



# On Feature Combination for Music Classification

Zhouyu Fu | Guojun Lu | Kai Ming Ting | Dengsheng Zhang

Gippsland School of Information Technology, Monash University, Australia

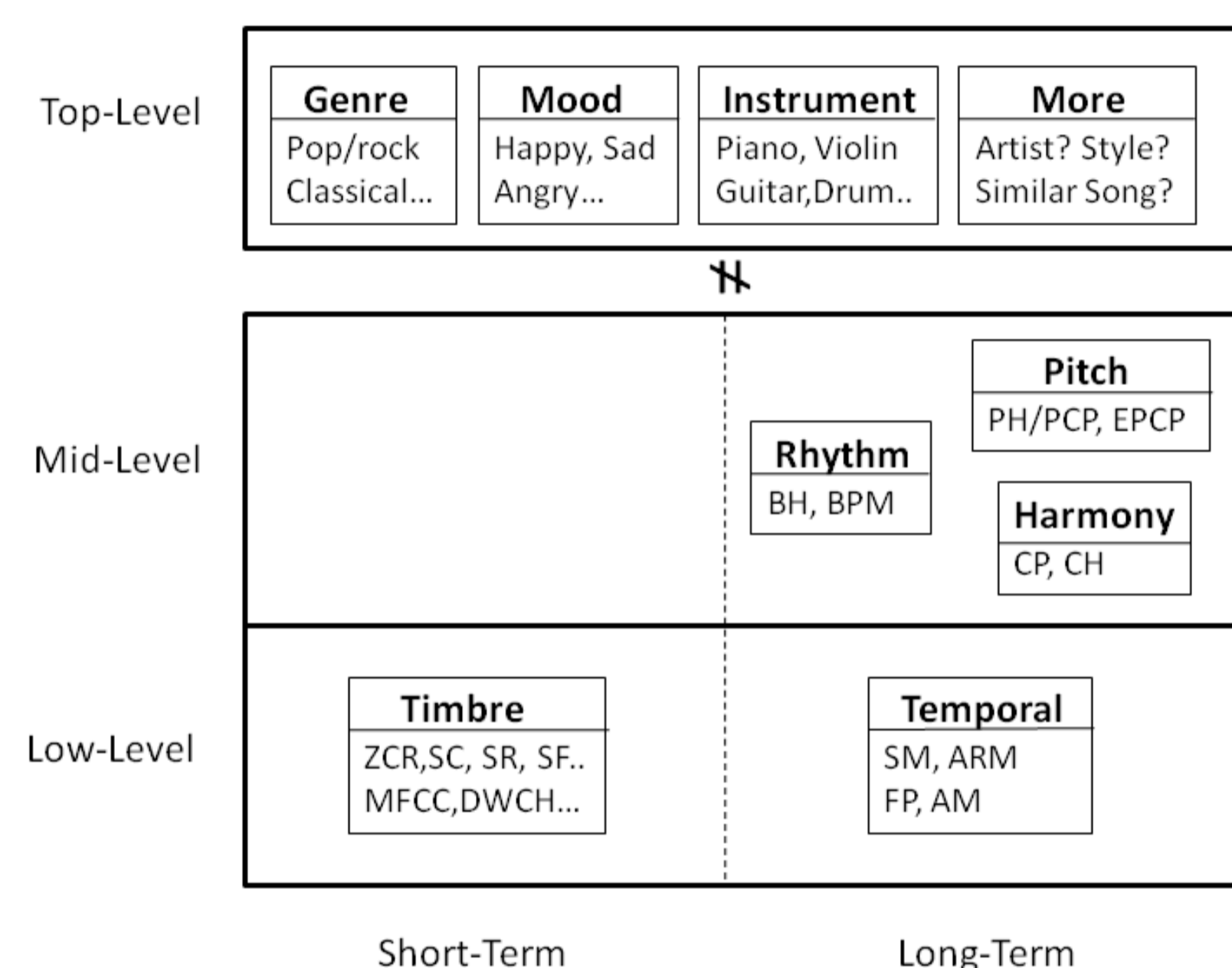
## Problem

**TASK:** combining multiple types of features for music classification from raw audio signals

- **Whether-** we need multiple features?
- **How-** to combine different features?
- **What-** is the best feature and combination scheme?

## Features

Taxonomy of audio features



- Timbre features capture the quality of the sound and has much to do with the instrumentation of the music.
- Temporal features capture the long-term variation of timbre and spectral features over time.
- Mid-level features are extracted on top of low-level features and more interpretable to human listeners.

## Methods

### Overview of Feature Combination

#### Problem Definition

- Input - data set  $\{(\mathbf{x}_i^1, \dots, \mathbf{x}_i^M), y_i\}_{i=1, \dots, N}$ , with  $\mathbf{x}_i^m \in \mathbb{R}^{d_m}$  and  $y_i \in \{1, \dots, K\}$
- Output - classification rule  $f: \mathcal{X} \rightarrow \{1, \dots, K\}$  with  $\mathcal{X} \subset \mathbb{R}^{d_1} \times \dots \times \mathbb{R}^{d_m}$

#### Overview of Methods

- **Decision level** - classifiers trained on individual features and fusion rules applied to the output of individual classifiers, e.g. **majority voting**, **sum rule**, **stacked generalization**, etc.
- **Feature level** - composite feature vector/similarity constructed from individual features, e.g. **feature concatenation**, **multiple kernel learning**, etc.

### Decision Level Fusion

#### Principle

$$\{\mathbf{x}^1, \dots, \mathbf{x}^M\} \xrightarrow{(1)} \{\mathbf{f}_1(\mathbf{x}^1), \dots, \mathbf{f}_M(\mathbf{x}^M)\} \xrightarrow{(2)} f: \mathbf{f}_1 \times \dots \times \mathbf{f}_M \rightarrow \{1, \dots, K\}$$

- (1) Train a separate classifier  $\mathbf{f}_m$  for each individual feature type  $\mathbf{x}^m$
- (2) Combine the decisions returned by individual classifiers

**Assumption:** **decision scores** are returned by each individual classifier  $\mathbf{f}_m = [f_m^1, \dots, f_m^K]$ , where  $f_m^k$  is the score for  $k$ th class returned by the classifier trained on  $m$ th feature type

#### Methods

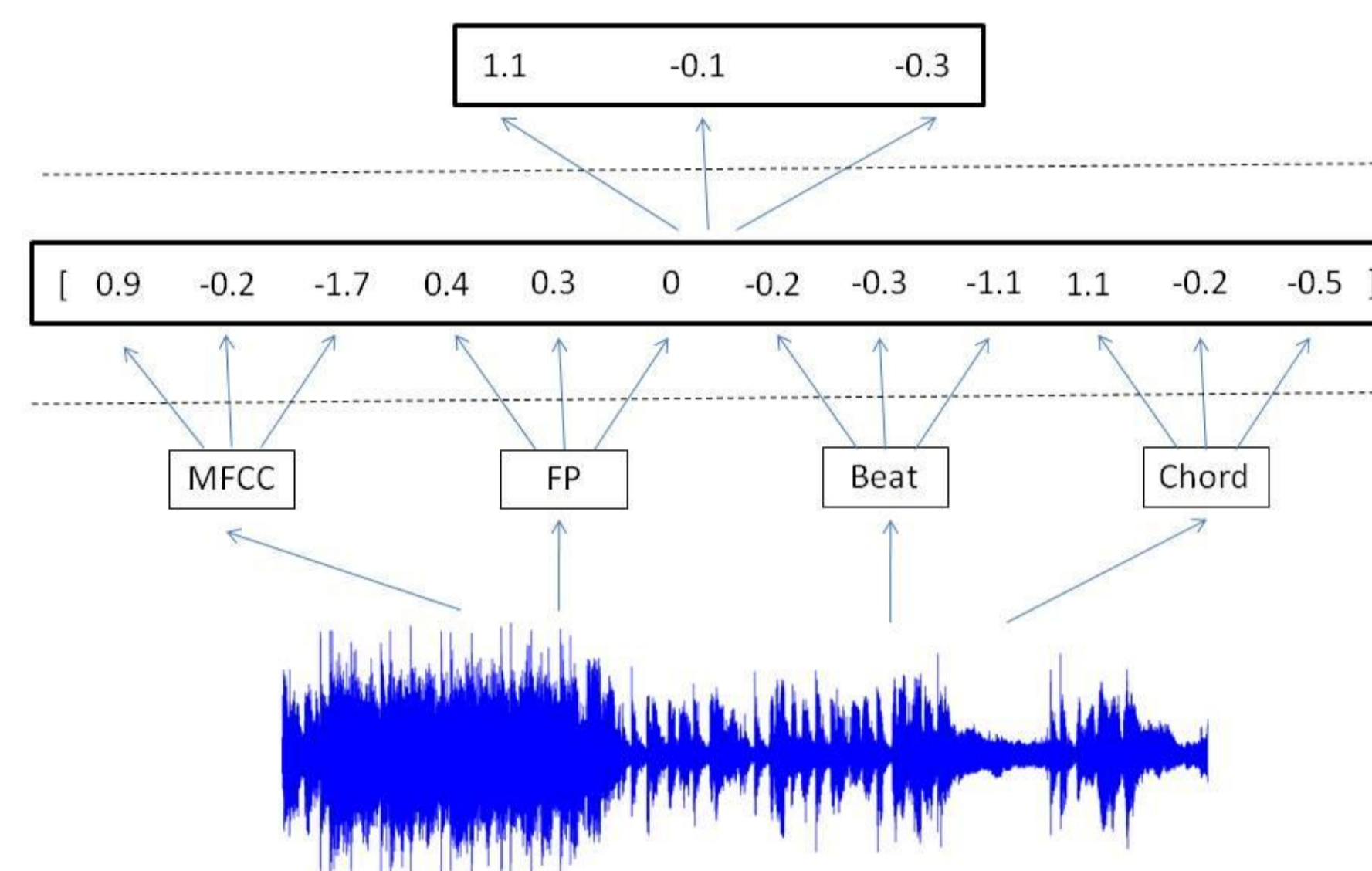
##### • Voting

$$f: \arg \max_{k=1, \dots, K} \sum_{m=1, \dots, M} f_m^k$$

- **Majority Voting**: provides a winner-takes-all voting scheme  
 $f_m^k = 1$  if  $m$ th classifier votes for class  $k$  and  $f_m^k = 0$  otherwise
- **Sum Rule**: provides a weighted voting scheme  
 $f_m^k \in \mathbb{R}$  encodes the confidence value for/against class  $k$  returned by  $m$ th classifier

##### • Stacked Generalization

- Stack the decision values returned by individual classifiers into a score vector
- Train a classifier using the score vectors as new input features
- **Classifier on Classifier**



### Feature-Level Combination

#### Principle

$$\{\mathbf{x}^1, \dots, \mathbf{x}^M\} \xrightarrow{(1)} \psi(\mathbf{x}^1, \dots, \mathbf{x}^M) \xrightarrow{(2)} f(\psi)$$

- (1) A composite feature (vector) is obtained by aggregating the individual features
- (2) A single classifier is trained on the composite feature

Choice of feature-level combination methods is class specific, here we focus on **SVM**

#### Methods

##### • Feature Concatenation

– **Input Space**  $\psi(\mathbf{x}) = [\mathbf{x}^1T, \dots, \mathbf{x}^MT]T$

– **Embedded Space (RKHS)**  $\psi(\mathbf{x}) = \frac{1}{\sqrt{M}}[\phi(\mathbf{x}^1)T, \dots, \phi(\mathbf{x}^M)T]T$

From kernel point of view, this is equivalent to **averaging** individual kernels

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = \frac{1}{M} \sum_{m=1}^M \langle \phi(\mathbf{x}_i^m), \phi(\mathbf{x}_j^m) \rangle = \frac{1}{M} \sum_{m=1}^M K(\mathbf{x}_i^m, \mathbf{x}_j^m)$$

##### • Multiple Kernel Learning (MKL)

- weighted feature concatenation in RKHS  $\psi(\mathbf{x}) = [\sqrt{\beta_1}\phi(\mathbf{x}^1)T, \dots, \sqrt{\beta_m}\phi(\mathbf{x}^M)T]T$
- equivalent to weighted composite kernel  $K_\beta(\mathbf{x}_i, \mathbf{x}_j) = \sum_m \beta_m K_m(\mathbf{x}_i^m, \mathbf{x}_j^m)$
- kernel weights and SVM classifier weights are learned simultaneously

$$\min_{\beta \geq 0} \max_{0 \leq \alpha_i \leq \lambda} \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j K_\beta(\mathbf{x}_i, \mathbf{x}_j)$$

Kernel Classifier

## Results

### Experimental Setup

- 1000 song clips equally distributed in 10 genres
- AU format with 22050Hz sampling rate, roughly 30 seconds in length for each clip
- **Eight** features are extracted from each clip, including 3 timbre features (SMFCC, SASE, SOSOC), 3 temporal features (TMFCC, TASE, TOSC) and 2 mid-level features of beat and chord.

### Performance of Individual Features

	SMFCC	SASE	SOSOC	TMFCC	TASE	TOSC	Beat	Chord
Blues	75.90	64.40	76.90	73.20	72.00	78.80	18.60	83.20
Classical	93.00	91.50	94.80	95.80	92.10	94.30	29.40	90.20
Disco	63.30	56.60	63.00	63.20	69.10	66.20	71.60	54.10
Hiphop	68.90	65.20	72.40	73.80	77.50	74.90	27.10	96.60
Metal	77.10	75.60	71.40	66.10	71.80	74.80	17.90	77.90
...								
10-class	73.55	68.00	73.10	73.81	73.20	75.00	24.66	78.92

- Chord is the best individual feature overall, and beat is the weakest
- No single type of feature performs consistently well for each individual class

### Performance of Feature Combination

	Best	Vote	FC	AvgK	Sum	MKL	SG
Blues	83.20	86.30	89.60	94.20	91.70	93.70	95.70
Classical	95.80	96.80	97.00	97.20	96.60	97.50	97.00
Disco	71.60	77.60	83.00	83.70	86.10	86.30	86.60
Hiphop	96.60	85.90	86.60	91.90	93.40	93.00	93.30
Metal	77.90	78.60	80.60	88.00	87.80	89.70	87.80
...							
10-class	78.92	84.29	84.75	89.08	89.80	90.38	90.85

- Feature combination is effective in enhancing classification performance
- Classes for which all individual features perform weakly benefit more from feature combination
- Learning-based combination methods (SG, MKL) perform better than heuristics-based methods

### Comparison with other methods

References	Accuracy	References	Accuracy
<b>Feature Combination</b>	<b>90.9 ± 1.02%</b>	T. Li <i>et al.</i>	78.5 ± 4.07%
Bergstra <i>et al.</i>	81%	I. Panagakis (2008) <i>et al.</i>	78.2 ± 3.82%
Lee <i>et al.</i>	90.6 ± 3.06%	T. Lidy <i>et al.</i>	74.9%
Panagakis <i>et al.</i>	92.7 ± 2%	G. Tzanetakis <i>et al.</i>	61.0%