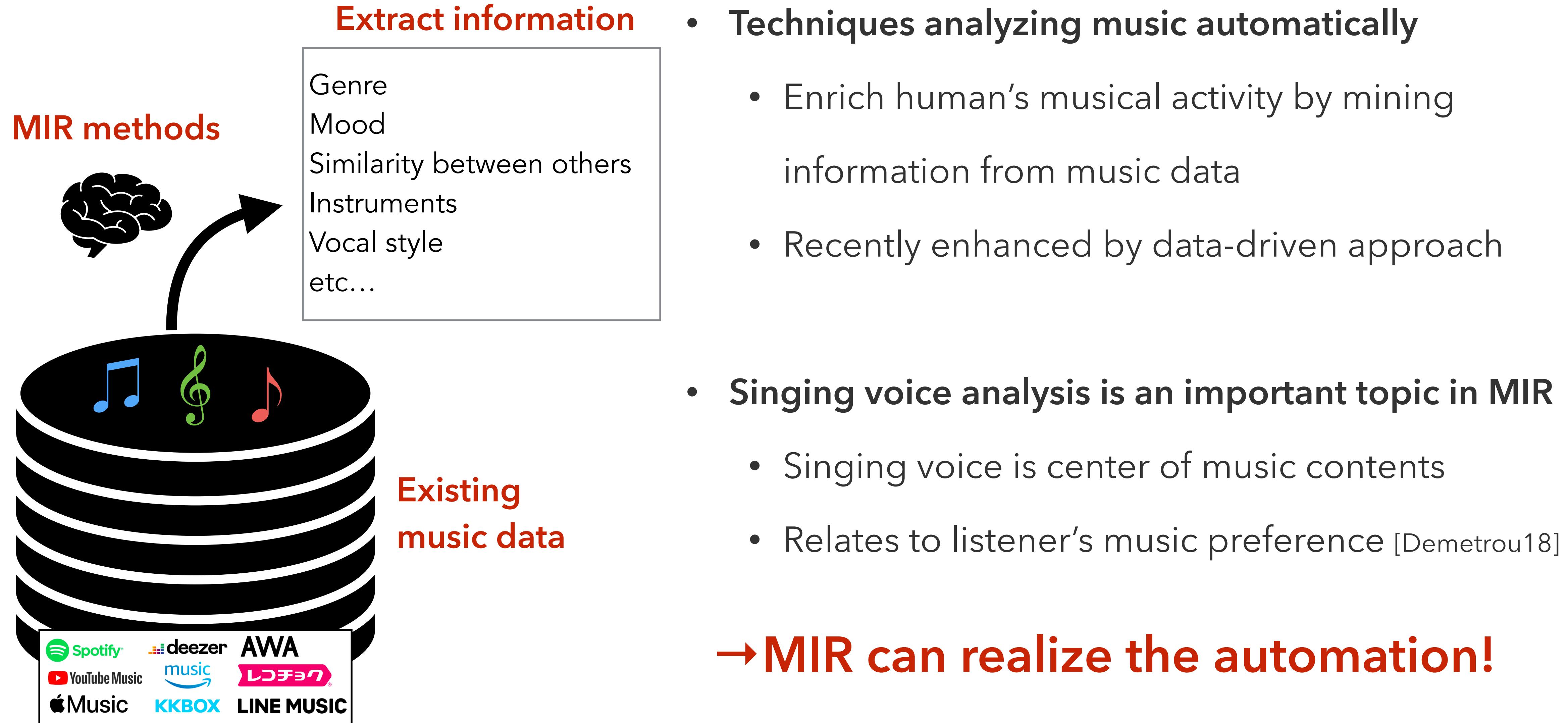


- Advantages
 - 1. Accelerates the singing technique analysis by automation
 - 2. Combines to various computational applications (e.g., analysis, pedagogy, discovery, creation, etc.)

- Manually analysis of singing styles to clarify the characteristics
 - Recording & Acoustic parameter analysis (e.g., Vibrato parameter [Seashore 37])
 - Qualitative song reading (e.g., “Kobushi” in Ken Hirai [Nakazato 09])
- **Valuable, yet still less applicable in a real-world scenario**
 - Lacks overview and objectivity
 - Needs expertise and time-consuming
 - Few samples and limited scope

Research field: Music information retrieval (MIR)

11



Related works ② Automatic vocal analysis in MIR

12

component



role in performance

What?

**MIR tasks/
research field**

**Singing note- / lyric-
transcription**

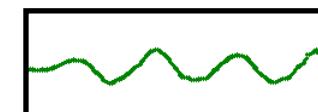
**Conventional
works**

[Mauch 15],
[Nishikimi 17],
[Gupta 20],
[Hsu 21],
[Wang 21],
[Demirel 21] and more.

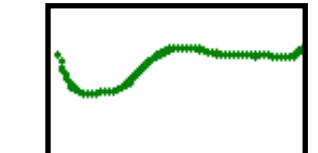
Vocalist



Expression



Vibrato at each "go"



Glissando at first "let"

Who?

**Singer
identification**

[Fujihara 10],
[Kroher 14],
[Wang 18]
[Nachmani 19],
[Lee 19],
[Hsieh 20] and more.

How?

**Singing technique
analysis**

- Indian raga analysis,
- 4 pitch techniques [Miryala 13]
- Vocal style transfer,
- 4 pitch techniques [Ikemiya 14]
- Track-wise identification,
- 10 various techniques [Wilkins 18]
- Heavy metal scream detection,
- 3 different scream [Kalbag 22]

Lacking in-song, various techniques, and with timestamp annotation

Dataset	singing songs	annotation with timestamp	monophonic (a capella)	genre/purpose	kinds	numbers
Phonation Modes [Proutskova 13]				Classic/ Phonation mode classification	4	763
VocalSet [Wilkins 18]	(Simple phrase e.g., scale, arpeggio)			Classic and popular/ Singer technique classification	10	3560
KVT Dataset [Kim 20]				K-POP/ Singing voice tagging	7	446
MVD Dataset [Kalbag 22]				Heavy metal/ Scream detection	3	57
SVQTD [Xu 22]				Classic/singing attribute classification	7	4000

[Proutskova 13] Breathy, Resonant, Pressed – Automatic Detection of Phonation Mode from Audio Recordings of Singing, P. Proutskova et al., Journal of New Music Research, 2013

[Wilkins 18] VocalSet: A Singing Voice Dataset, J.Wilkins et al., ISMIR 2018

[Kim 20] Semantic Tagging of Singing Voices in Popular Music Recordings, K. L Kim et al., TASLP 2020.

[Kalbag 22] Scream Detection in Heavy Metal Music, V. Kalbag et al. SMC 2022

[Xu 22] Paralinguistic singing attribute recognition using supervised machine learning for describing the classical tenor solo singing voice in vocal pedagogy, Y. Xu et al. EURASIP JASM 2022.

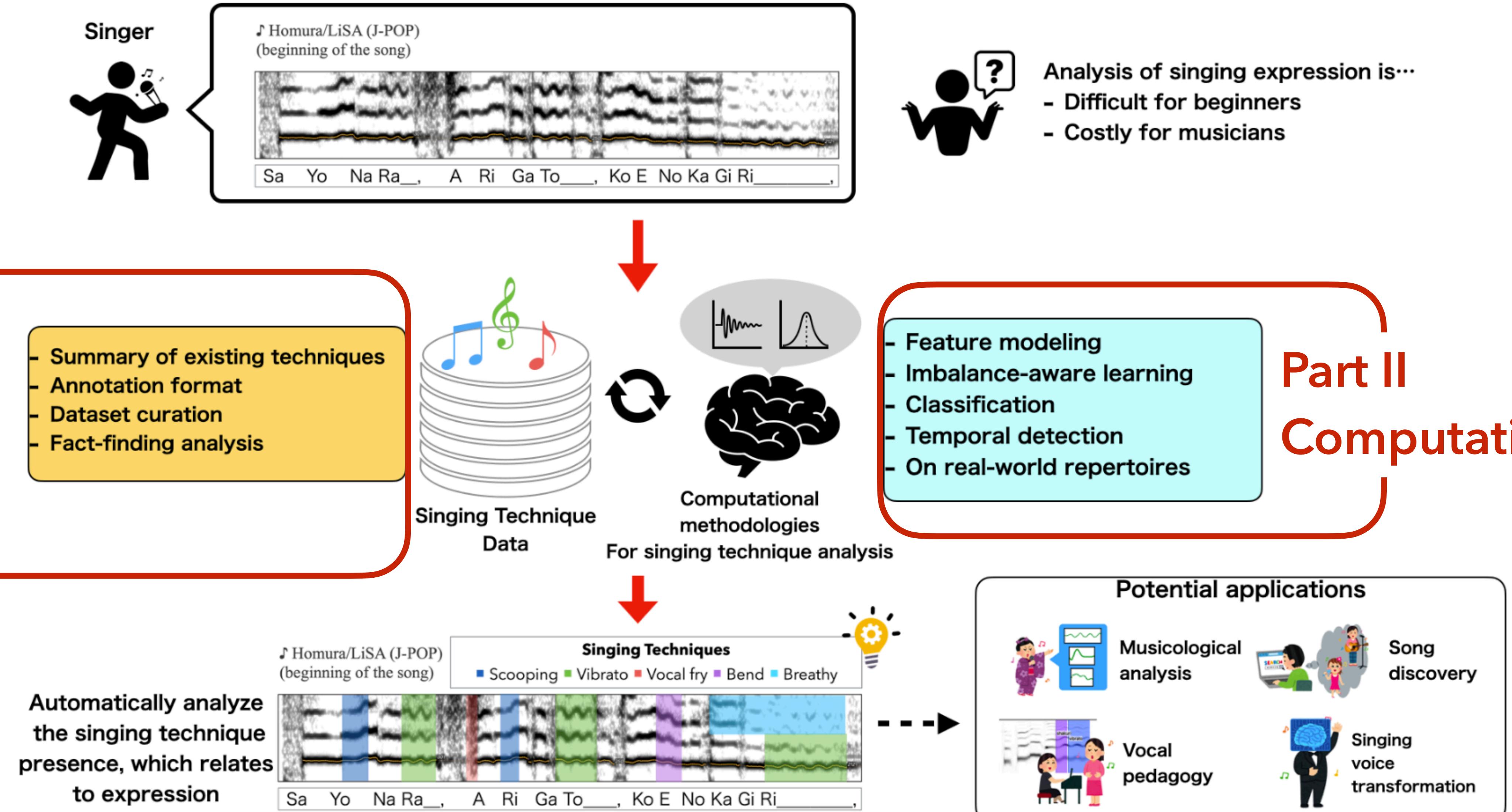
The lack on framework; dataset and computational methodology

- **1. Absence of data, annotation and its characteristics**
 - What singing techniques should be annotated?
 - Need datasets, but how to annotate?
 - What is their specific characteristics?
- **2. Less established identification methods**
 - How to design the automatic model?
 - Can we detect techniques from real-world singing tracks?

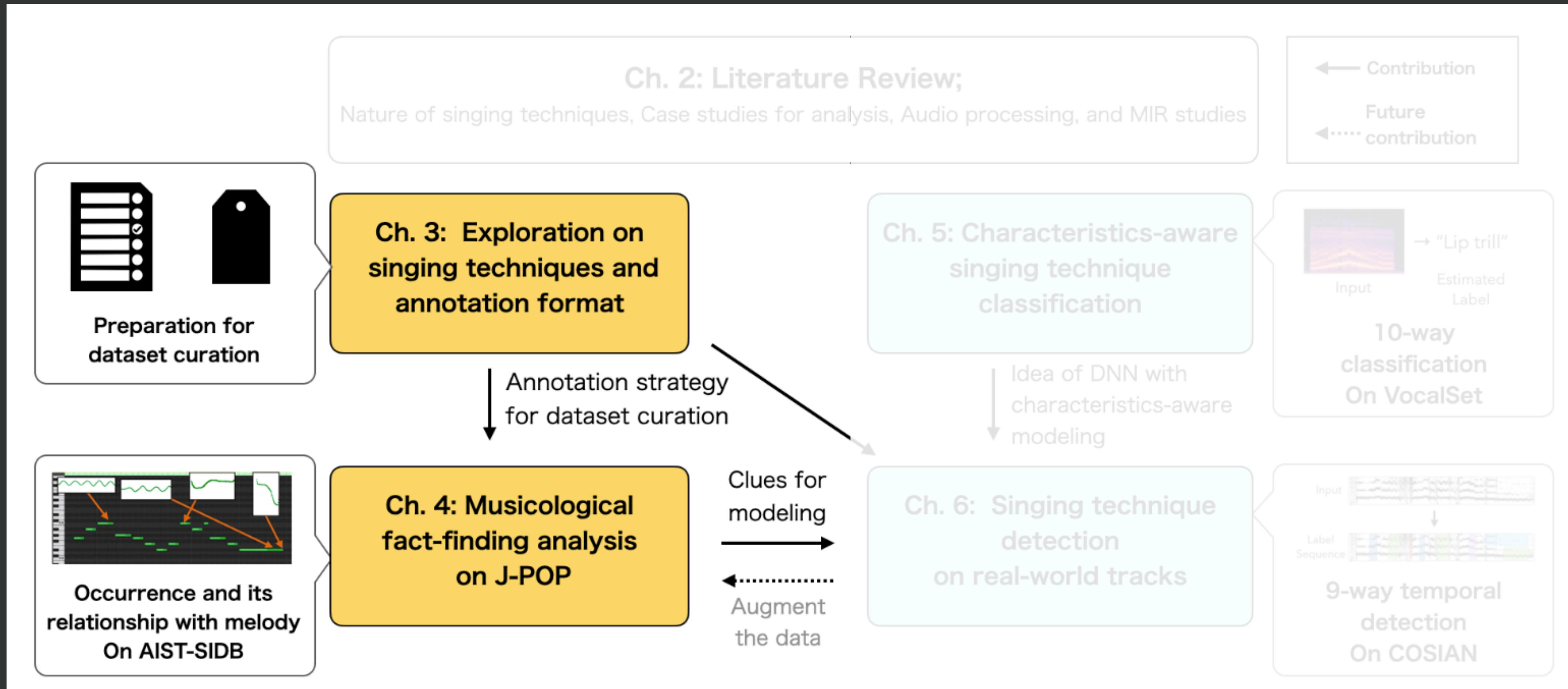
Proposal: Computational foundations for singing technique analysis

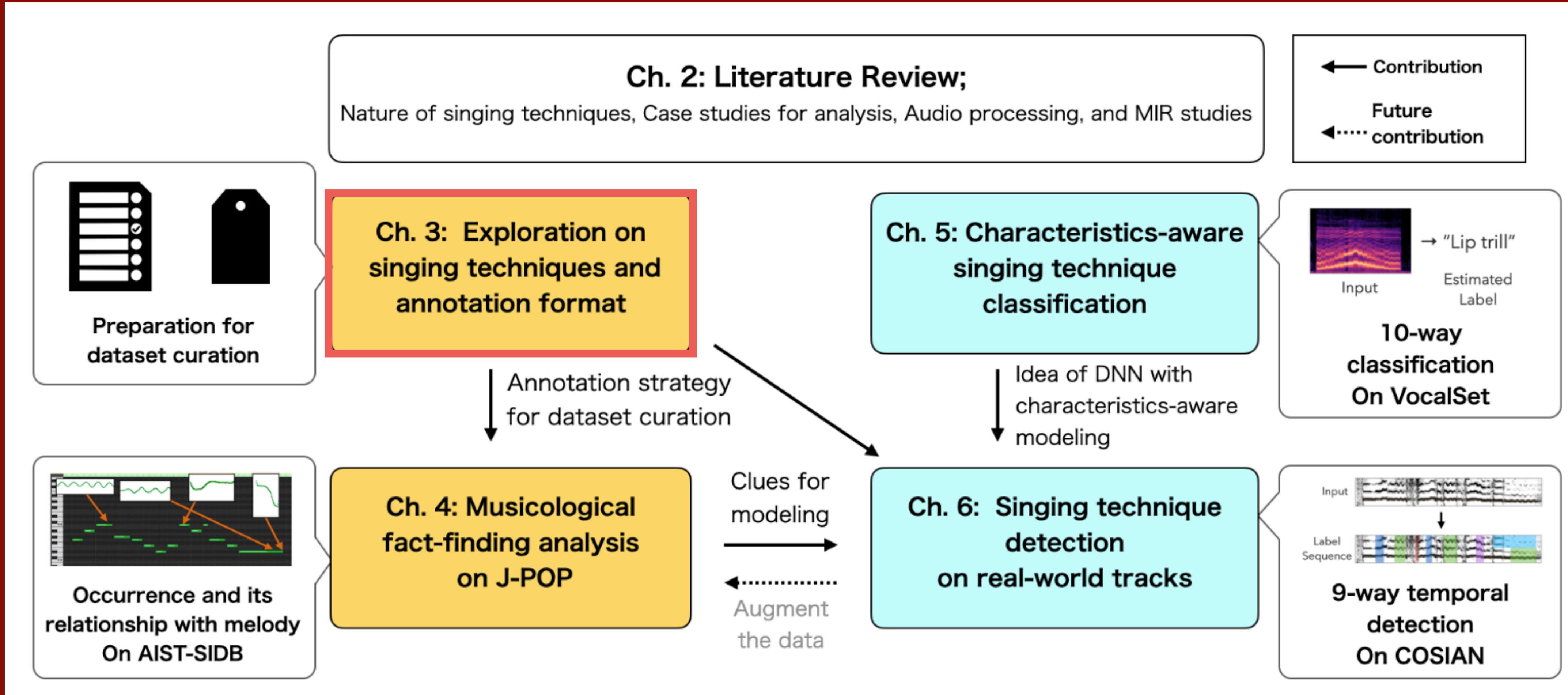
15

Goal of the Thesis



Part I. Data exploration



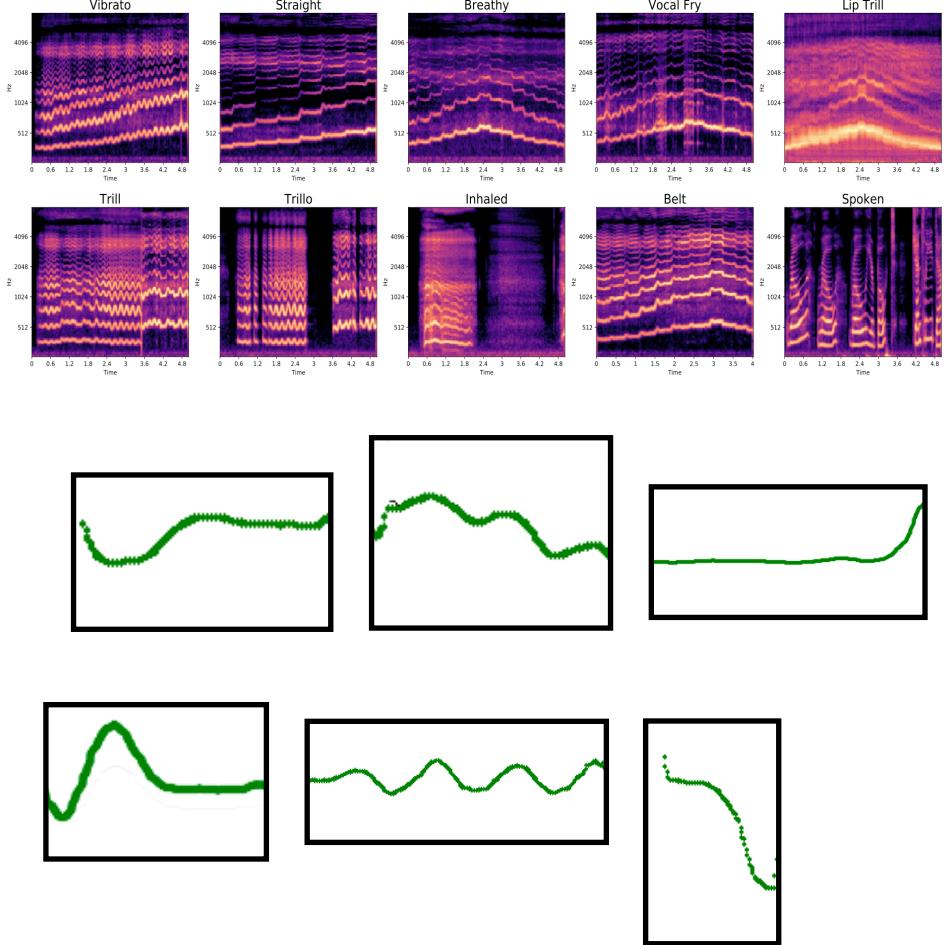


Chapter 3

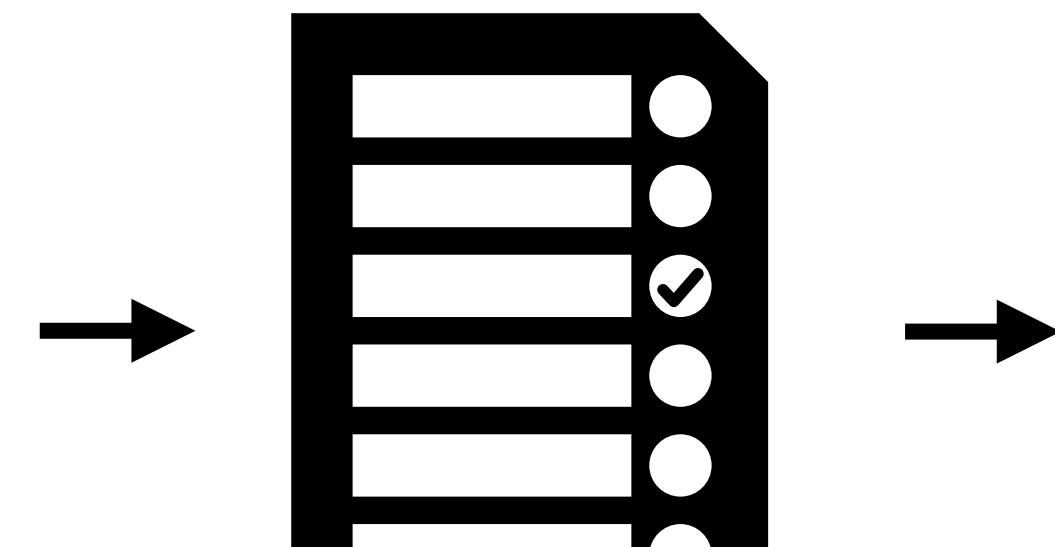
Exploration of Singing Techniques for Dataset Creation

Exploration of existing singing techniques to build datasets

Existing
singing techniques

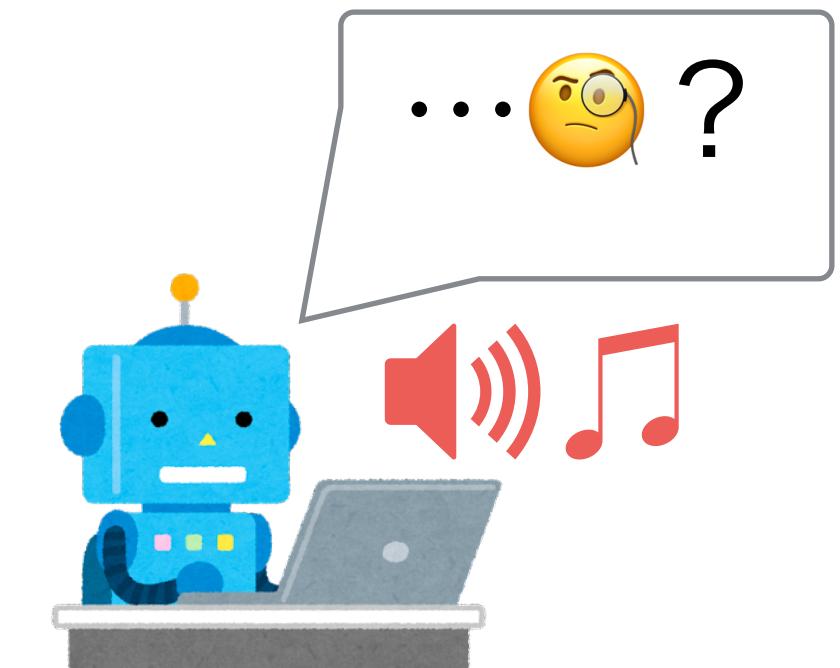


1. Making the
categorized list of
singing techniques



RQ1. What should be
considered?

2. Determine how to
represent in computer



RQ2. Which format is suitable for
singing technique analysis?

Goal: Establish annotation
strategy of datasets

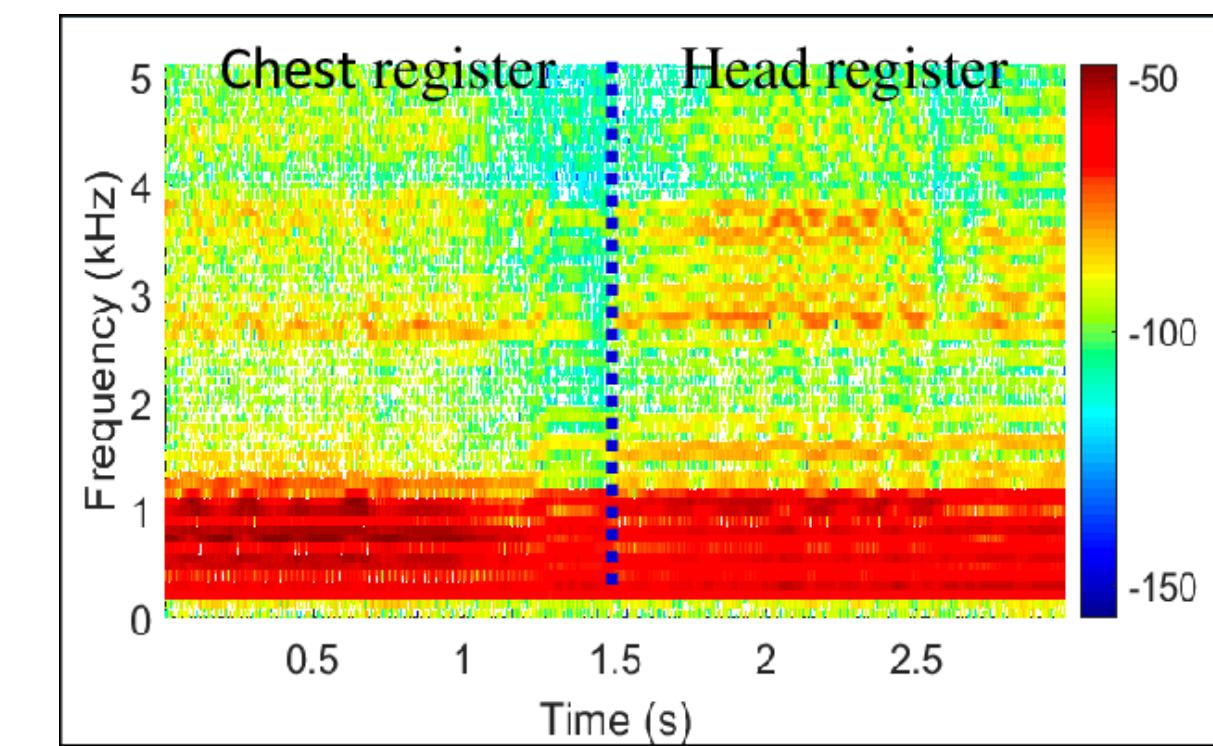
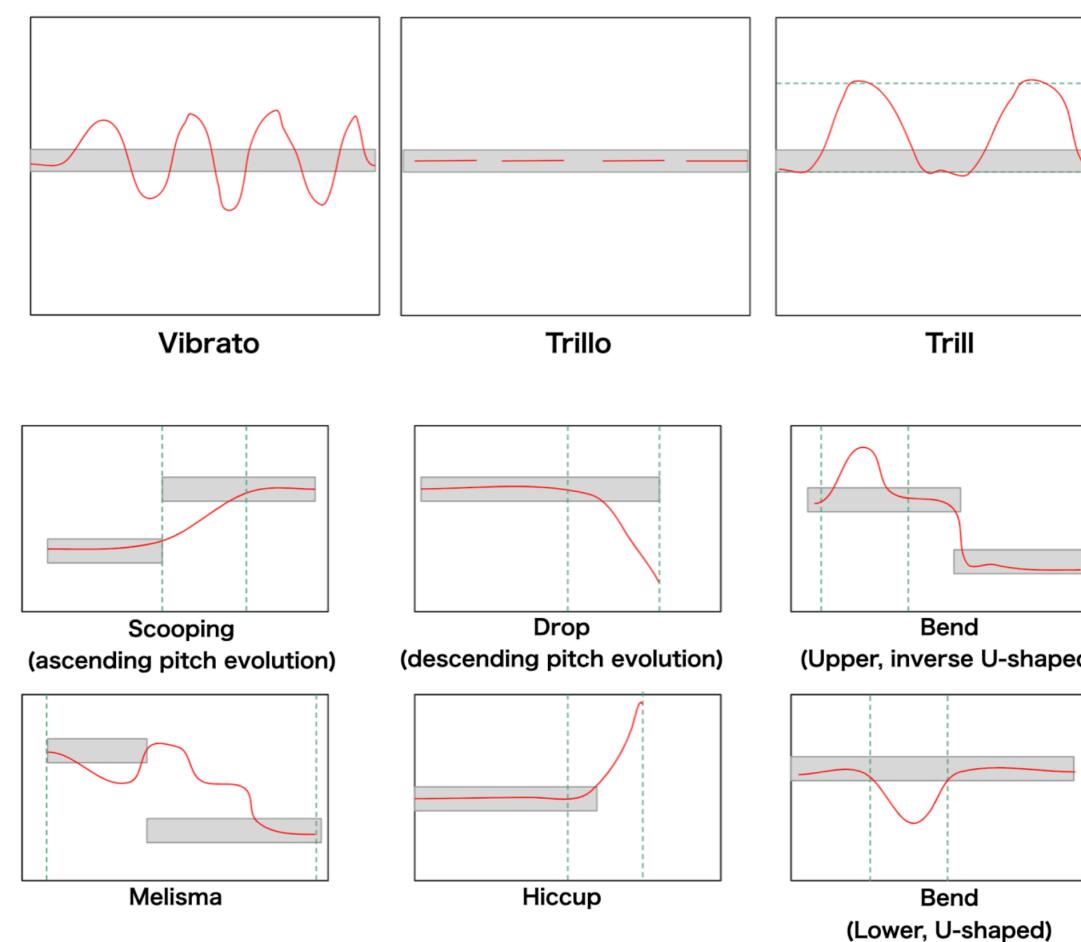


- **Contributions:**

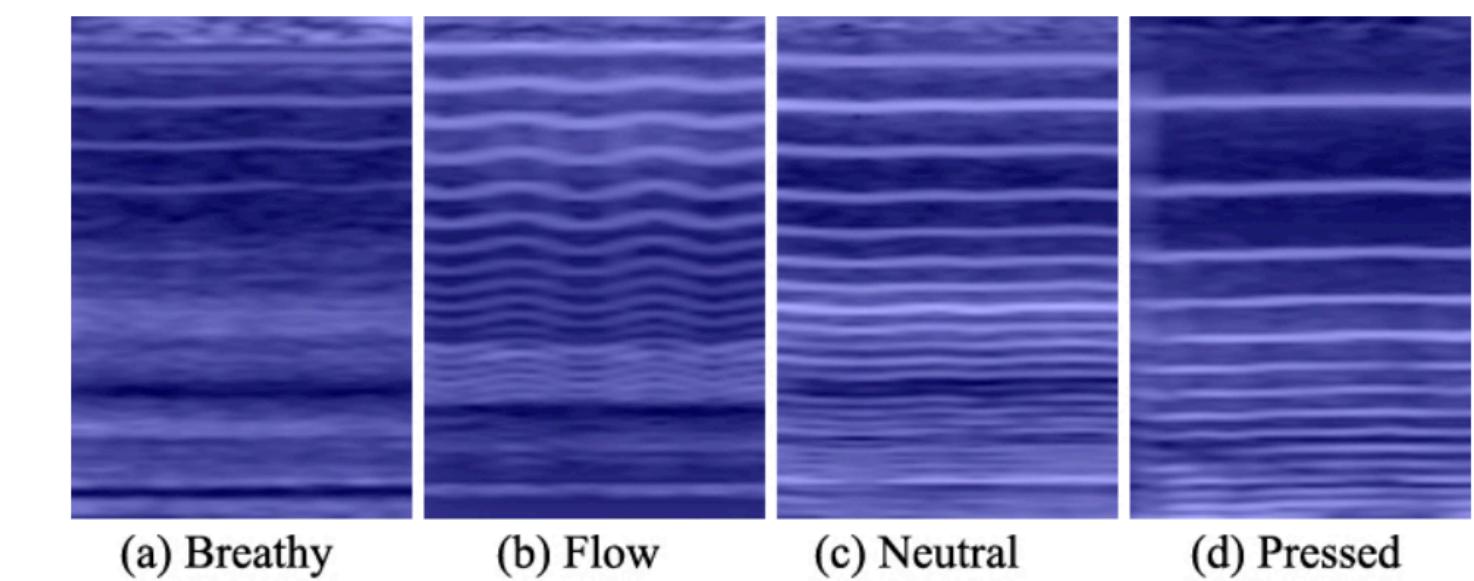
- 1. Listed up named singing techniques from various musical aspects (e.g., pitch, timbre, etc.) -> RQ1
- 2. Adopt region label (name+ start & end time) for annotation -> RQ2

Listed named techniques by literature survey both on **academic** and **non-academic**

- Explored on various named techniques in the world
- Categorize coarsely by what is the fluctuated components
 - Pitch: modulation, portamento (continuously changing), bending.
 - Timbre: register, phonation mode, extreme effects, etc.
 - Others



R. Elbarougy, Acoustic Analysis for Chest-to-Head Register Transition in Singing Voice International Journal of Computer Applications 177(10):11-16, 2019



X. Sun et al. RESIDUAL ATTENTION BASED NETWORK FOR AUTOMATIC CLASSIFICATION OF PHONATION MODES, arxiv 2021.

Candidates of singing techniques

20

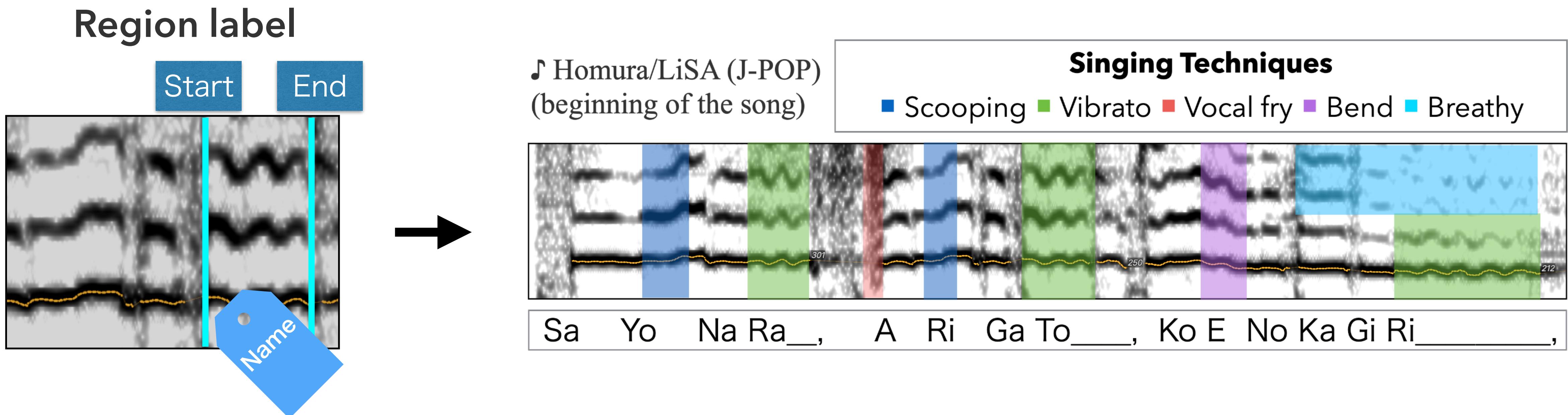


For more detail:
Refer to the thesis!

Annotation strategy

21

- 1. Represent by **region labels** (name+ start&end time boundaries)
- 2. Based on **observable singing techniques** in data



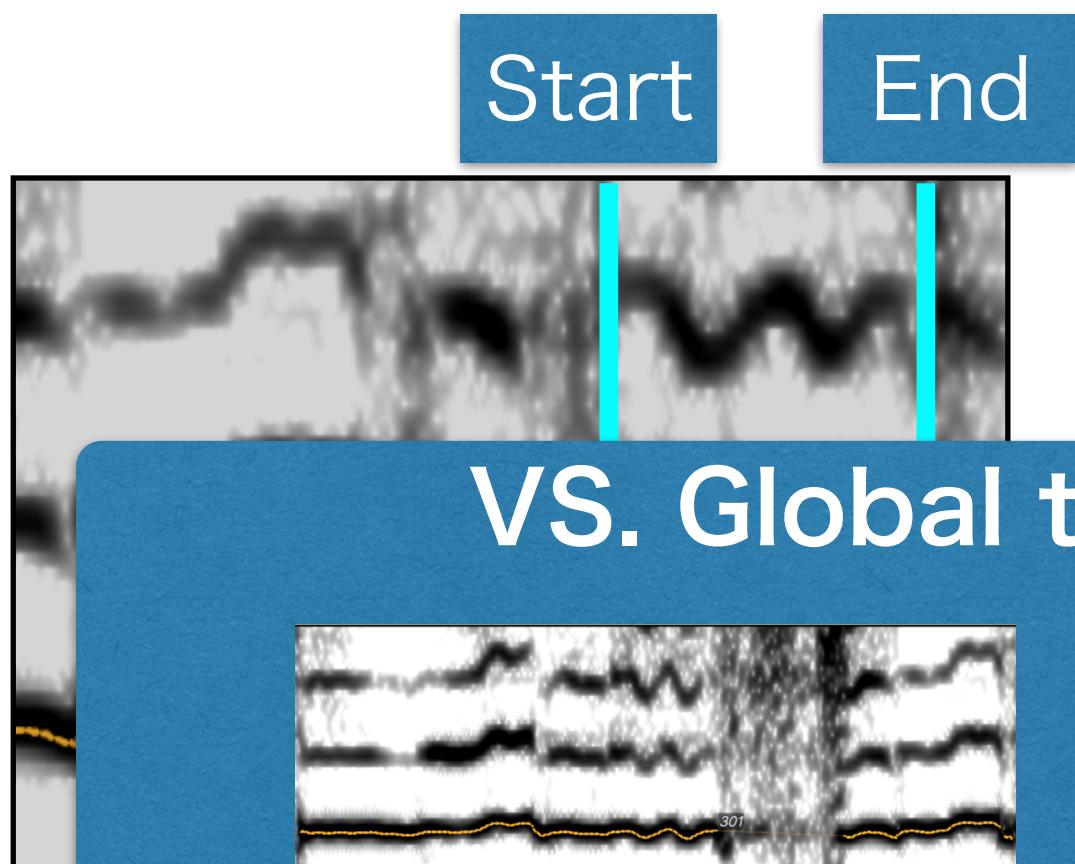
- Advantages: ① Can represents local phenomena, ② Intuitive, ③ Note-information-free

Annotation strategy

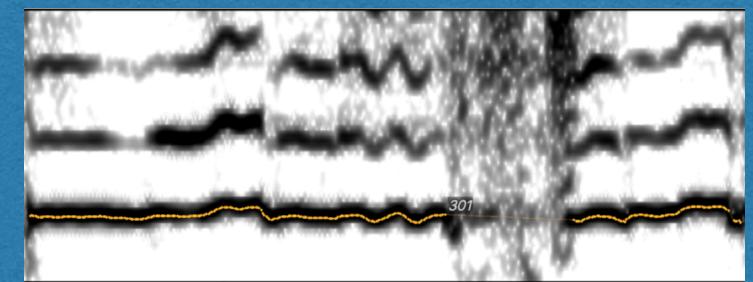
22

- 1. Represent by **region labels** (name+ start&end time boundaries)
- 2. Based on **observable singing techniques** in data

Region label



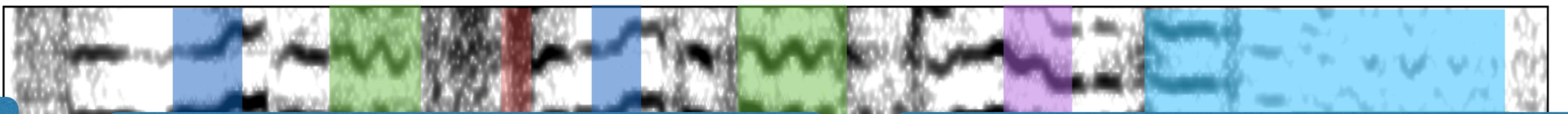
VS. Global tag



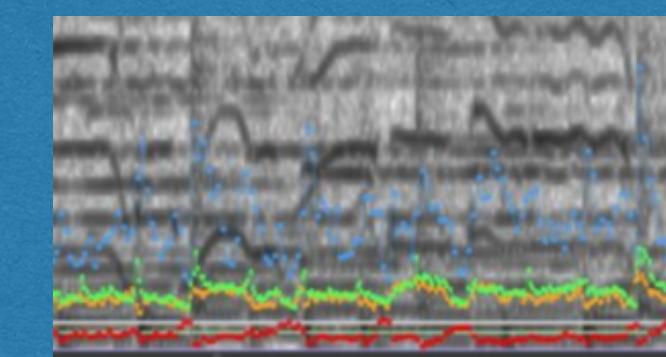
-> vibrato,
Scooping

Can not represents location

♪ Homura/LiSA (J-POP)
(beginning of the song)



VS. Numeric parameters



Less intuitive for humans

Singing Techniques

■ Scooping ■ Vibrato ■ Vocal fry ■ Bend ■ Breathy

VS. Note-wise label



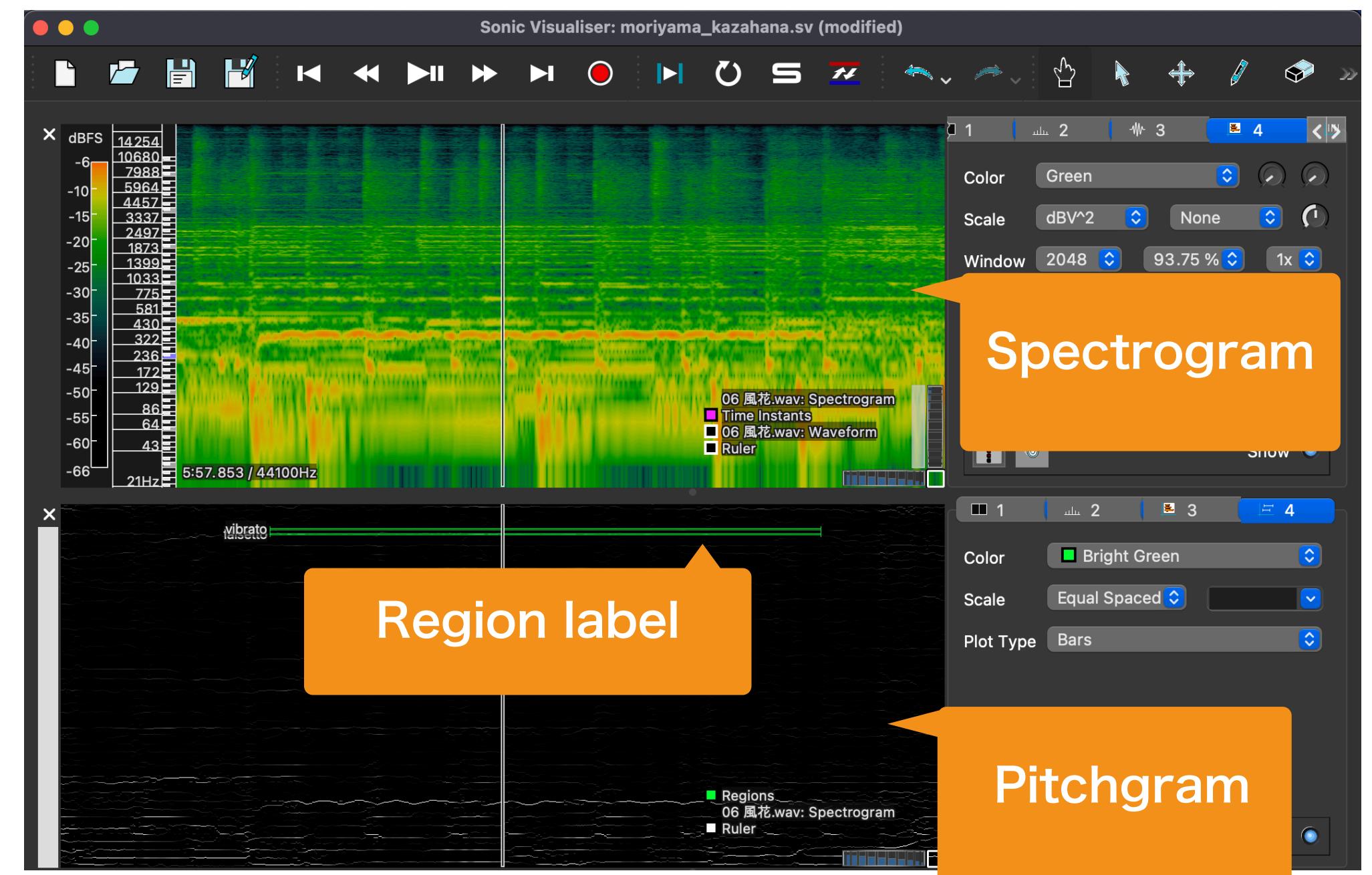
Needs additional note annotation

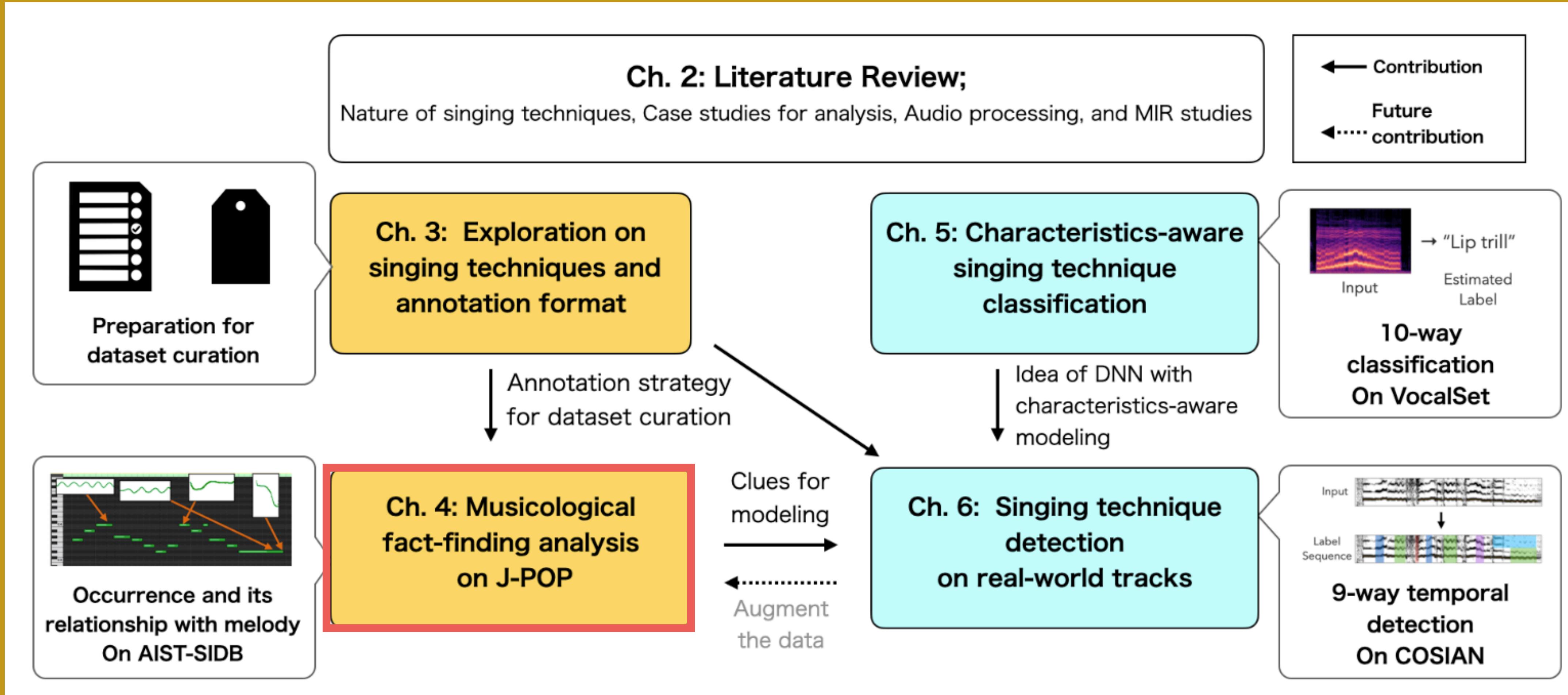
- Advantages: ① Can represents local phenomena, ② Intuitive, ③ Note-information-free

Annotation process

23

- Manually annotated
 - Annotator: author
 - amateur (no academic degree on music)
 - 9 years of popular vocal, 4 years of chorus (tenor),
 - has relative pitch
 - Software: Sonic visualiser [Cannam 10]
 - visualizing spectrogram and pitchgram
 - both aid of visual & audio feedback
 - set region label on pitchgram



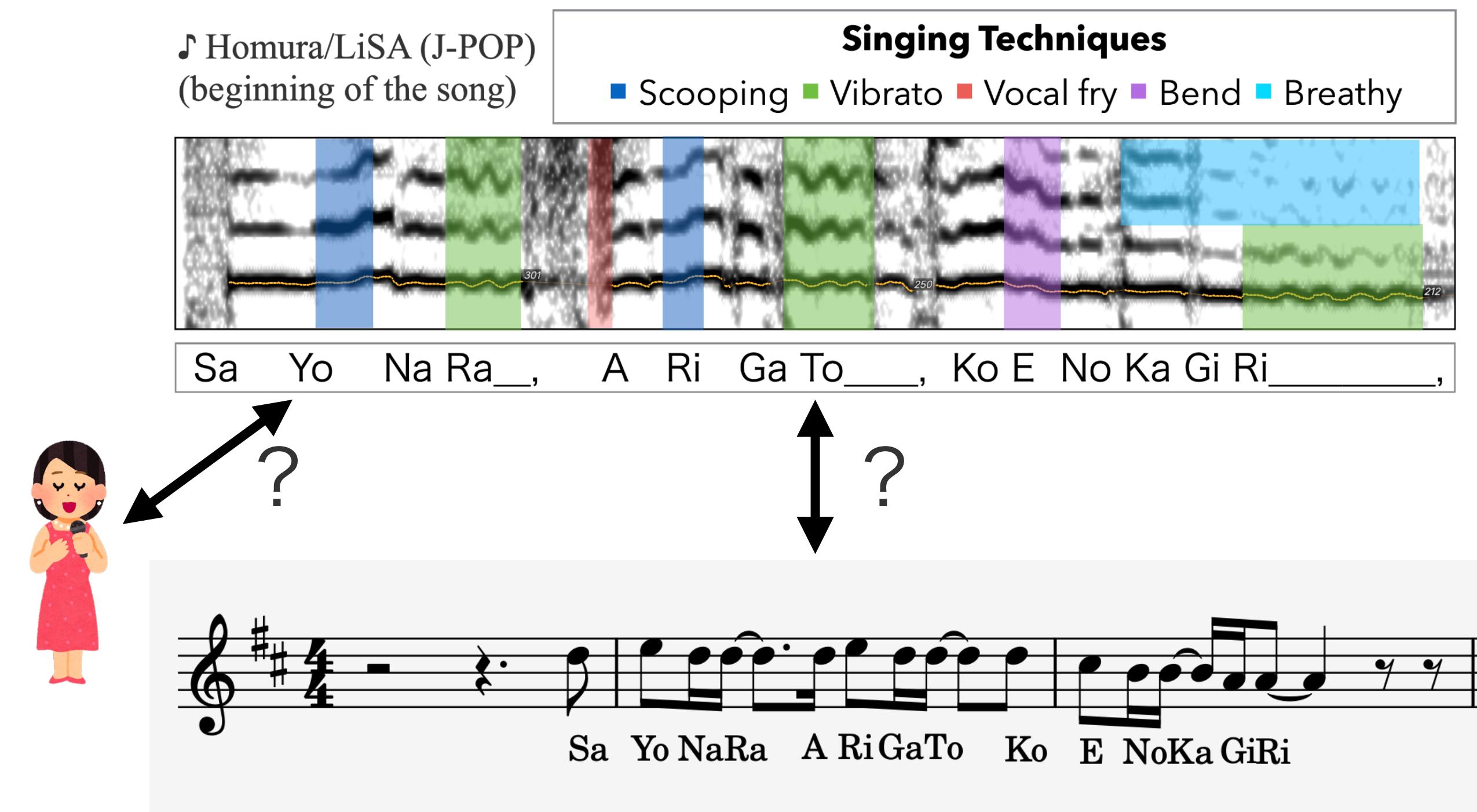


Chapter 4

Musicological Analysis of Singing Techniques with Correspondence to Musical Score on Imitative Singing

Investigate the relationship between **singing technique** and **song**

- **1. Occurrence frequency**
(RQ: How many on each singer?)
- **2. Vibrato parameter**
(RQ: How to produce, where?)
- **3. Occurrence locations**
(RQ: How many where in the song?)
- **Goal: Better understanding singing techniques**
- **Contributions**
 - New discoveries about tendency of singing technique appearance
 - Enhanced utility of singing techniques annotation and their analysis by showing the relationship with singer and song



- **Imitation of J-POP famous singers**

- A Cappella, studio quality
- Original : 24 (12 for each gender)
- Imitator: professional singer (7 F/M)
- 48 tracks (two imitator per song)
- **Private database, possessed by AIST**

- **Examples**

Keisuke Kuwata/ Katte ni Sinbat

Imitator1 **Imitator2**

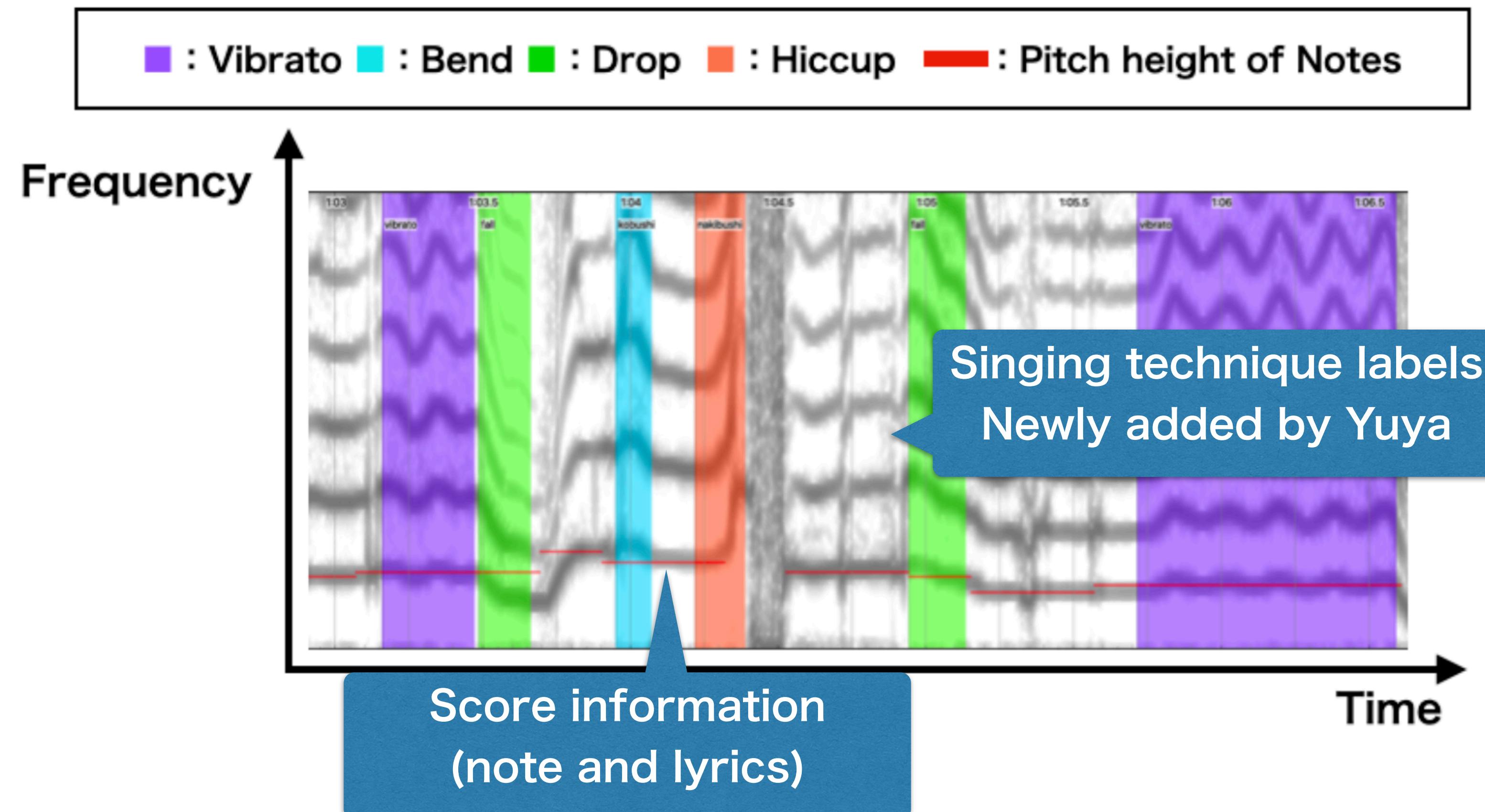
♪ Su Na Ma Ji Ri No Chi Ga Sa Ki
Hi To Mo Na Mi Mo Ki E Te…



Original singer	Song name	Gender	Imitator1	Imitator2
玉置浩二 (Tamaki)	出逢い	M	M03	M04
小田和正 (oda)	キラキラ	M	M02	M06
Gackt (gackt)	ありったけの愛で	M	M01	M05
桑田佳祐 (sazan)	勝手にシンドバット	M	M06	M07
チバユウスケ (skapara)	カナリヤ鳴く空	M	M05	M01
西川貴教 (tmr)	Heat Capacity	M	M05	M01
hyde (larcenciel)	Lies and Truth	M	M03	M04
平井堅 (hirai)	瞳を閉じて	M	M01	M05
福山雅治 (fukuyama)	桜坂	M	M04	M03
楳原敬之 (makihara)	桃	M	M04	M03
森山直太朗 (moriyama)	さくら (独唱)	M	M02	M06
山崎まさよし (yamazaki)	未完成	M	M06	M07
aiko (aiko)	ボーイフレンド	F	F01	F06
絢香 (ayaka)	三日月	F	F05	F02
宇多田ヒカル (utada)	Can You Keep A Secret?	F	F03	F04
鬼束ちひろ (onitsuka)	月光	F	F03	F04
倖田來未 (koda)	夢のうた	F	F06	F01
小柳ゆき (koyanagi)	愛情	F	F04	F07
chara (chara)	大切をきずくもの	F	F02	F05
浜崎あゆみ (hamasaki)	seasons	F	F06	F01
一青窈 (hitotoyo)	ハナミズキ	F	F05	F02
平原綾香 (hirahara)	明日	F	F03	F04
松浦亜弥 (matsuura)	♡ 桃色片思い ♡	F	F01	F06
YUKI (jam)	motto	F	F02 (key-1)	F05

Annotation for AIST-SIDB

27



Annotated techniques

28

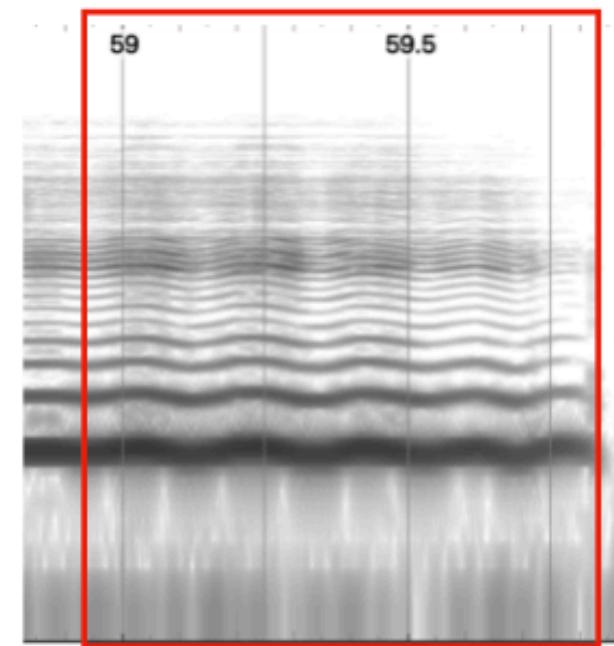
- Annotated 13 singing techniques
- Based on the survey in Chapter 3 and observation of data

Table 4.2: Definitions, synonyms, and comparable techniques for singing techniques.

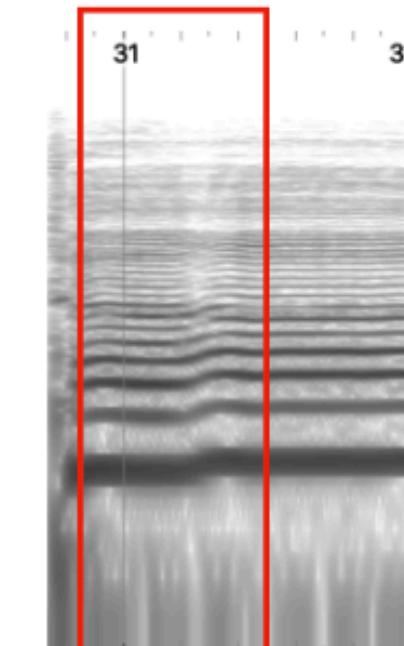
Techniques	What is modulated	How modulates	Discrepancies	Similar techniques
Vibrato	pitch, loudness	Singing with a wavering effect, introducing periodic oscillation	Shake	Tremolo, Trill
Scooping	pitch	Continuously changing pitch upward	Glissando, Portamento, Scoop, Scoop up	
Bend	pitch	Continuously changing pitch in a U-shaped or inverted U-shaped pattern	Bend, Tremolo	
Drop	pitch	Continuously changing pitch downward	Drop, Scoop down	
Hiccup	pitch, timbre	Producing a momentary falsetto or tightened throat singing voice	Cry, Sob, Vocal break	Yodel
Melisma	pitch	Assigning multiple pitches to a single syllable	Fake	
Vocal fry	timbre	Producing a raspy sound	Edge voice, Creaky voice	Growl
Falsetto	timbre	Singing in the falsetto range	Head voice, Falsetto	
Breathy	timbre	Mixing in breathy sounds		
Whisper	timbre	Singing in a whispering manner		
Shout	—	Shouting		
Spoken	—	Singing in a spoken manner		Rap
Tongue trill	—	Using rolled tongue	Tongue roll, Rolled tongue	Lip roll

Spectrograms of each technique

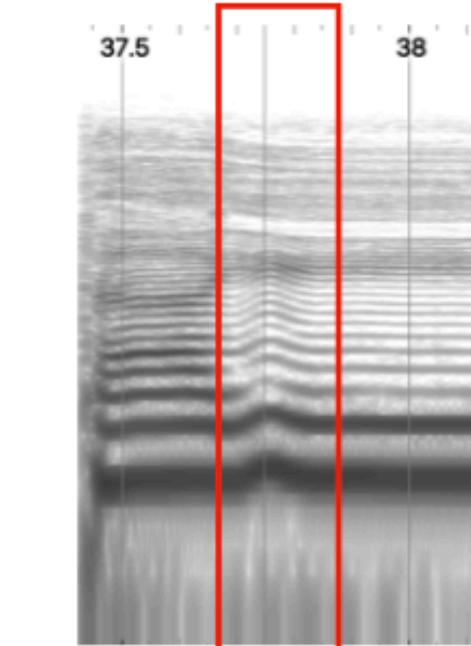
29



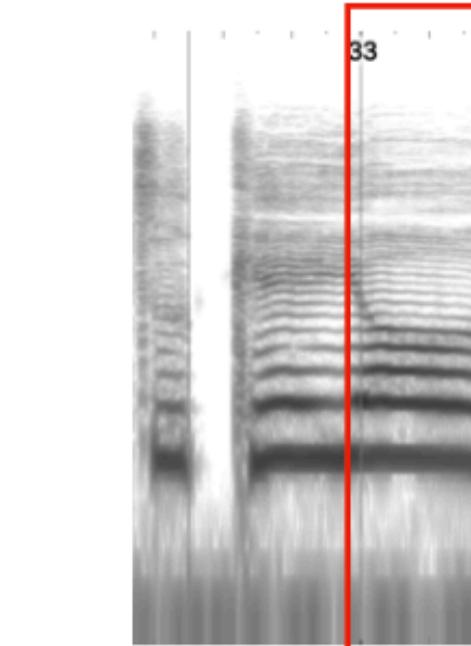
Vibrato
ビブラート



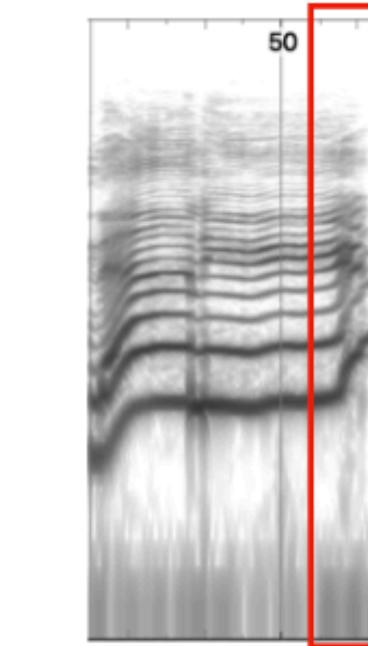
Scooping
しゃくり



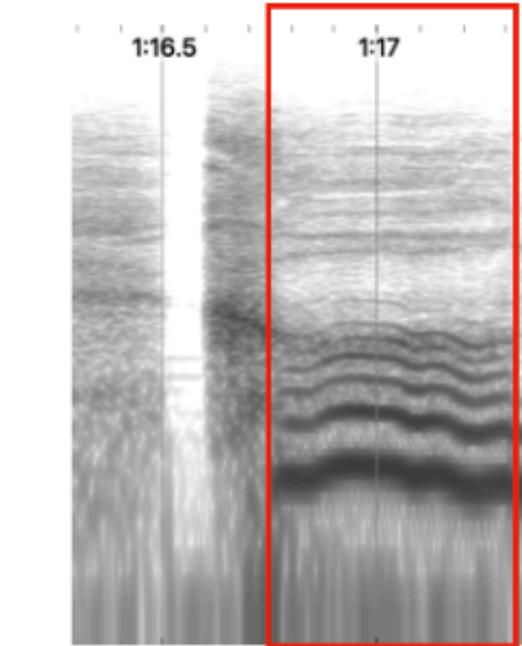
Bend
こぶし



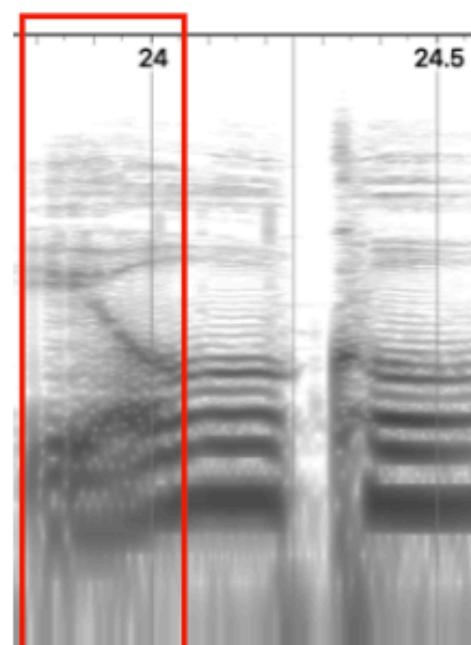
Drop
フォール



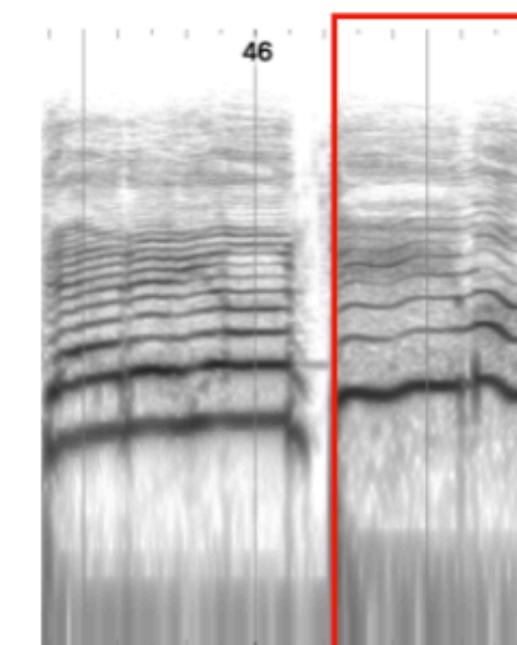
Hiccup
ヒーカップ



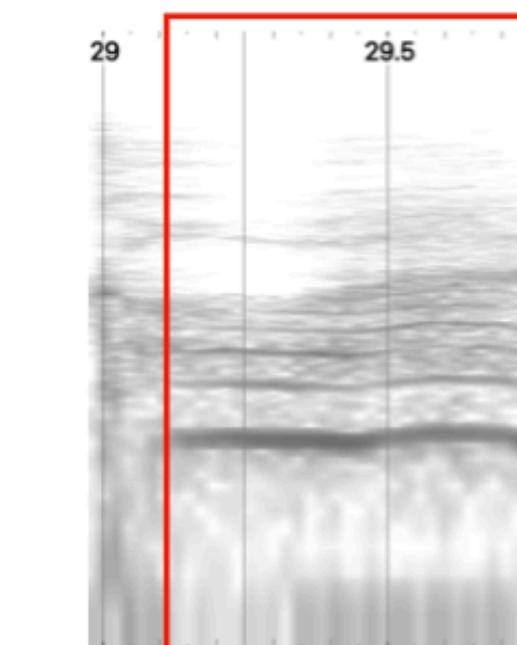
Melisma
メリスマ



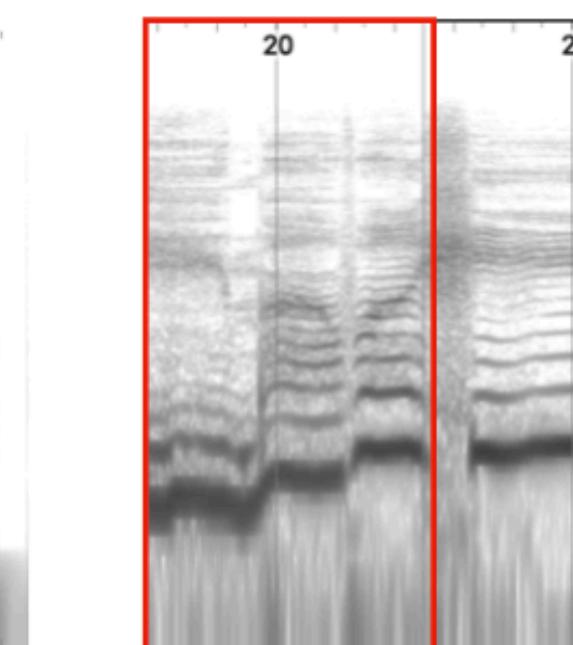
Vocal fry
ボーカル
フライ



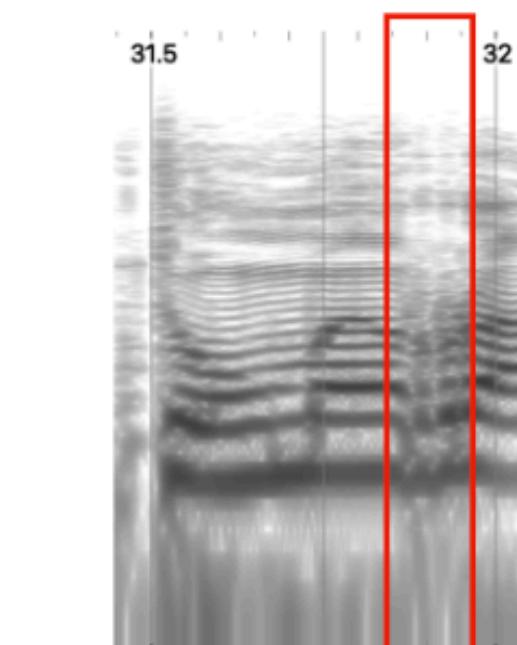
Falsetto
ファル
セット



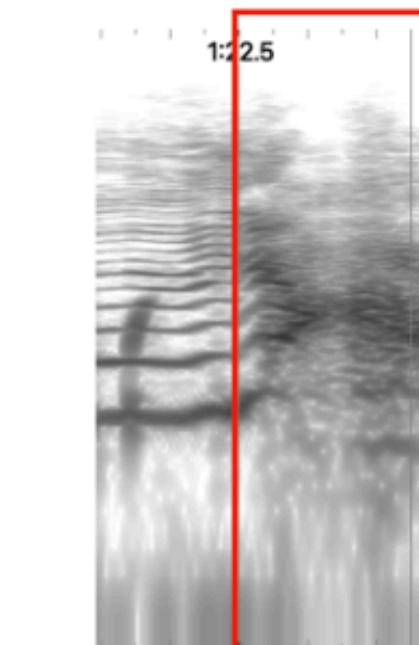
Whisper
ウィスパー



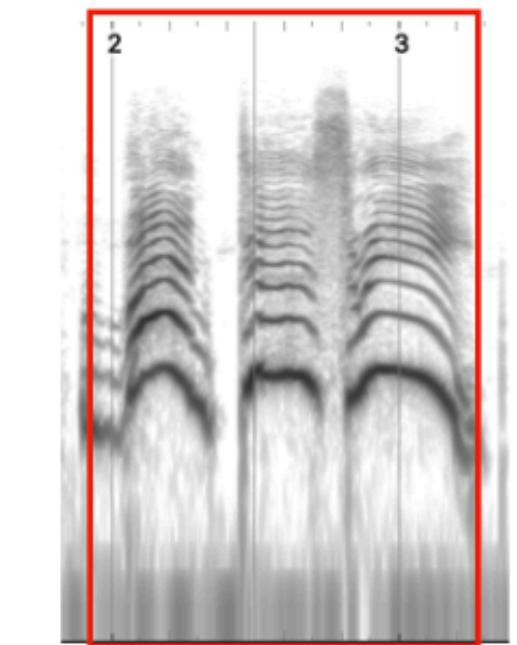
Breathy
ブレシー



Tongue trill
タング
トリル



Shout
シャウト



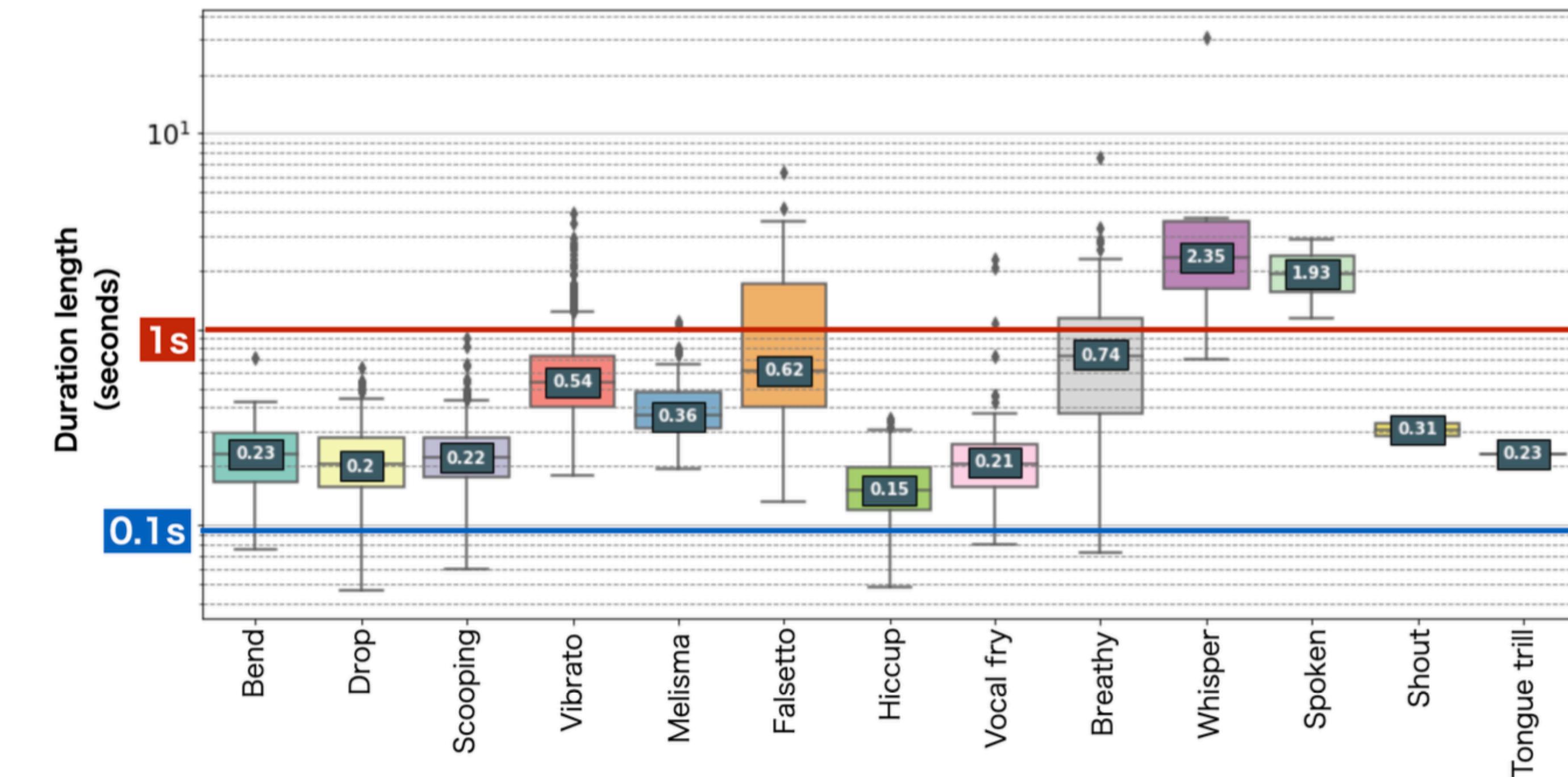
Spoken
セリフ語り

Analysis 1: Occurrence of techniques

30

Whole statistics

Techniques	Number of labels	Total duration [s]	Average duration [s]
Vibrato	717	448.57	0.63
Scooping	528	118.40	0.24
Bend	144	33.14	0.23
Drop	140	31.25	0.23
Hiccup	126	20.35	0.16
Melisma	38	16.36	0.44
Whisper	11	54.5	4.95
Falsetto	86	96.16	1.14
Breathy	52	41.57	1.03
Vocal fry	82	21.73	0.28
Tongue trill	1	0.36	0.23
Shout	2	1.16	0.39
Spoken	4	13.71	1.98

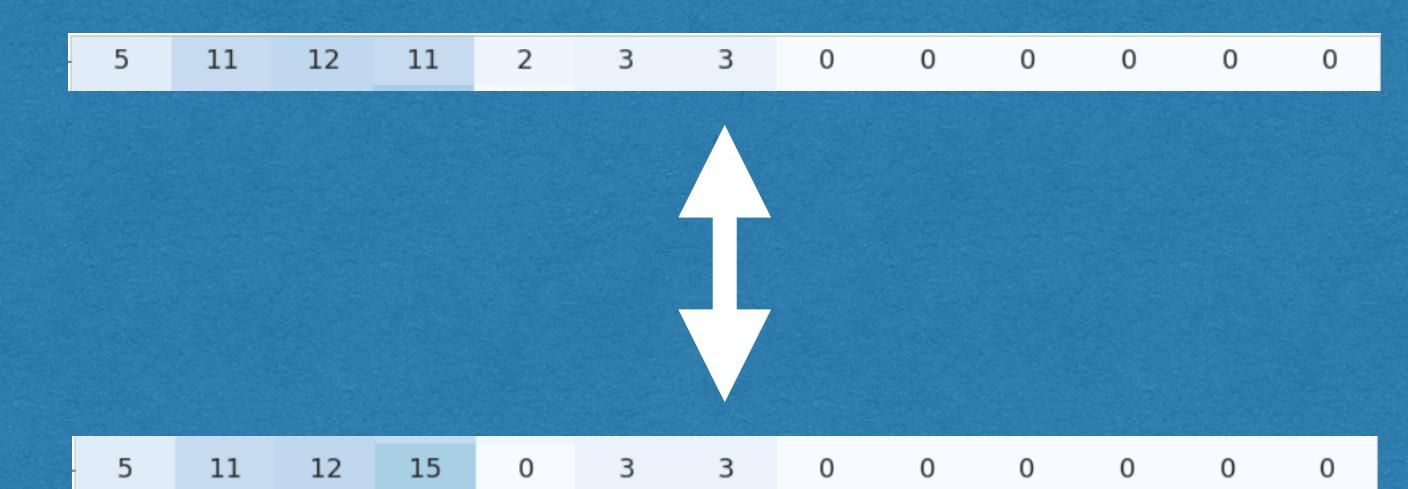


- Pitch techniques (vibrato, scooping, bend, drop) are frequent
- Most of techniques are short length (0.1s - 1s)

Analysis 1: Occurrence of techniques

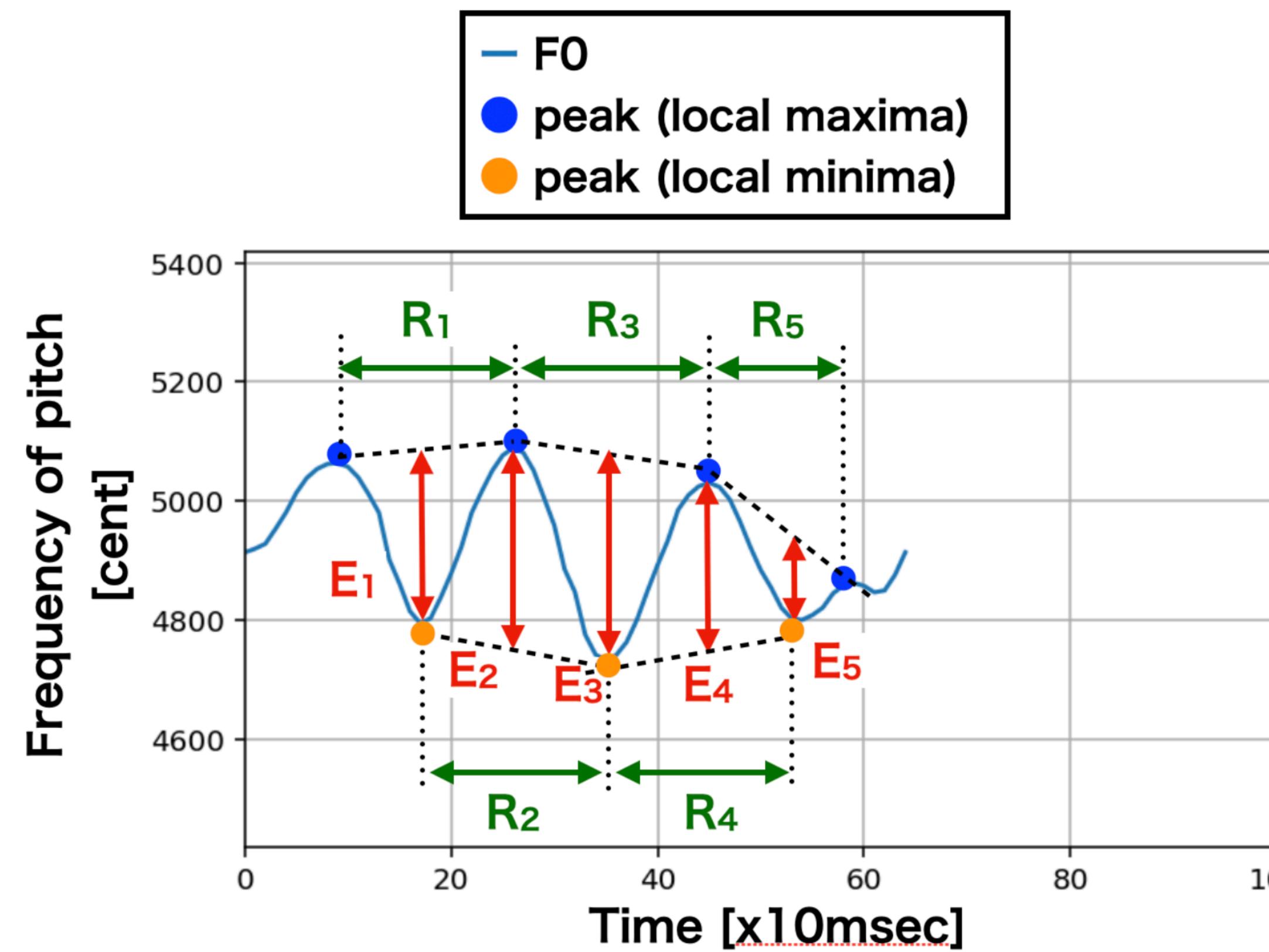
- The distribution is different by singer
 - The distributions are similar between same or similar style original singer

When imitating the same singer,
occurrence of singing techniques is
also likely to be similar ?



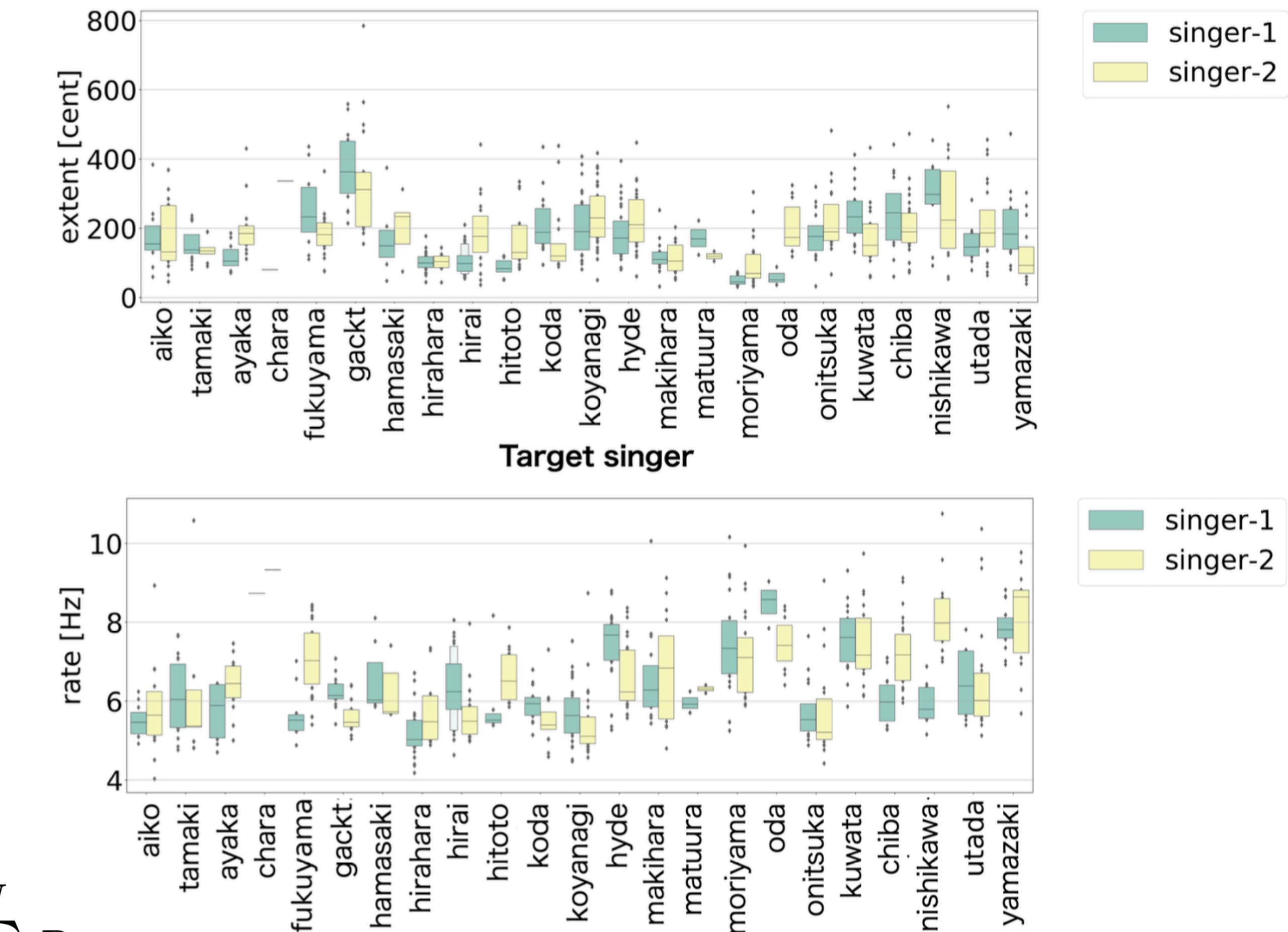
Cosine similarity
Between the occurrences
Avg. -> Intra: 0.83, Inter: 0.7

Picked up the vibrato labels, calculates the vibrato parameters (extent and rate)



Calculation

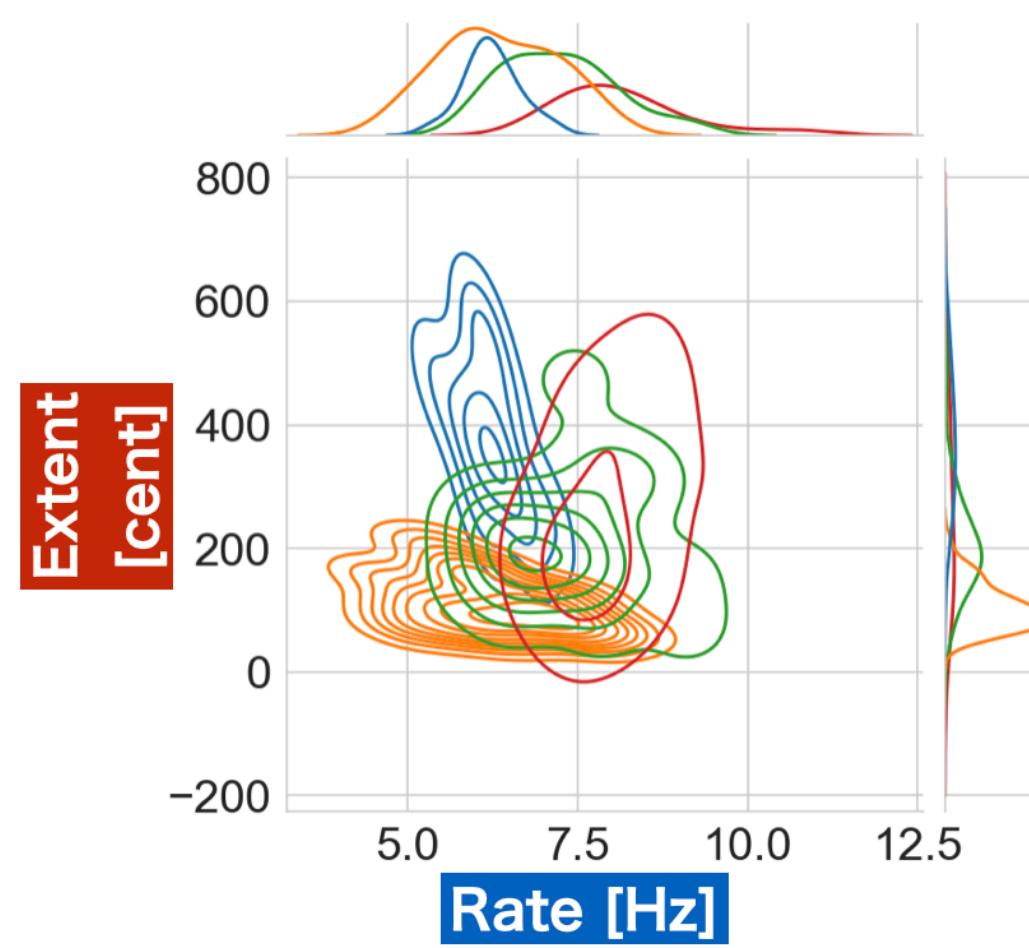
$$\text{extent} = \frac{1}{N} \cdot \sum_{n=1}^N E_n, \quad \text{rate} = \frac{1}{N} \cdot \sum_{n=1}^N R_n$$



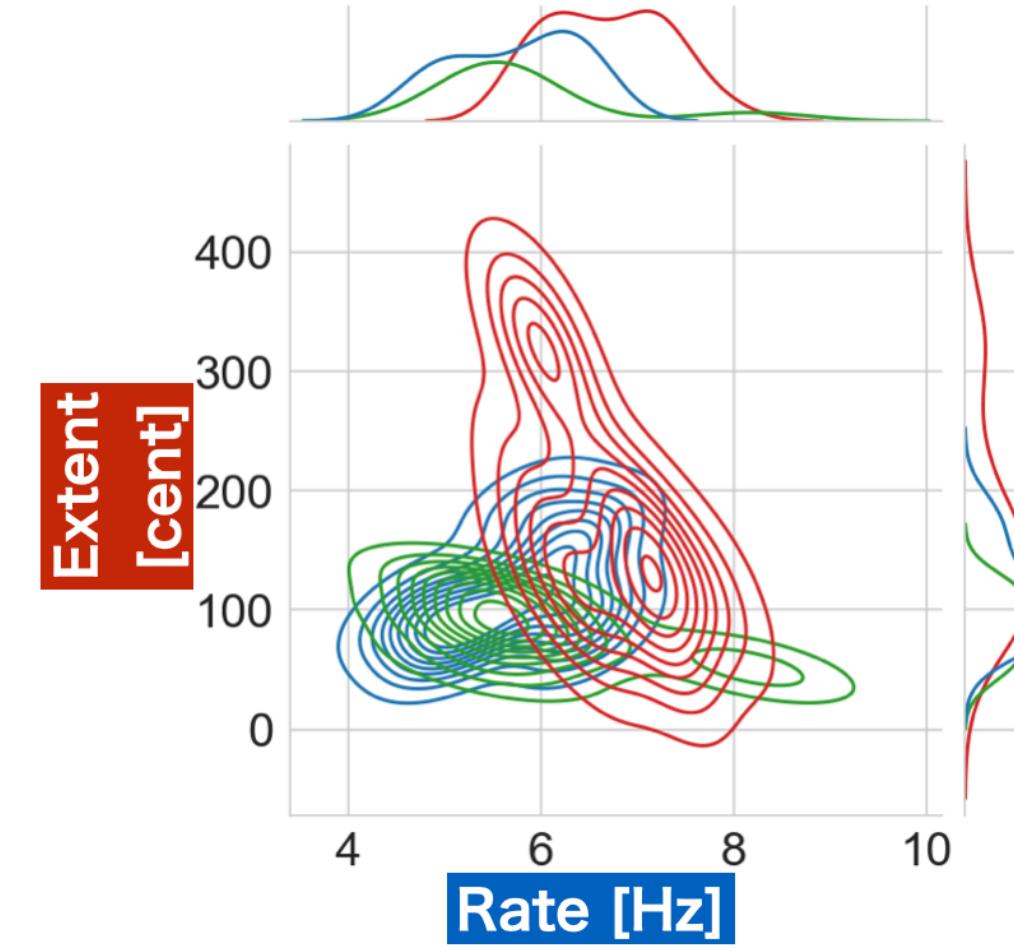
Summary of discovery

**Kernel distribution estimation (KDE)
of rate-extent plot**

Imitation singer ID: M01



Imitation singer ID: F05



**Pearson's correlation coefficients
with pitch and duration**

	Pitch (female)	Pitch (male)	Pitch (Normalized)	Label duration
Extent	-0.424	-0.336	-0.31	-0.127
Rate	0.225	0.169	0.025	-0.268

Negative correlation between...

- Vibrato extent and pitch height (deep \propto low)
- Vibrato rate and label duration (slow \propto long)

The shape changes when the imitated subject is different

-> professional may be able to imitate the vibrato, in terms of the extent and rate \leftrightarrow difficult for amateur singers [Saitou 11]

- Take statistics of co-occurrence of each singing technique
 - Associates the technique labels to each notes
 - by aggregating the label amounts

Ex. 夜に駆ける/YOASOBI

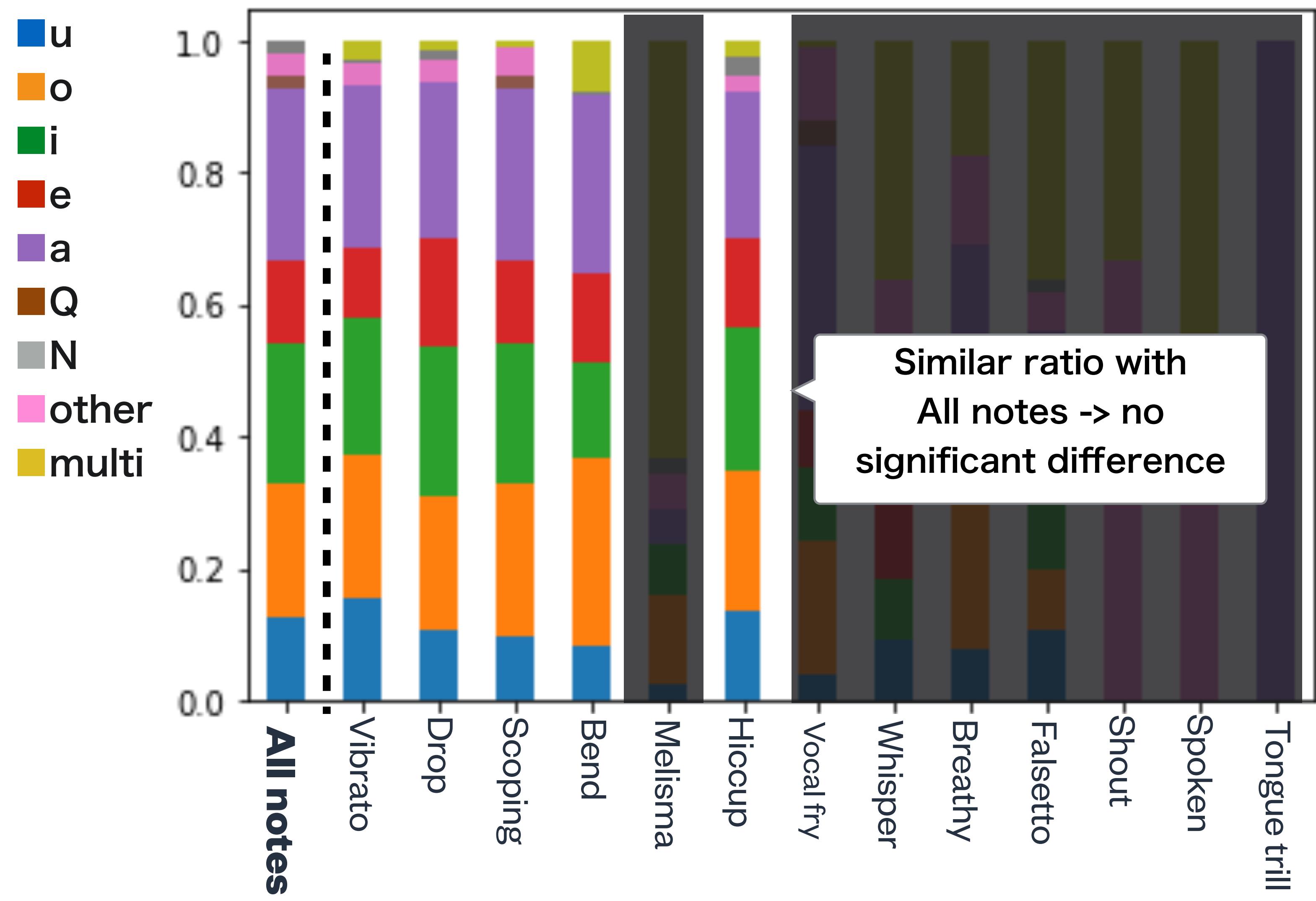
Shi Zu Mu Yo oh Ni To Ke Te Yu KuYo oh Ni

	Lyrics' vowel	Note heights	Note intervals	Duration	Phrase position
★ Vibrato	i x 2, u	C5, F4, Es4	+2, -1, -2	3/8, 1/8, 1/2	mid x 2, tail
★ Bend	o	As	As	3/8	mid

- **Lyric** : No salient correlation
- **Note heights** : More falsetto and scooping on high notes.
- **Note interval**
 - when pitch rose: more falsetto and scooping
 - when next pitch will fall: more falsetto and drop
- **Duration**
 - Many vocal fry, hiccup, bend and drop, on short notes
 - Many vibrato and scooping, on long notes.
- **Phrase position** : More vibrato on phrase end, Less Bend on phrase beginning

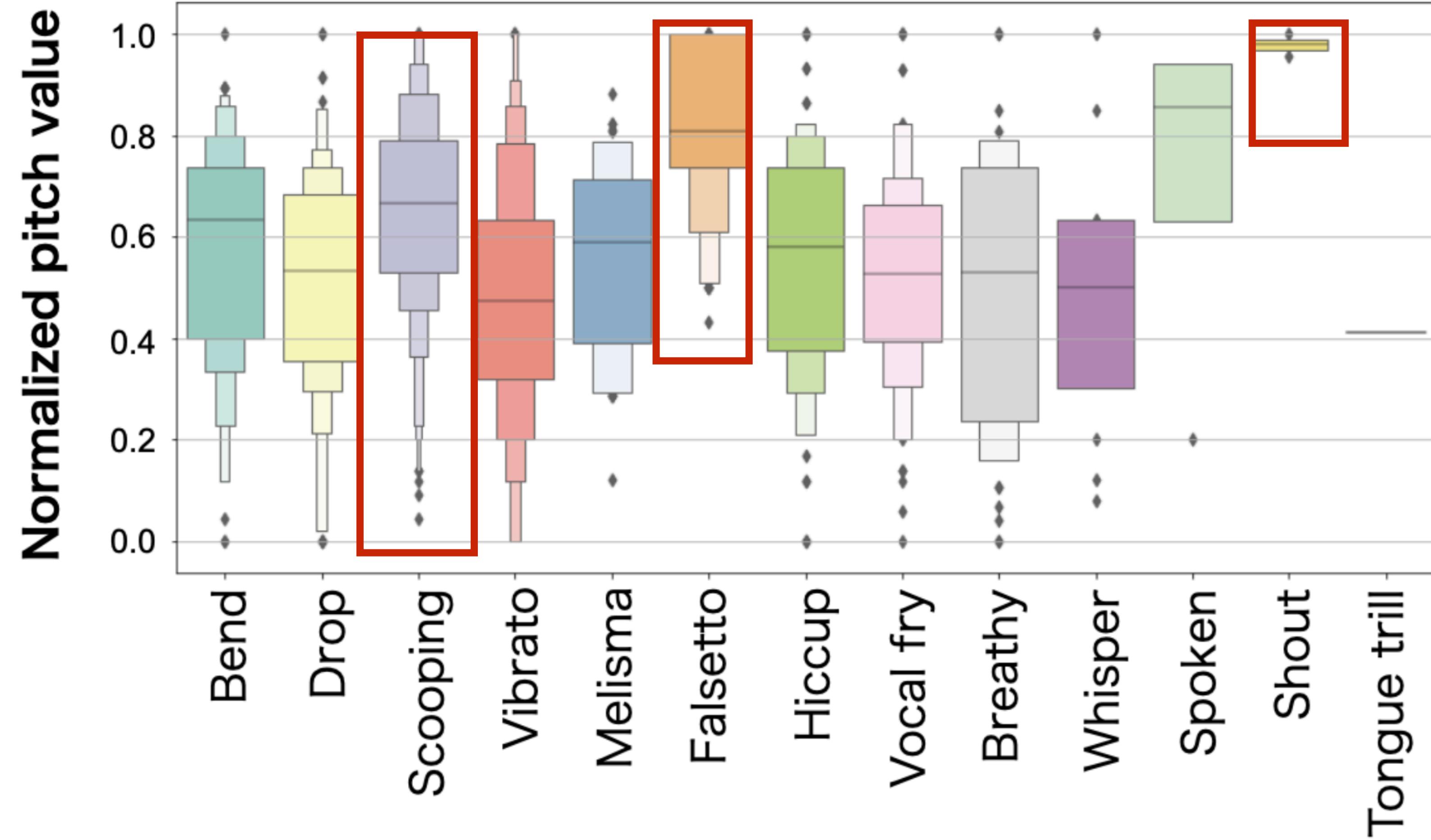
Vowels of note's lyrics

Omitted in presentation due to the time limit



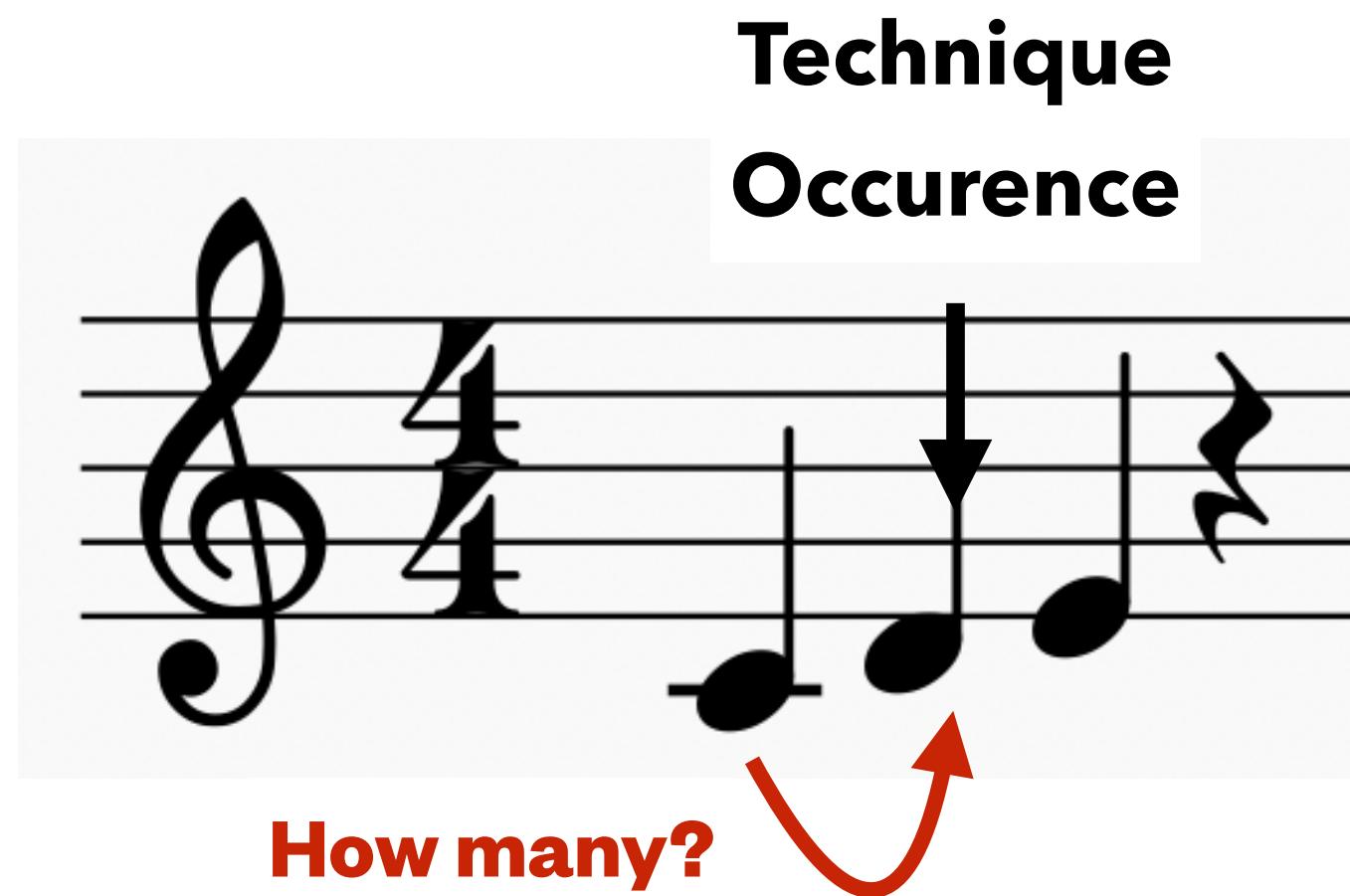
Note heights (after adjustment)

Omitted in presentation
due to the time limit

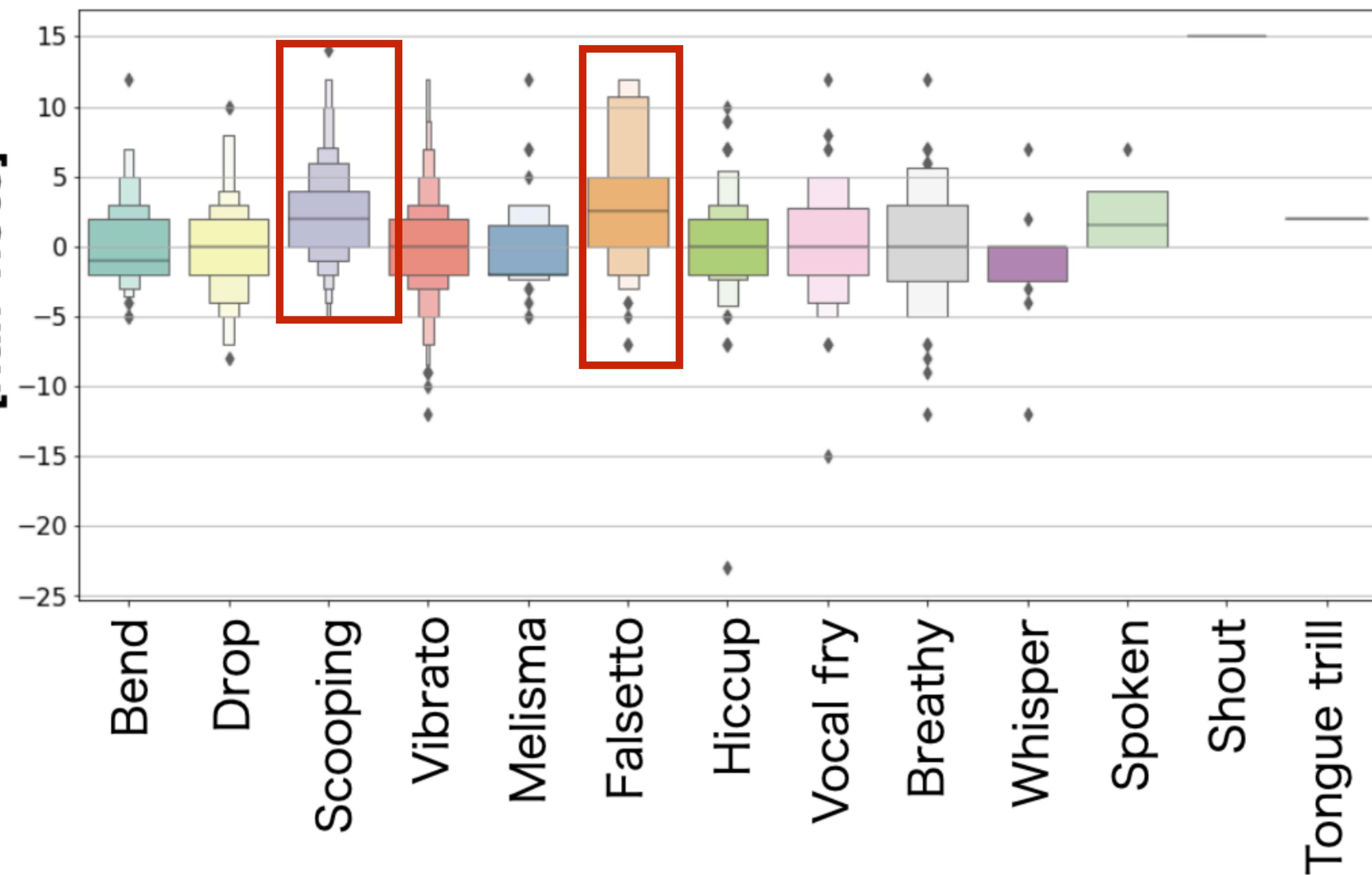


Preceding note interval

Omitted in presentation
due to the time limit

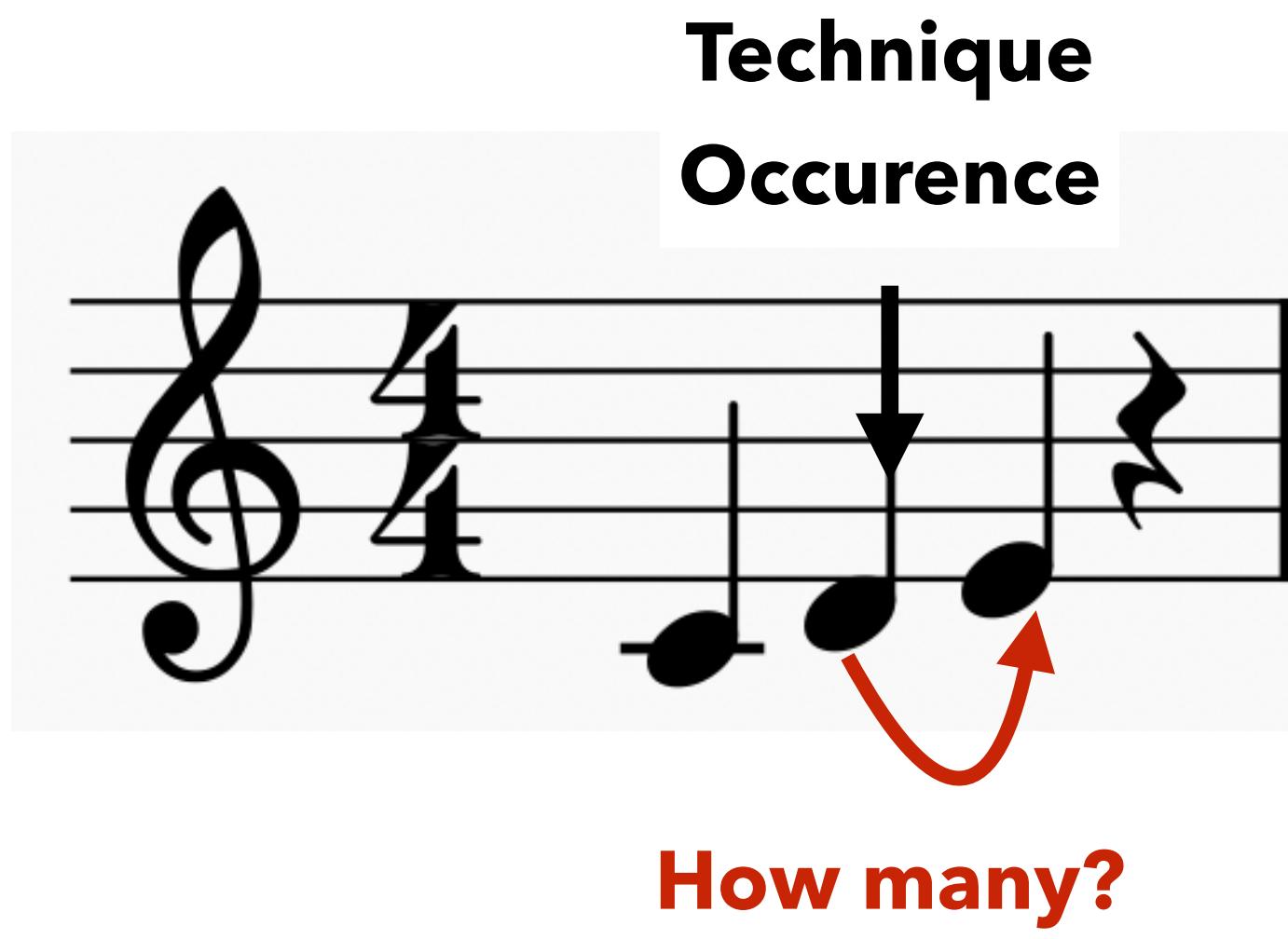


Intervals with the preceding note [half-note]

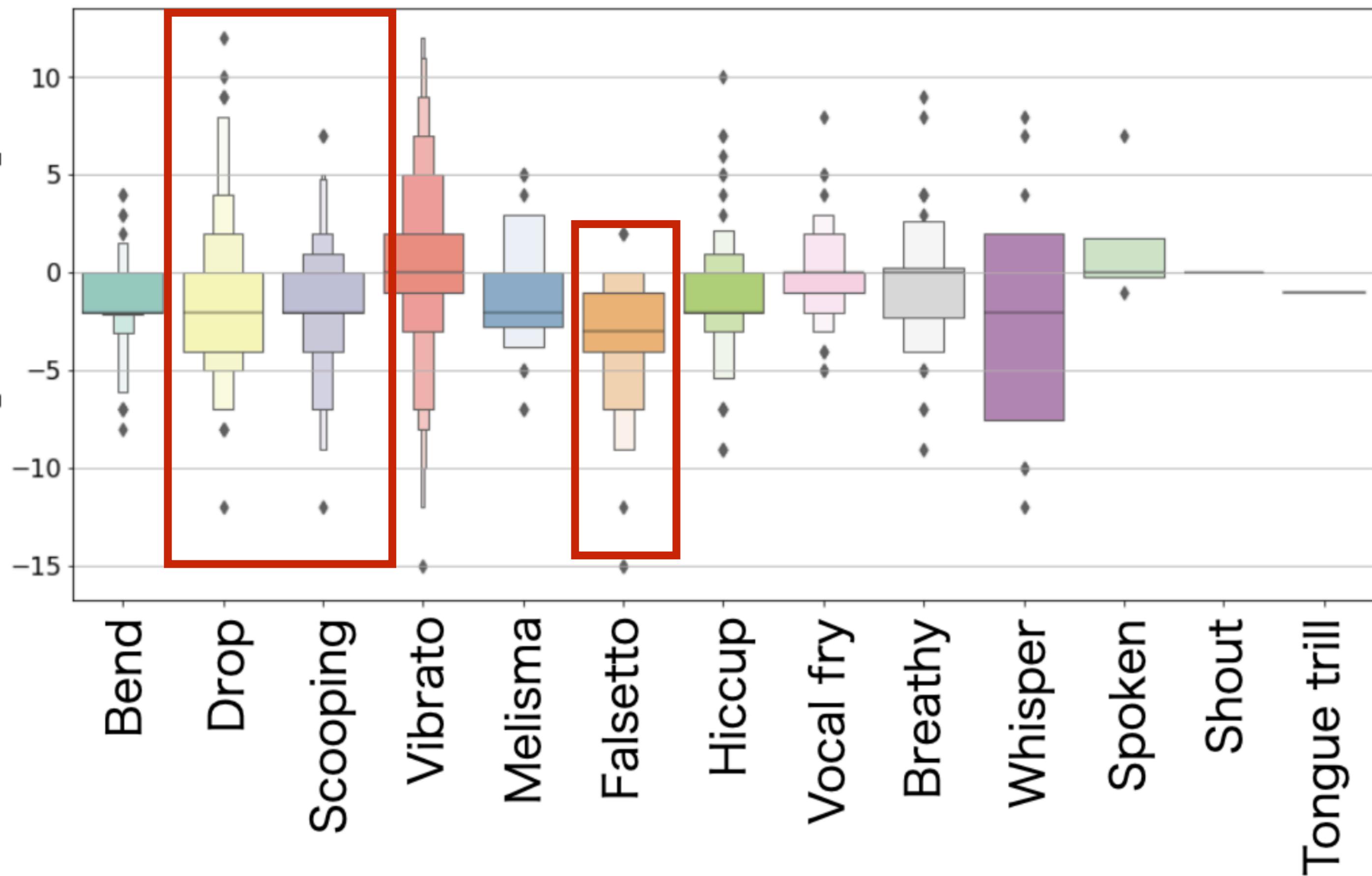


Following note interval

Omitted in presentation
due to the time limit

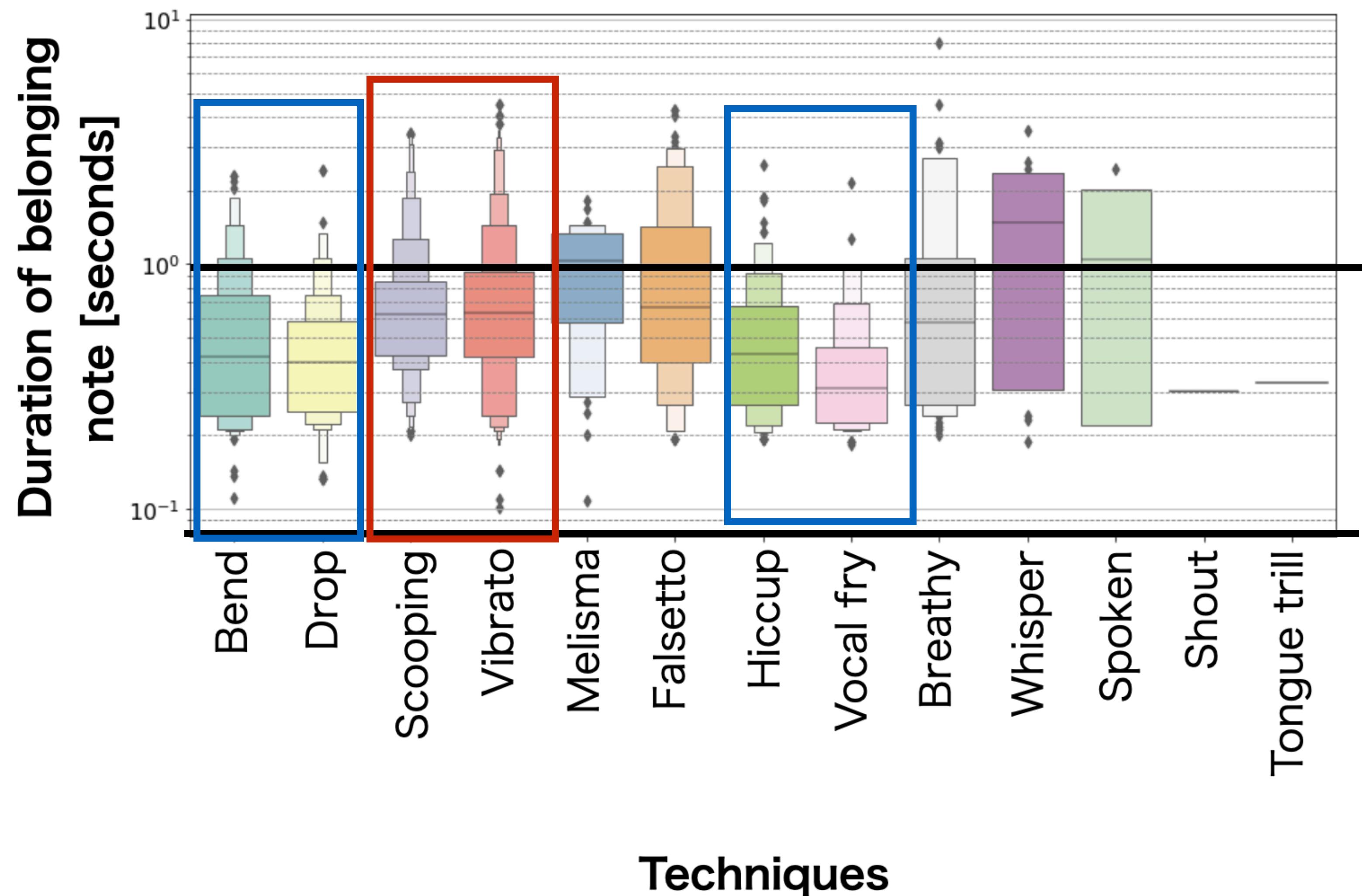


Intervals with the following note [half-note]



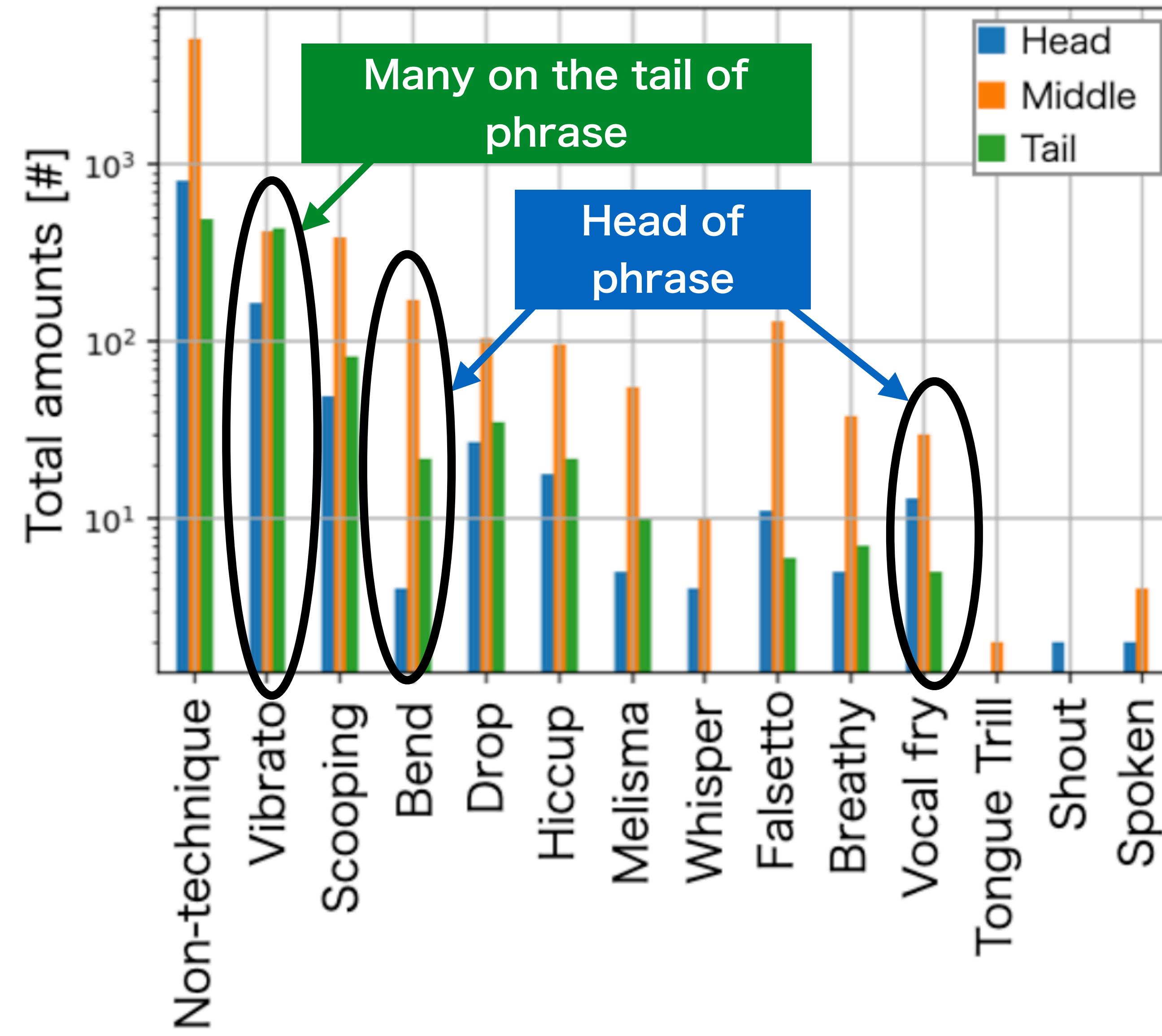
Duration

Omitted in presentation
due to the time limit



Phrase position

Omitted in presentation
due to the time limit



There are certain relationship between musical context and techniques

• Occurrence : What and how often appeared?

- Whole analysis
- Track-wise analysis

- Indicated occurrence and duration distribution of each technique
- Each singer has different distribution
- **The distributions are similar between same or similar style original singer**

• Vibrato parameter : How to realize vibrato?

- Analysis of vibrato parameters (depth and speed)
- Differences of vibrato parameters in same imitator

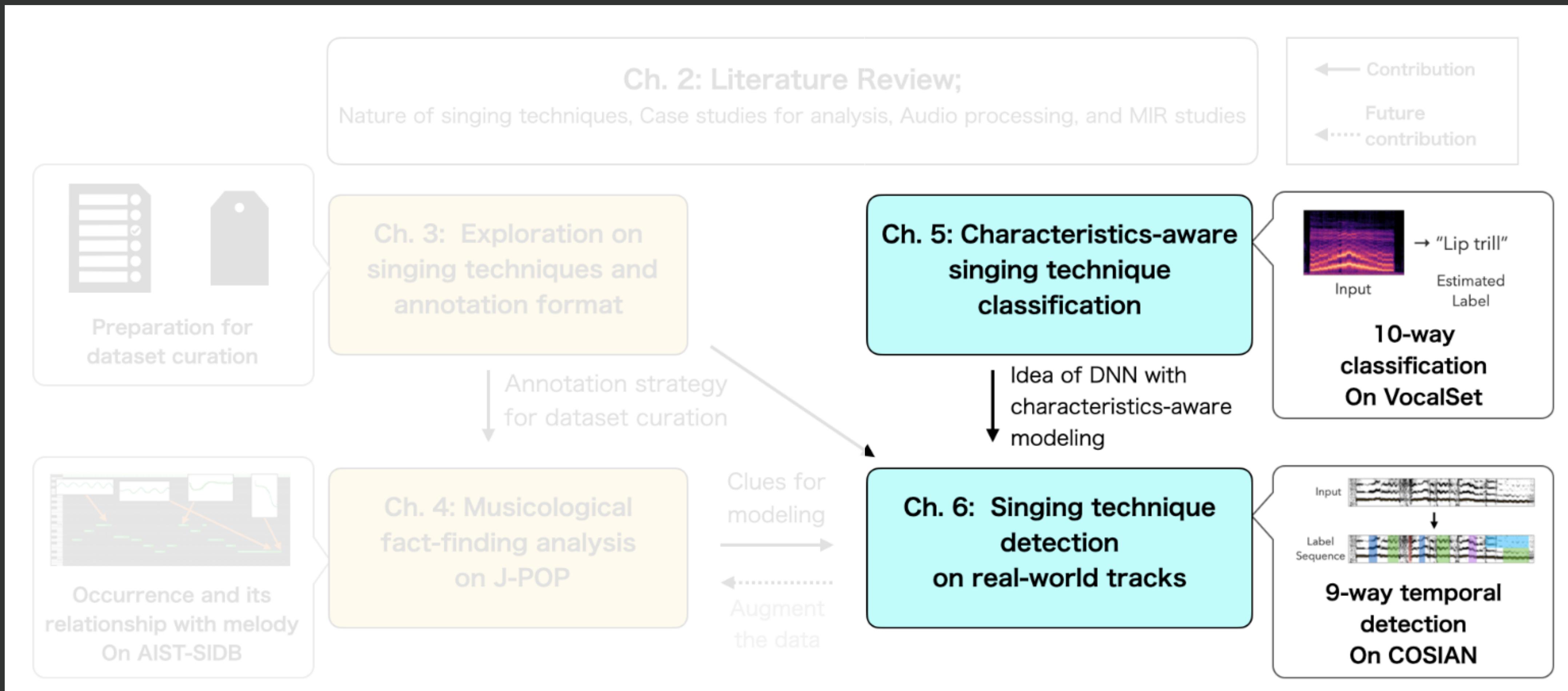
- The average depth is 181.1 cent (\approx 2 semitones)
- The average speed is 6.53 Hz
- **Observed the phenomena that the singers modify the parameters based on target style**

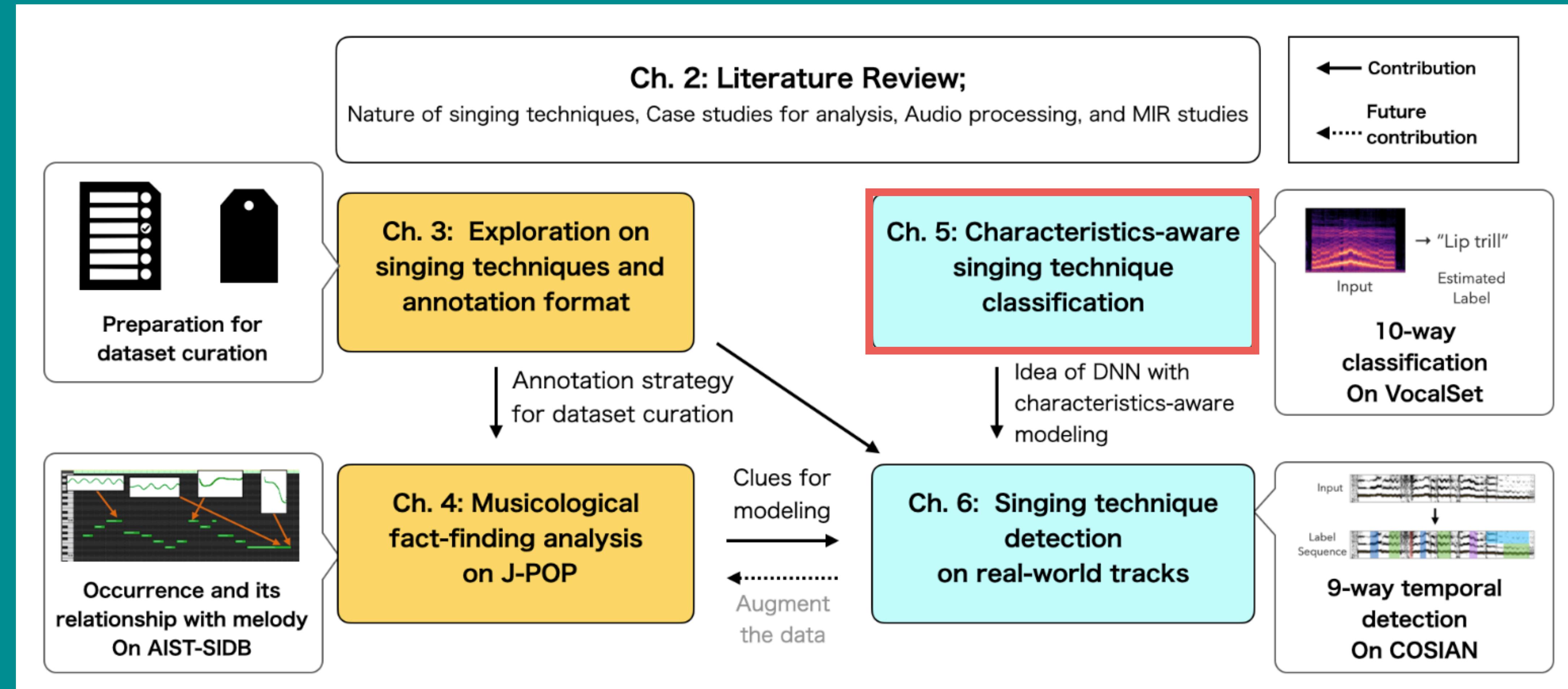
• Location : When techniques occur?

- Correlation between lyrics, previous and next pitch, note duration and position in the phrase
- Correlation between occurrence location and vibrato parameters

- lyric : No salient correlation
- pitch : More falsetto and scooping on high notes.
- pitch interval : when pitch rises: more falsetto and scooping, when next pitch will fall: more falsetto and drop
- duration : **Many vocal fry, hiccup, bend and drop, on short notes. Many vibrato and scooping, on long notes.**
- position : More vibrato on phrase end, Less Bend on phrase beginning
- parameters : **The higher the note, the shallower the vibrato, and the longer the vibrato, the slower it tends to be**

Part II. Computation



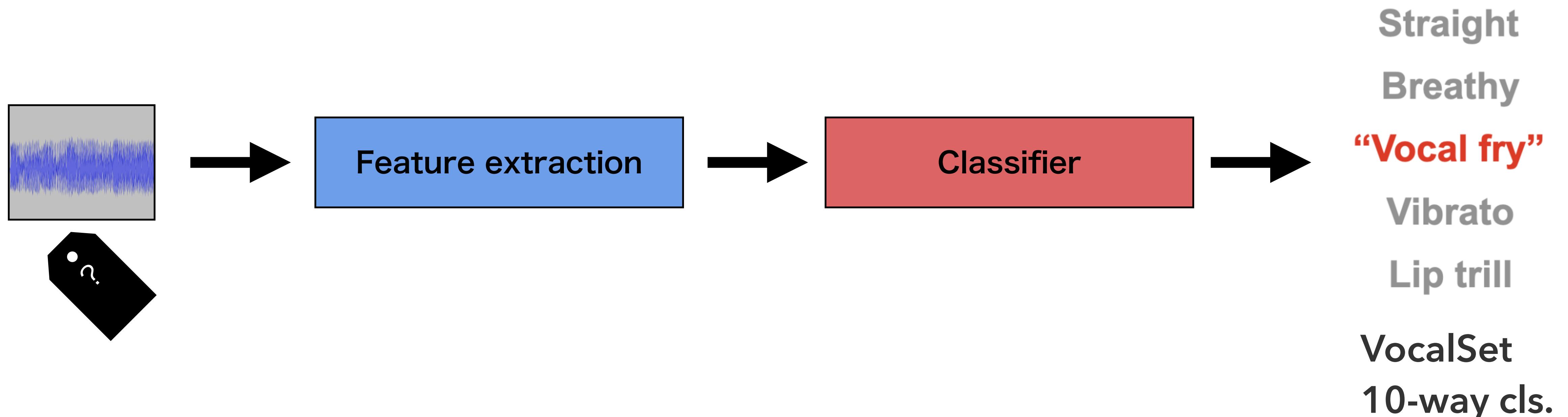


Chapter 5

Singing Technique Classification

considering Feature Extraction and Imbalance-aware Learning

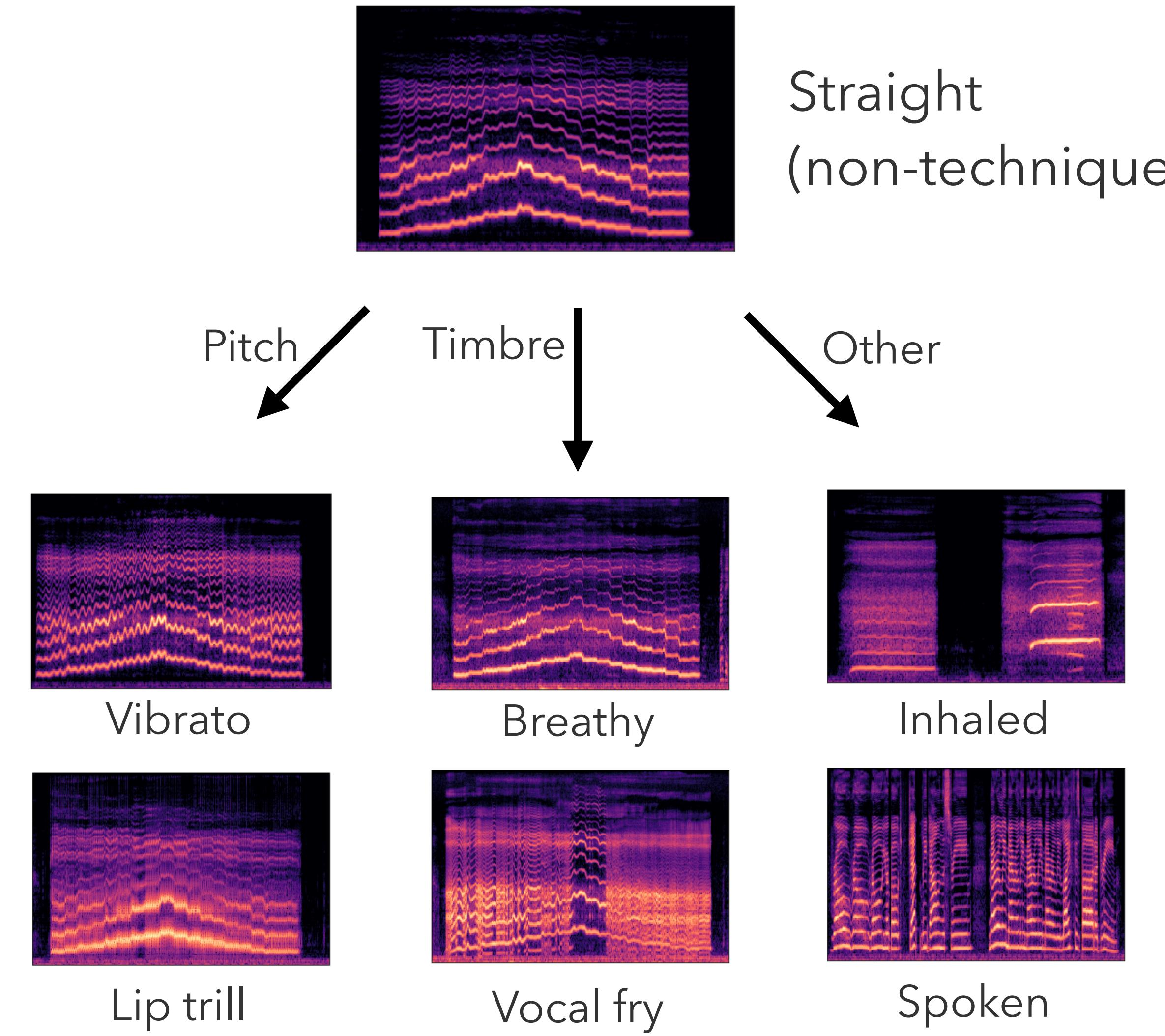
Exploration of singing technique classification methods



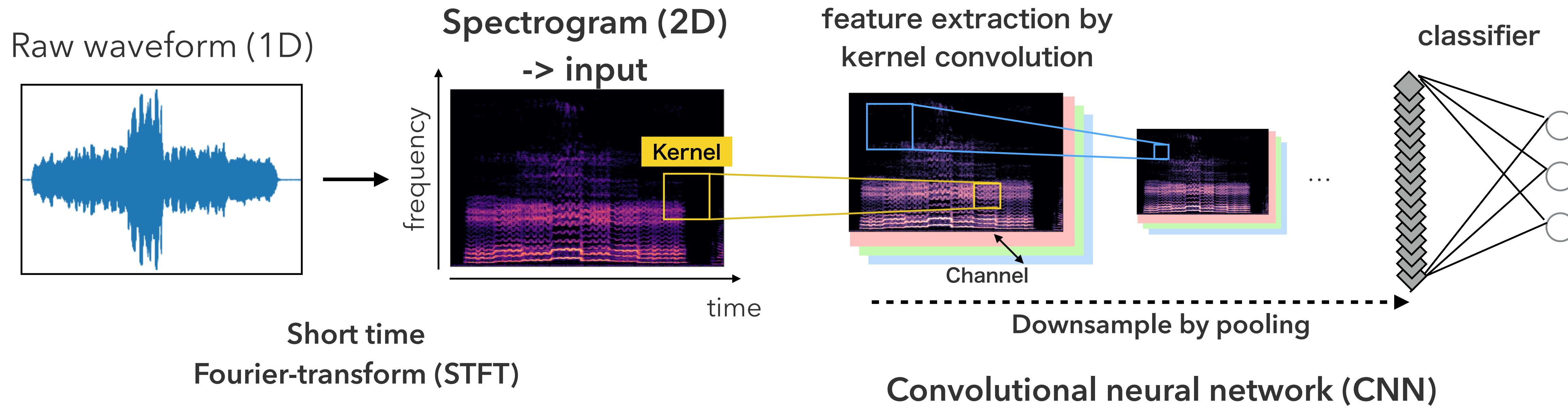
- Challenges:
 1. Feature extraction -> Which approach is better to model singing techniques?
 2. Imbalanced data -> How to mitigate bad effect from label imbalance?

Task: singing technique classification

- **Single label classification [Wilkins 18]**
 - An emerging task in music classification
 - Given audio, identify which singing techniques
 - Dataset: VocalSet
- **Challenges**
 - **1. Many fluctuation components**
 - Pitch - vibrato, lip trill etc.
 - Timbre - breathy, vocal fry etc.
 - Other - inhaled, spoken etc.
 - **2. Imbalance data**
 - Some techniques are few, as real-world



- Making spectrogram -> process as 2D image and apply CNN
 - Now, most common method for DNN-based audio classification [Purwins 19]

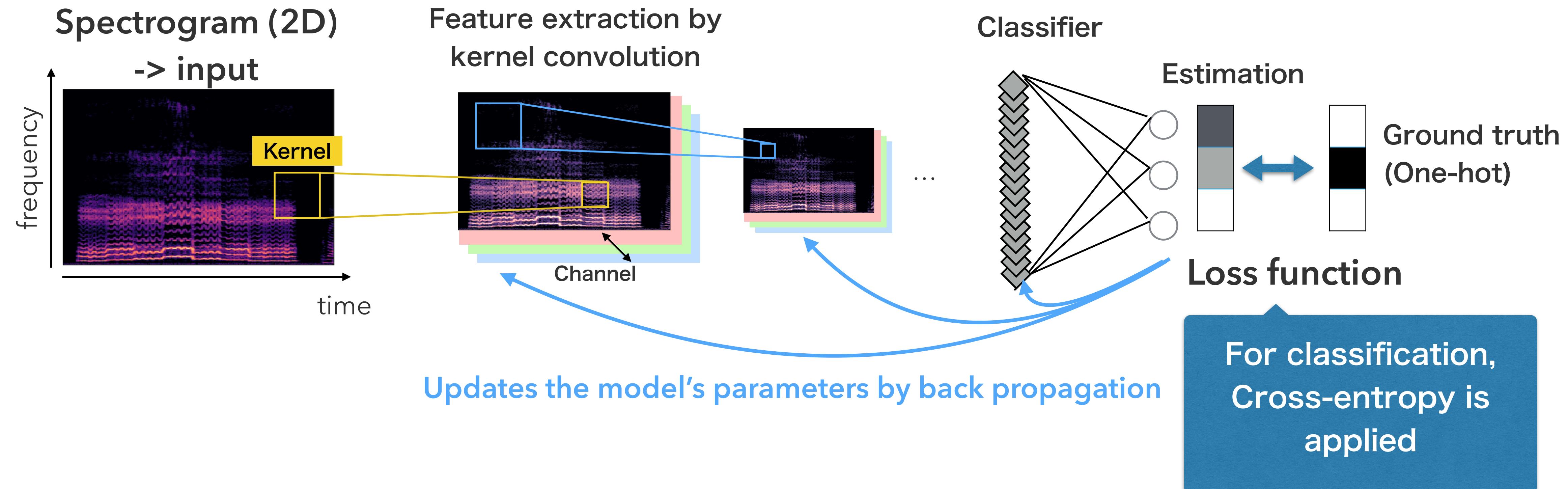


-> Expects powerful performance on feature extraction

Training of Neural network

48

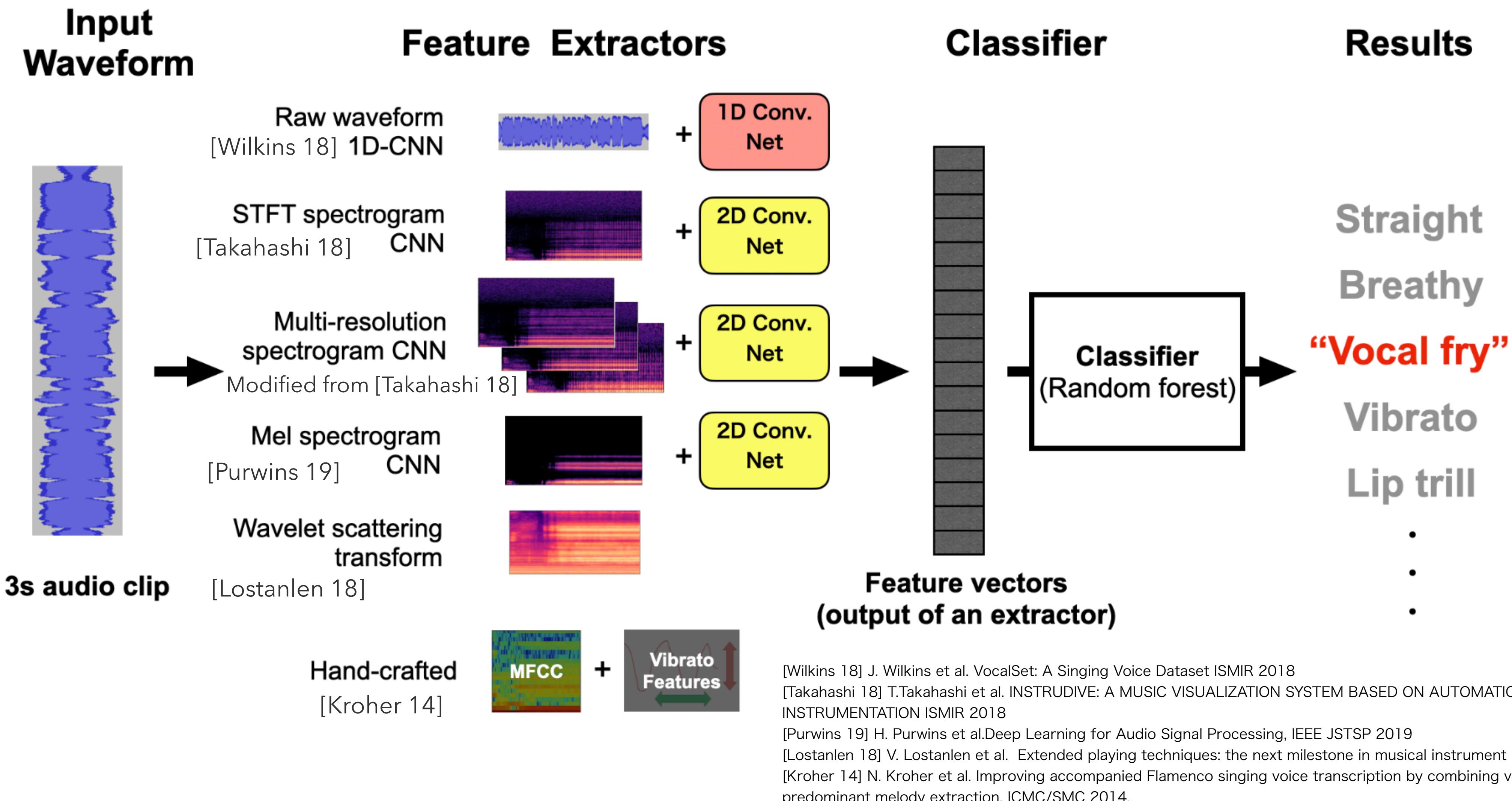
- Based on the minimization of loss function, updates the parameters



① Which audio representation is better?

49

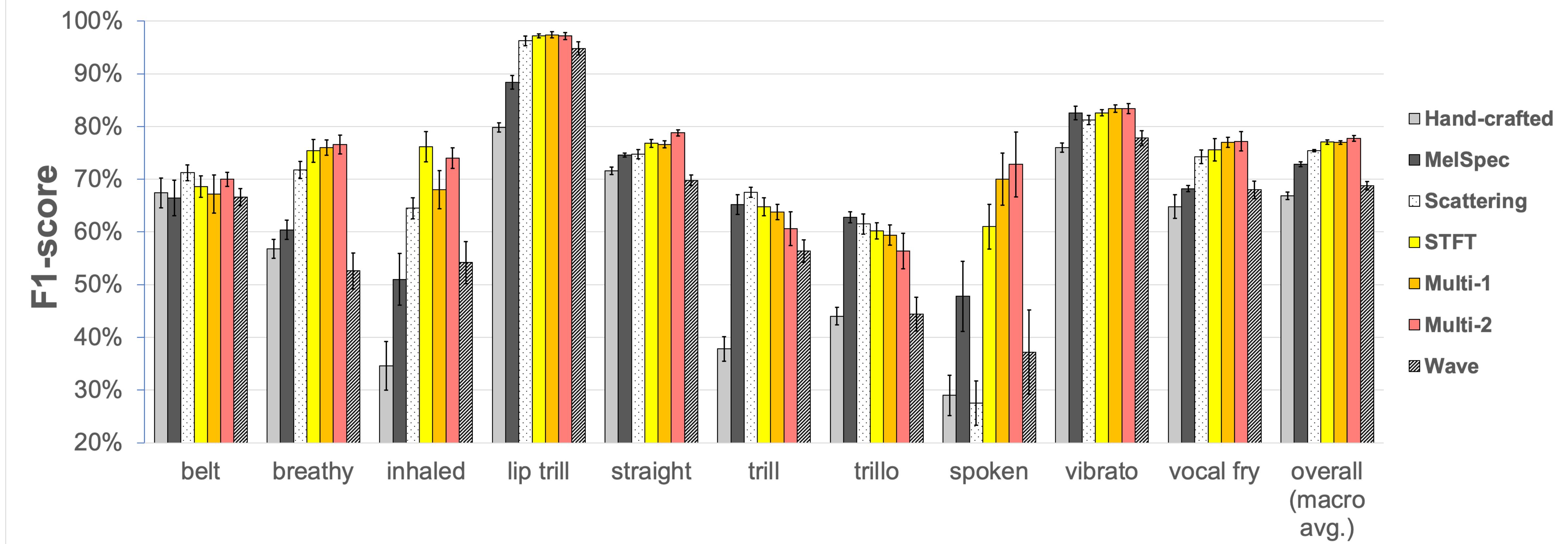
Compared various setting including hand-crafted feature and CNN-based learning



① Which feature extraction is better?

50

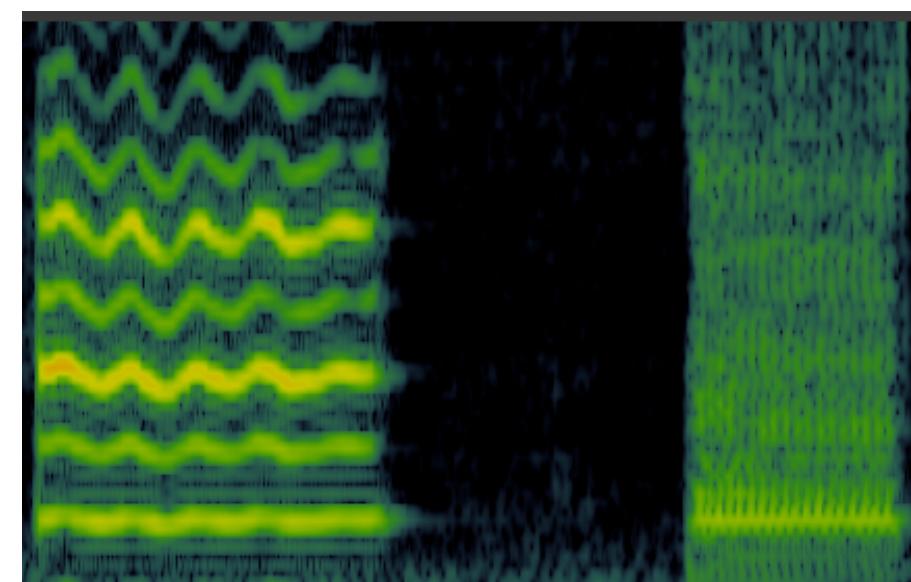
Multi-resolution spectrogram + 2D Oblong-kernel CNN achieved the best performance



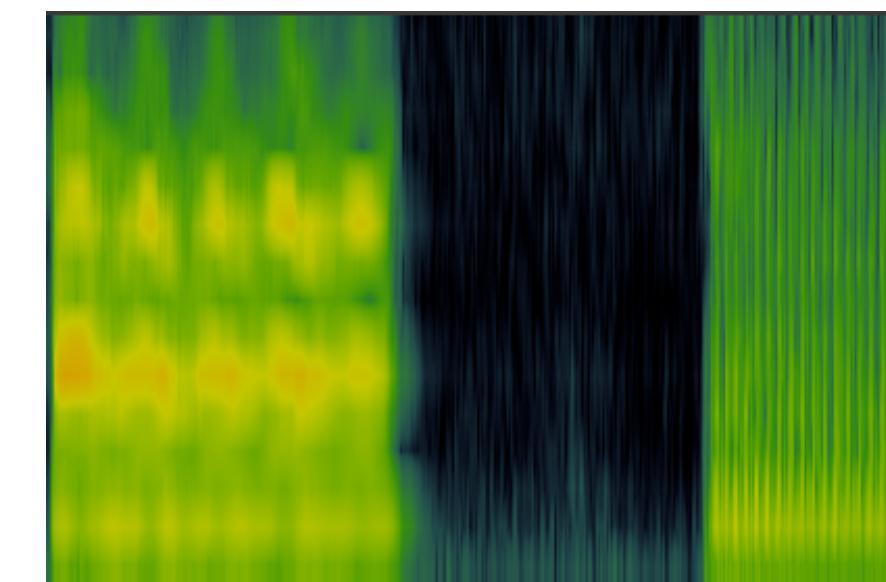
Modification for Input and CNN

51

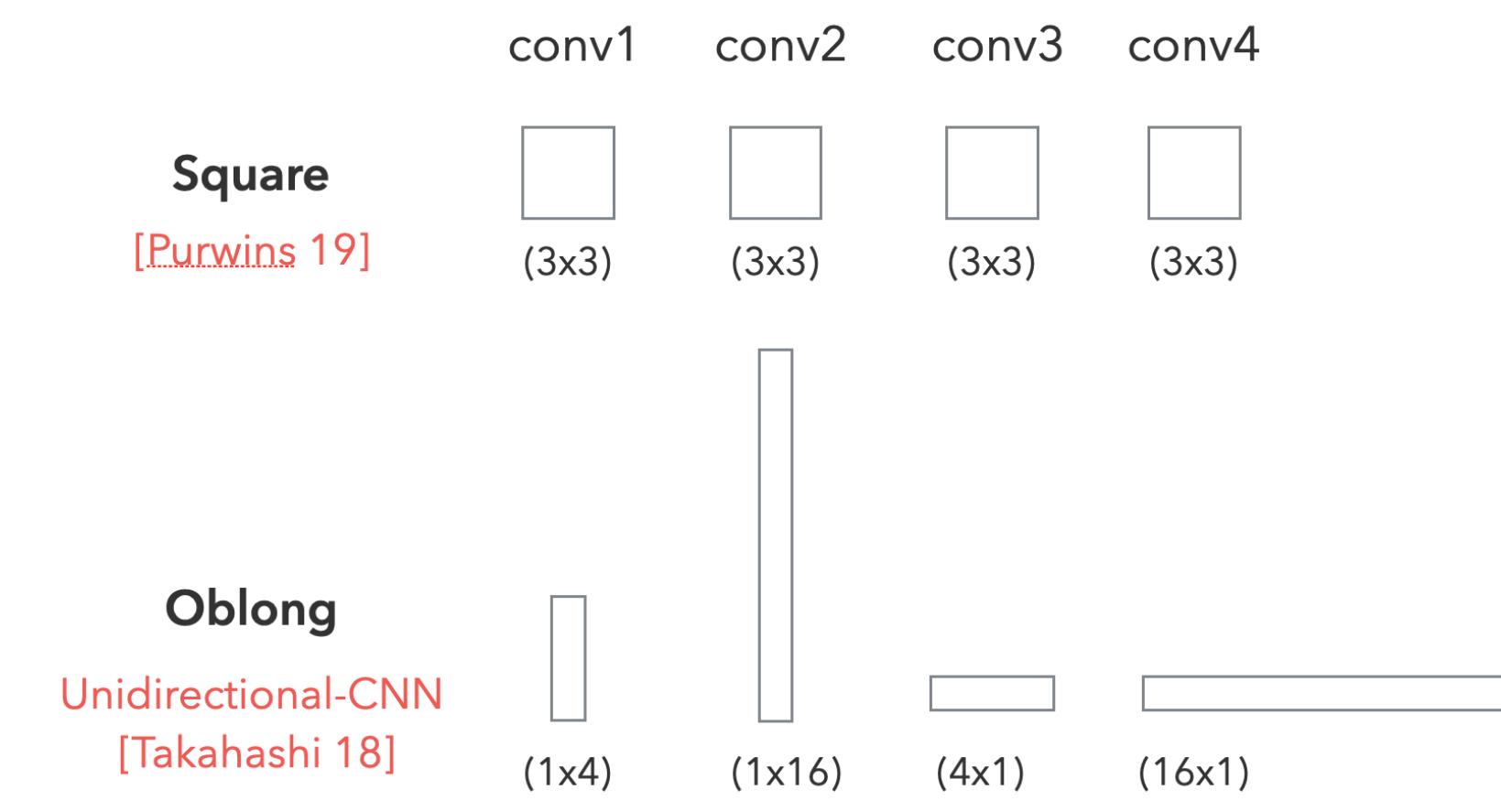
- **Input representation -> Multi-resolution**
 - Stacked 3 different resolution spectrograms on channel axis, by differentiating FFT-size at STFT
 - Expects to adapt various fluctuation
- **CNN -> Adapts the kernel shape**
 - Modified oblong-shaped kernel, from 3x3 square-shaped kernel, which is widely used in CNN
 - Expects to capture more meaningful locality-contexts
 - Vertical long -> timbral feature, Horizontal long -> temporal feature



FFT-length: long (2048)
frequency: **high**
time: **low**



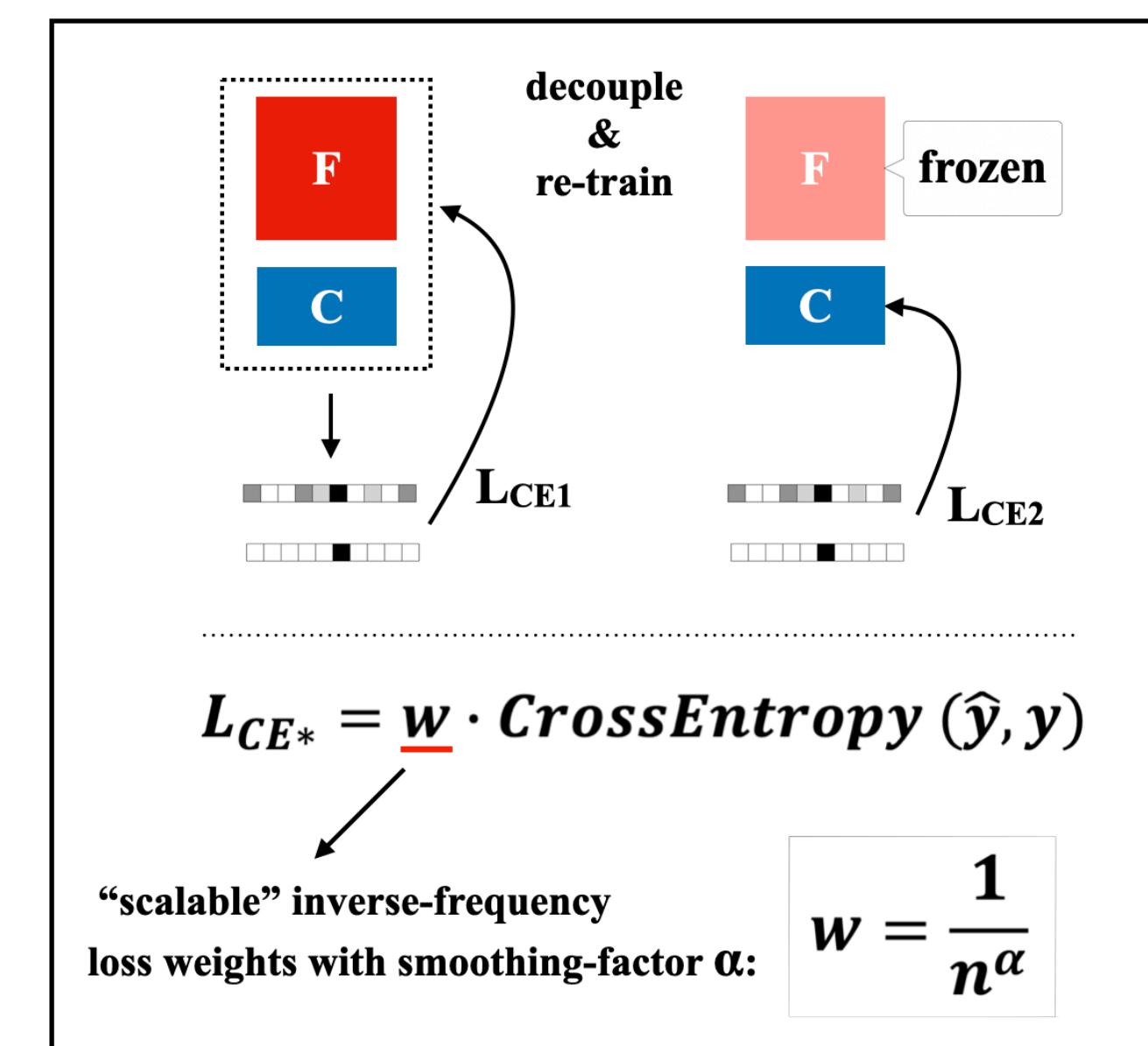
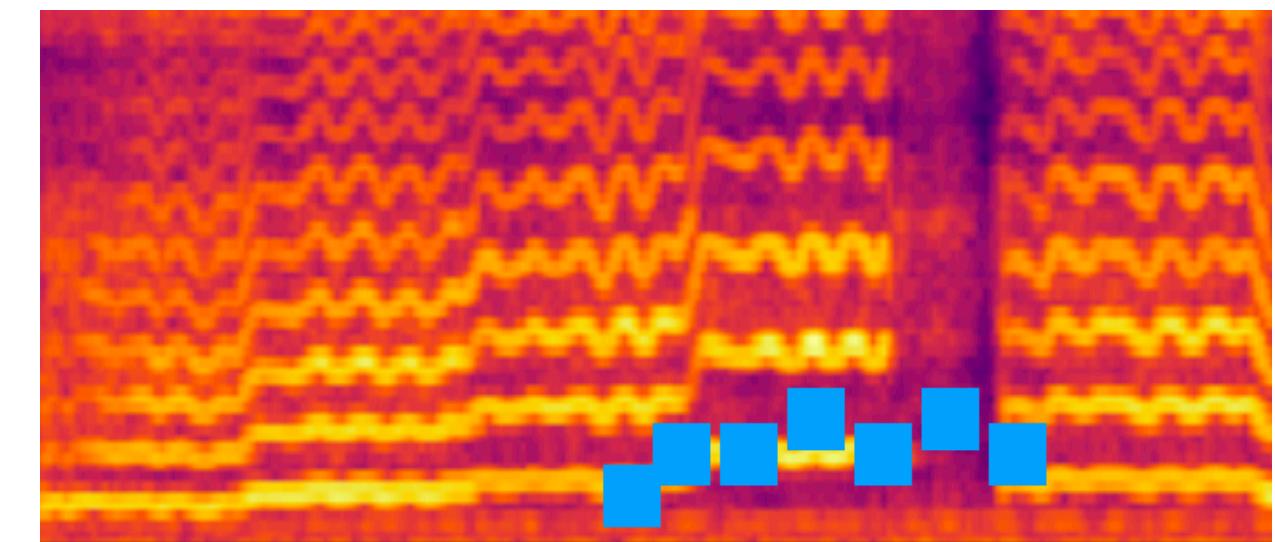
FFT-length: short (512)
frequency: **low**
time: **high**



② Further characteristics-aware modeling

52

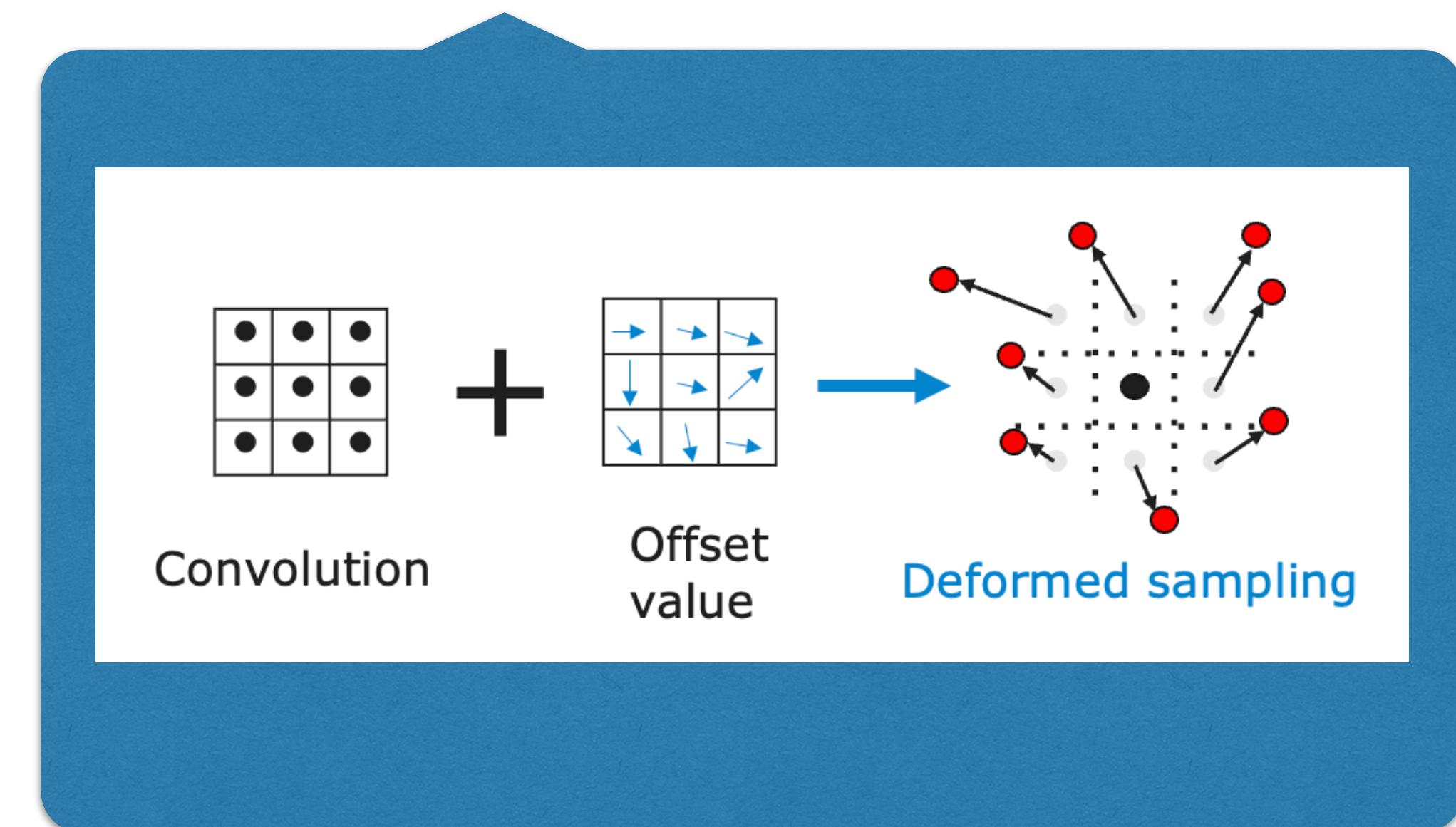
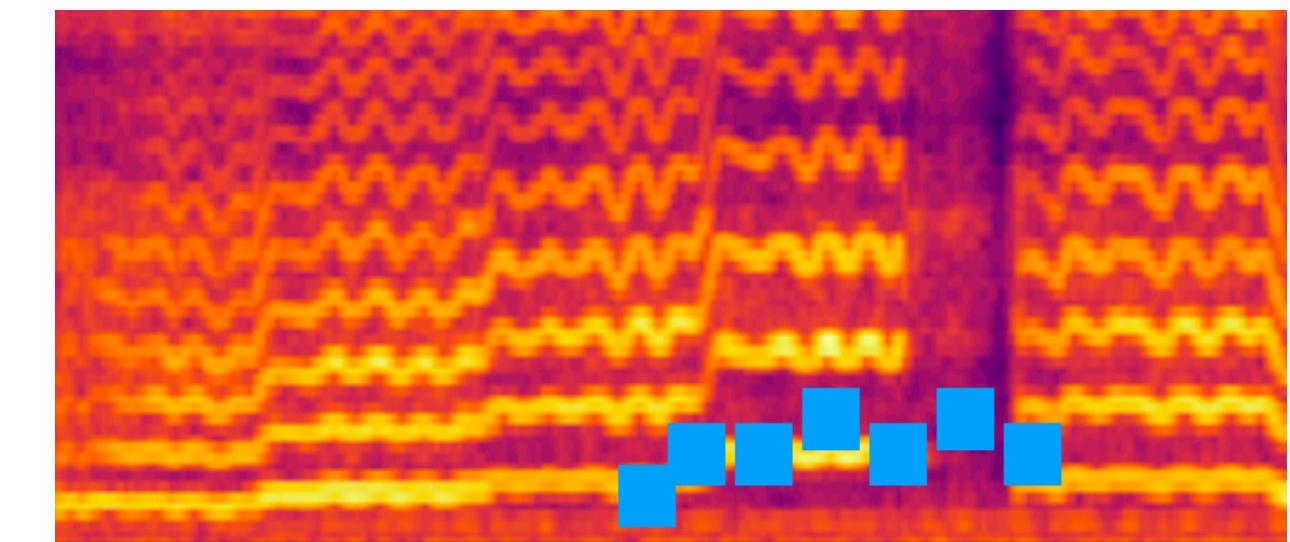
- **1. Adapt Deformable convolution [Dai 17]**
 - **Dynamically** determined the kernel shape
→ Expects to adapt more on various fluctuation
 - “Singing technique exhibits the geometric patterns on the spectrogram”
- **2. Classifier retraining (cRT) with inverse frequency weight [Kang 20]**
 - 1st: train entire part, 2nd: only **retrain the classifier** part
 - Adopts **inverse-frequency-weighted** cross entropy for loss function. (Scaling factor $\alpha=0.2$)



② Further characteristics-aware modeling

53

- **1. Adapt Deformable convolution [Dai 17]**
 - Dynamically determined the kernel shape
→ Expects to adapt more on various fluctuation
 - “Singing technique exhibits the geometric patterns on the spectrogram”
- **2. Classifier retraining (cRT) with inverse frequency weight [Kang 20]**
 - 1st: train entire part, 2nd: only retrain the classifier part
 - Adopts inverse-frequency-weighted cross entropy for loss function. (Scaling factor $\alpha=0.2$)



$$L_{CE^*} = \underline{w} \cdot \text{CrossEntropy}(\hat{y}, y)$$

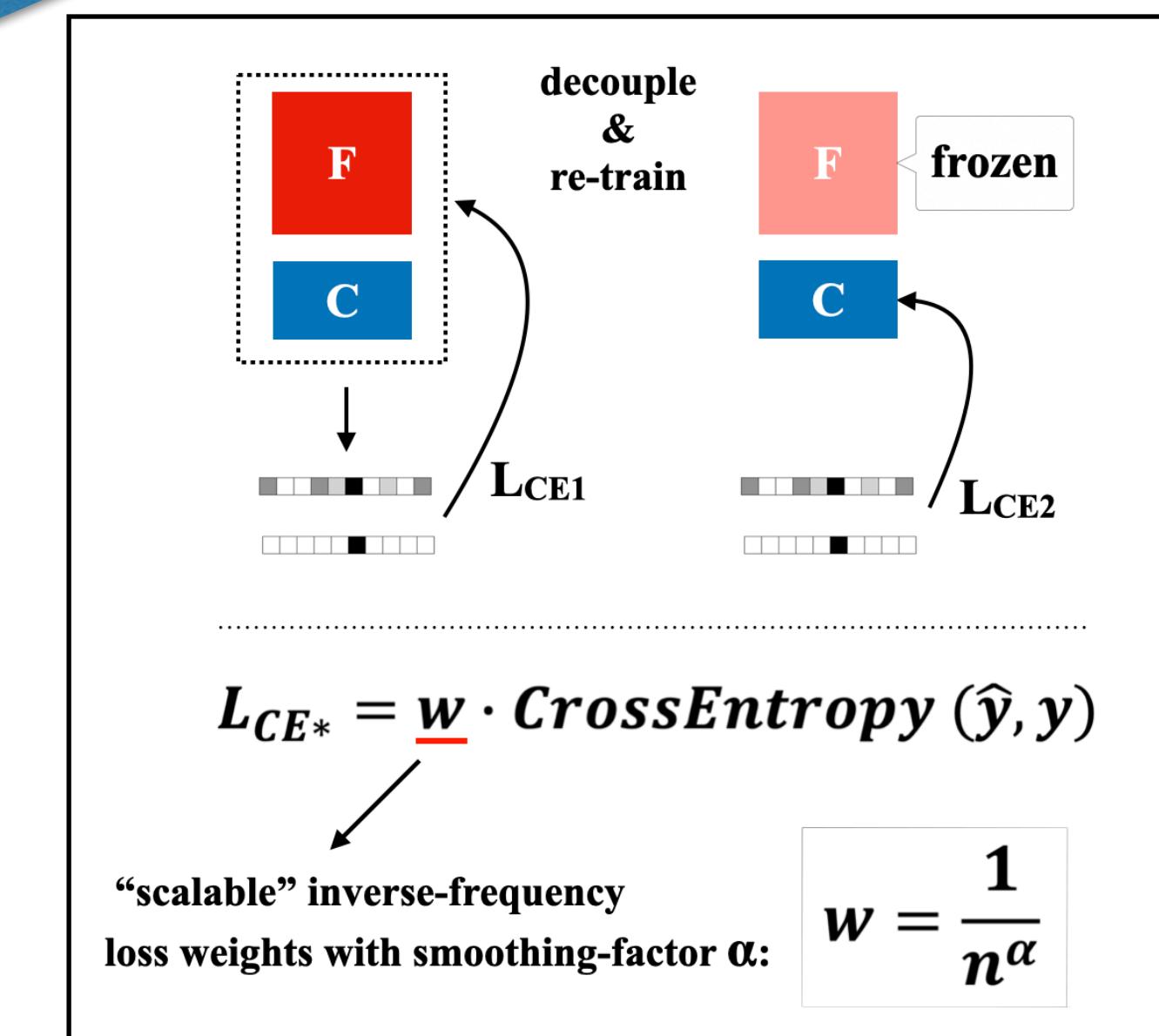
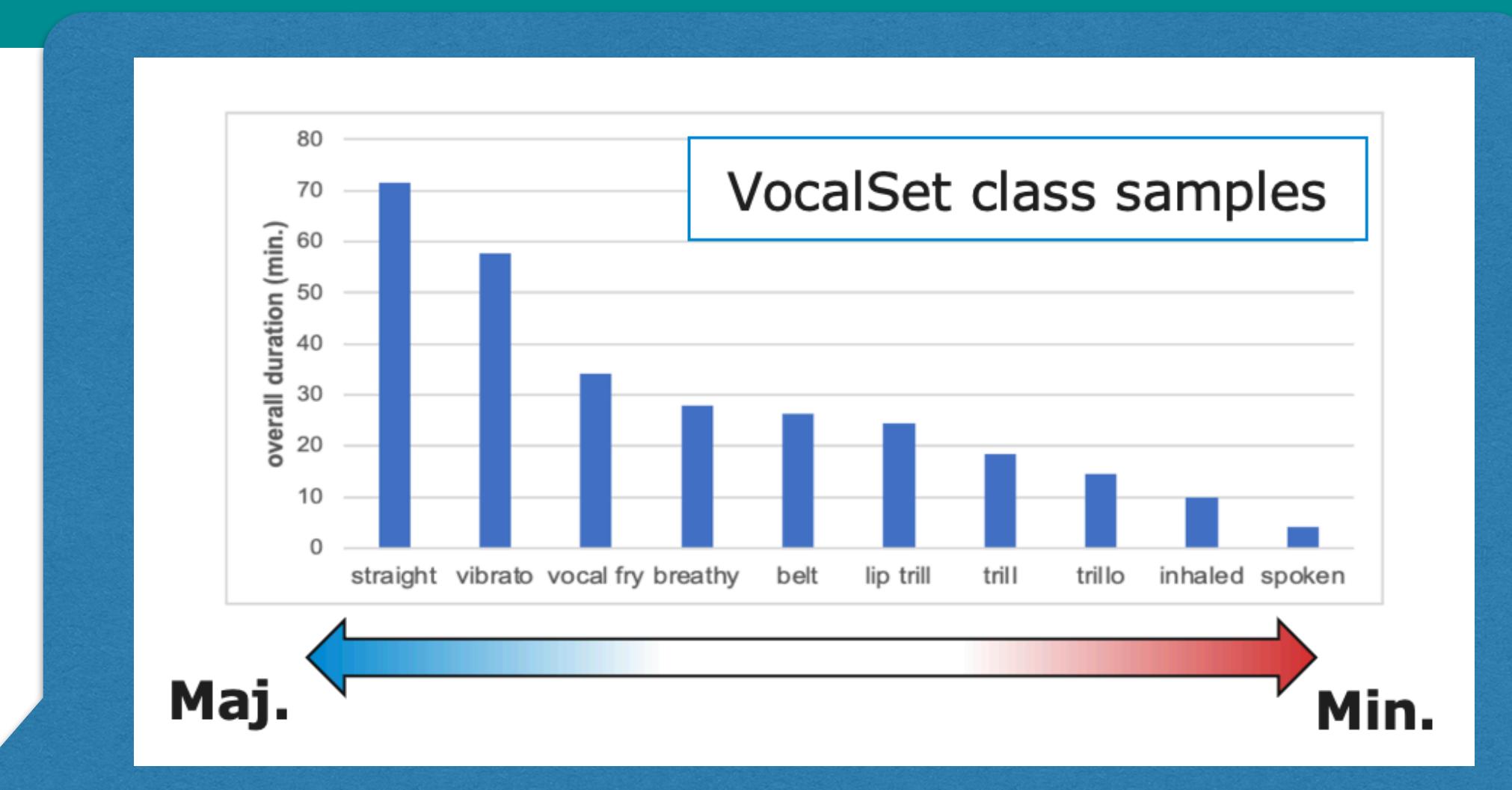
“scalable” inverse-frequency loss weights with smoothing-factor α :

$$w = \frac{1}{n^\alpha}$$

② Further characteristics-aware modeling

54

- 1. Adapt Deformable convolution [Dai 17]
 - Dynamically determined the kernel shape
→ Expects to adapt more on various fluctuation
 - “Singing technique exhibits the geometric patterns on the spectrogram”
- 2. Classifier retraining (cRT) with inverse frequency weight [Kang 20]
 - 1st: train entire part, 2nd: only retrain the classifier part
 - Adopts inverse-frequency-weighted cross entropy for loss function. (Scaling factor $\alpha=0.2$)



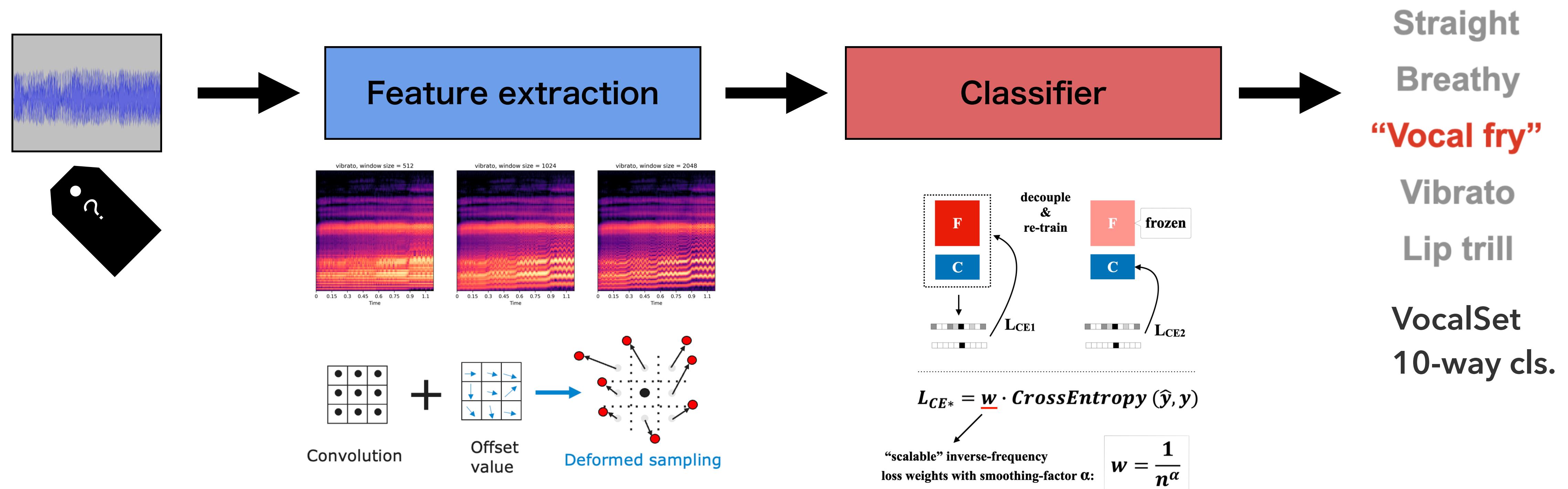
Results of the full model

*unlike the comparison part, the split is by singer. Since the official split is found after the comparison

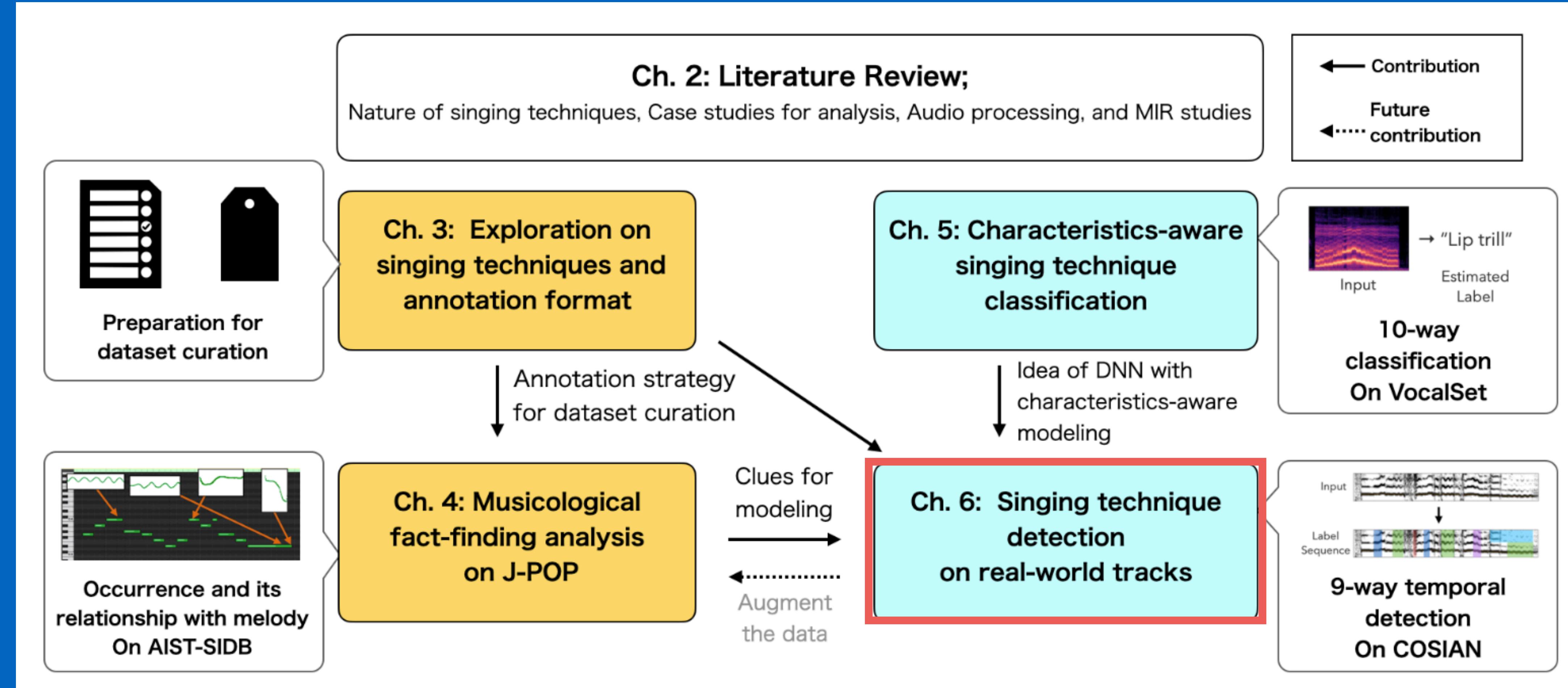
- + Weight loss
- + Deformable convolution
(On 3rd and 4th layer)
- + cRT, only applied weight loss on retraining phase

Methods	F1-score	Accuracy	Balanced Accuracy
Multi-resolution Oblong CNN	0.404	0.492	0.472
Multi-resolution Oblong CNN w/ weight loss ($\alpha=0.2$)	0.513	0.554	0.575
Deformable CNN w/ weight loss ($\alpha=0.2$)	0.559	0.610	0.635
Deformable CNN w/ weight loss on cRT ($\alpha=0.2$)	0.620	0.656	0.655

Characteristics-aware CNN model for singing technique classification



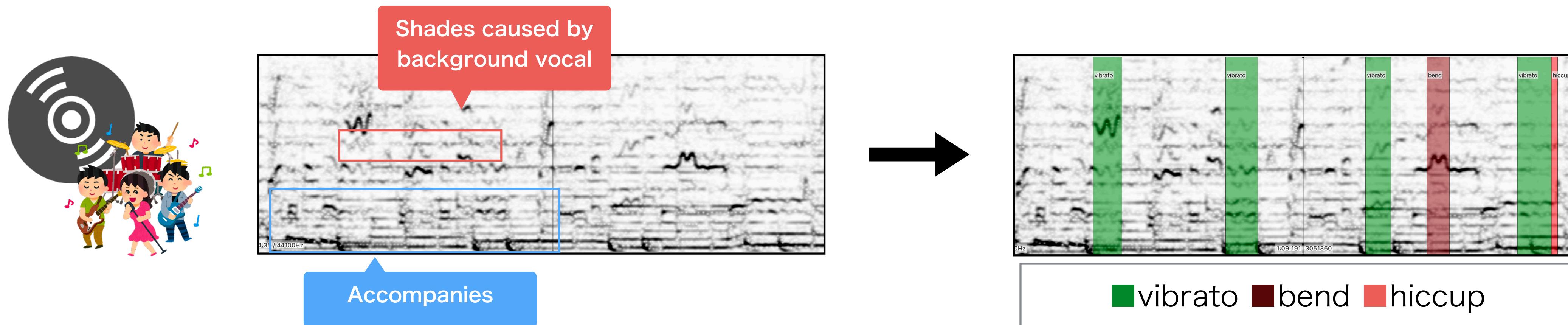
- Contributions:
 - Confirmed the effectiveness of **Multi-resolution** spectrogram and **CNN-kernel modification**
 - Proposed the model based on **Deformable CNN** and **Imbalance-aware learning**



Chapter 6

Singing Technique Detection from Real-world Popular Music

Extracting appearance of singing techniques on real-world vocal tracks

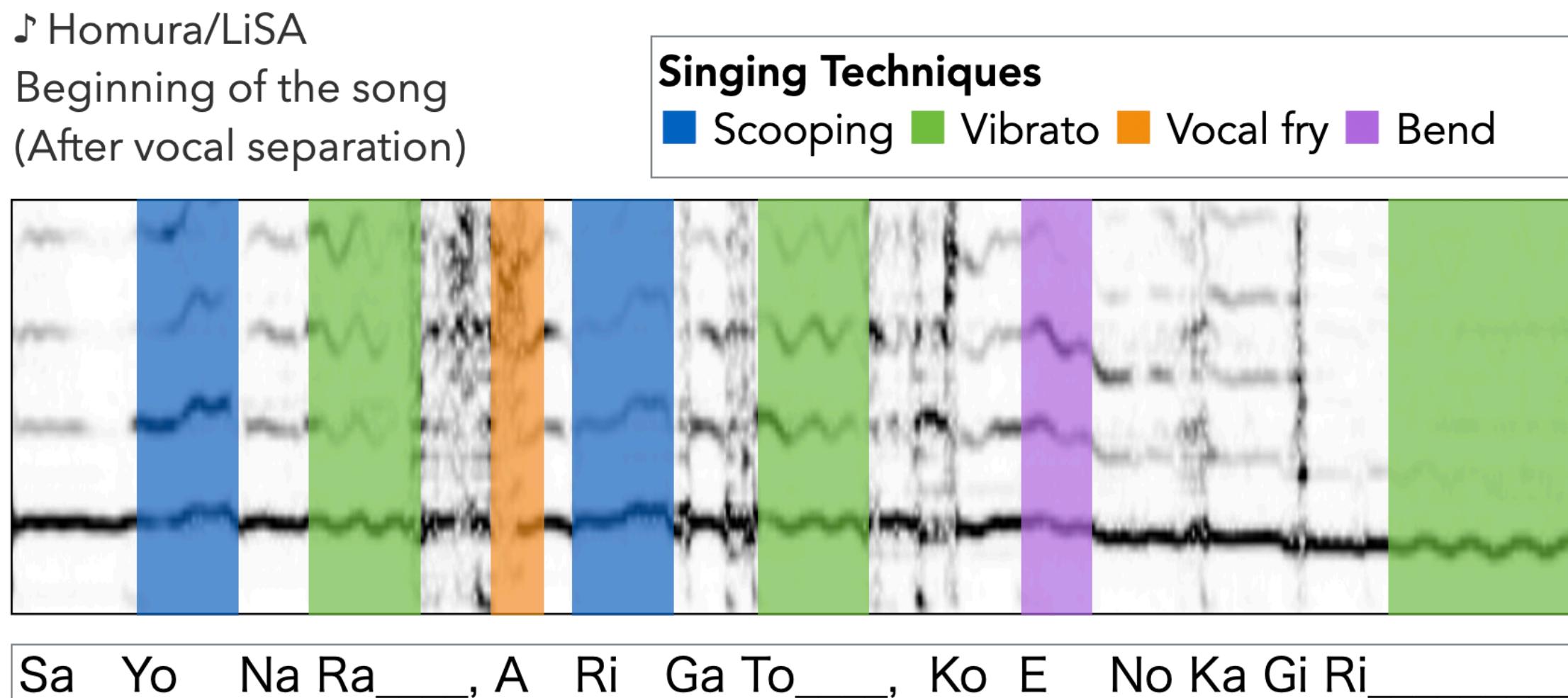
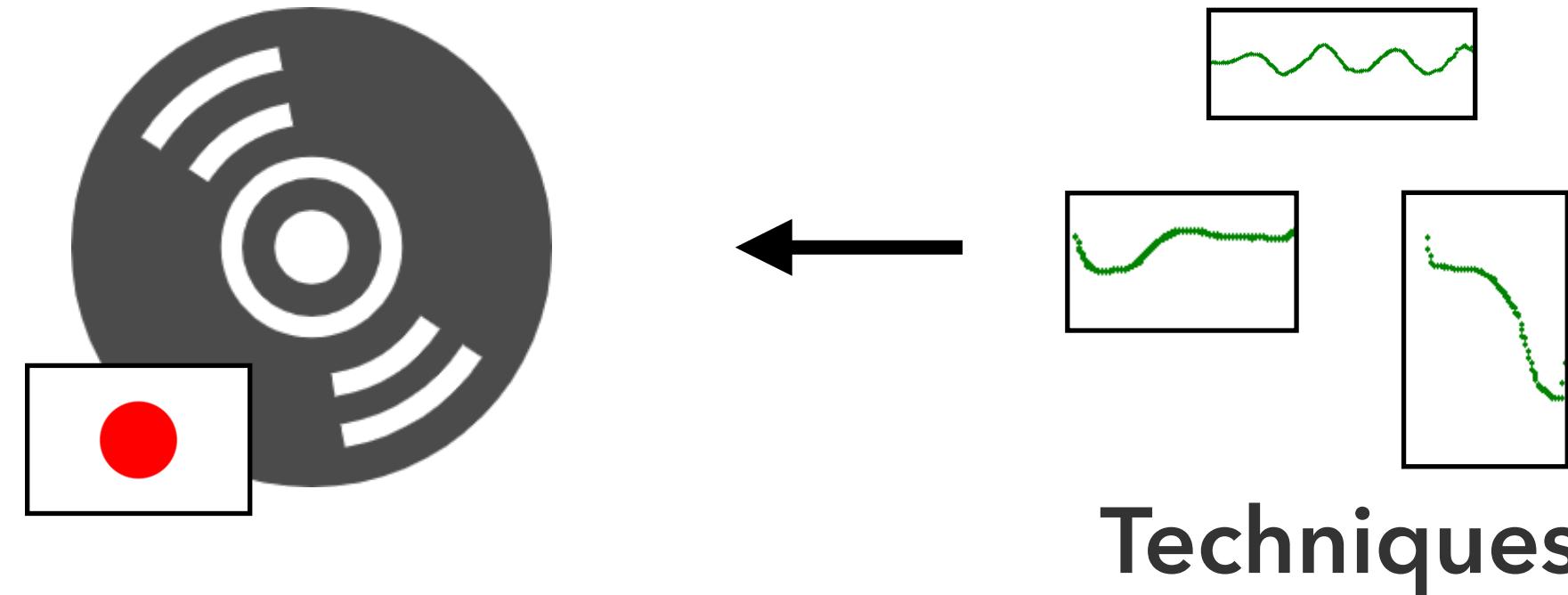


Singing voices from commercial musics
(i.e., convolved effects & backgrounds)

Goal: Automatically detects
singing technique appearance

- Challenges
 - Task/Dataset creation -> How to implement a computational task?
 - Detection model -> How to handle temporarily appearing techniques and real-world tracks?

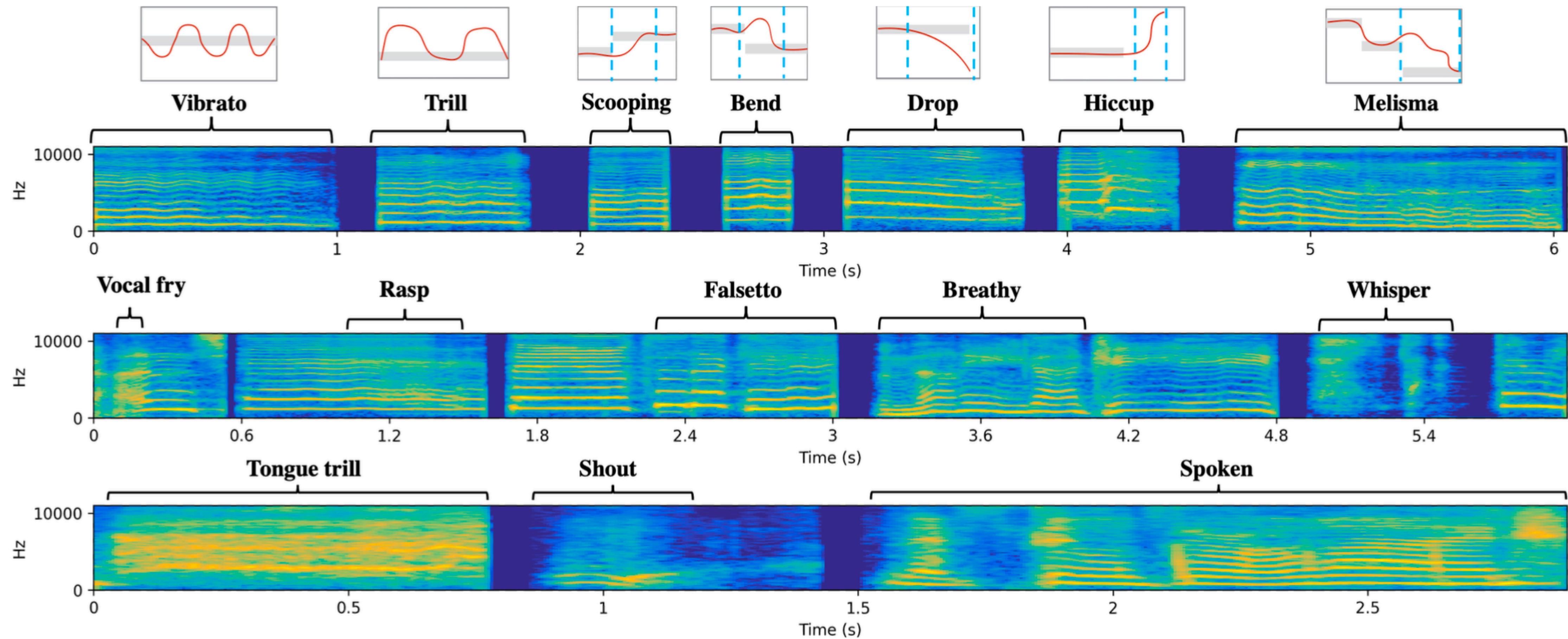
Built a dataset to enable training model and evaluation



- **168 commercial J-POP songs' first section**
 - 21 male, 21 female, various types
 - 4 songs from each vocalist
 - **4h 47m 39s**, in total
 - Voice is separated by Demucs v3 [Defossez 21]

- **Annotation**
 - Region labels of singing techniques
 - Pitch f0 value

Various 15 techniques



States the detailed criteria for each techniques (p.69-70)

Technique	Sketch	Beginning	End	Difference with...	Samples from audio (mel spectrogram)
Vibrato		Visible beginning of the pitch change	Visible end of the pitch change	w/ NA: has visible sinusoid and periodicity w/ Trill: does not have the target pitch of the edge of pitch endpoint	
Scooping		Visible beginning of the preparation	Visible end of the overshoot	w/ NA: has hearable pitch change w/ Hiccup: occurs on the attack and does not have abrupt higher pitch change	
Bend		Visible beginning of the preparation	Visible end of the unstable pitch	w/ NA: has hearable pitch change w/ Vibrato: <1 roundtrip of the pitch w/ Hiccup: not so abrupt pitch change	
Drop		Visible beginning of the pitch dropping	Visible end of the pitch dropping	w/Bend: occurs on the release	
Hiccup		Visible beginning of pitch rising	Visible end of pitch region	w/Bend: has extreme pitch rising (> 4 semitones) w/Falsetto: has instantaneous region of high pitch	
Melisma		Visible beginning of pitch change	Visible beginning of stable pitch region	w/ NA: only has one syllable and fast pitch change w/ Bend: > 1 stable pitch targets	
Trill		Visible beginning of pitch change	Visible end of pitch change	w/ NA: only has one syllable w/Vibrato: has the target pitch of the edge of pitch endpoint	

For pitch: p.69

For timbre: p.70

States the detailed criteria for each techniques (p.69-70)

Technique	Difference with...	Samples from audio (mel spectrogram)			
Breathy	w/ NA: has higher breathiness and more frequent noisy components compared to ordinary voice region w/ Whisper: its pitch component is not so missing, relatively w/ Falsetto: its vocal register is not falsetto (mixed or modal, etc.)				
Falsetto	w/ NA: accompanied by high vocal note and is in different register as ordinary w/ Breathy: its vocal register is falsetto w/ Hiccup: the region sung by falsetto register is not instantaneous				
Whisper	w/ Breathy: its pitch component is relatively missing				
Rasp	w/ NA: has distorted timbre, with visible subharmonics on spectrogram w/ Vocal fry: has main accompanied pitch				
Vocal fry	w/ NA: has creaky sound, with visible pulse pattern on spectrogram w/ Rasp: more instantaneous, not accompanied main pitch and tend to be used in attack				

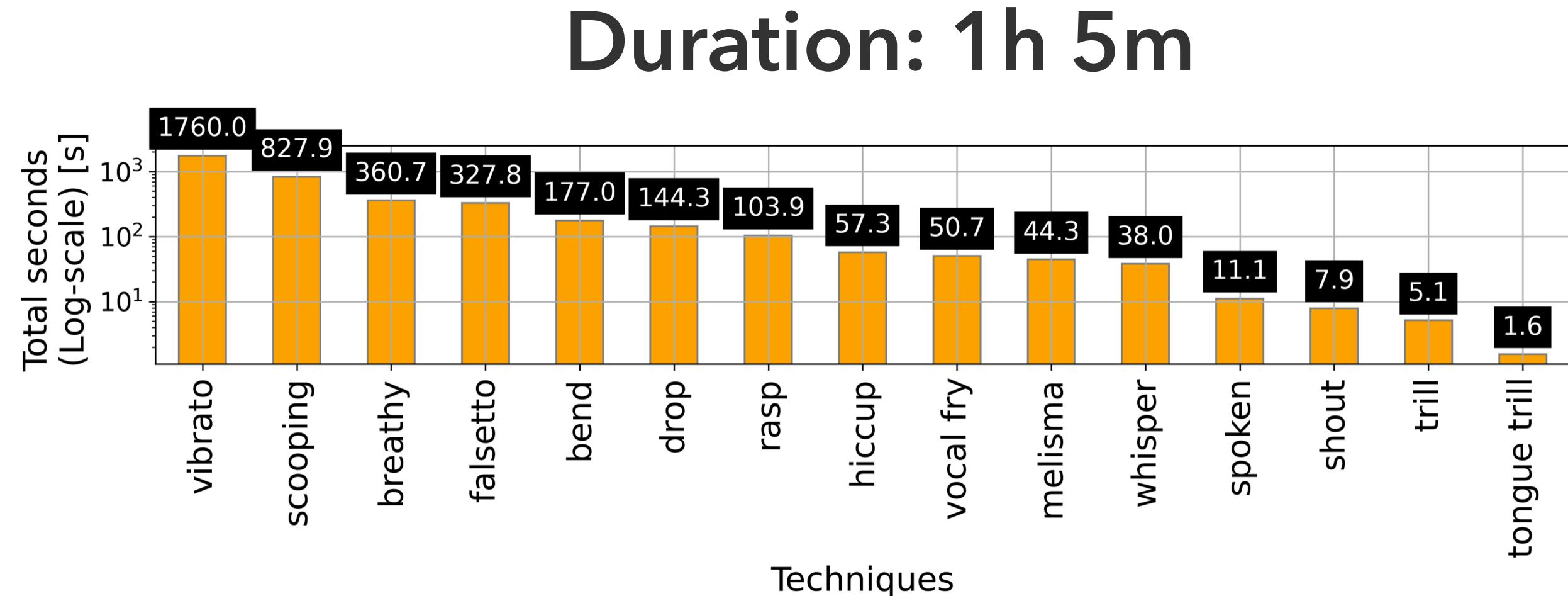
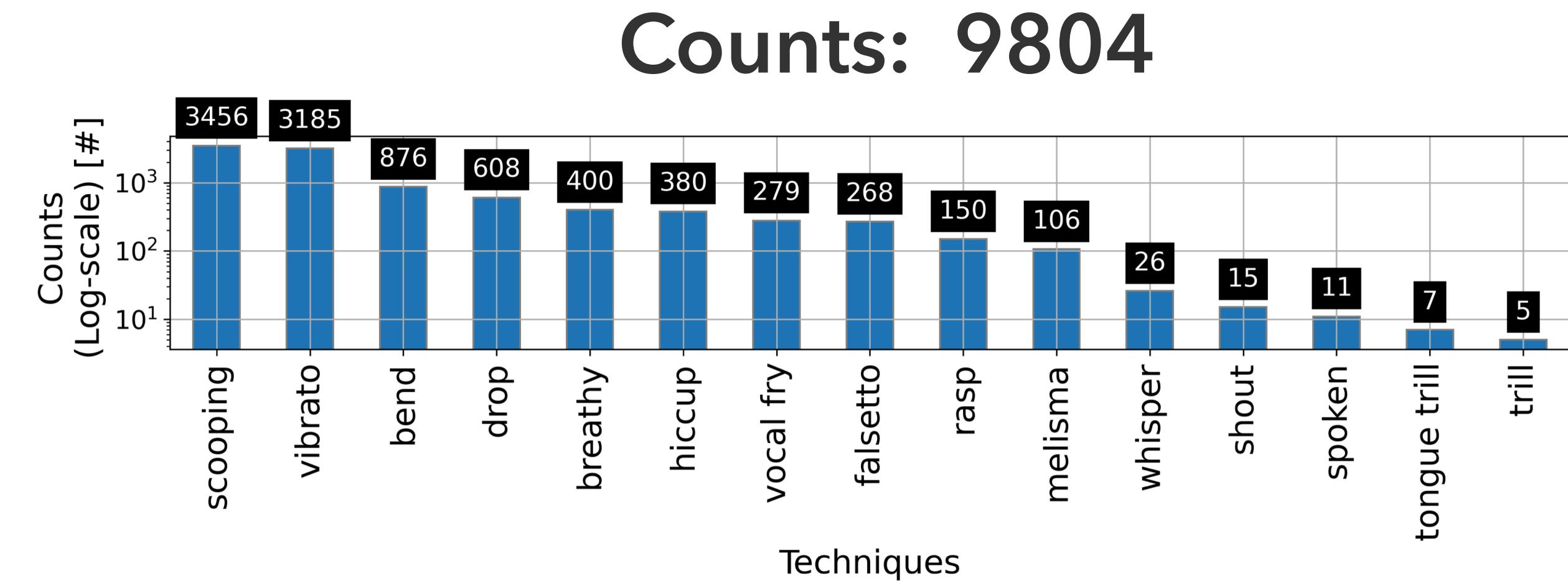
For pitch: p.69

For timbre: p.70

Distribution of labels

63

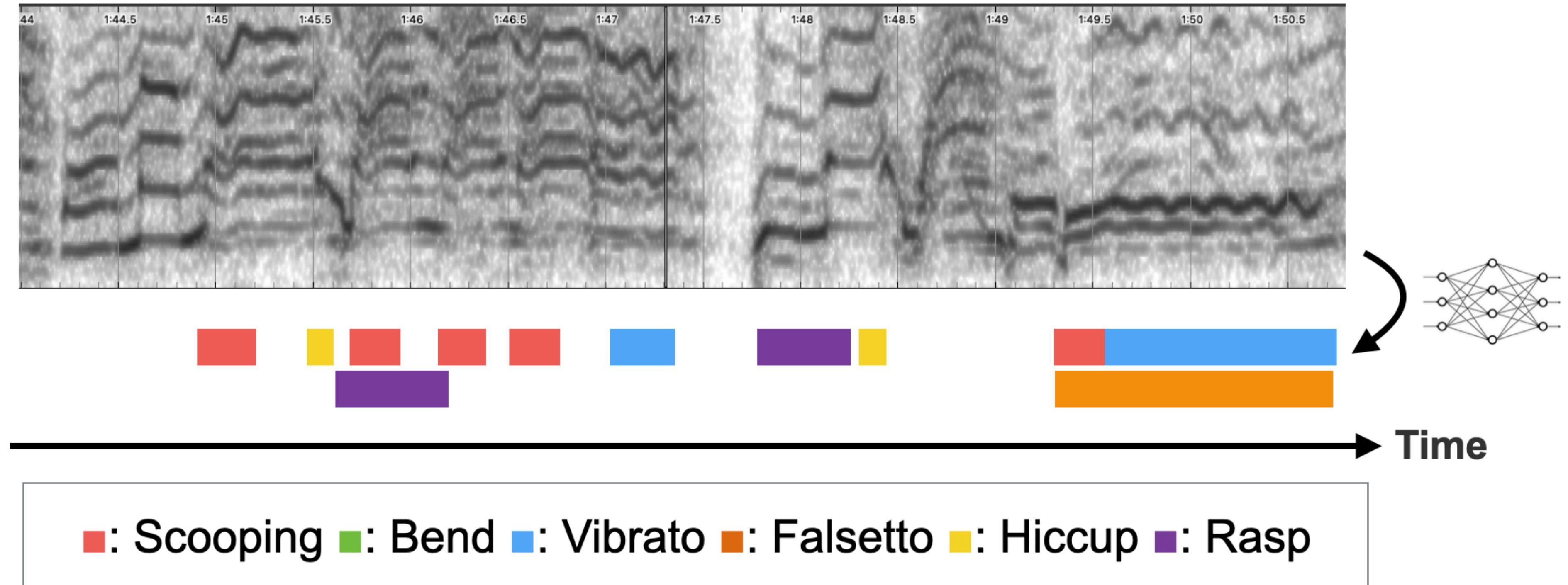
- **Technique ratio**
 - vs audio length: **22.8%**
 - vs vocal length: **38.1%**
- **Vibrato and Scooping are most frequent techniques**



Task: Singing technique detection

64

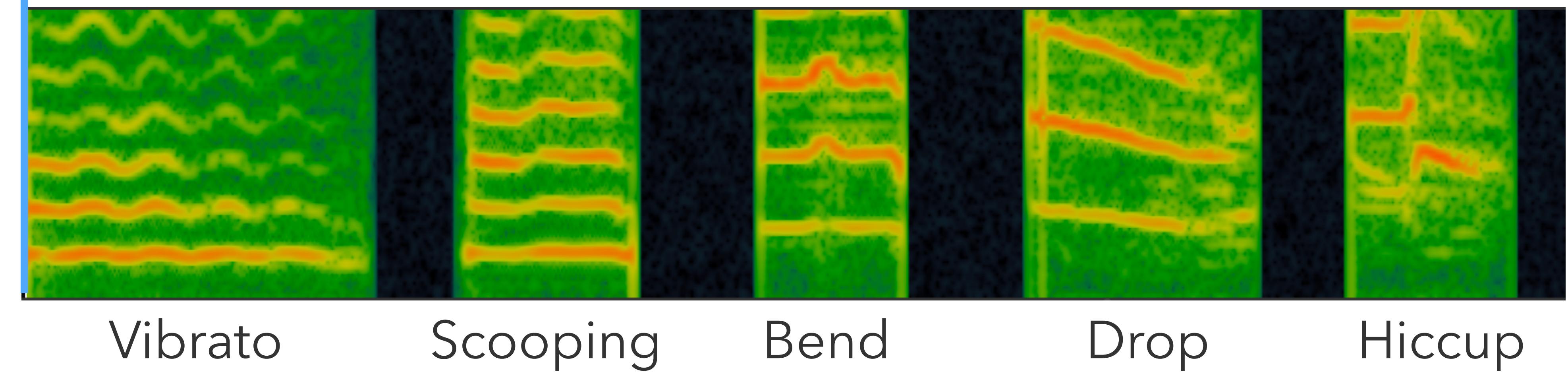
**Sung voice
audio
(separated)**



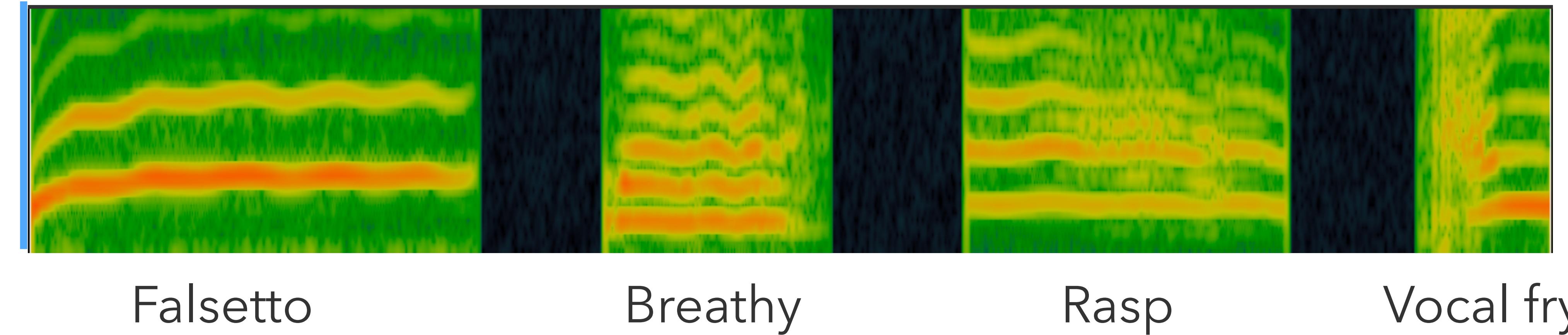
Target: 9 singing techniques

65

Pitch
techniques



Timbre
techniques

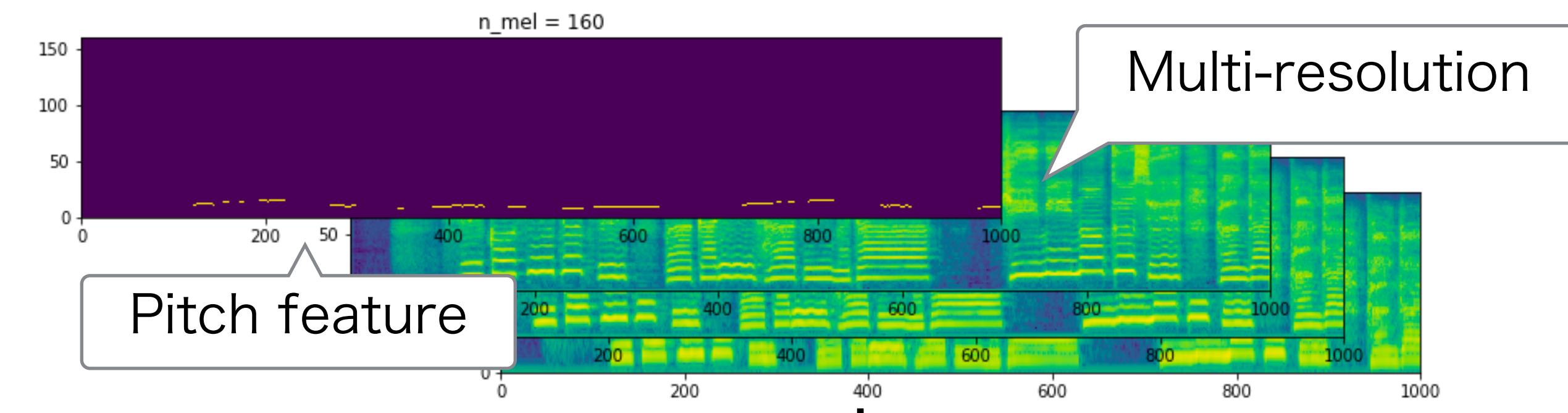


PrimaDNN': Singing technique detection CRNN

66

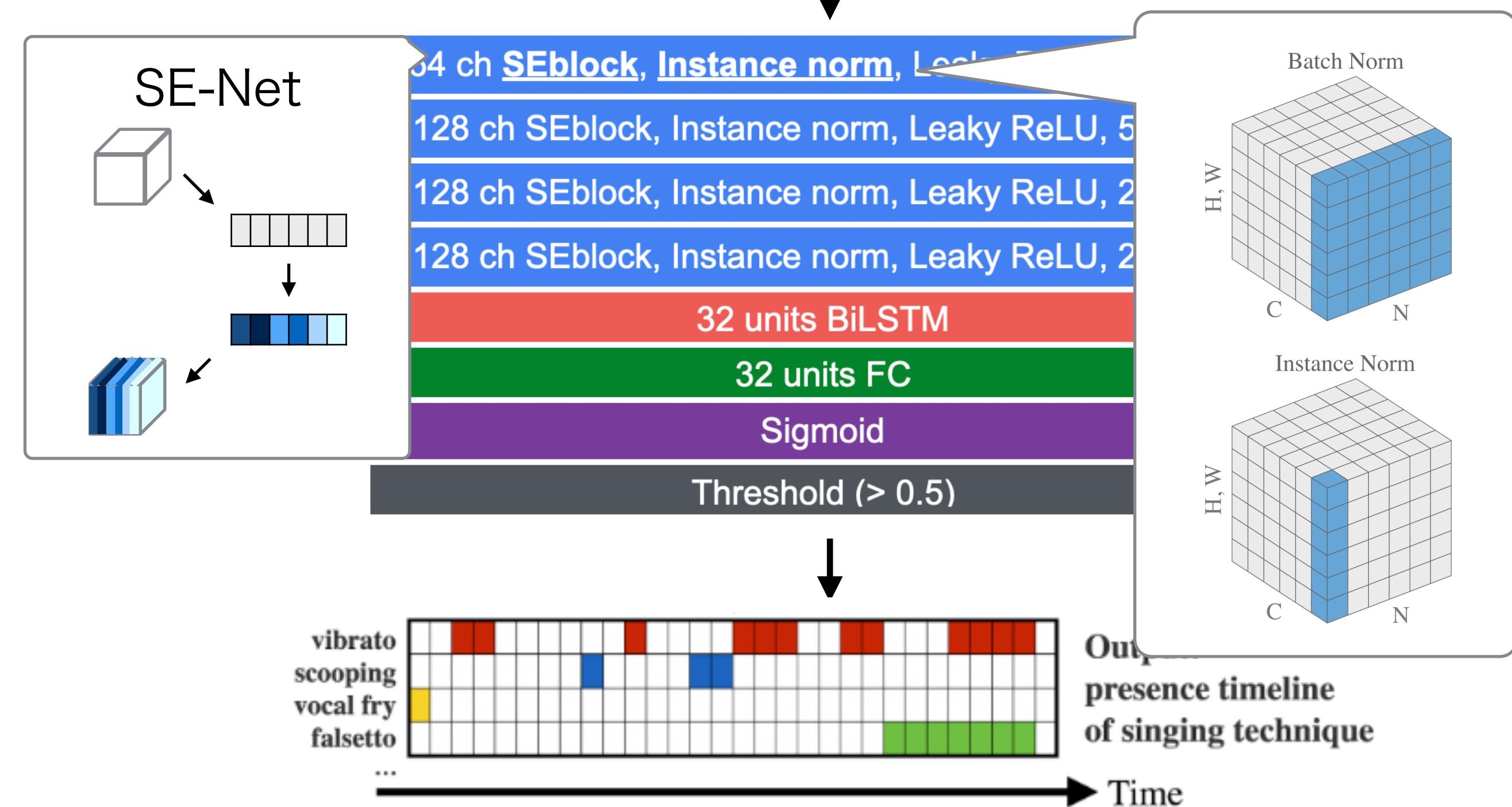
• Input:

- Multi-resolution mel-spectrogram
- Pitch feature -> estimated from CREPE [Kim 18]



• Model

- CNN based feature extraction
 - Squeeze and Excitation (SE-Net) [Hu 18]
 - Instance normalization
- RNN (LSTM) based temporal model



[Kim 18]: J.W.Kim et al. CREPE: A CONVOLUTIONAL REPRESENTATION FOR PITCH ESTIMATION, ICASSP 2018.

[Hu 18] J. Hu et al. Squeeze-and-Excitation Networks. CVPR 2018

[Lin 17] T. Lin et al. Focal loss for Dense Object Detection ICCV 2017

$$Focal = -\alpha_t(1 - p_t)^\gamma \log(p_t)$$

Brief description about each modification of PrimaDNN'

67

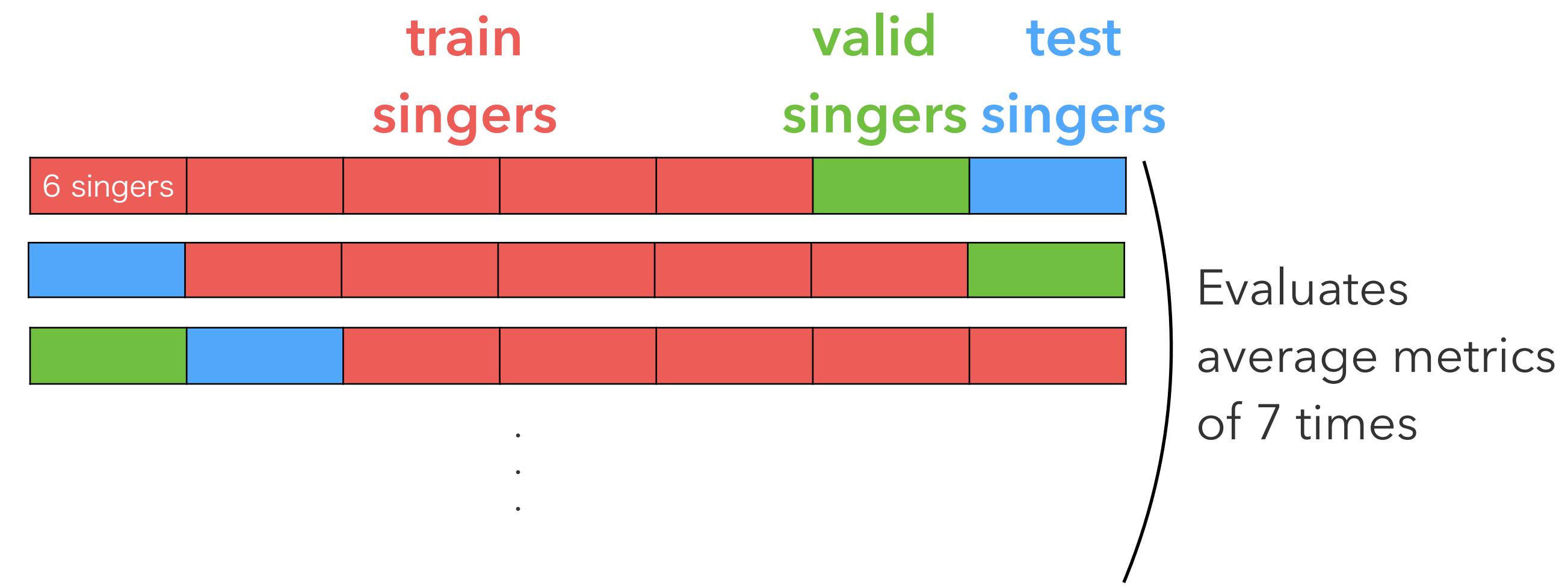
Modifications	Brief explanation	For what?
Pitch feature	Highlights the pitch heights	To take a hint of the melody
Focal loss	Up-weights the hard-samples, Applied for imbalance problems	To enhance the detection of technique segments (sparse)
Multi-resolution	Stacking three different resolutional Mel-spectrograms	To capture a wide type of fluctuation (Similar to Chapter 5)
Squeeze & Excitation (SE)	Weights the importance of feature maps along channel axis	To pick up the more important feature maps -> adapts variation more
Instance Normalization	Calculates normalization at instance level, not batch level	To suppress the effects from Non-targets (singer, effect, etc)

Experimental conditions

68

- **Singer-wise 7-fold cross validation**

- In order to evaluate unseen situation
- Validation set are used for controlling training time (early stopping)



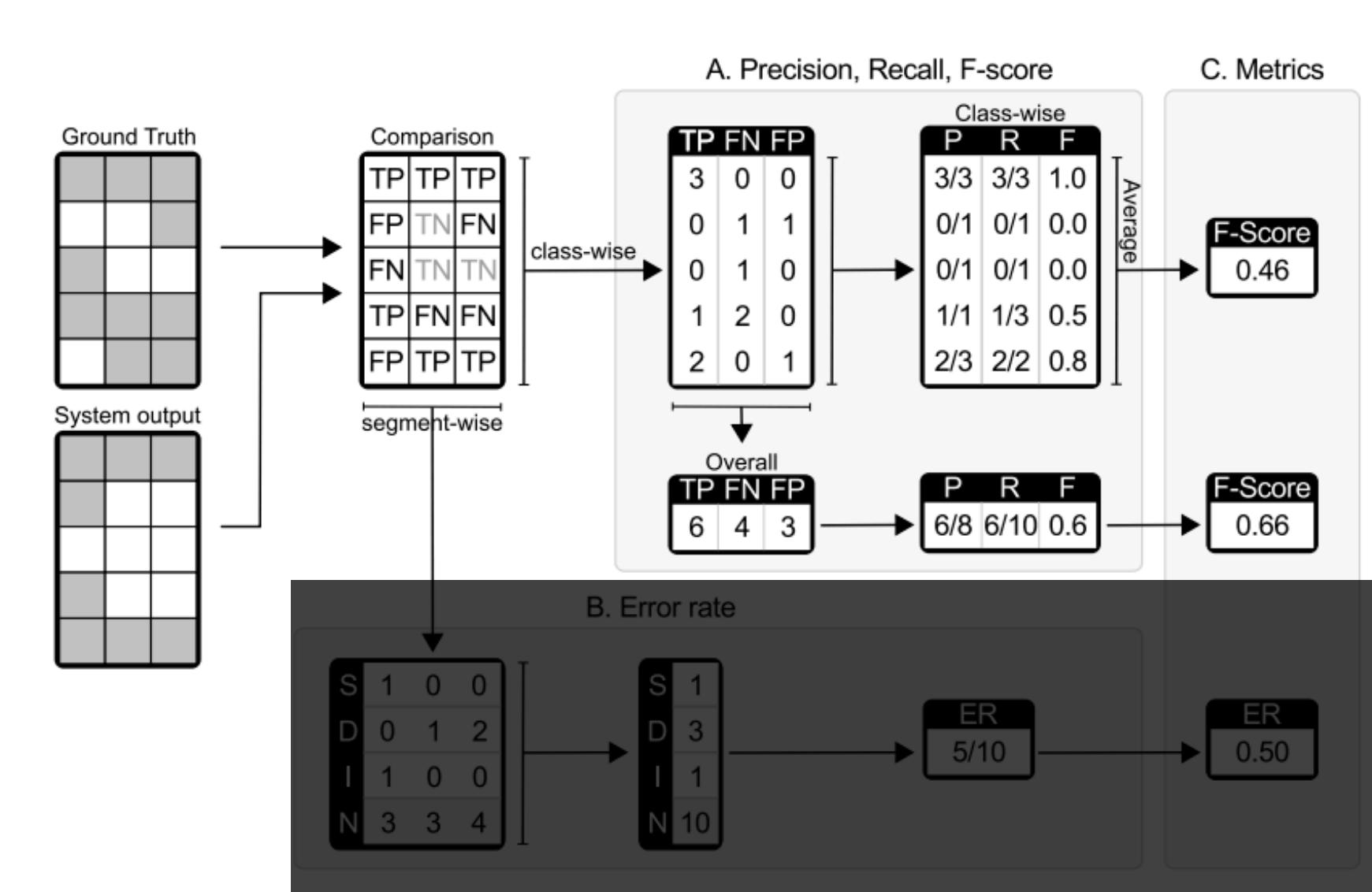
- **Evaluation: Segment-based metrics**

- Precision (P), Recall (R)
- **Macro-F**: class-wise average of F-measure
 - Equally consider every classes
- Micro-F: overall average of F-measure
 - Emphasis more on majority classes
- Segment length: 50 [ms]

$$P = \frac{TP}{TP + FP}$$

$$R = \frac{TP}{TP + FN}$$

$$F = \frac{2 \cdot P \cdot R}{P + R}$$



PrimaDNN' achieved the best performance

Methods	Macro-F	Micro-F	Precision	Recall
eGeMAPS+LSTM [Eyben 15]	0.092	0.063	0.113	0.016
CRNN (versatile) [Imoto 21]	0.377	0.563	0.422	0.392
CRNN-PitchFocal [ISMIR 22] (ours)	0.402	0.551	0.377	0.480
PrimaDNN' [EUSIPCO 23] (ours)	0.449	0.606	0.438	0.483
CNN Self-attention [Imoto 21]	0.420	0.593	0.434	0.477

Hand-crafted feature to DNN

+ Pitch feature & Focal loss

+ Scale-up & Multi-res spec.

& SE-Net & Instance norm.

[Eyben 15] F.Eyben et al. The Geneva Minimalistic Acoustic Parameter Set (GeMAPS) for Voice Research and Affective Computing, IEEE Trans. affective computing 2015

[Imoto 21] K. Imoto et al. Impact of Sound Duration and Inactive frames on Sound Event Detection, ICASSP 2021

[ISMIR 22] Y. Yamamoto et al. Analysis and Detection of Singing Techniques in Repertoires of J-POP Solo Singers, ISMIR 2022

[EUSIPCO 23] Y.Yamamoto et al. PrimaDNN': A Characteristics-aware DNN Customization for Singing Technique Detection, EUSIPCO 2023

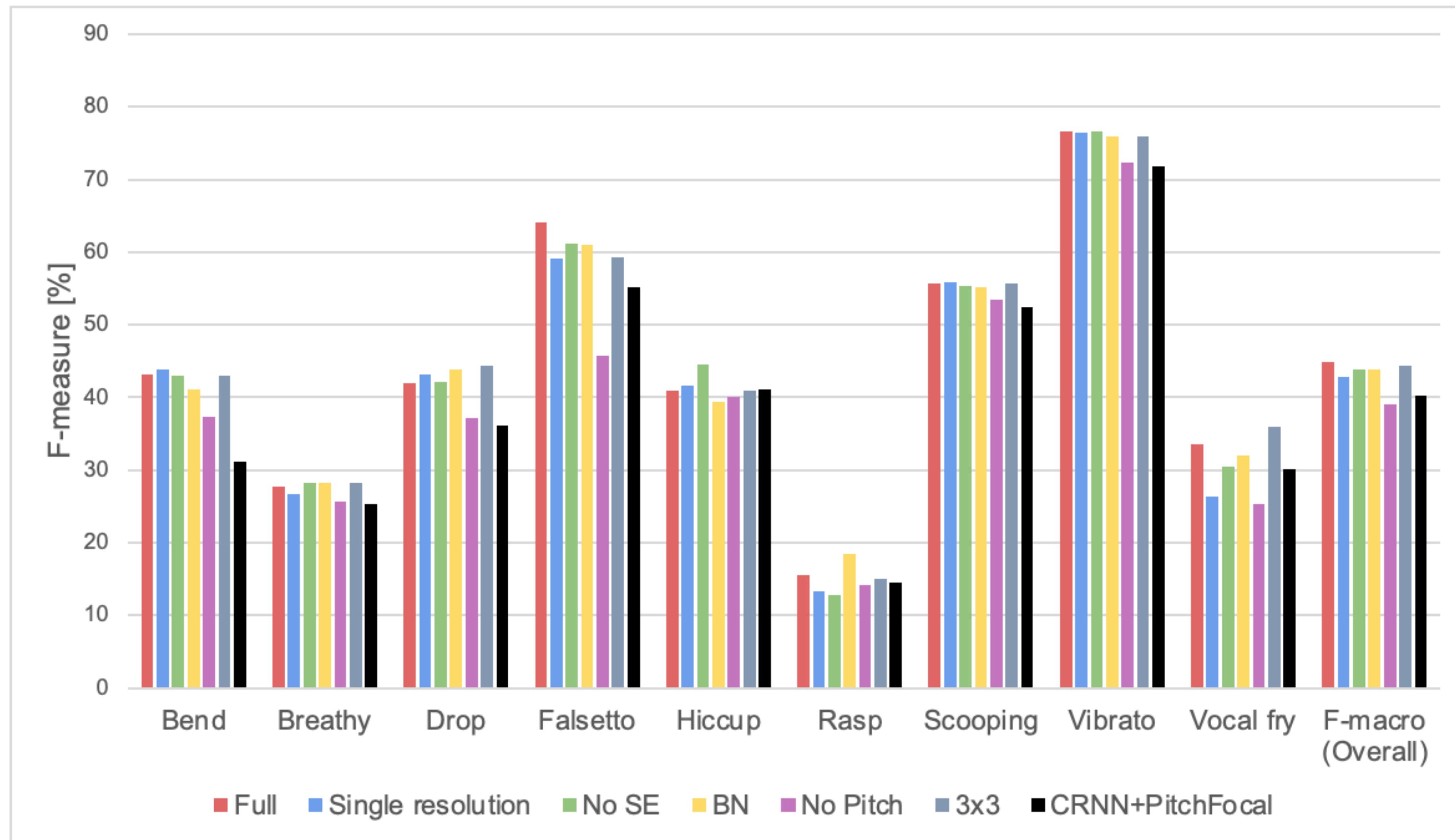
Ablation study

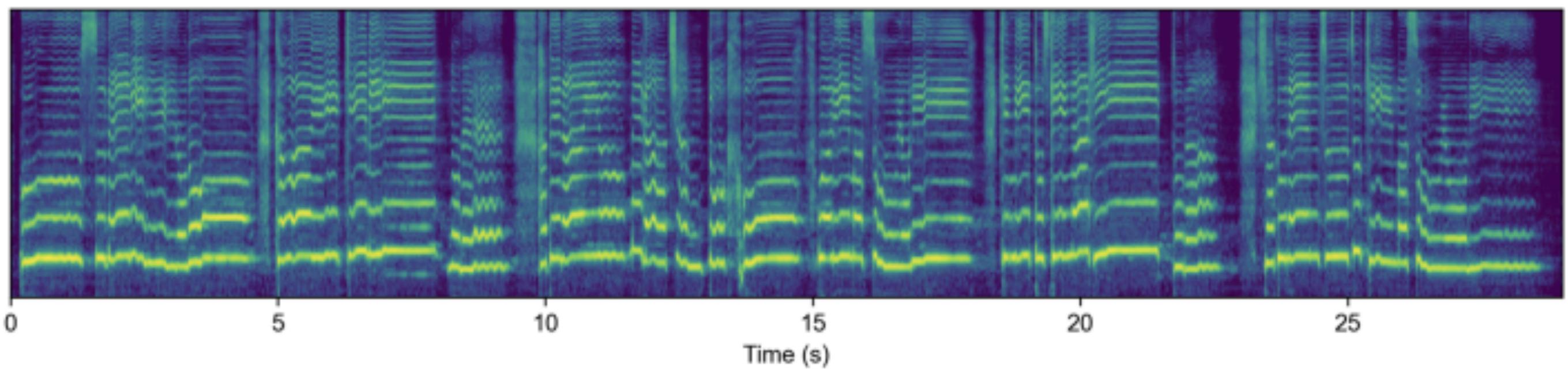
70

Methods	Macro-F	Micro-F	Precision	Recall
Full model	0.449	0.606	0.438	0.483
w/o pitch	0.390	0.548	0.366	0.473
w/o multi-resolution (Only use single resolution)	0.429	0.602	0.441	0.466
W/o Squeeze & Excitation	0.438	0.603	0.430	0.481
W/o Instance norm. (Use batch norm)	0.439	0.596	0.446	0.481
W/o wide kernel (3x3 kernel)	0.443	0.600	0.432	0.488

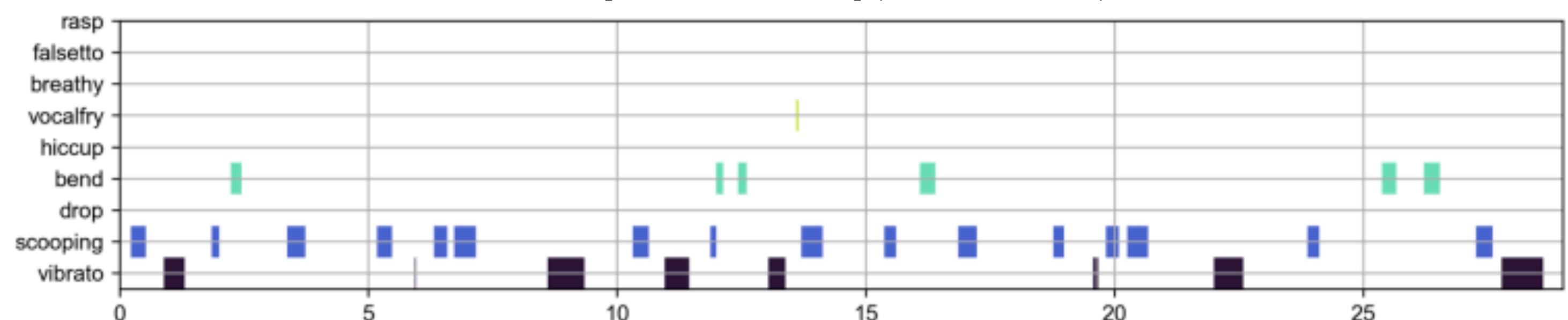
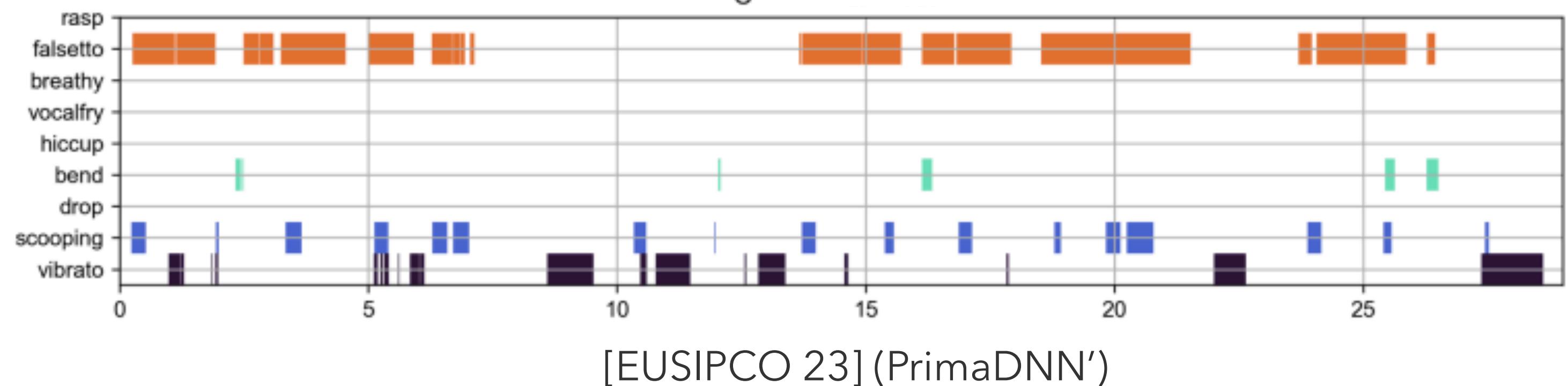
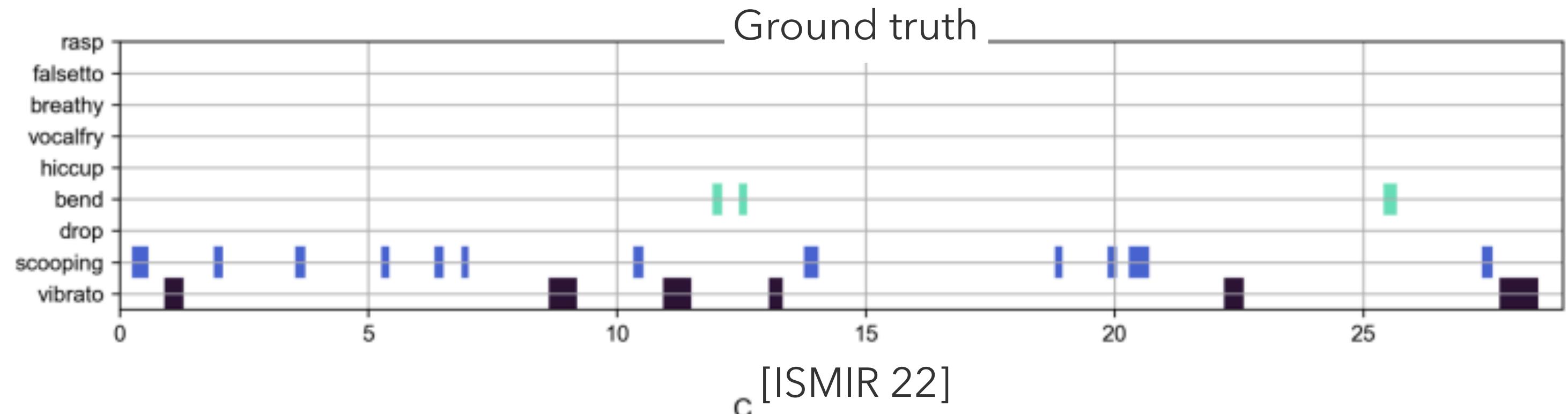
Ablation study

71





♪ Yo Hitoto, Hanamizuki



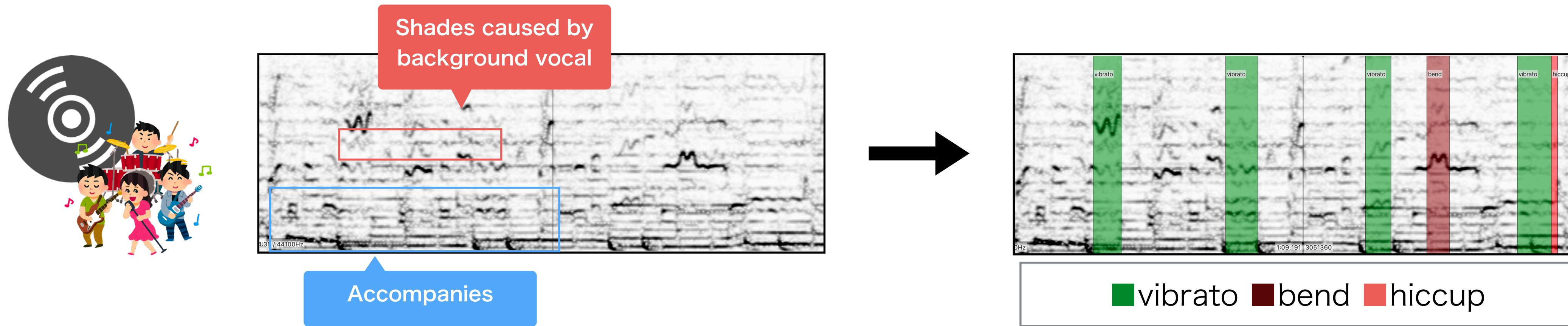
← More example

2: PrimaDNN' is more robust on the localization of pitch techniques

1: CRNN+PitchFocal misidentified falsetto, which is not appeared in the track while PrimaDNN' accurately identifies.

Notable improvements

Extracting appearance of singing techniques on real-world vocal tracks



Singing voices from commercial musics
(i.e., convolved effects & backgrounds)

Goal: Automatically detects
singing technique appearance

- Contributions
 - Build a new MIR task with the first temporally annotated datasets on J-POP tracks
 - PrimaDNN', A SoTA model for 9-way singing technique detection considering characteristics of data (i.e., Input modification, CNN modification, and Focal loss)

Conclusion

- **1. Absence of data and its characteristics**

- What singing techniques should be annotated?
- Need datasets, but how to annotate?
- What is their specific characteristics?

Summary of Techniques
(Ch.3)

Adopt region labels (Ch.3)
→ Demonstrated analysis (Ch.4)

Find some discoveries (Ch.4)

- **2. Less established identification methods**

- How to design/evaluate the automatic model?
- Can we detect techniques from real-world singing tracks?

Modified DNN with spectrogram
input seems better (Ch.5/6)
Data imbalance-aware training is
effective (Ch.5/6)

Built a dataset to evaluate it (Ch.6)
Nine-common techniques can be
detected by customized model at
44.9% on Macro-F(Ch.6)

- **Dataset expansion**
 - Data amount
 - More data and annotation is needed
 - Data quality
 - Associating with other data format (parameter, note-wise annotation, text, etc.)
 - Annotating other information such as vocal notes, lyrics, etc.

- **Further improvement on detection**
 - Transfer learning (e.g., self-supervised speech models)
 - Auxiliary task (e.g., onset detection, activity detection, etc.)
- **The range of detection/scope of techniques**
 - There are still undetectable techniques -> Few-shot learning, Anomaly detection, etc.
 - Wider genres -> Cross genre study
- **Association with other application**
 - Combine with singing transcription or other tasks
 - Leverage the detection results for singing voice analysis, synthesis, etc.

Establishment of computational foundation for singing technique analysis

- **Chapter 3:** Summarized various singing techniques and set the annotation strategy (i.e., annotate observable techniques by region label) to create datasets
- **Chapter 4:** Conducted singing technique analysis using annotation to investigate the relationship with song and singer.
- **Chapter 5:** Explored singing technique classification models and Proposed DNN models based on characteristics-aware feature extraction and imbalance-aware learning
- **Chapter 6:** Established nine-way singing technique detection on real-world pieces with a new dataset and DNN models with characteristics-aware customization

Achievements and Acknowledgements

Publications related to the thesis

80

- **Core papers**
 - [IPSJ 23] Yuya Yamamoto, Tomoyasu Nakano, Masataka Goto, Hiroko Terasawa. **Singing technique analysis with correspondence to musical score on imitative singing of popular music.** IPSJ Journal Vol. 64, No.10, 2023 (in Japanese), (Referred, Prized IPSJ Journal Specially Selected Paper (Equal to top-10%)) **Core paper 1**
 - [ISMIR 22] Yuya Yamamoto, Juhan Nam, Hiroko Terasawa. **Analysis and Detection of Singing Techniques in Repertoires of J-POP Solo Singers.** In Proceedings of the 23rd International Society for Music Information Retrieval Conference (ISMIR), 2022 (Referred, acceptance rate: 43%, oral and poster) **Core paper 2**
- **Others**
 - [APSIPA 21] Yuya Yamamoto, Juhan Nam, Hiroko Terasawa, Yuzuru Hiraga. **Investigating Time- Frequency Representations for Audio Feature Extraction in Singing Technique Classification,** In Proceedings of the 2021 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC), 2021 (Referred, poster)
 - [INTERSPEECH 22] Yuya Yamamoto, Juhan Nam, Hiroko Terasawa. **Deformable CNN and Imbalance-aware Feature Learning for Singing Technique Classification.** In Proceedings of the 23rd Annual Conference of the International Speech Communication Association (INTERSPEECH), 2022 (Referred, acceptance rate: 51%, oral)
 - [EUSIPCO 23] Yuya Yamamoto, Juhan Nam, Hiroko Terasawa. **PrimaDNN': A Characteristics-aware DNN Customization for Singing Technique Detection.** Proceedings of the 31st European Signal Processing Conference (EUSIPCO), 2023 (Referred, poster)
 - [SLIS 21] Yuya Yamamoto, **Establishing foundations for automatic singing technique detection.** (in Japanese), Master dissertation, University of Tsukuba, 2021

7 referred papers, including top-conference (ISMIR, INTERSPEECH) and awarded journal paper

- [APSIPA 23] Yuya Yamamoto, **Toward Leveraging Pre-Trained Self-Supervised Frontends for Automatic Singing Voice Understanding Tasks: Three Case Studies**, In Proceedings of the 2021 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC), 2023 (Referred, oral)
- [EUSIPCO 23] Yuya Yamamoto, Juhani Nam, Hiroko Terasawa. **PrimaDNN': A Characteristics-aware DNN Customization for Singing Technique Detection**. Proceedings of the 31st European Signal Processing Conference (EUSIPCO), 2023 (Referred, poster)
- [IPSJ 23] Yuya Yamamoto, Tomoyasu Nakano, Masataka Goto, Hiroko Terasawa. **Singing technique analysis with correspondence to musical score on imitative singing of popular music**. IPSJ Journal Vol. 64, No.10, 2023 (in Japanese), (Referred, Prized IPSJ Journal Specially Selected Paper (Equal to top-10%))
- [ISMIR 22] Yuya Yamamoto, Juhani Nam, Hiroko Terasawa. **Analysis and Detection of Singing Techniques in Repertoires of J-POP Solo Singers**. In Proceedings of the 23rd International Society for Music Information Retrieval Conference (ISMIR), 2022 (Referred, acceptance rate: 43%, oral and poster)
- [APSIPA 22] Yoshiteru Matsumoto, Hiroyoshi Ito, Hiroko Terasawa, Yuya Yamamoto, Yuzuru Hiraga, Masaki Matsubara. **Human-In-The-Loop Chord Progression Generator With Generative Adversarial Network**, In Proceedings of the 2021 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC), 2022 (Referred, oral)
- [INTERSPEECH 22] Yuya Yamamoto, Juhani Nam, Hiroko Terasawa. **Deformable CNN and Imbalance-aware Feature Learning for Singing Technique Classification**. In Proceedings of the 23rd Annual Conference of the International Speech Communication Association (INTERSPEECH), 2022 (Referred, acceptance rate: 51%, oral)
- [APSIPA 21] Yuya Yamamoto, Juhani Nam, Hiroko Terasawa, Yuzuru Hiraga. **Investigating Time- Frequency Representations for Audio Feature Extraction in Singing Technique Classification**, In Proceedings of the 2021 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC), 2021 (Referred, poster)

Reviewing on top journal, many guest talks, educational activities, etc.

- **Reviewer:** IEEE/ACM Transaction on Audio, Speech, and Language Processing (**IF: 5.4 at 2023**): 2023, 2024
- **Guest talk:**
 - Lightning talk on Music Analysis Meetup (MUANA) 2021, 2022, 2023
 - Komagata junior high school (Career development), 2021
 - **Guest talk in SIGMUS 136** (report on ISMIR 2022), 2023
 - Tsukuba University of Technology, 2024
- **Educational activity:**
 - **Teaching fellow** on Music and Acoustic Information Processing, College of Media Arts, Science and Technology, University of Tsukuba, 2021
 - **Organizing paper reading meetup** of ISMIR 2022, 2023
 - **Organizing weekly lecture on Music and audio X Deep learning**, 2022 (<https://github.com/yamathcy/music-deeplearning-japanese>)
 - Curated list of music and audio processing, 2021 (<https://github.com/yamathcy/awesome-music-informatics>, 150+ stars on 2024 Jan.)

Won 6 awards and Got 3 grants, during the graduate school

- Awards
 - **IPSJ Journal Specially Selected Paper (equal to top-10%)**, from IPSJ, 2023
 - **Sound Symposium Student Excellence Presentation Award**, from IPSJ SIGMUS and SIGSLP, 2023, as a co-author (First author: Tsugumasa Yutani)
 - **IPSJ Yamashita SIG Research Award (equal to the annual best paper)**, from IPSJ, 2023 paper title: Analysis of frequency, acoustic characteristics, and occurrence location of singing techniques using imitated j-pop singing voice (at SIGMUS 132, 2021.)
 - **Best Presentation Award (Best research)**, from IPSJ SIGMUS, 2021
 - **Dean's Award** of University of Tsukuba, 2021
 - **Student Award**, from IPSJ SIGMUS, 2019
- Grants
 - **JST SPRING, tier1 (top-25%)**
 - **Travel Grant of The Telecommunications Advancement Foundation (JPY 190,000)**, 2022
 - **ISMIR student author grant 100 % wavier**, 2022

Acknowledgements

84

- **Dr. Hiroko Terasawa, Dr. Nobutaka Suzuki, Dr. Hiroyoshi Ito:** Supervisors
- **Dr. Atsushi Toshimori, Dr. Shuichi Moritsugu:** The committee members
- **Dr. Juhan Nam:** Co-advisor and The guest committee
- **Dr. Masataka Goto, Dr. Tomoyasu Nakano:** Collaboration and Mentor for research works
- **Dr. Yuzuru Hiraga:** Supervisor (B.S. and M.S.)
- **The members of LSPC**
- **The members of MACLab@KAIST**
- **Researchers/Students whom I met in the conferences**
- **JST SPRING grant program:** Financial support for research and living expenses
- **Vocalists in the world**
- **My family**
- **Yuya Yamamoto:** Myself

Thank you!!

A Computational Approach to Analysis and Detection of Singing Techniques

January 29th, 2024 Ph.D. Defense
Yuya Yamamoto

Supervisor:
Nobutaka Suzuki
Hiroko Terasawa
Hiroyoshi Ito

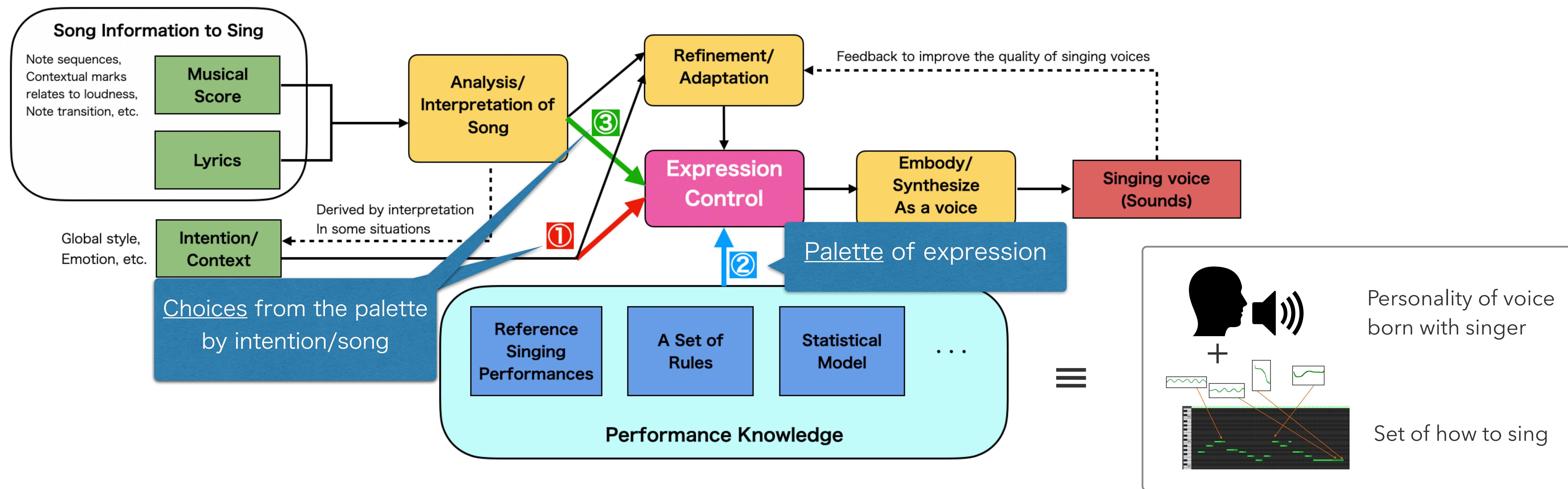
Committee:
Nobutaka Suzuki
Hiroko Terasawa
Atsushi Toshimori
Shuichi Moritsugu
Juhan Nam

Appendix slides

Singing technique is one of the ways to embody the expression

1. The intent of the singer, including their motivations and messages
2. The singing style that the singer possesses
3. The intentions derived from the song (via the sheet music and lyrics)

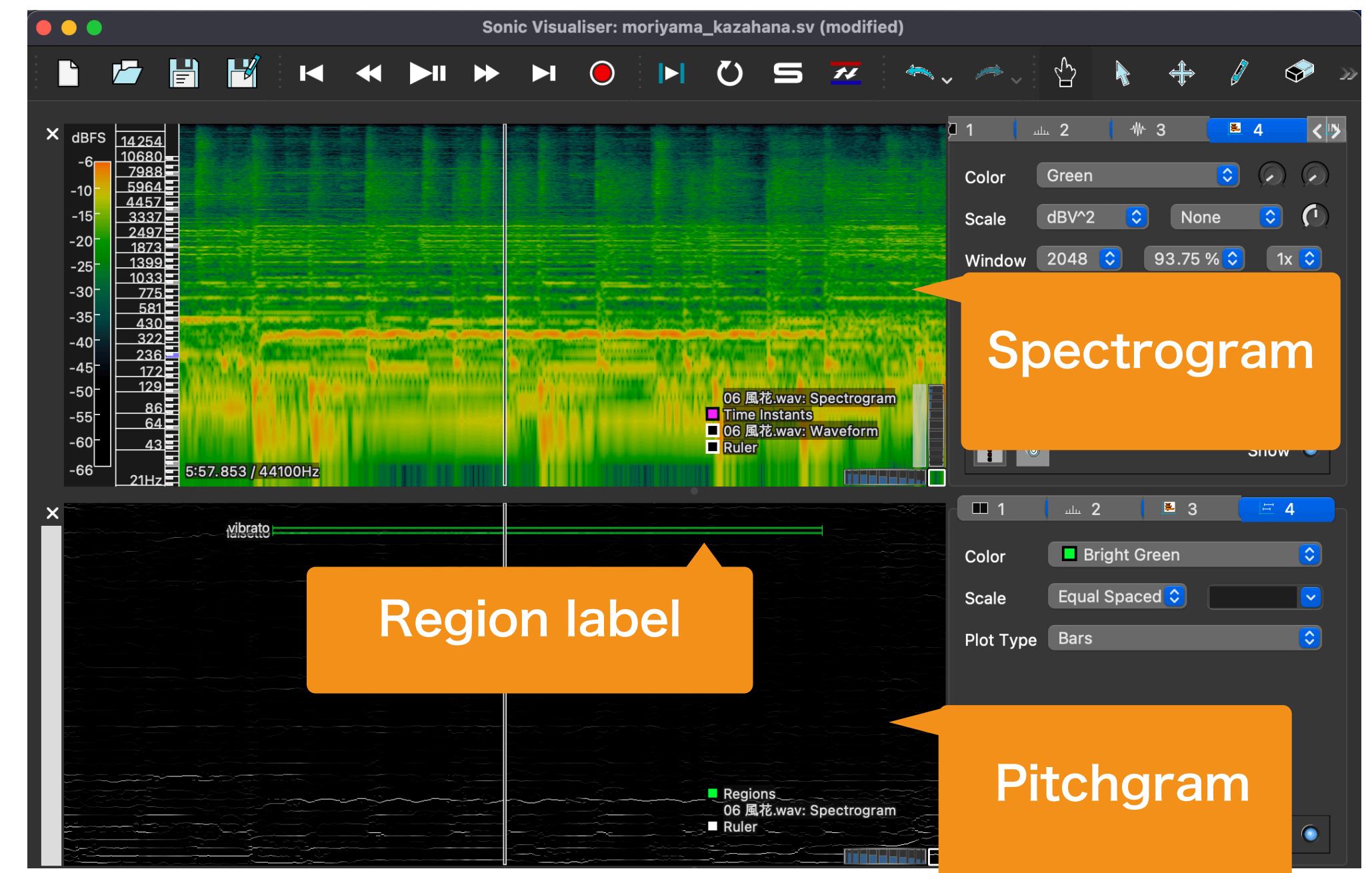
Specifically what?



Annotation process

88

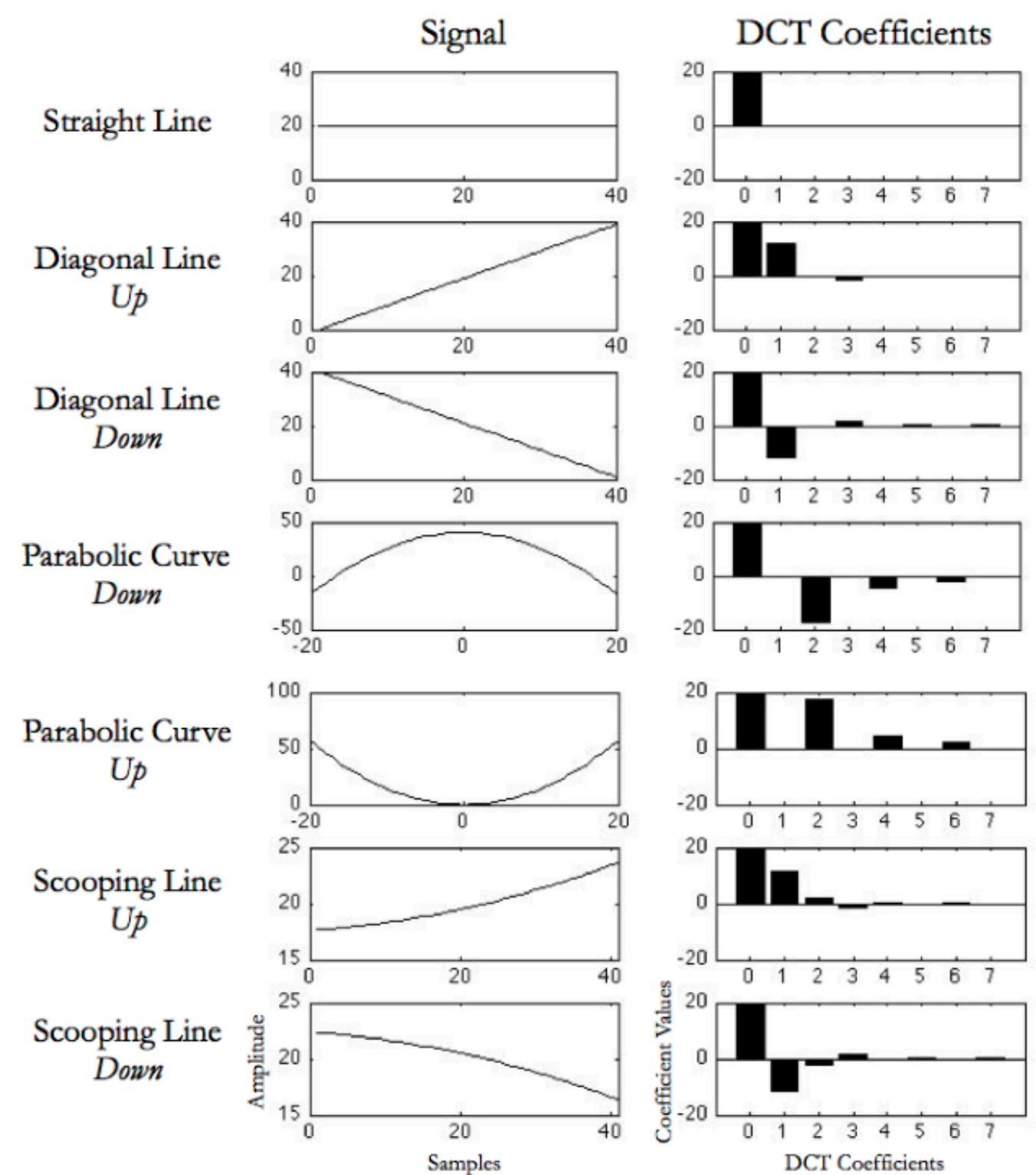
- Manually annotated
 - Annotator: author
 - amateur (no academic degree on music)
 - 9 years of popular vocal, 4 years of chorus (tenor),
 - has relative pitch
 - Software: Sonic visualiser [Cannam 10]
 - visualizing spectrogram and pitchgram
 - both aid of visual & audio feedback
 - set region label on pitchgram



Numeric parameter

89

- “Intuitive” -> in terms of what?
 - Not interpretable for amateur
 - Ex: pitch trajectory representation by discrete cosine transform

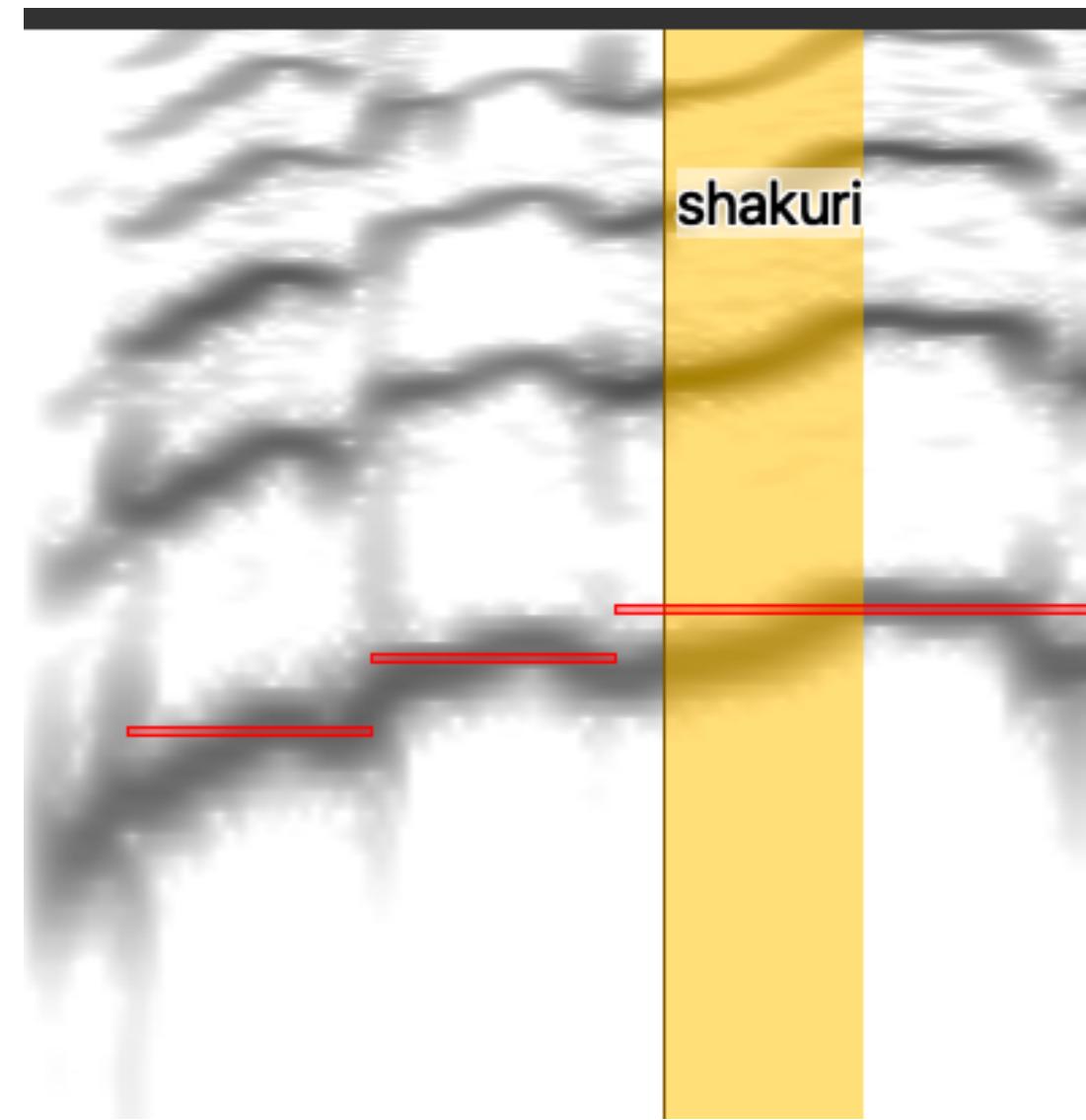


- **Imitation of J-POP famous singers**
 - A Cappella,
 - Original : 24 (12 for each gender)
 - Imitator: professional singer (7 F/M)
 - 48 tracks (two imitator per song)
 - Private database, possessed by AIST
- **Examples**

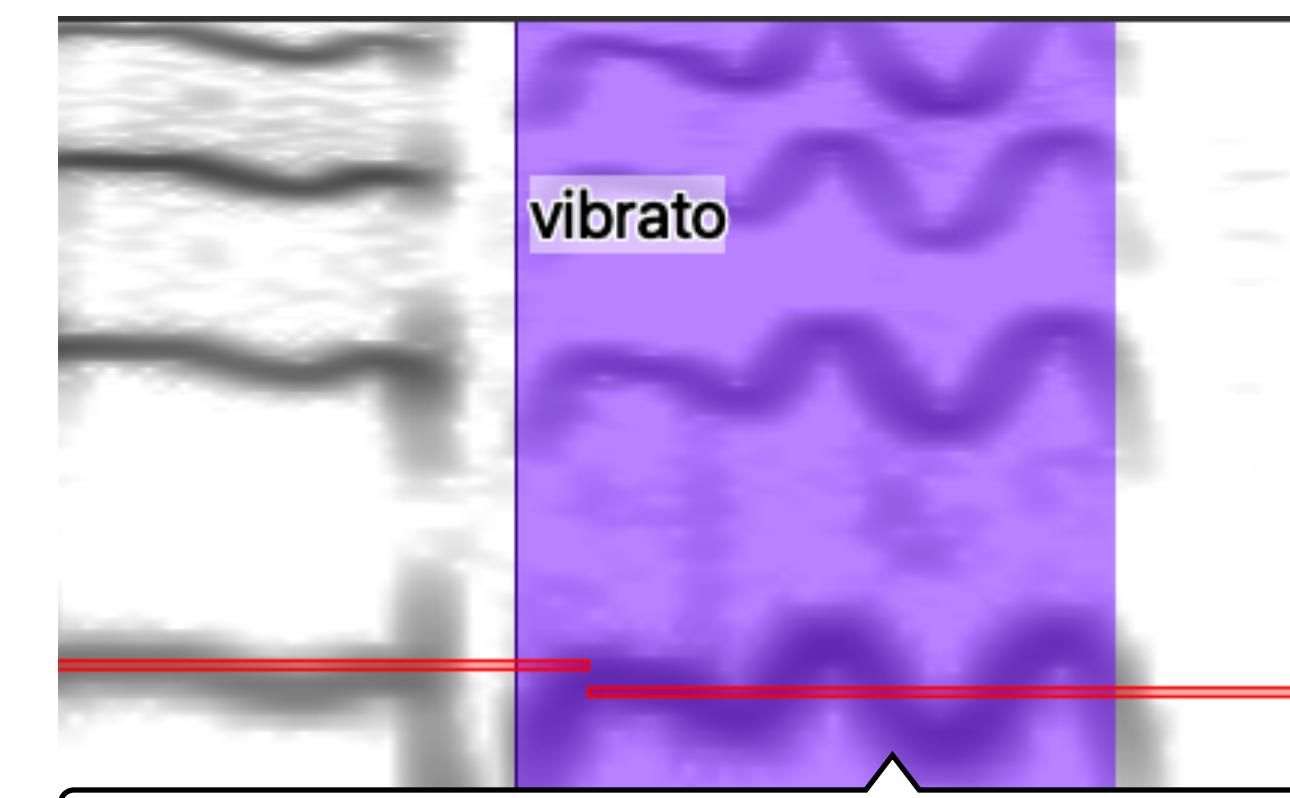
Keisuke Kuwata/ Katte ni Sinbat
 Imitator1 **Imitator2**
 ♪ Su Na Ma Ji Ri No Chi Ga Sa Ki
 Hi To Mo Na Mi Mo Ki E Te…

Original singer	Song name	Gender	Imitator1	Imitator2
玉置浩二 (Tamaki)	出逢い	M	M03	M04
小田和正 (oda)	キラキラ	M	M02	M06
Gackt (gackt)	ありったけの愛で	M	M01	M05
桑田佳祐 (sazan)	勝手にシンドバット	M	M06	M07
チバユウスケ (skapara)	カナリヤ鳴く空	M	M05	M01
西川貴教 (tmr)	Heat Capacity	M	M05	M01
hyde (larcenciel)	Lies and Truth	M	M03	M04
平井堅 (hirai)	瞳を閉じて	M	M01	M05
福山雅治 (fukuyama)	桜坂	M	M04	M03
楳原敬之 (makihara)	桃	M	M04	M03
森山直太朗 (moriyama)	さくら (独唱)	M	M02	M06
山崎まさよし (yamazaki)	未完成	M	M06	M07
aiko (aiko)	ボーイフレンド	F	F01	F06
絢香 (ayaka)	三日月	F	F05	F02
宇多田ヒカル (utada)	Can You Keep A Secret?	F	F03	F04
鬼束ちひろ (onitsuka)	月光	F	F03	F04
倖田來未 (koda)	夢のうた	F	F06	F01
小柳ゆき (koyanagi)	愛情	F	F04	F07
chara (chara)	大切をきずくもの	F	F02	F05
浜崎あゆみ (hamasaki)	seasons	F	F06	F01
一青窈 (hitotoyo)	ハナミズキ	F	F05	F02
平原綾香 (hirahara)	明日	F	F03	F04
松浦亜弥 (matsuura)	♡ 桃色片思い ♡	F	F01	F06
YUKI (jam)	motto	F	F02 (key-1)	F05

Defining assignment rule

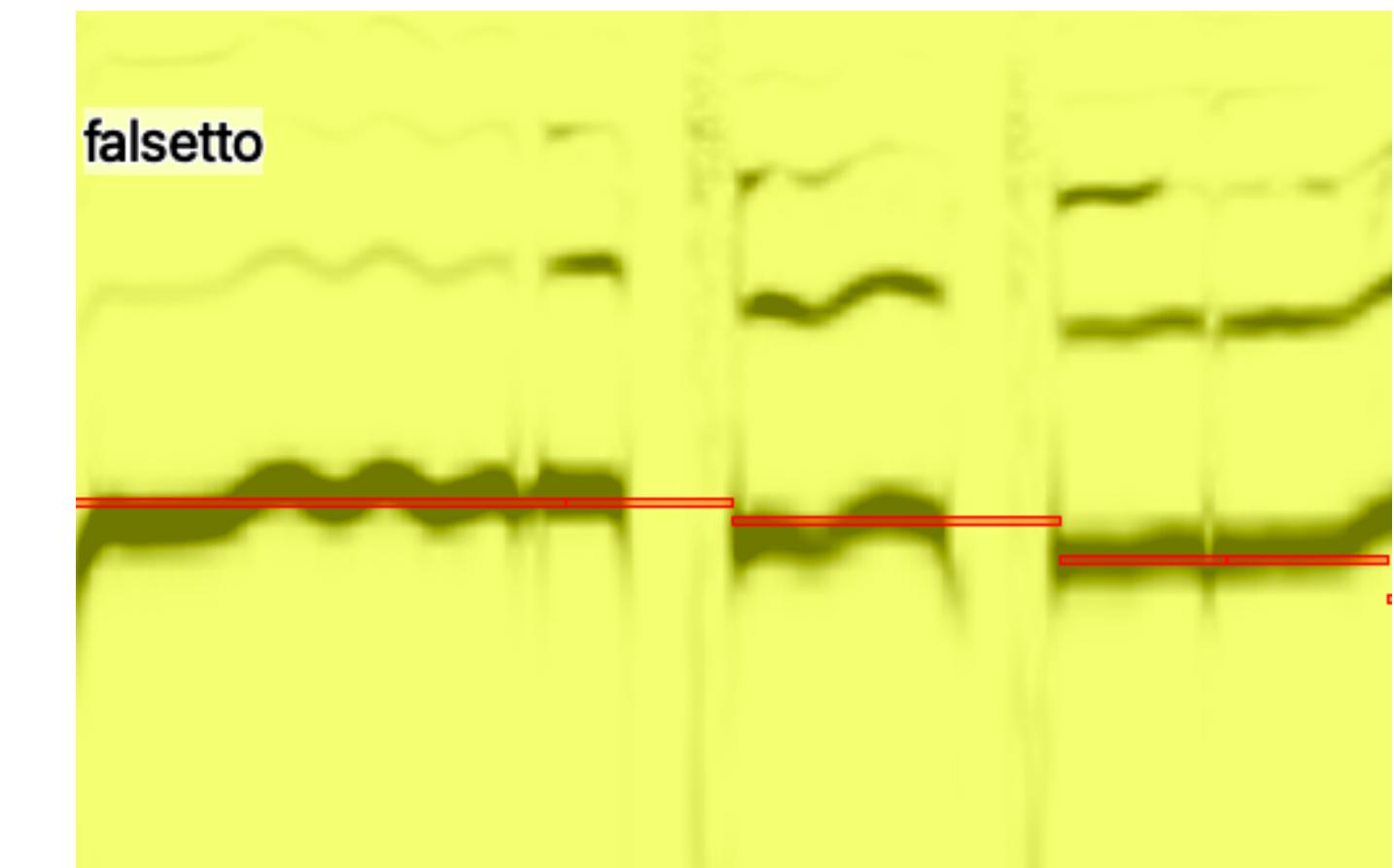


Case1: in-note, single
-> just assign the note



The latter note has longer overlap, vibrato will be assigned the note

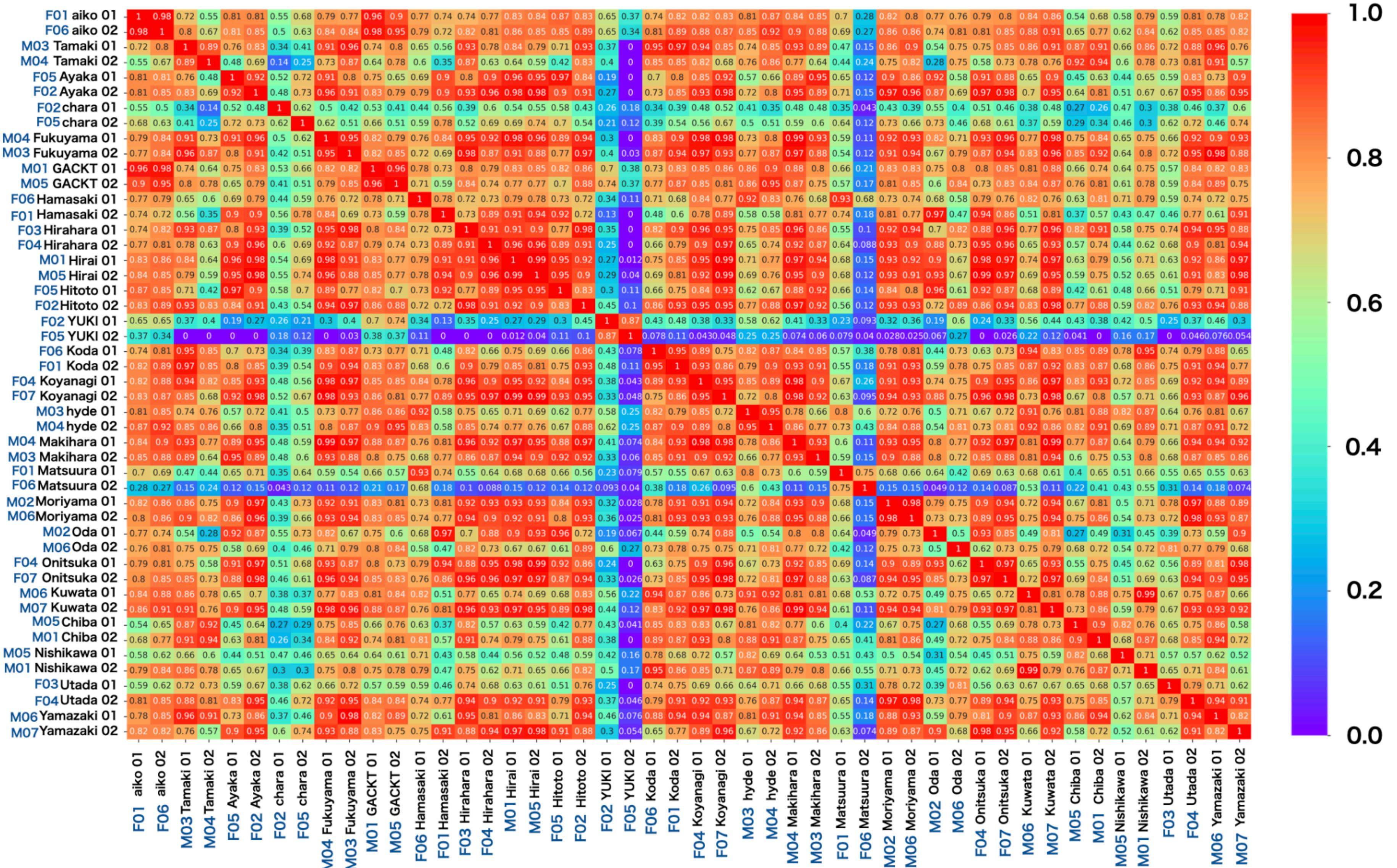
Case2: spanning across two notes
-> assign the note with longer overlap



Case3: spanning across more than three notes
-> assign all notes

Cosine-similarity of every singer

92



Detailed results by techniques (pitch)

Technique	Discovery indicated from the analysis
Vibrato	<ul style="list-style-type: none">• The most frequent technique• Depends on each, in terms both of extent and rate (100-300cent, 5-8Hz)• Professional vocalist may be able to imitate the characteristics of vibrato by other singer (↔ Not for amateur [Saitou 11])• Frequent on… : long notes, constant height portion, and phrase tails• Negative correlation between vibrato extent and pitch height (deep-low)• Negative correlation between vibrato rate and pitch duration (slow-long)
Scooping	<ul style="list-style-type: none">• Frequent technique• Frequent on… : long notes, ascending note heights, high notes
Bend	<ul style="list-style-type: none">• Depends on singer (Frequent on imitators of Yo Hitoto, Kumi Koda, and Ken Hirai, in the data)• Frequent on… : short notes, constant height portion
Drop	<ul style="list-style-type: none">• Depends on singer (Frequent on imitators of hyde, GACKT, aiko)• Frequent on… : short notes, constant height portion, the note to descend in next
Hiccup	<ul style="list-style-type: none">• Depends on singer (Frequent on imitators of hyde, Aya Matsuura , Takanori Nishikawa)• Frequent on… : short notes, constant height portion
Melisma	<ul style="list-style-type: none">• Depends on singer (Frequent on imitators of Hikaru Utada, Masayoshi Yamaguchi)• Depend on song

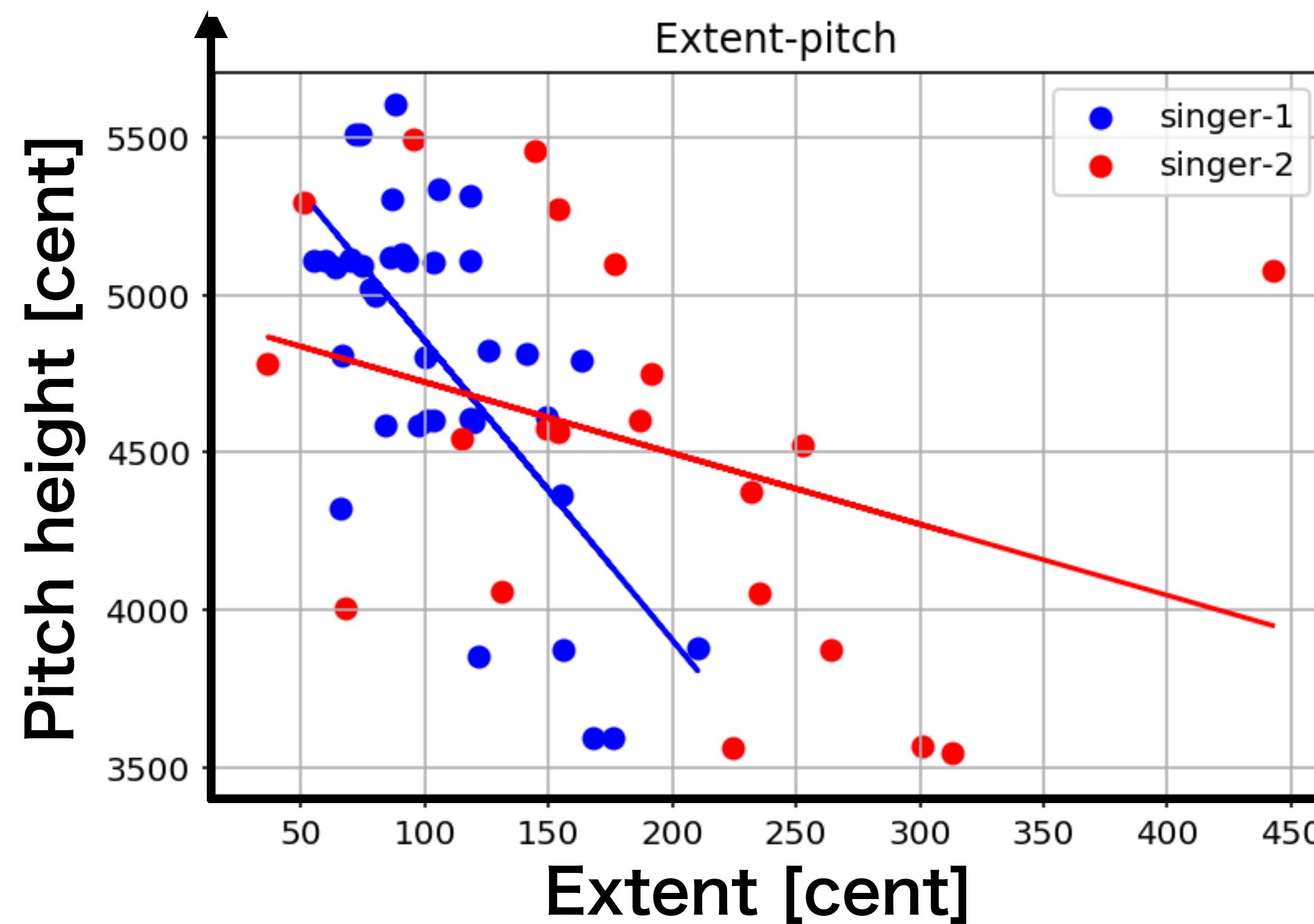
Detailed results by techniques (timbre and others)

Techniques	Discovery indicated from the analysis
Vocal fry	<ul style="list-style-type: none">Depends on the singer and song (Frequent on imitators of Yuki Koyanagi and Ken Hirai)Frequent on… : Short notes, constant height portion, The head of phrases
Falsetto	<ul style="list-style-type: none">Common occurrence and between two imitatorsFrequent on… : High notes, Ascending note heights, The note to descend in next
Breathy	<ul style="list-style-type: none">Depends on the singer and song (Frequent on imitators of Ayaka Hirahara and Hikaru Utada)Long, Occurring across multiple notes
Whisper	<ul style="list-style-type: none">Depends on the singer and song (Frequent on imitators of Chara)Long, Occurring across multiple notes
Shout	<ul style="list-style-type: none">Depends on the song
Spoken	<ul style="list-style-type: none">Depends on the song
Tongue trill	<ul style="list-style-type: none">Depends on the song, only on an imitator of Keisuke Kuwata

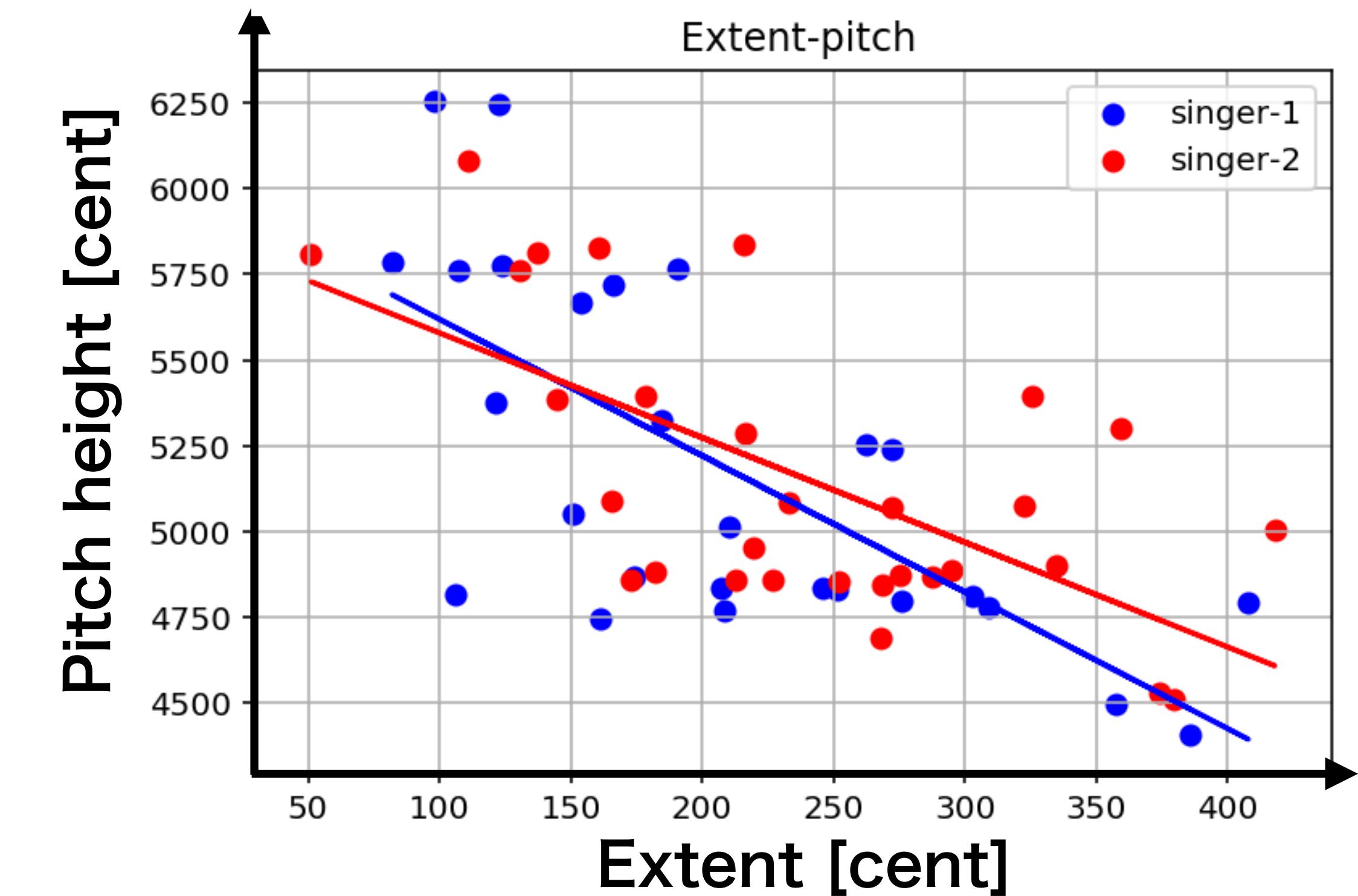
Extent-Pitch (negative correlation)

95

Two imitation singers of “Hitomi wo tojite”
By Ken Hirai



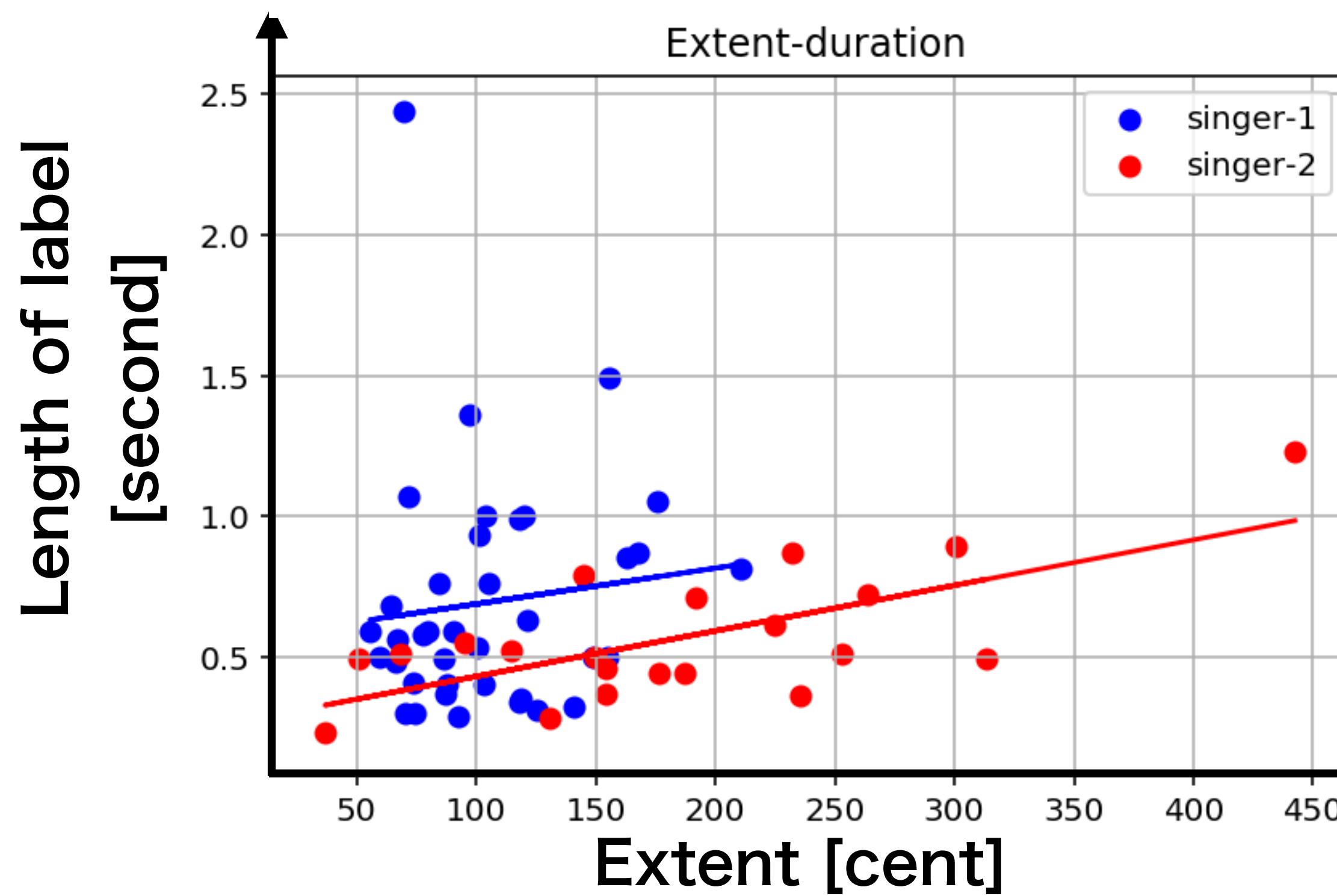
Two imitation singers of “Aijyou”
By Yuki Koyanagi



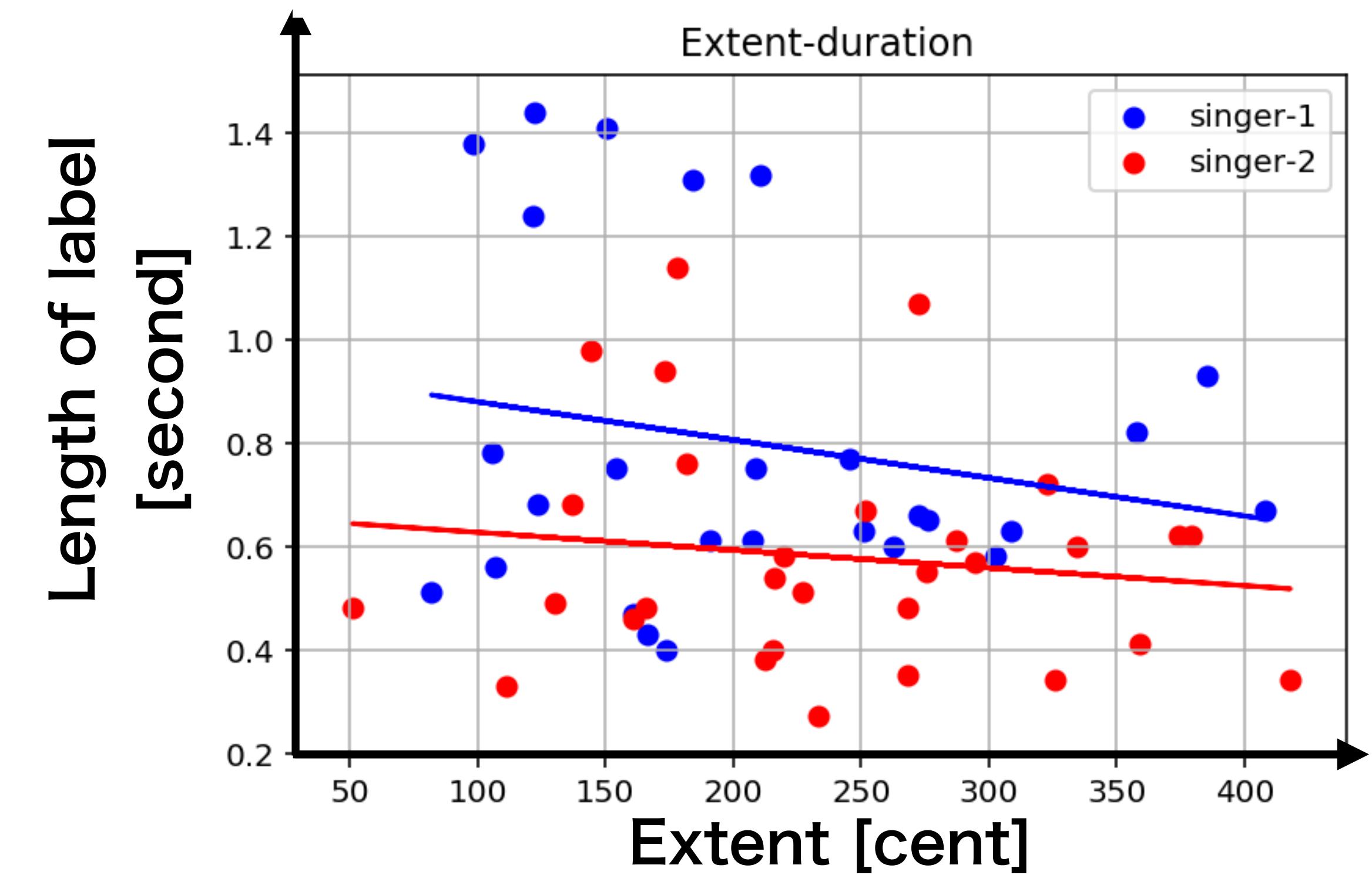
Extent-Duration

96

Two imitation singers of “Hitomi wo tojite”
By Ken Hirai



Two imitation singers of “Aijyou”
By Yuki Koyanagi

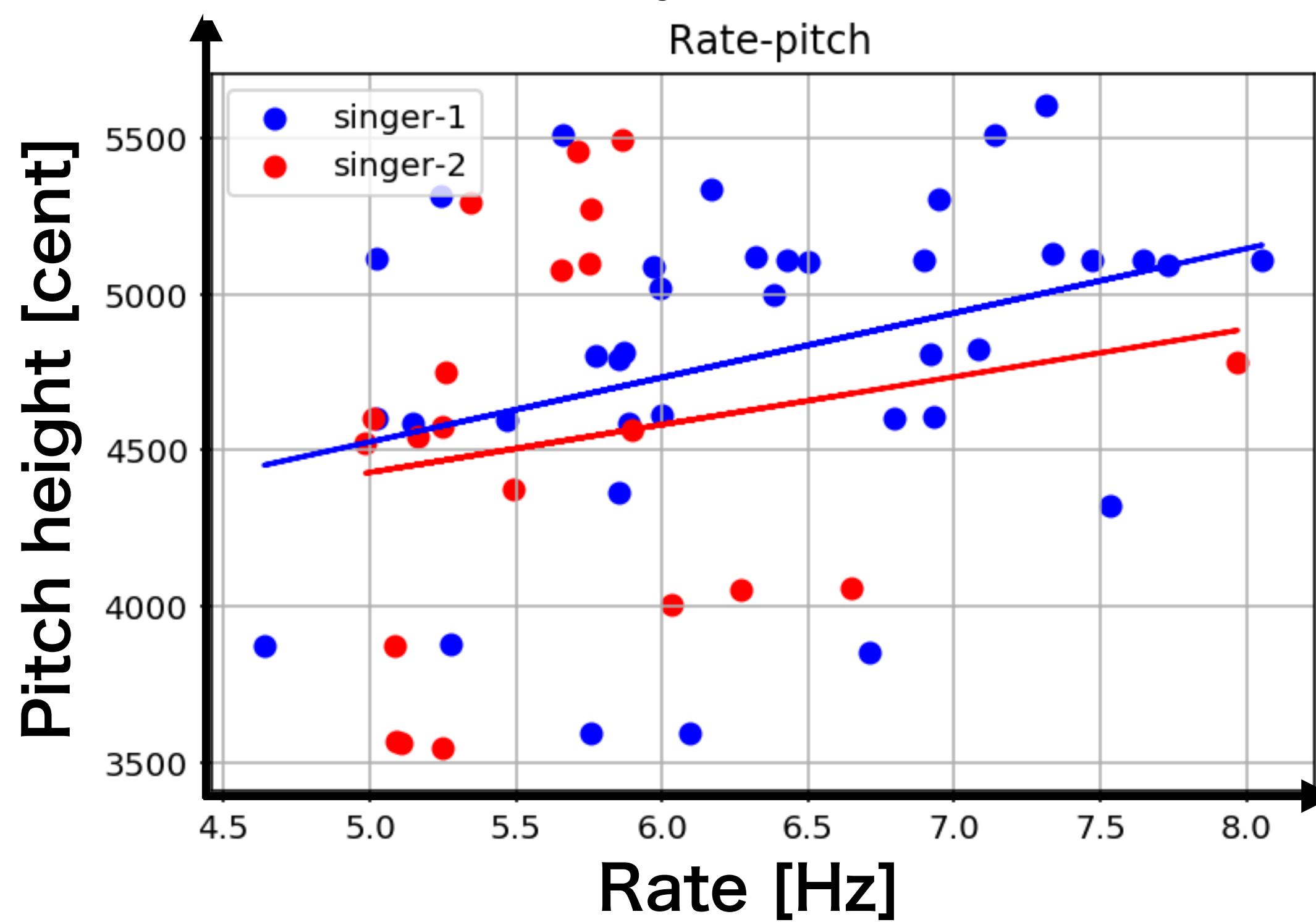


Rate-Pitch

97

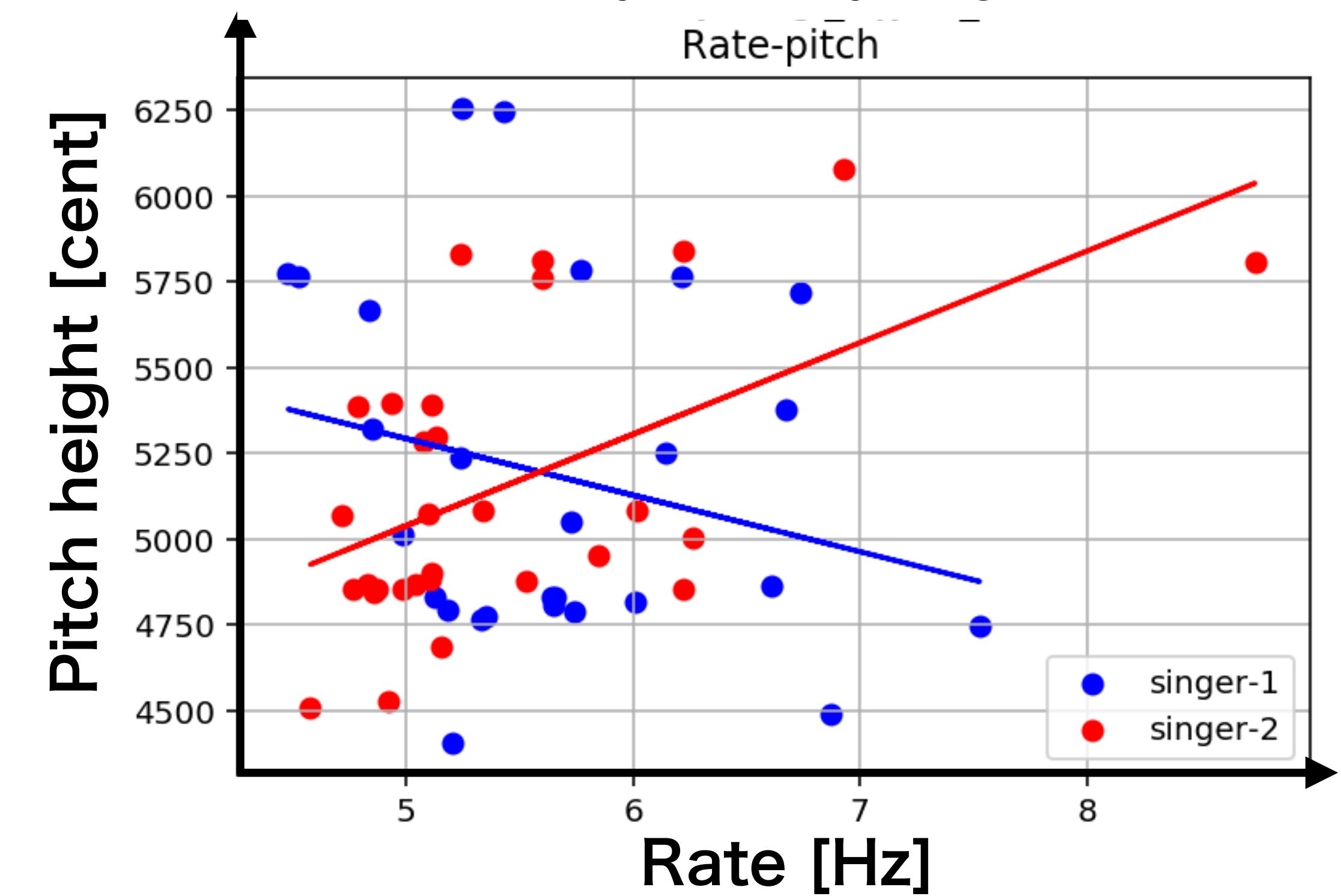
Two imitation singers of “Hitomi wo tojite”

By Ken Hirai



Two imitation singers of “Aijyou”

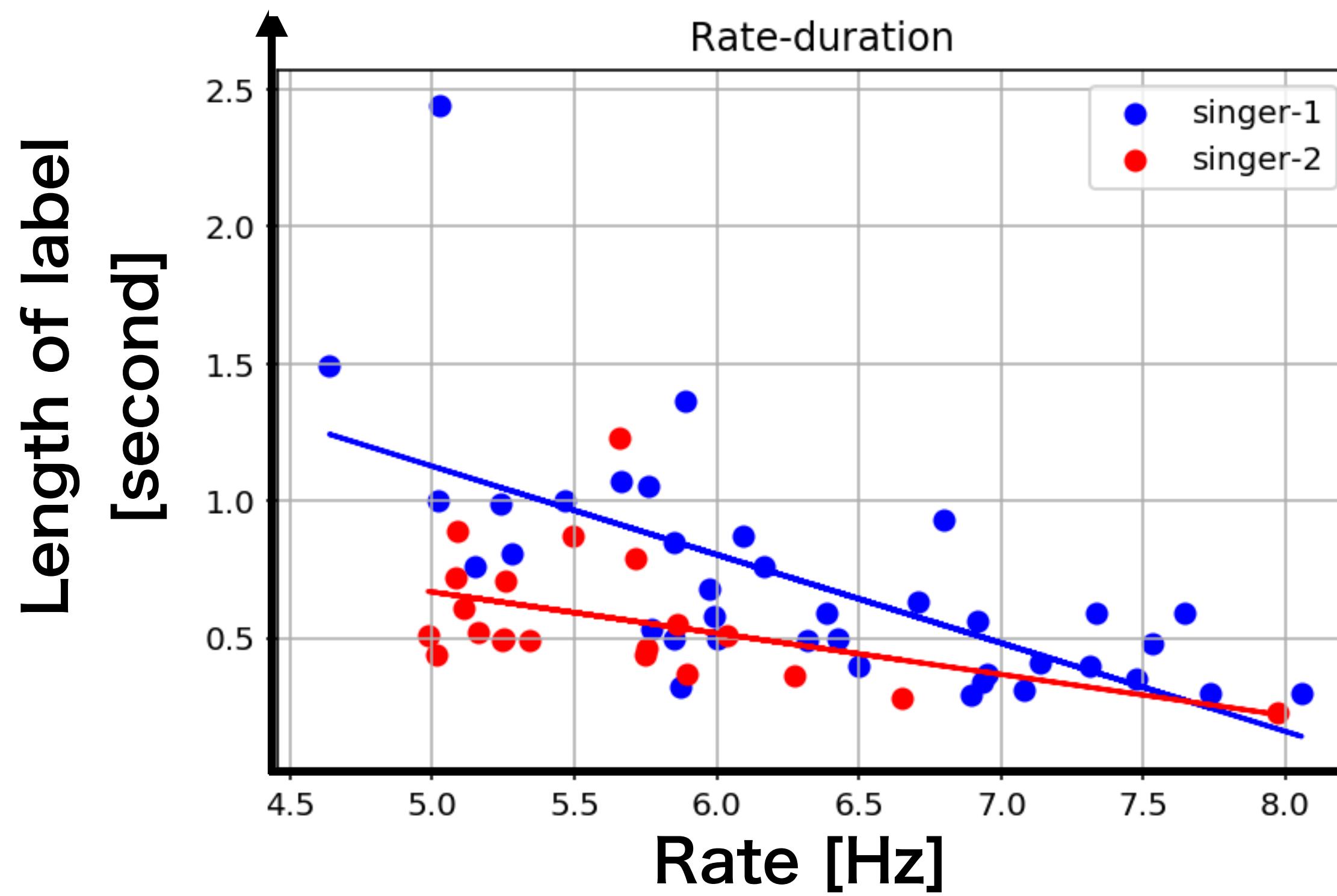
By Yuki Koyanagi



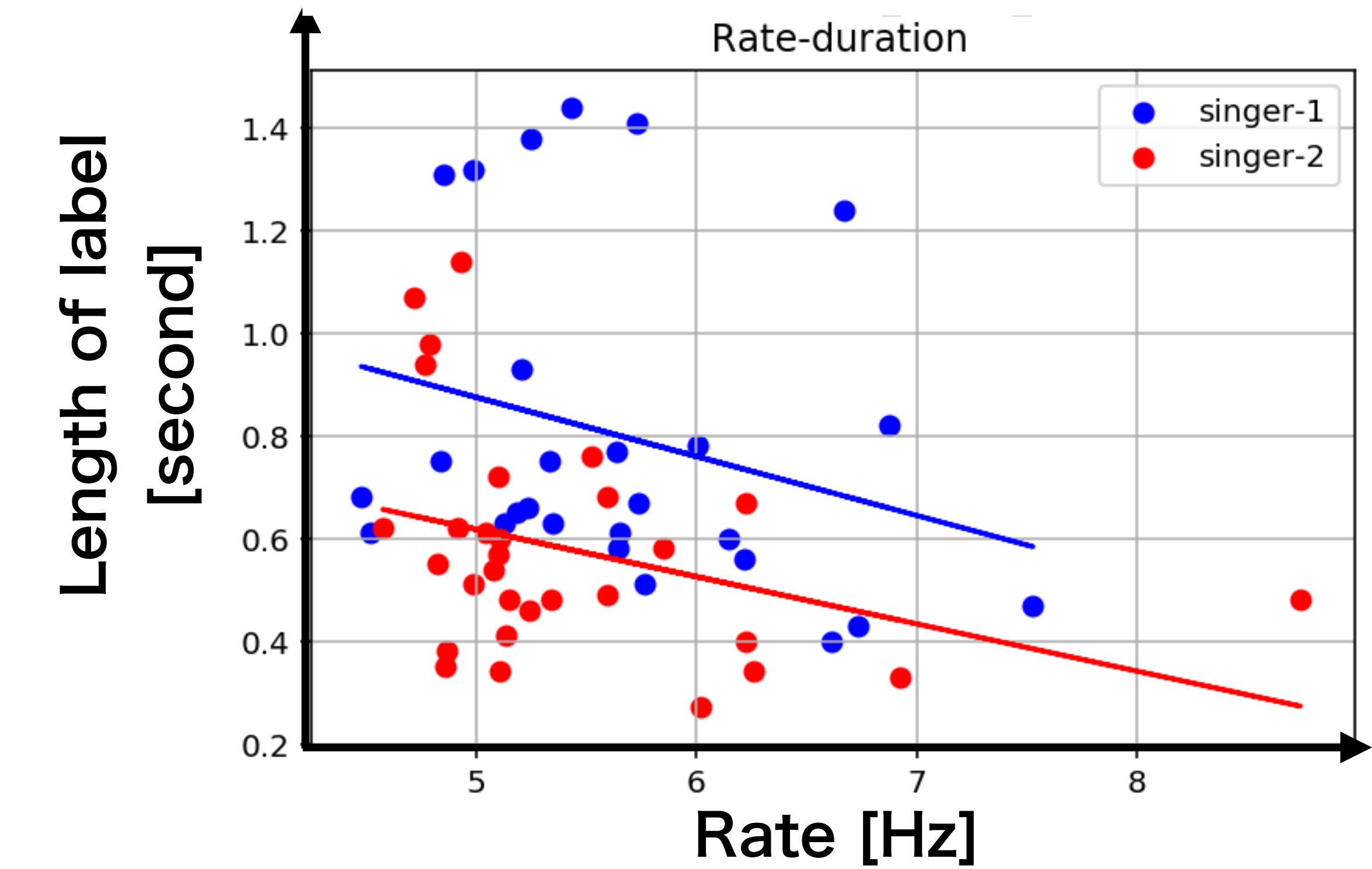
Rate-Duration (negative correlation)

98

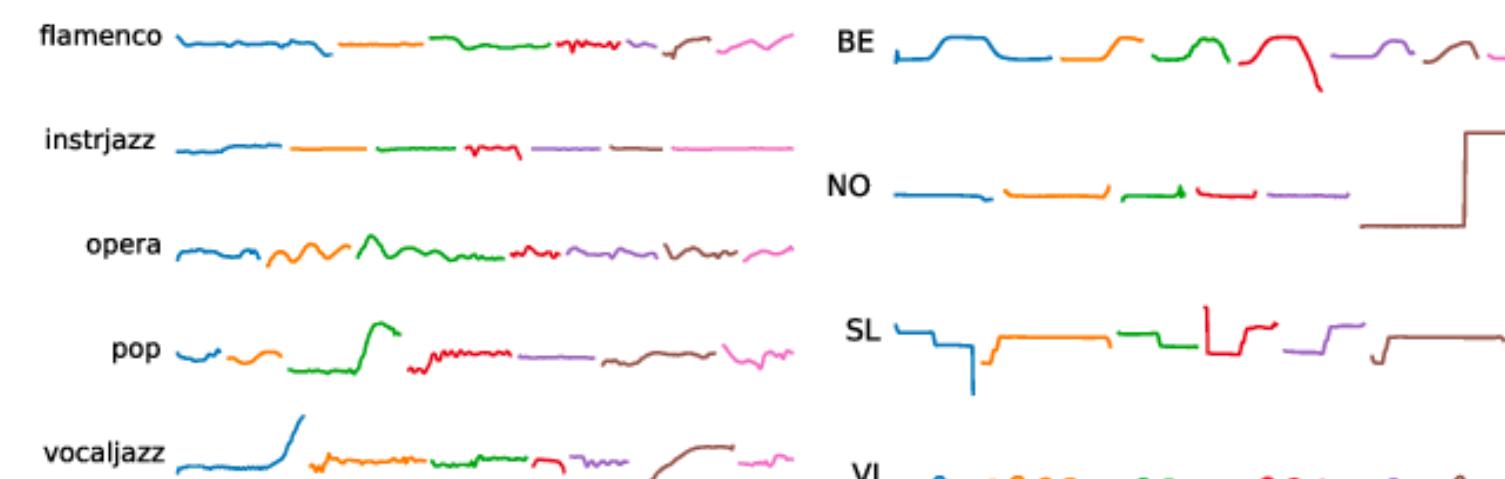
Two imitation singers of “Hitomi wo tojite”
By Ken Hirai



Two imitation singers of “Aijyou”
By Yuki Koyanagi

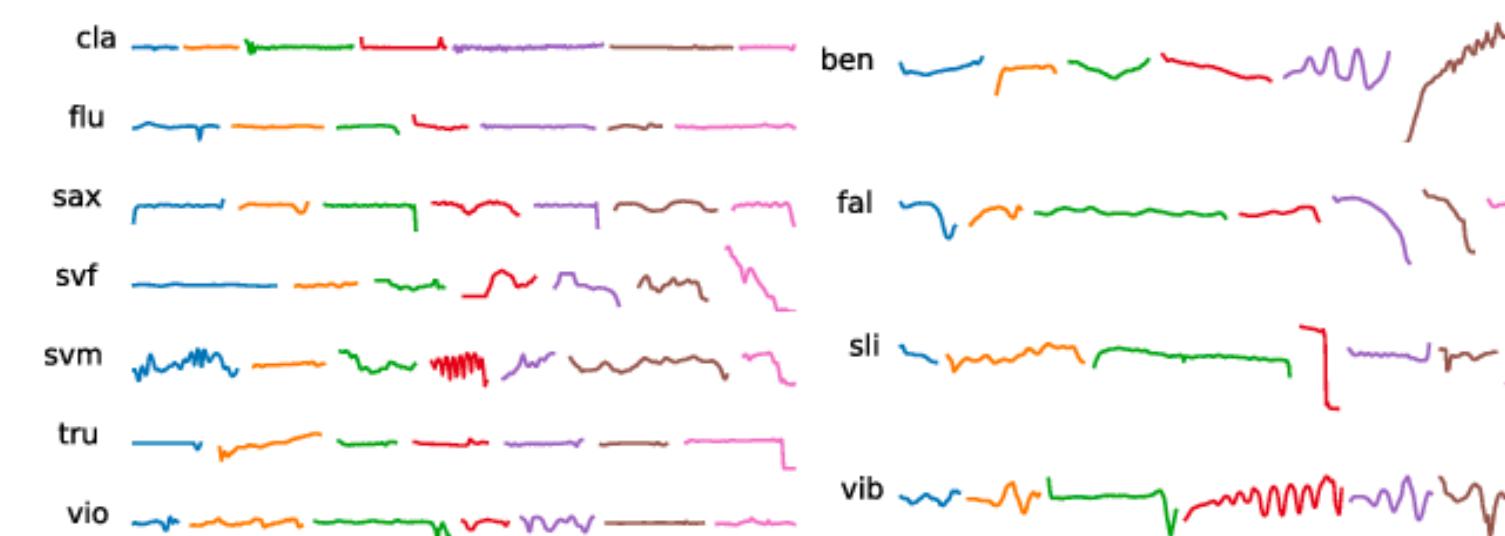


- A comparison on hand-crafted feature and CNN learnt feature on musical playing technique identification (only pitch technique)
- CNN can achieve equivalent or higher performance -> how about singing techniques?



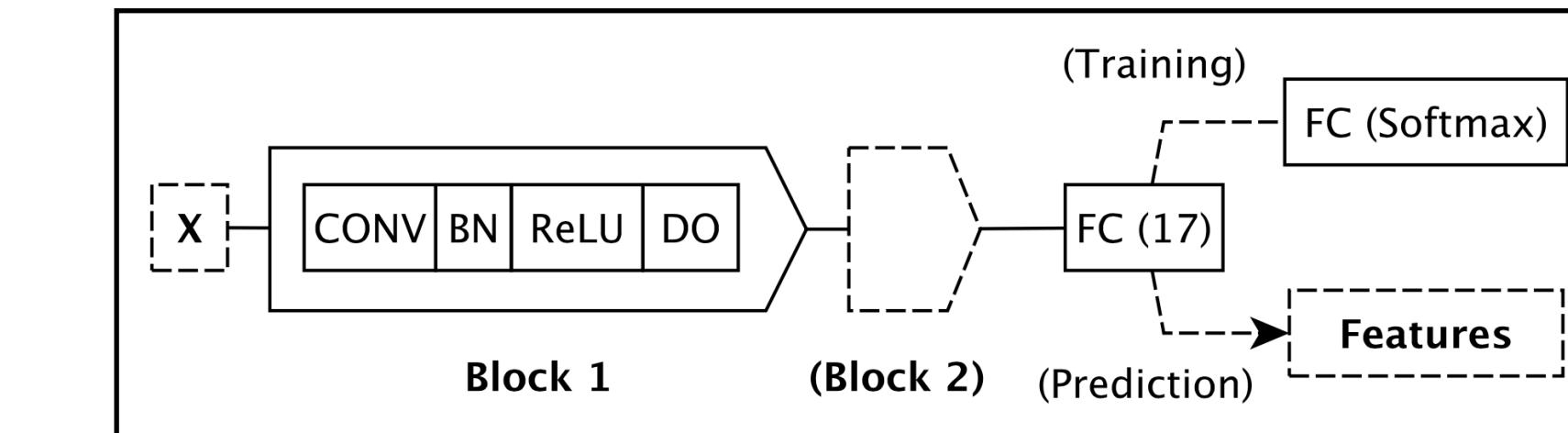
(a) GENRE

(b) GUITAR



(c) INST

(d) WJD

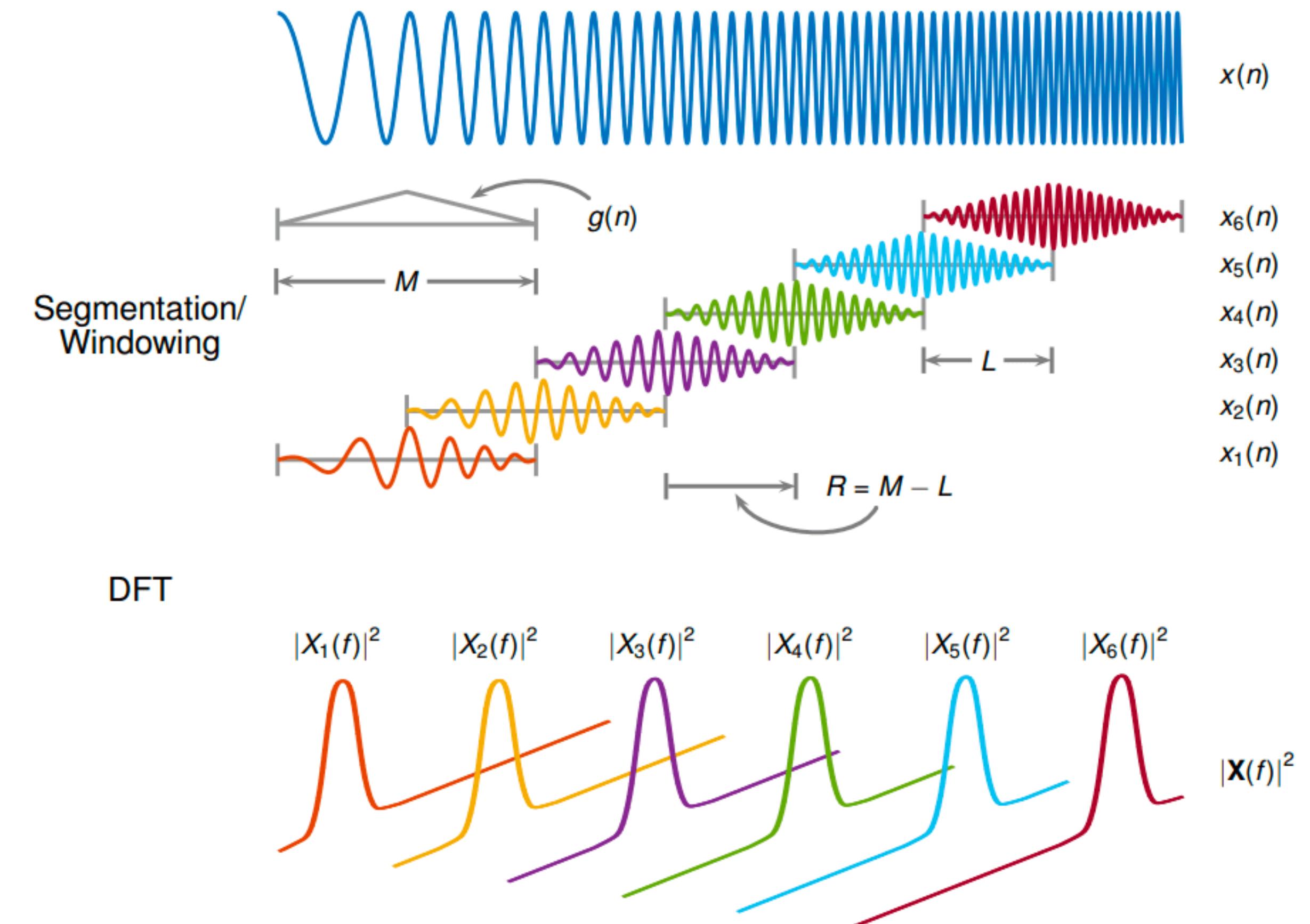


Extractor	Aggr. Dataset	Contour-Level			File-Level	
		GENRE	GUITAR	INST	WJD	GENRE
BITTELI	C	0.51	0.97	0.38	0.87	0.73
	SC	0.54	0.96	0.43	0.82	0.76
PYMUS	C	0.53	0.98	0.35	0.87	0.79
	SC	0.55	0.97	0.31	0.83	0.85
CNN-1	SC	0.54	0.95	0.34	0.83	0.85
CNN-2	SC	0.63	0.96	0.43	0.84	0.94
						0.67

Short-time Fourier transform

100

- Given input signal $x(n)$,
 - Chunked by window $g(n)$, that has size of M
 - Apply the Discrete Fourier transform on the chunked signal
 - Then, derived $X(f)$ is a complex value, whose magnitude is amplitude, and the argument is phase
 - Typically, the amplitude spectrogram is only used for the representation



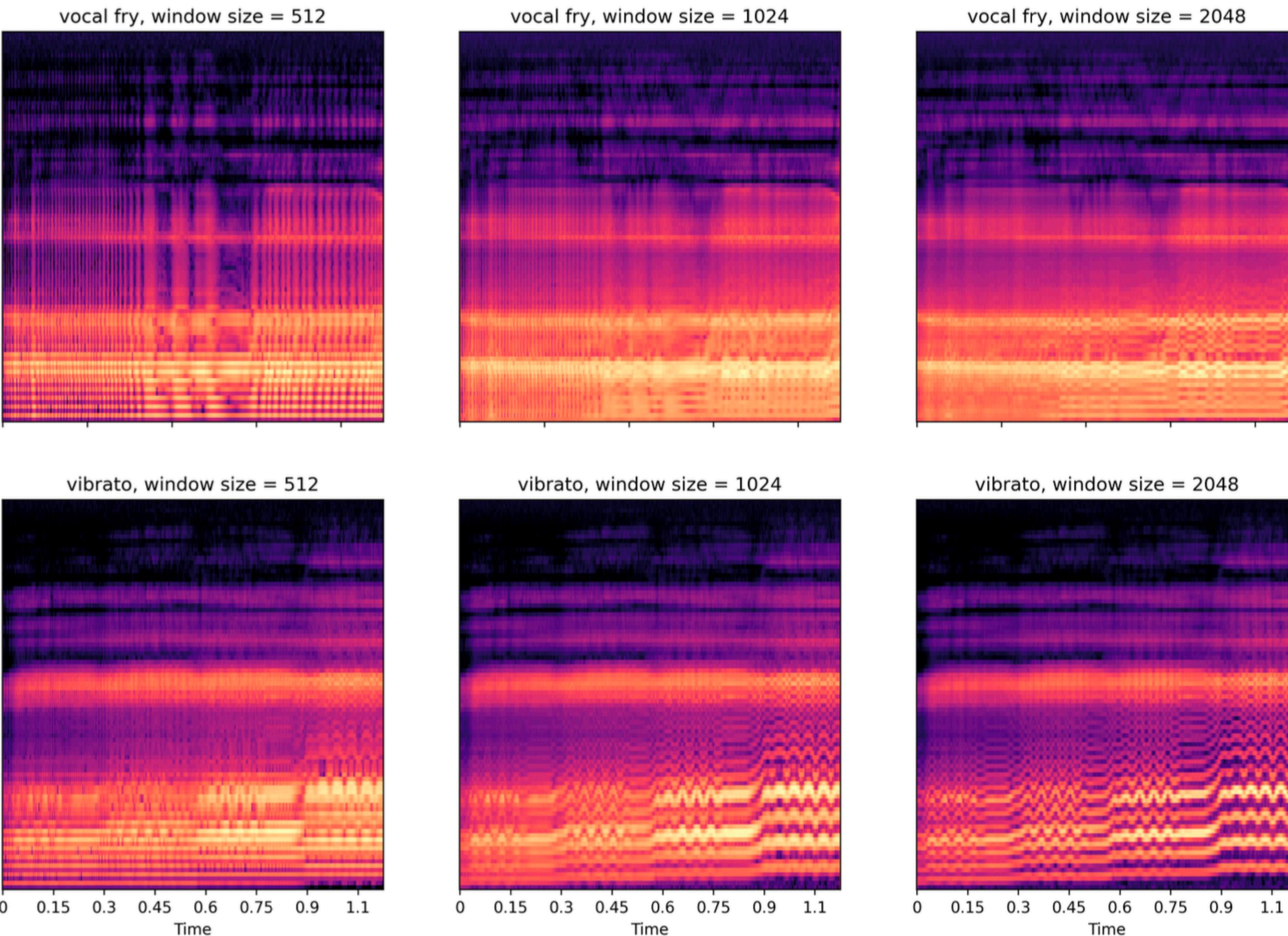
[https://jp.mathworks.com/
help/signal/ref/stft.html](https://jp.mathworks.com/help/signal/ref/stft.html)

Resolution trade-off

101

Spectrograms

Vertical: Frequency
Horizontal: Time



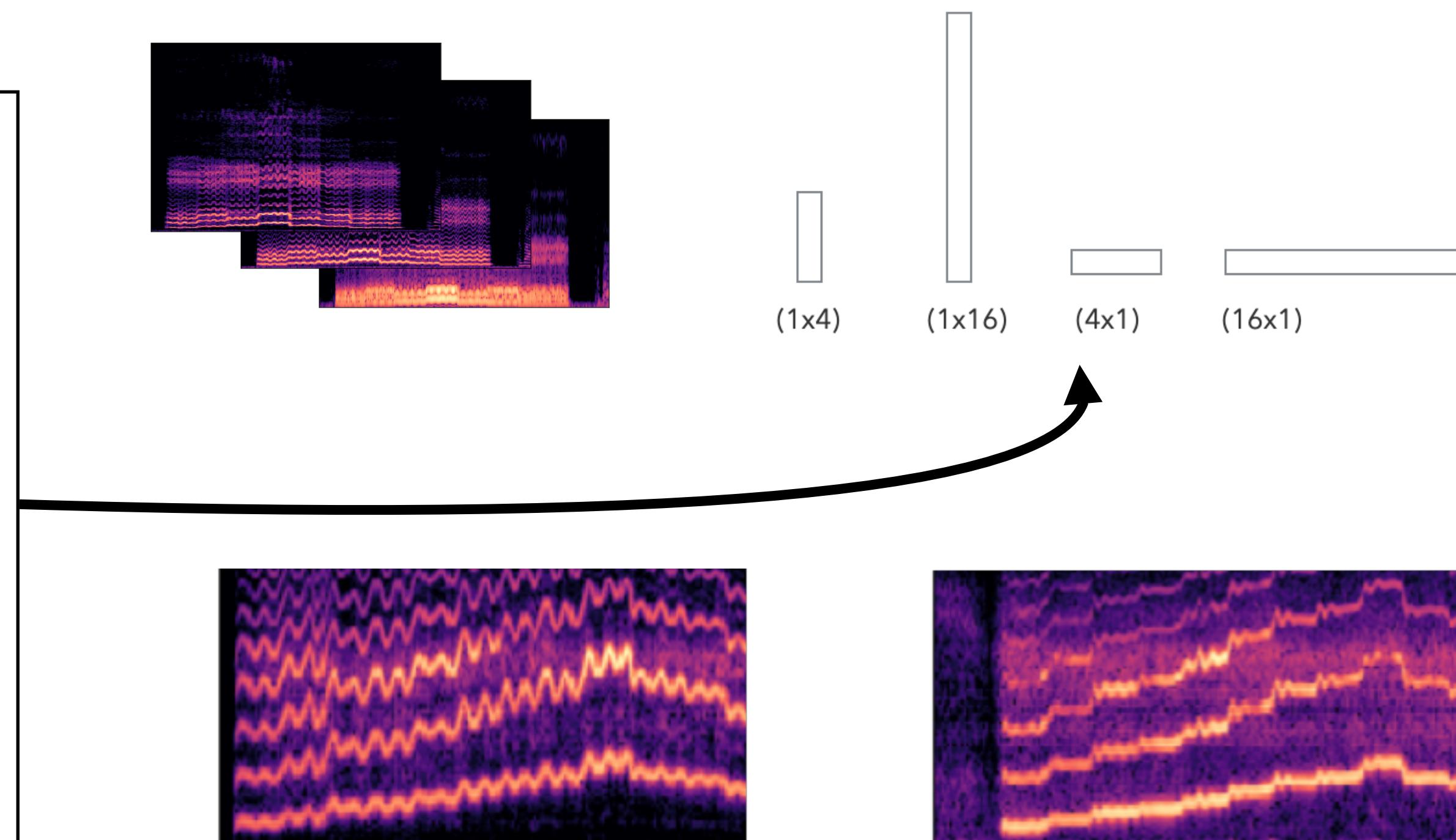
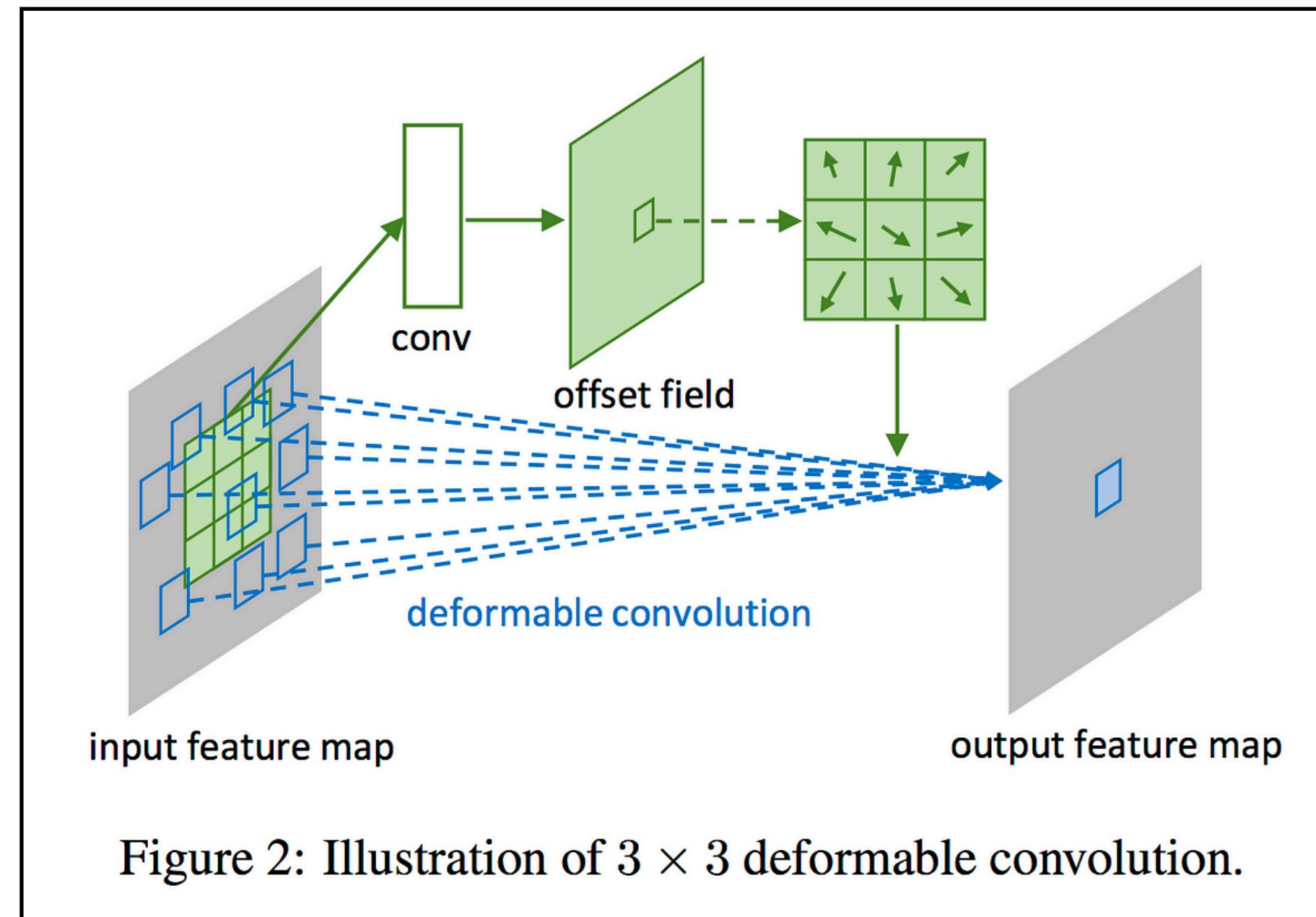
Length of
window

512

1024

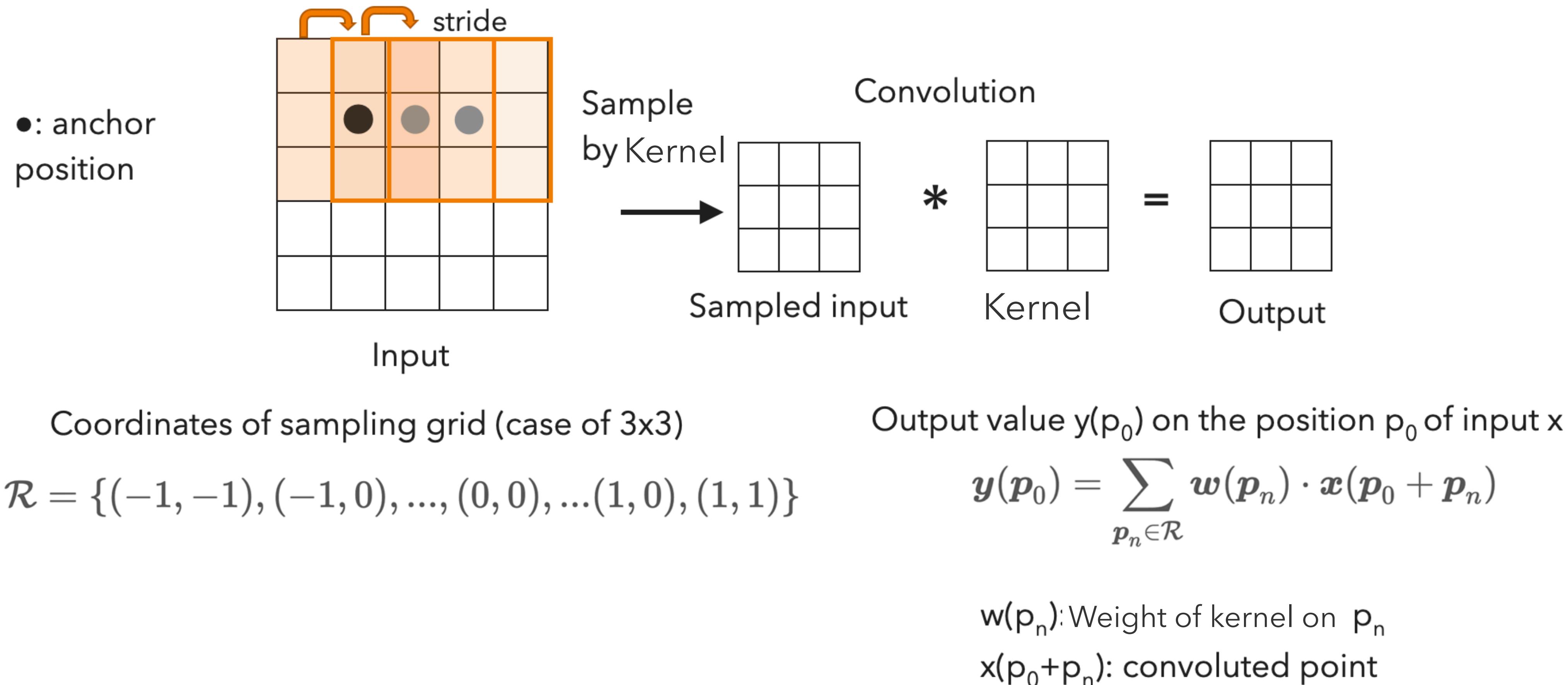
2048

Flexible kernel, that adapts the geometric meaningful pattern



Focused on singing techniques exhibit
geometrical patterns on spectrograms

Normal convolution



About deformable convolution (1/2)

104

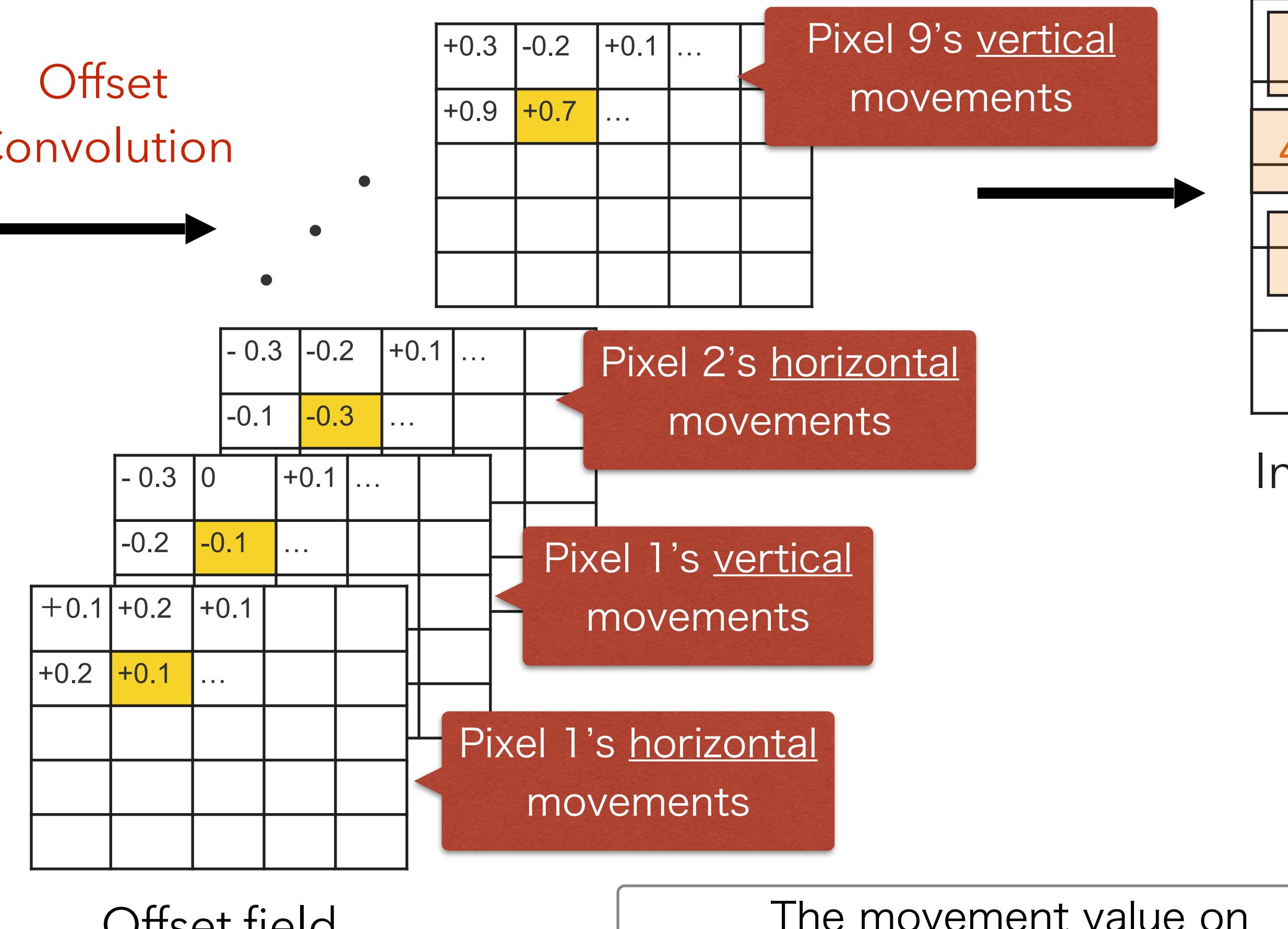
Apply “Offset convolution” to calculate where to move on each pixel

●: anchor position

1	2	3			
4	5	6			
7	8	9			

Input
(N, H, W)

Offset
Convolution



1	2	3			
4	5	6			
7	8	9			

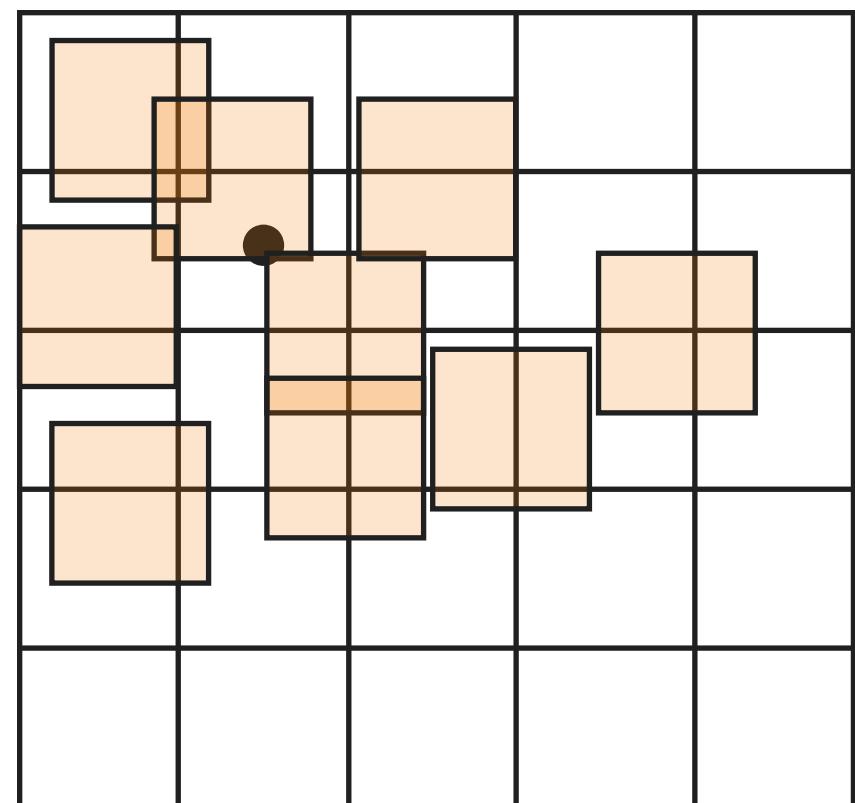
Input (with offset)

About deformable convolution (2/2)

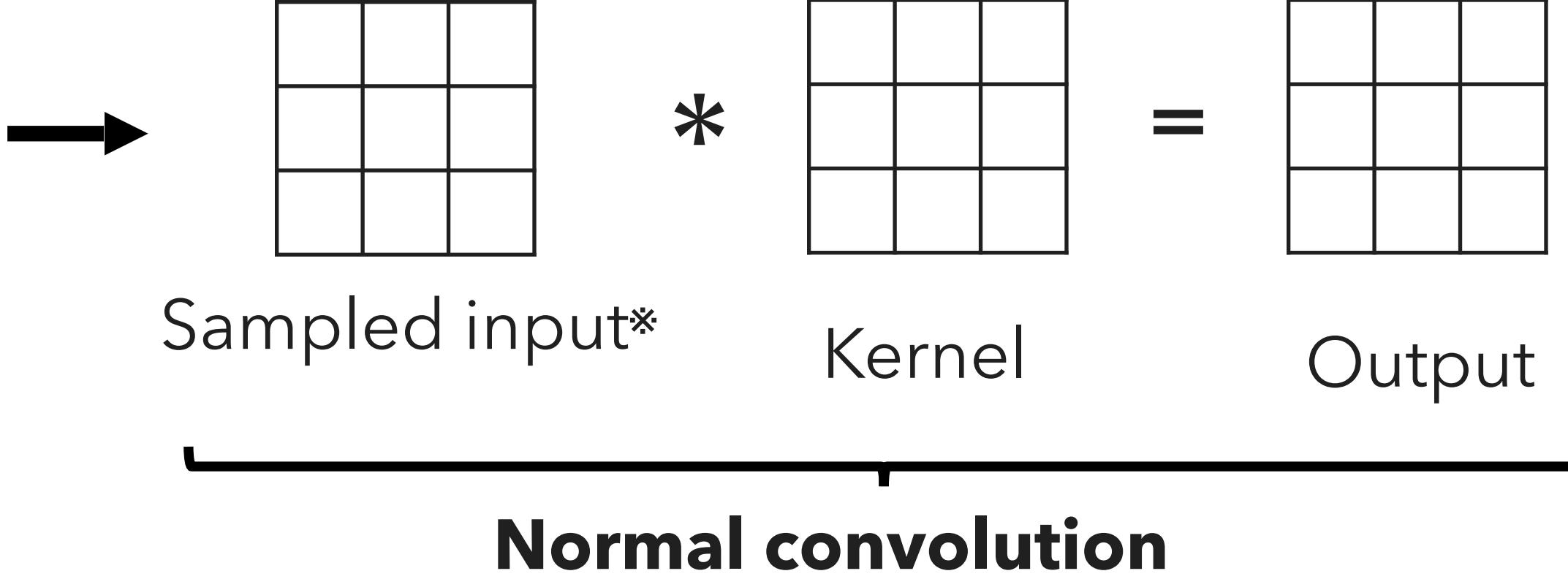
105

Convolutional operation is applied on input with offsets

●: anchor position



Input (with offset)



* Actual sampled values are interpolated by **bilinear interpolation** on nearest grid points

Coordinates of sampling grid (case of 3x3)

$$\mathcal{R} = \{(-1, -1), (-1, 0), \dots, (0, 0), \dots, (1, 0), (1, 1)\}$$

Output value $y(p_0)$ on the position p_0 of input x

$$y(p_0) = \sum_{p_n \in \mathcal{R}} w(p_n) \cdot x(p_0 + p_n + \Delta p_n)$$

$w(p_n)$: Weight of kernel on p_n

$x(p_0 + p_n + \Delta p_n)$: convoluted point

Δp_n : offset of each point

Evaluation metrics on classification

106

- Accuracy
- Balanced-accuracy
- F1-score

$$\text{Accuracy} = \frac{TP + TN}{TP + FN + TN + FP}$$

$$\text{Balanced Accuracy} = \frac{1}{2} \left(\frac{TP}{TP + FN} + \frac{TN}{TN + FP} \right)$$

$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \left(\text{Precision} = \frac{TP}{TP + FP}, \text{Recall} = \frac{TP}{TP + FN} \right)$$

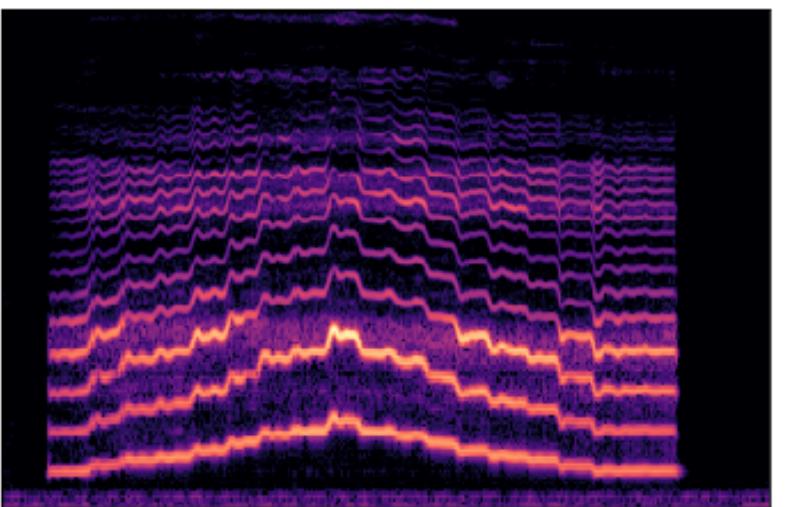
Examples of vocalset

107

- **Many fluctuation components**

- Pitch - vibrato, lip trill etc.
- Timbre - breathy, vocal fry etc.
- Other - inhaled, spoken etc.

-> Need a well-discriminative audio feature

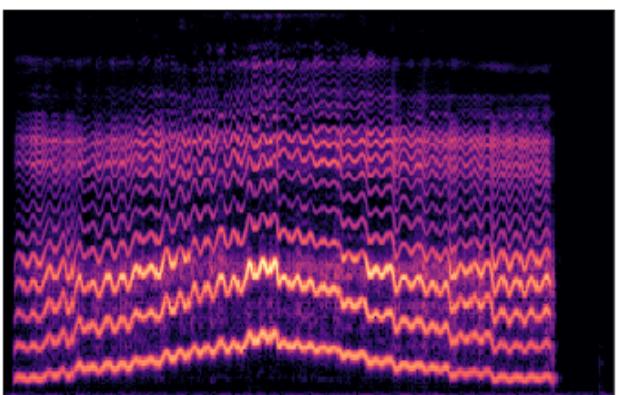


Straight
(non-technique)

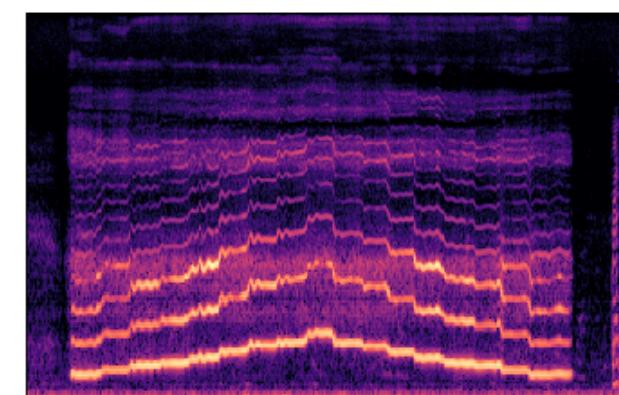
Pitch

Timbre

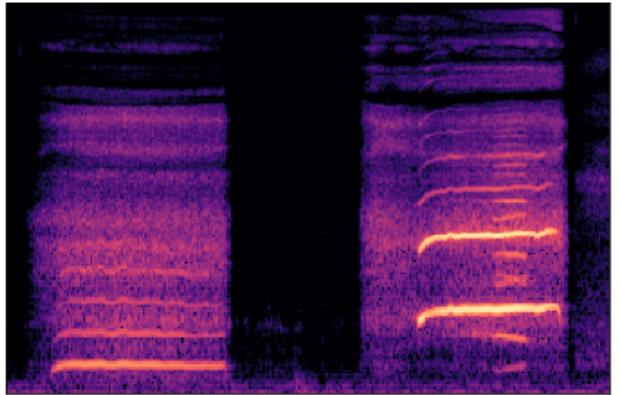
Other



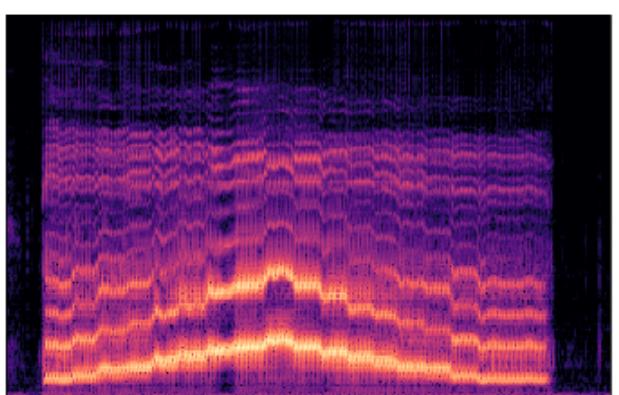
Vibrato



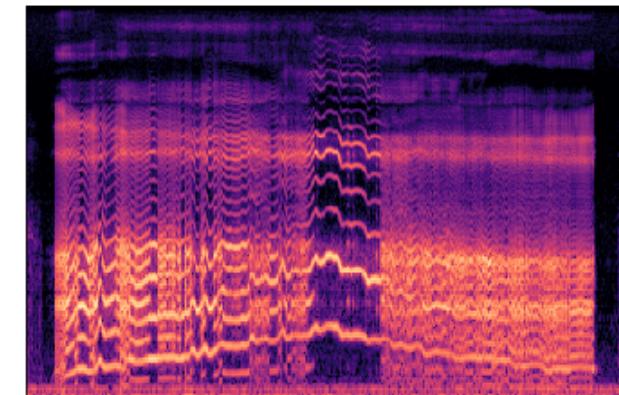
Breathy



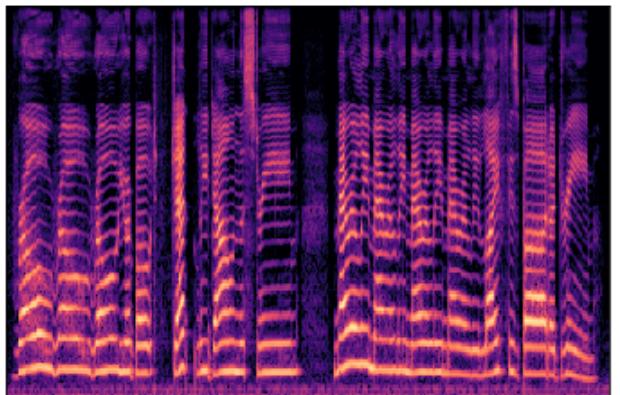
Inhaled



Lip trill



Vocal fry



Spoken

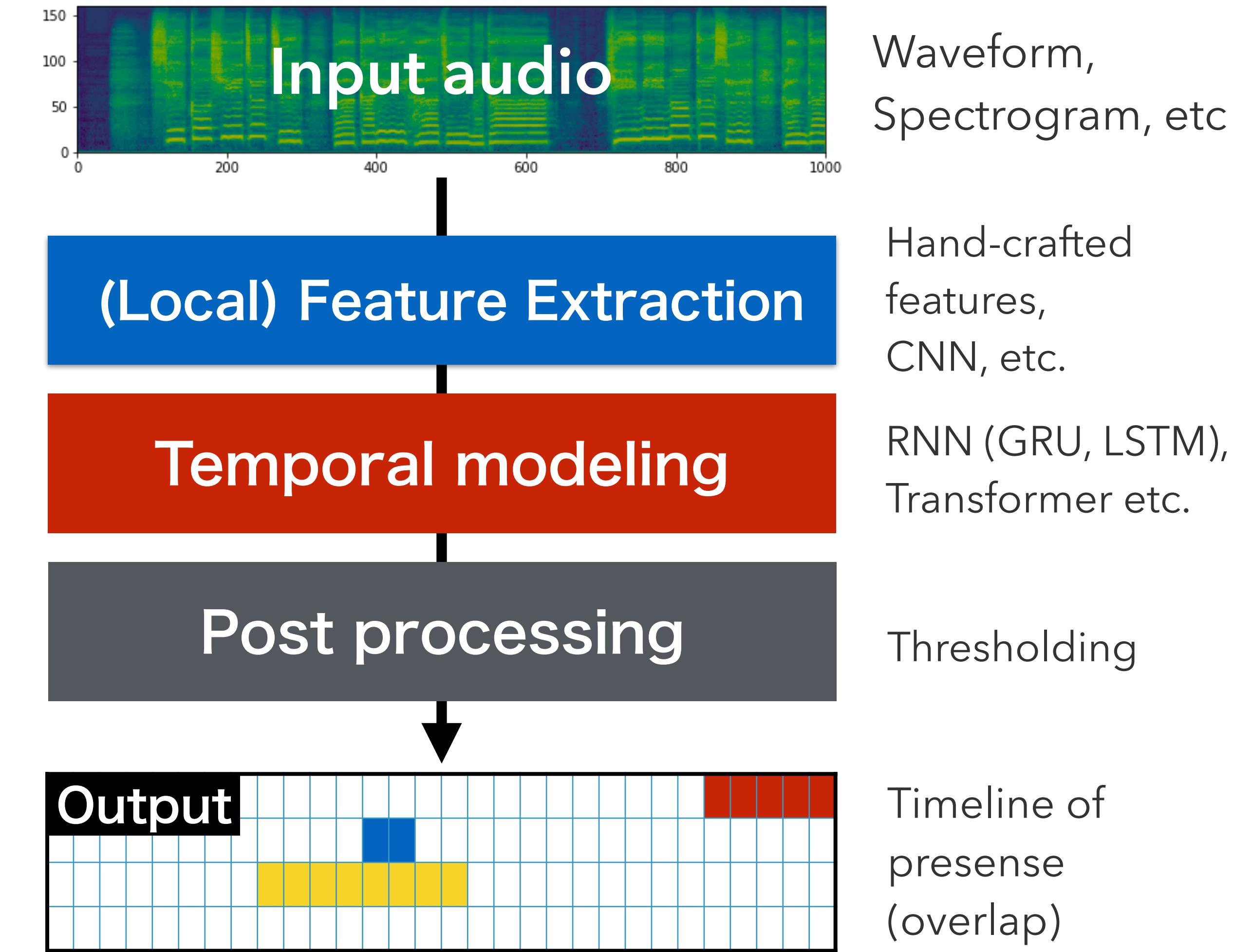
Detection models for singing technique detection

108

- CRNN model is common architecture
 - Especially used in audio tasks [Cakir 17, Imoto21, etc.]
 - CNN captures timbre characteristics
 - RNN captures their temporal fluctuation
- Treat as binary classification of each time, each technique
 - Binary cross entropy loss (BCE) is used

$$BCE = -\log(p_t)$$

$$p_t = \begin{cases} p & (y = 1) \\ 1 - p & (\text{otherwise}) \end{cases}$$



- Widely used in speech emotion recognition
- Feature sets of hand-crafted features
- z-standardization is applied

Frequency related parameters:

- **Pitch**, logarithmic F_0 on a semitone frequency scale, starting at 27.5 Hz (semitone 0).
- **Jitter**, deviations in individual consecutive F_0 period lengths.
- **Formant 1, 2, and 3 frequency**, centre frequency of first, second, and third formant
- **Formant 1**, bandwidth of first formant.

Energy/Amplitude related parameters:

- **Shimmer**, difference of the peak amplitudes of consecutive F_0 periods.
- **Loudness**, estimate of perceived signal intensity from an auditory spectrum.
- **Harmonics-to-Noise Ratio (HNR)**, relation of energy in harmonic components to energy in noise-like components.

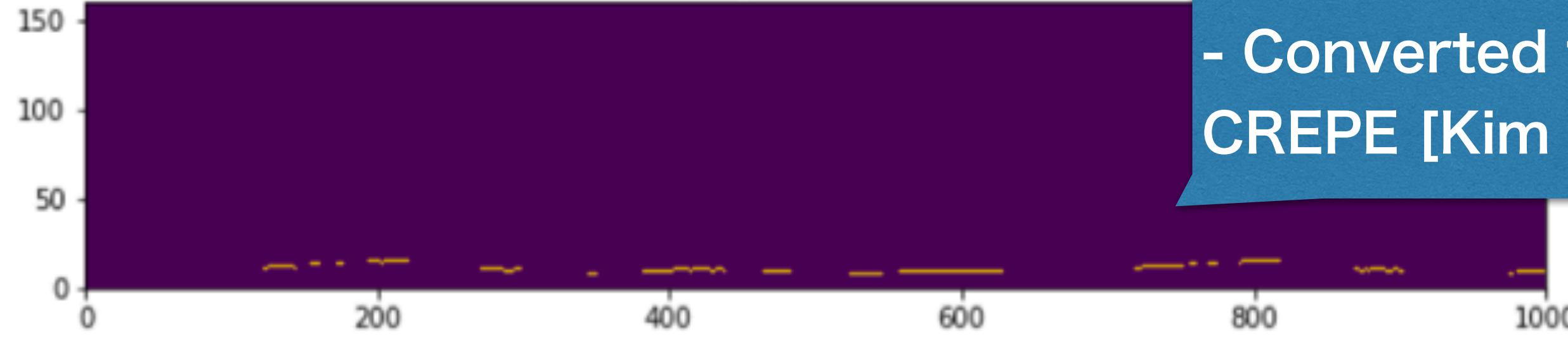
Spectral (balance) parameters:

- **Alpha Ratio**, ratio of the summed energy from 50–1000 Hz and 1–5 kHz
- **Hammarberg Index**, ratio of the strongest energy peak in the 0–2 kHz region to the strongest peak in the 2–5 kHz region.
- **Spectral Slope 0–500 Hz and 500–1500 Hz**, linear regression slope of the logarithmic power spectrum within the two given bands.
- **Formant 1, 2, and 3 relative energy**, as well as the ratio of the energy of the spectral harmonic peak at the first, second, third formant's centre frequency to the energy of the spectral peak at F_0 .
- **Harmonic difference H1–H2**, ratio of energy of the first F_0 harmonic (H1) to the energy of the second F_0 harmonic (H2).
- **Harmonic difference H1–A3**, ratio of energy of the first F_0 harmonic (H1) to the energy of the highest harmonic in the third formant range (A3).

PrimeDNN': input representation

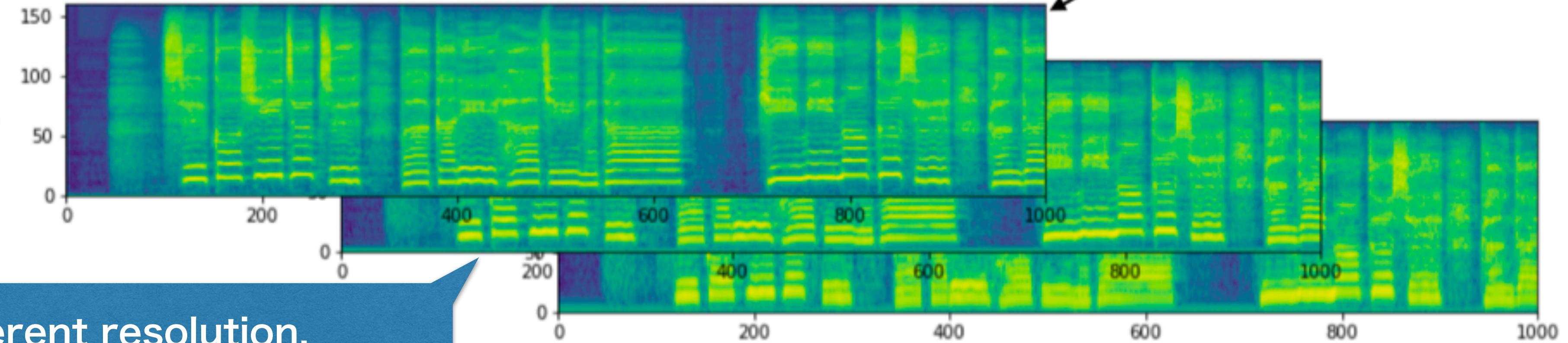
110

Mel-band Pitchgram



- Binary representation,
Indicates where the pitch height is
- Converted from estimation from
CREPE [Kim 18]

MMelSpecs

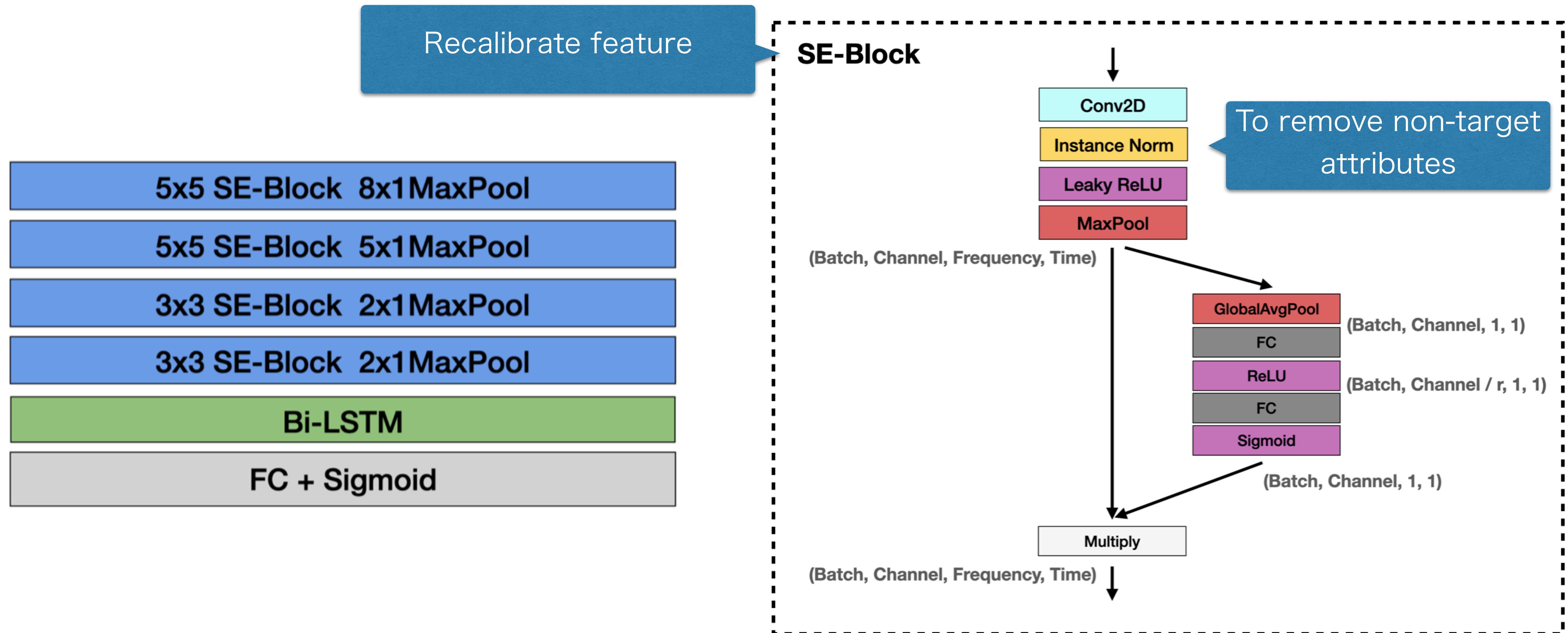


Different resolution,
by altering FFT length
(Shown effective in Chapter 5)

Channel-axis stack

PrimeDNN': architecture

111

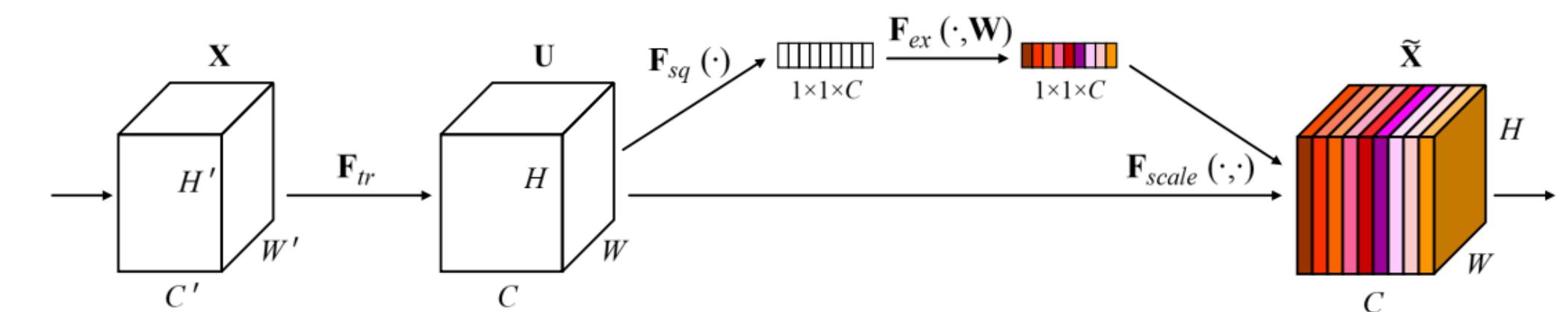


Squeeze and Excitation Network (SE-Net)

112

- Apply weights on output feature maps

- enhancing the important features
- Re-calibrates the feature maps



Conv.block

B, C, F, T

Global AvePool

B, C, 1, 1

Linear

B, C/r, 1, 1

ReLU

B, C/r, 1, 1

Linear

B, C, 1, 1

Sigmoid

multiply

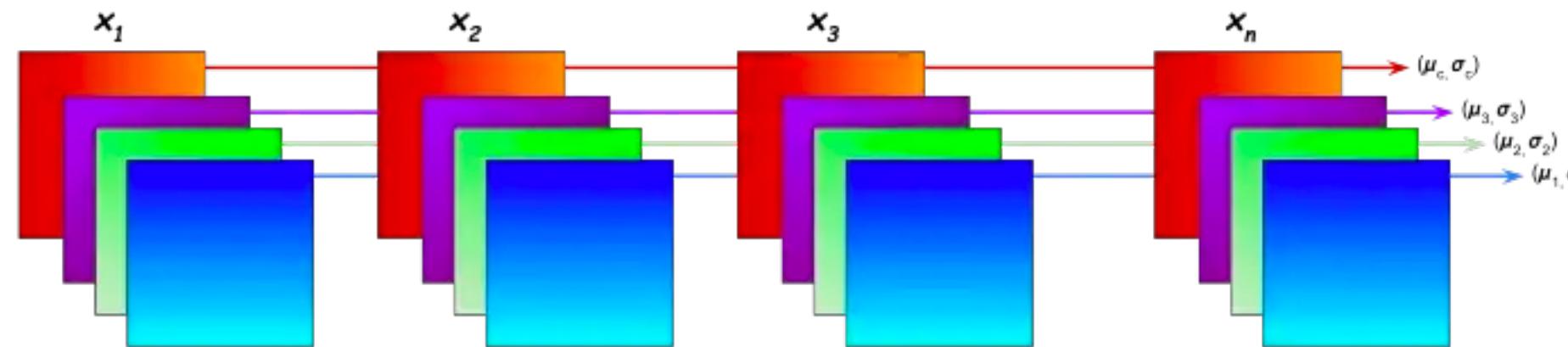
r: reduction ratio

(Empirically set r=2)

Instance normalization [Ulyanov 16]

113

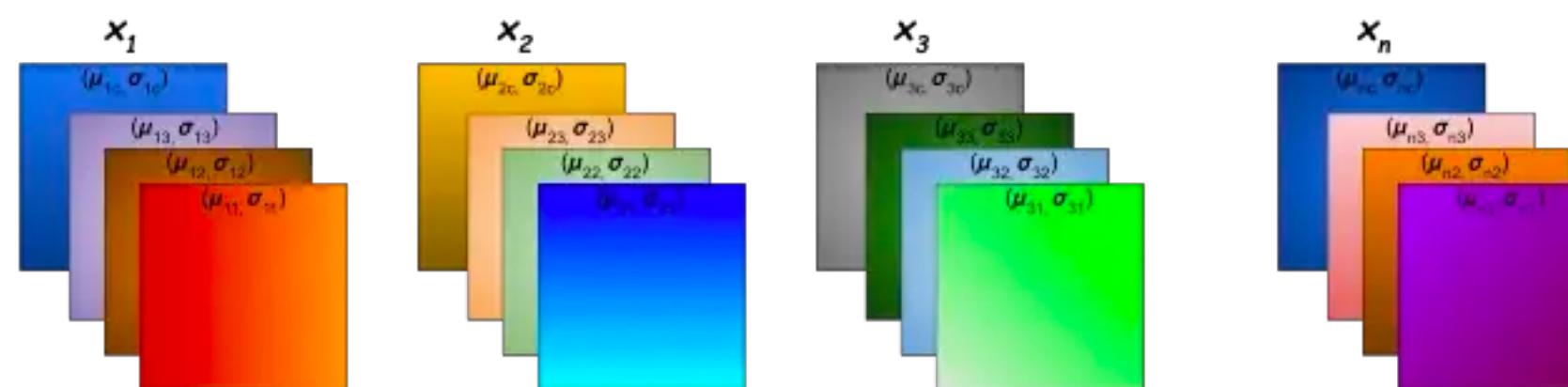
Batch normalization



$$y_{tijk} = \frac{x_{tijk} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}, \quad \mu_i = \frac{1}{HWT} \sum_{t=1}^T \sum_{l=1}^W \sum_{m=1}^H x_{tilm}, \quad \sigma_i^2 = \frac{1}{HWT} \sum_{t=1}^T \sum_{l=1}^W \sum_{m=1}^H (x_{tilm} - \mu_i)^2.$$

The problem of batch normalization:
The variability increases significantly per mini-batch,
leading to instability in the learning process

Instance normalization

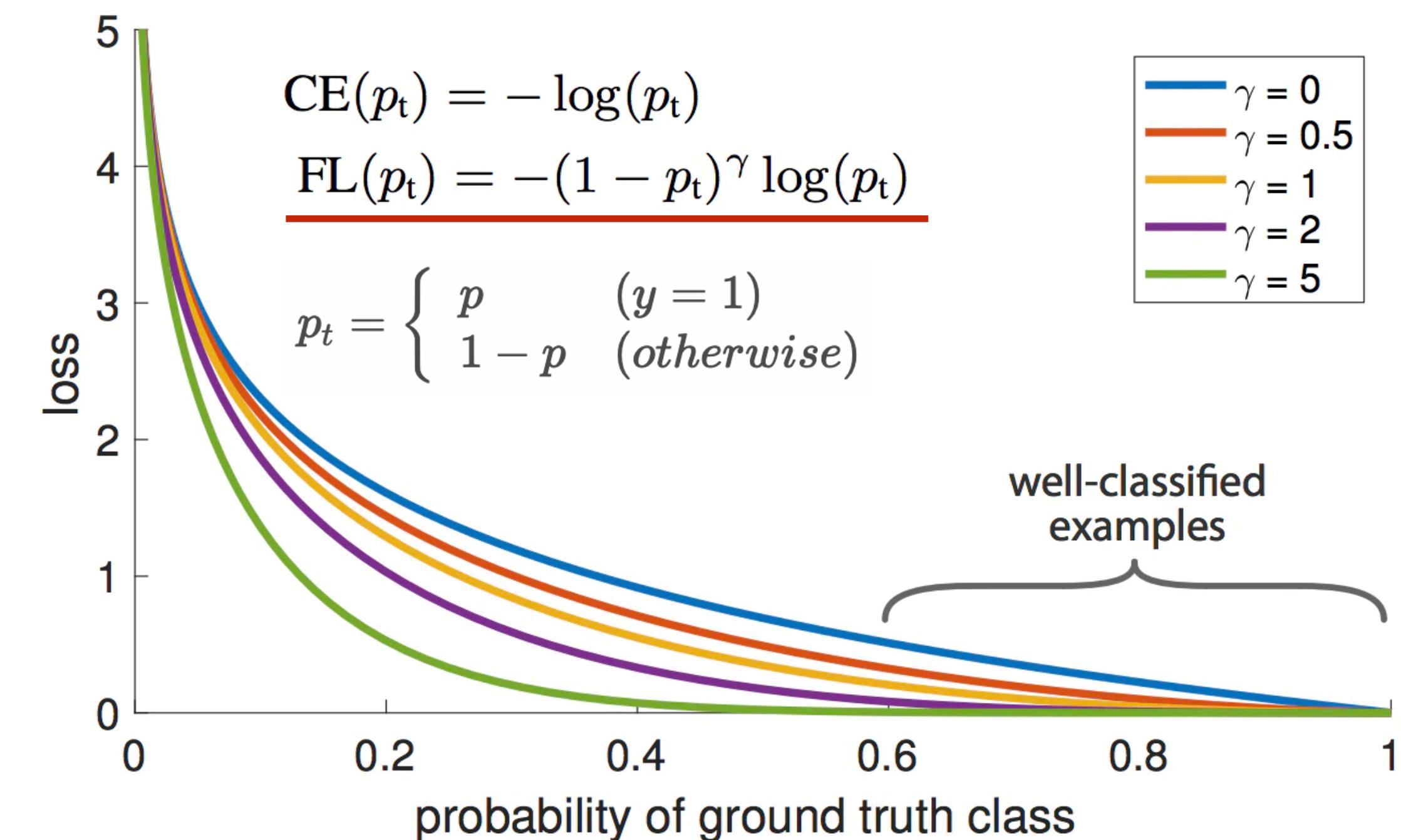
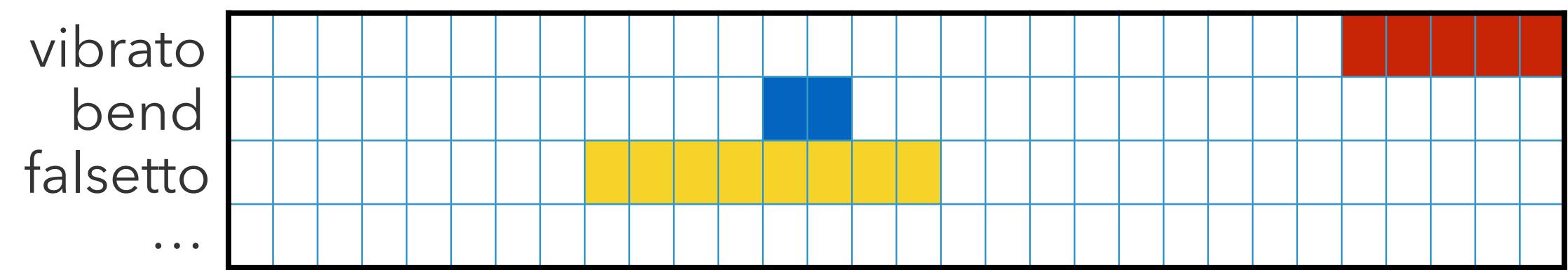


$$y_{tijk} = \frac{x_{tijk} - \mu_{ti}}{\sqrt{\sigma_{ti}^2 + \epsilon}}, \quad \mu_{ti} = \frac{1}{HW} \sum_{l=1}^W \sum_{m=1}^H x_{tilm}, \quad \sigma_{ti}^2 = \frac{1}{HW} \sum_{l=1}^W \sum_{m=1}^H (x_{tilm} - \mu_{ti})^2.$$

Instance normalization: contrast-invariant,
Can work at less batch size

Used Focal loss, considering label sparseness

- The presence of singing techniques is **sparse**-> most segment are non-techniques
- Focal loss down weights the importance of easy-samples (= non-technique segments)

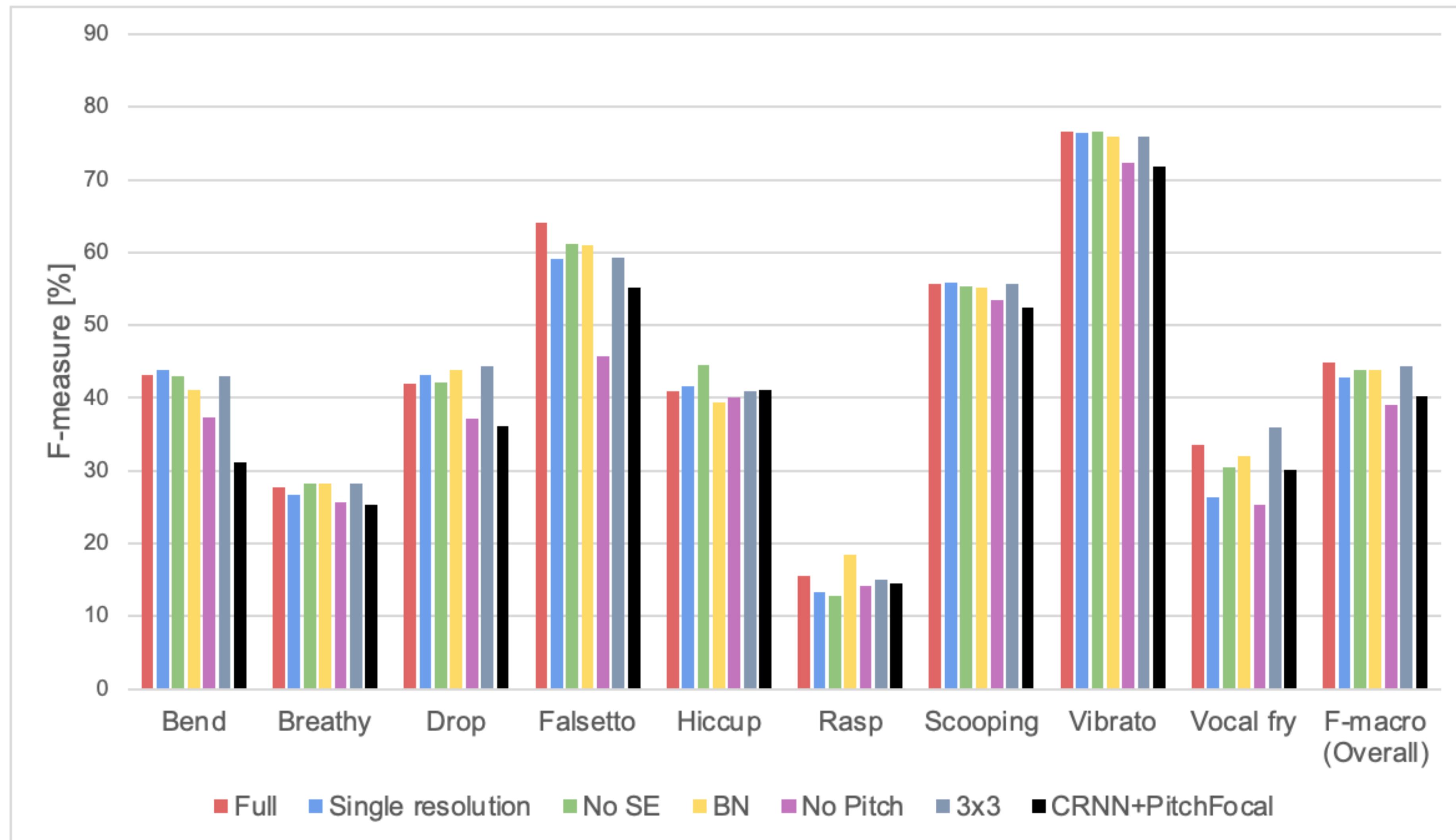


Applying “weight value” to prioritize minor samples

- Chapter 5: Inverse frequency of each class
 - The value become large on minority classes
- Chapter 6: Exponentiation of $(1-p)$ (probability of non-technique frame)
 - The training focusing on hard positive samples (i.e., low probability value on estimation)

Results: ablation study

116



- Easy samples
 - Pitch techniques -> relatively easy to capture, exhibiting geometric pattern on spectrogram
 - Non effect vocal -> vocal effect affects the quality of detection
- Hard samples
 - Timbre techniques -> observations of signal are different each other
 - Effect ed vocal -> make the detection hard due to blurring techniques
 - Short techniques -> less frame samples than long ones