

# Polyphonic Piano Note Transcription with Non-negative Matrix Factorization of Differential Spectrogram

Lufei Gao\*, Li Su†, Yi-Hsuan Yang†, Tan Lee\*

\*Department of Electronic Engineering, The Chinese University of Hong Kong

†Academia Sinica, Taiwan

## Highlights

- Note-level music transcription of pitched percussive instruments
- Investigate the idea of highlighting local energy increase in the TF representation;
- Propose algorithms based on existing NMF based methods;
- Validate the advantages of the differential spectrogram.

## Baseline Methods

- Standard NMF (NMF)
  - Approximate the STFT  $X_{ft}$  as the product of two non-negative matrices:

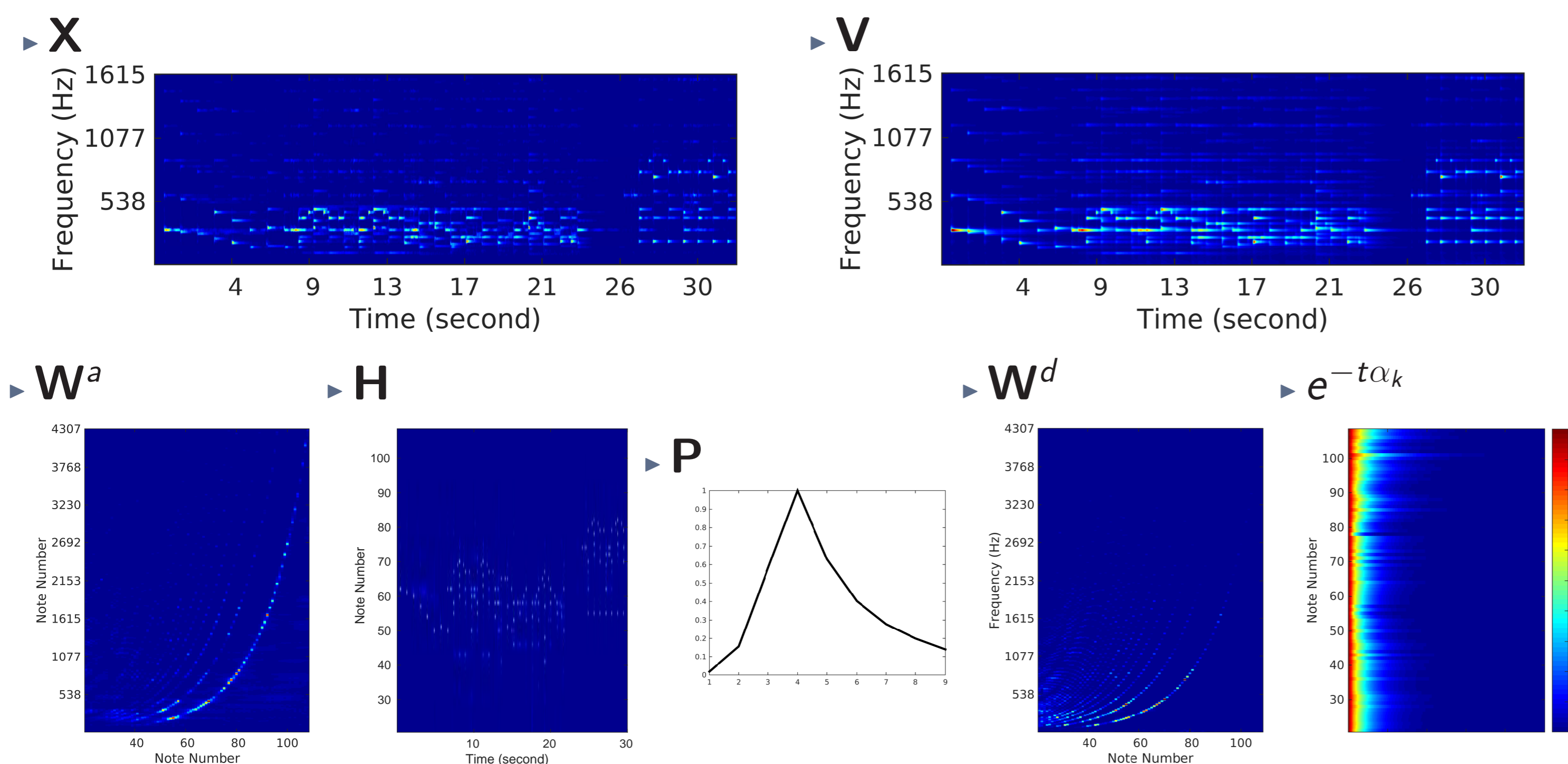
$$X_{ft} \simeq V_{ft} = \sum_{k=1}^K W_{fk} H_{kt}$$

$\mathbf{W}$  is the template of single-note spectra;  $\mathbf{H}$  is the time-varying activation.

- The distortion  $D(\mathbf{X}|\mathbf{V})$  is measured by the  $\beta$ -divergence.
- Attack/Decay Convolutional NMF (CNMF-AD)
  - $X_{ft}$  is assumed to be the summation of the attack phase and the decay phase. The model is defined as

$$V_{ft} = \sum_{k=1}^K W_{fk}^a \sum_{\tau=t-T_t}^{t+T_t} H_{k\tau} P(t-\tau) + \sum_{k=1}^K W_{fk}^d \sum_{\tau=1}^t H_{k\tau} e^{-(t-\tau)\alpha_k}$$

$\mathbf{W}^a$  is the percussive template for the attack phase;  $\mathbf{W}^d$  is the harmonic template for the decay phase;  $\mathbf{P}$  and  $\alpha_k$  are the transient pattern and the exponential decay rate, respectively.

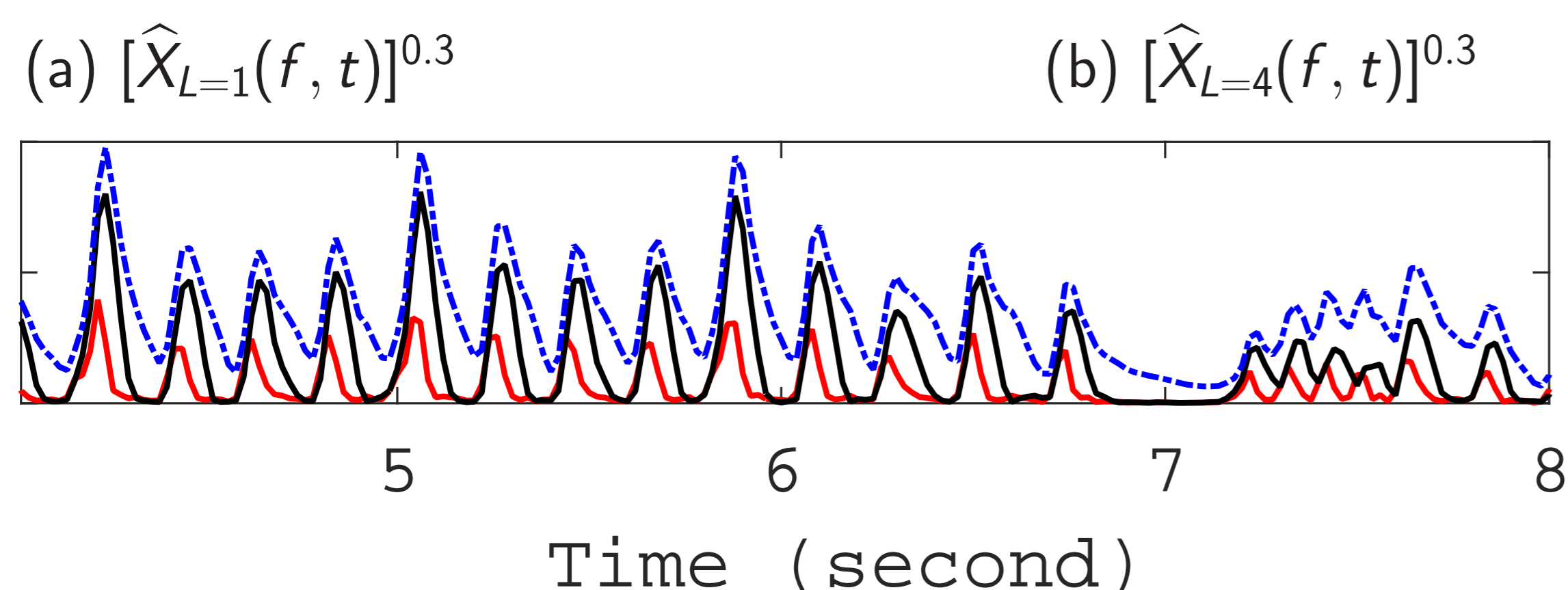
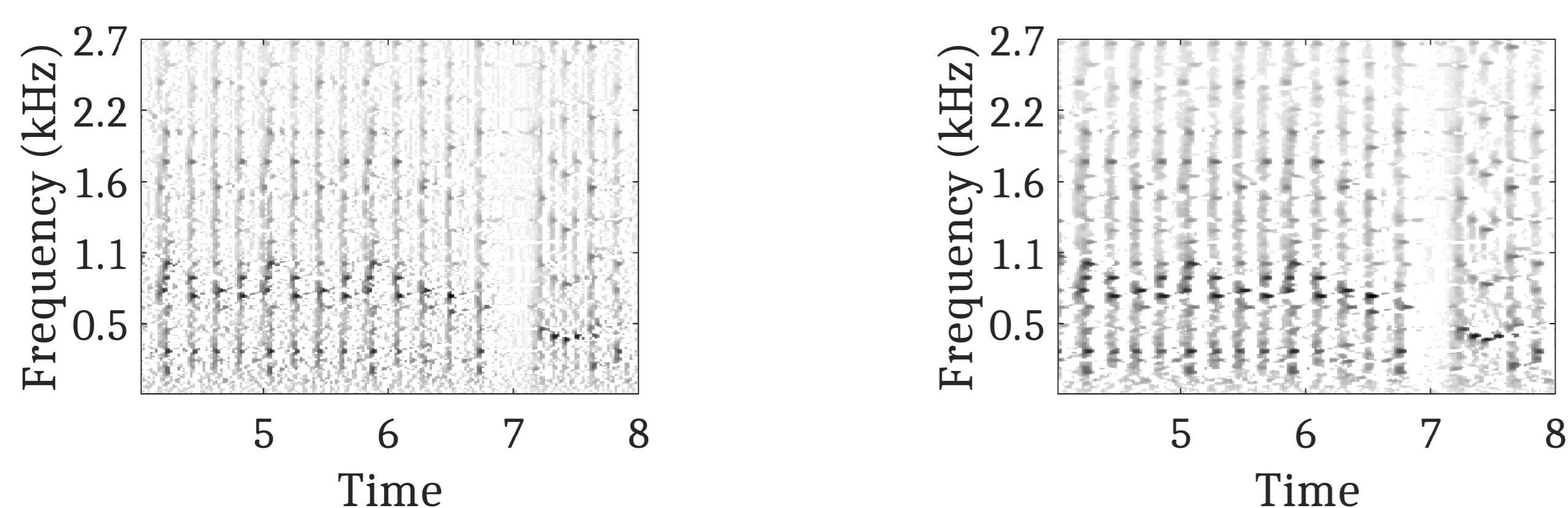


## Differential Spectrogram

- Assuming that the instrument exhibits harmonics with locally stable frequencies, the differential spectrogram  $\hat{X}_L(f, t)$  is defined as:

$$\hat{X}_L(f, t) = \text{HWR}(|X(f, t+L)| - |X(f, t)|)$$

HWR stands for the half-wave rectification ( $\text{HWR}(x) = \frac{x+|x|}{2}$ ).



## Model Adaptation with Differential Spectrogram

- Standard NMF adaptation (NMF- $\Delta$ ): The following feature is used to replace  $\mathbf{X}$  in the model:

$$\hat{X}_L(f, t) = c_1 X(f, t) + c_2 \hat{X}_L(f, t),$$

## Model Adaptation with Differential Spectrogram

- Convolutional NMF adaptation (CNMF- $\Delta$ ): The following model is utilized to estimate the note activation.

$$\hat{X}_L(f, t) \simeq \hat{V}_{ft} = \sum_{k=1}^K \hat{W}_{fk} \sum_{\tau=t-T_t}^{t+T_t} \hat{H}_{k\tau} \hat{P}(t-\tau).$$

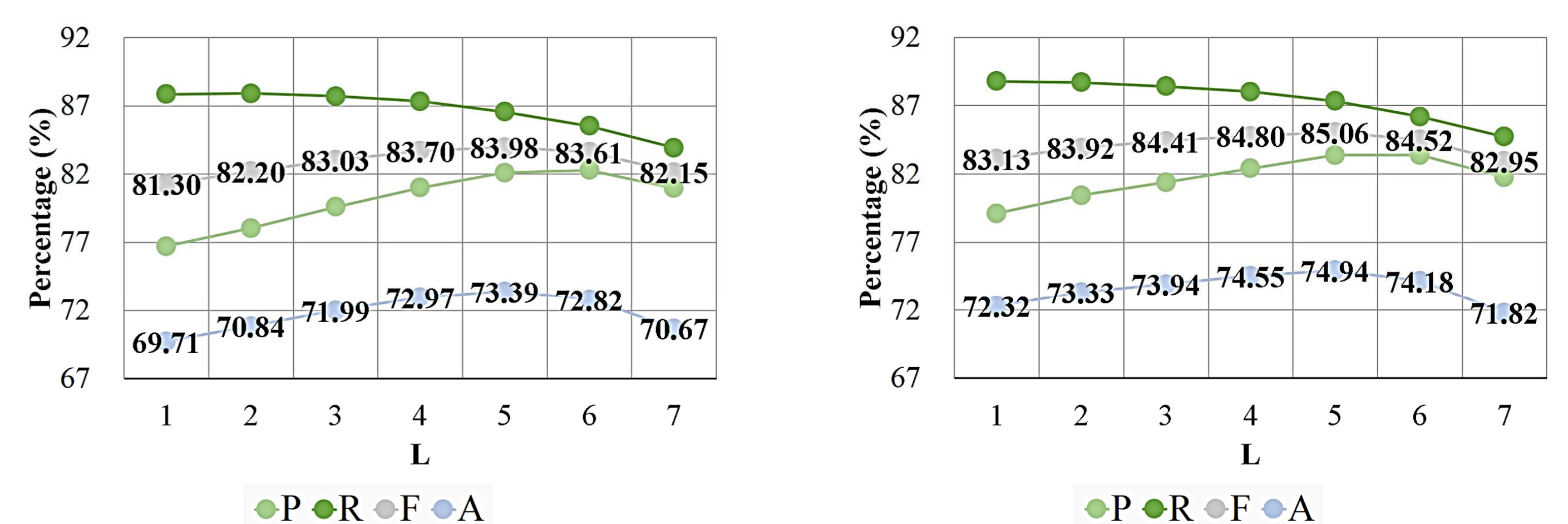
Convolving  $\hat{\mathbf{H}}$  with  $\hat{\mathbf{P}}$  yields the attack activation denoted by  $\hat{\mathbf{H}}^a$ .

- Model initialization (CNMF-AD- $\Delta$ ): Initialize  $\hat{\mathbf{H}}$  by  $\mathbf{H}$  which is estimated using CNMF-AD before the estimation with differential spectrogram.

## Experiment setting

- Databases
  - For dictionary construction
    - The 88 forte isolated note recordings in MAPS-ENSTDkCI.
  - For testing tasks
    - The first 30-second excerpt of the 30 music pieces from MAPS-ENSTDkCI.
- Metric: Precision  $\mathcal{P}$ , Recall  $\mathcal{R}$ , F-measure  $\mathcal{F}$ , Accuracy  $\mathcal{A}$ .

## Results: System settings



(a) CNMF- $\Delta$

(b) CNMF-AD- $\Delta$

## Results: Comparison

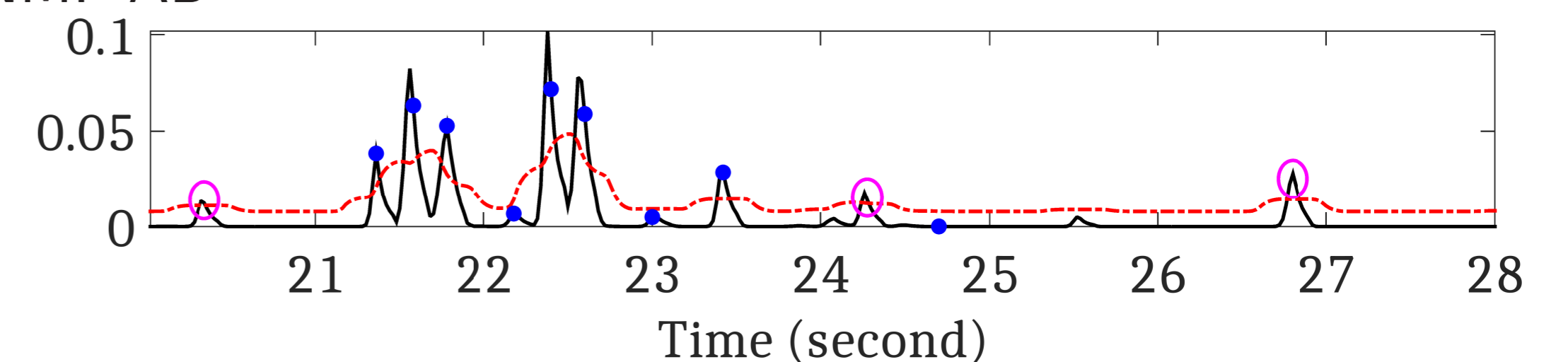
Table: Performance comparison on "ENSTDkCI"

Method	P	R	F	A
NMF ( $\beta = 0.5$ )	59.70	34.51	41.81	27.24
NMF- $\Delta$ ( $\beta = 0.5$ )	71.04	42.48	50.70	35.13
NMF ( $\beta = 2$ )	51.67	43.11	46.34	30.54
NMF- $\Delta$ ( $\beta = 2$ )	67.83	58.21	61.76	45.10
CNMF- $\Delta$	82.11	86.57	83.98	73.39
CNMF-AD- $\Delta$	83.38	<b>87.34</b>	<b>85.06</b>	<b>74.94</b>
CNMF-AD	<b>89.22</b>	78.35	82.91	71.55
Böck	–	–	–	68.70
Berg-Kirkpatrick	78.10	74.70	76.40	–

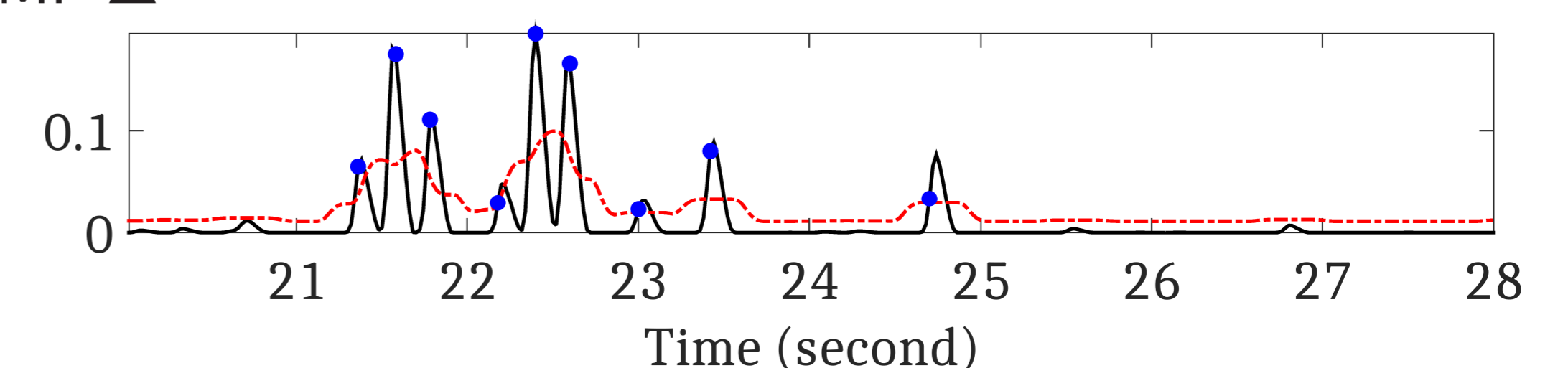
## Result Analysis

- The attack activations of note E4 in one test file.

- $\mathbf{H}^a$  of CNMF-AD



- $\hat{\mathbf{H}}^a$  of CNMF- $\Delta$



- The piano rolls of attack activations of our methods

