

Large Scale Polyphonic Music Transcription Using Randomized Matrix Decompositions

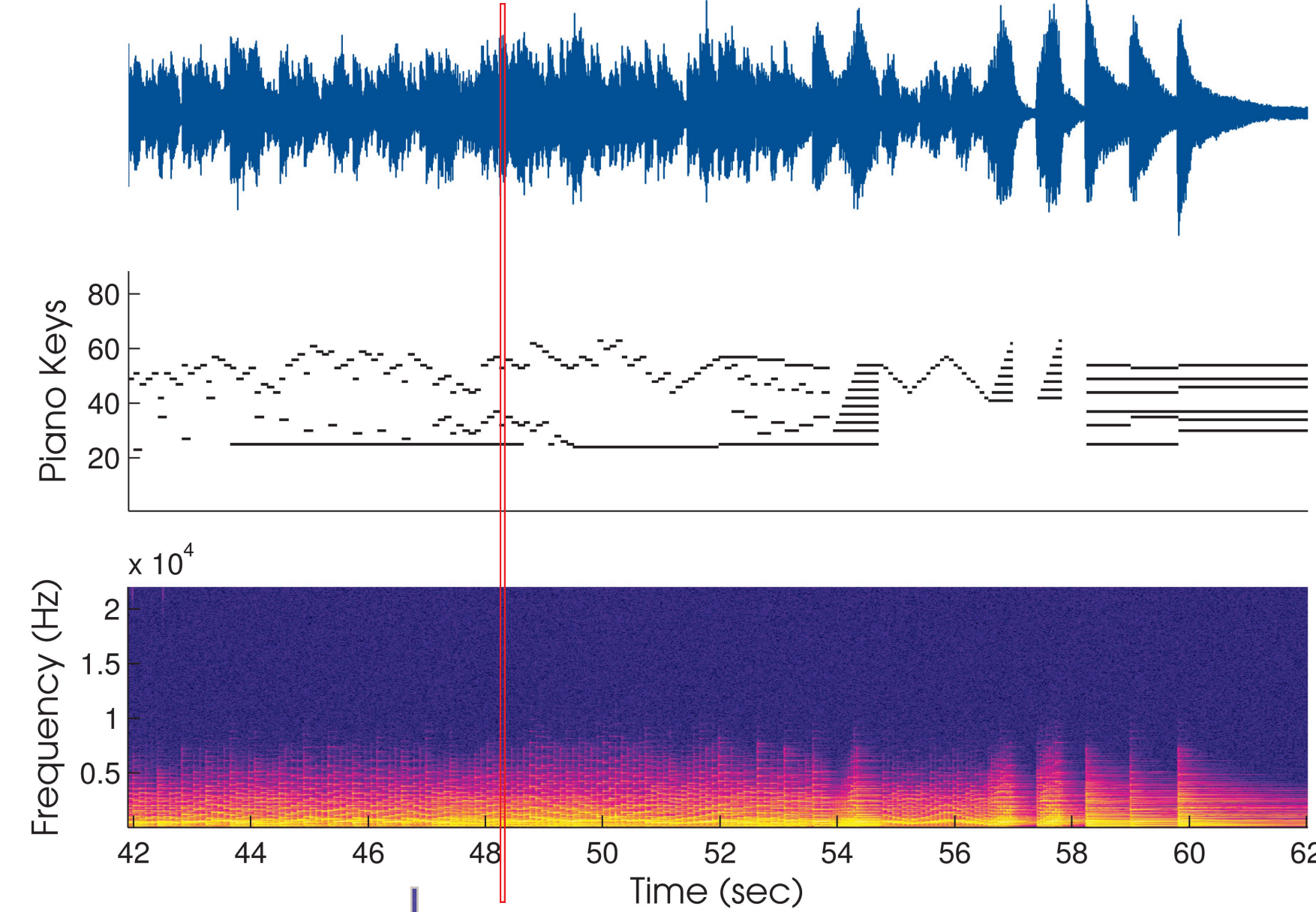
Ismail Arı, Umut Şimşekli,
Ali Taylan Cemgil, Lale Akarun



Dept. of Computer Engineering,
Bogazici University, TURKEY

Introduction

Objective: Polyphonic music transcription



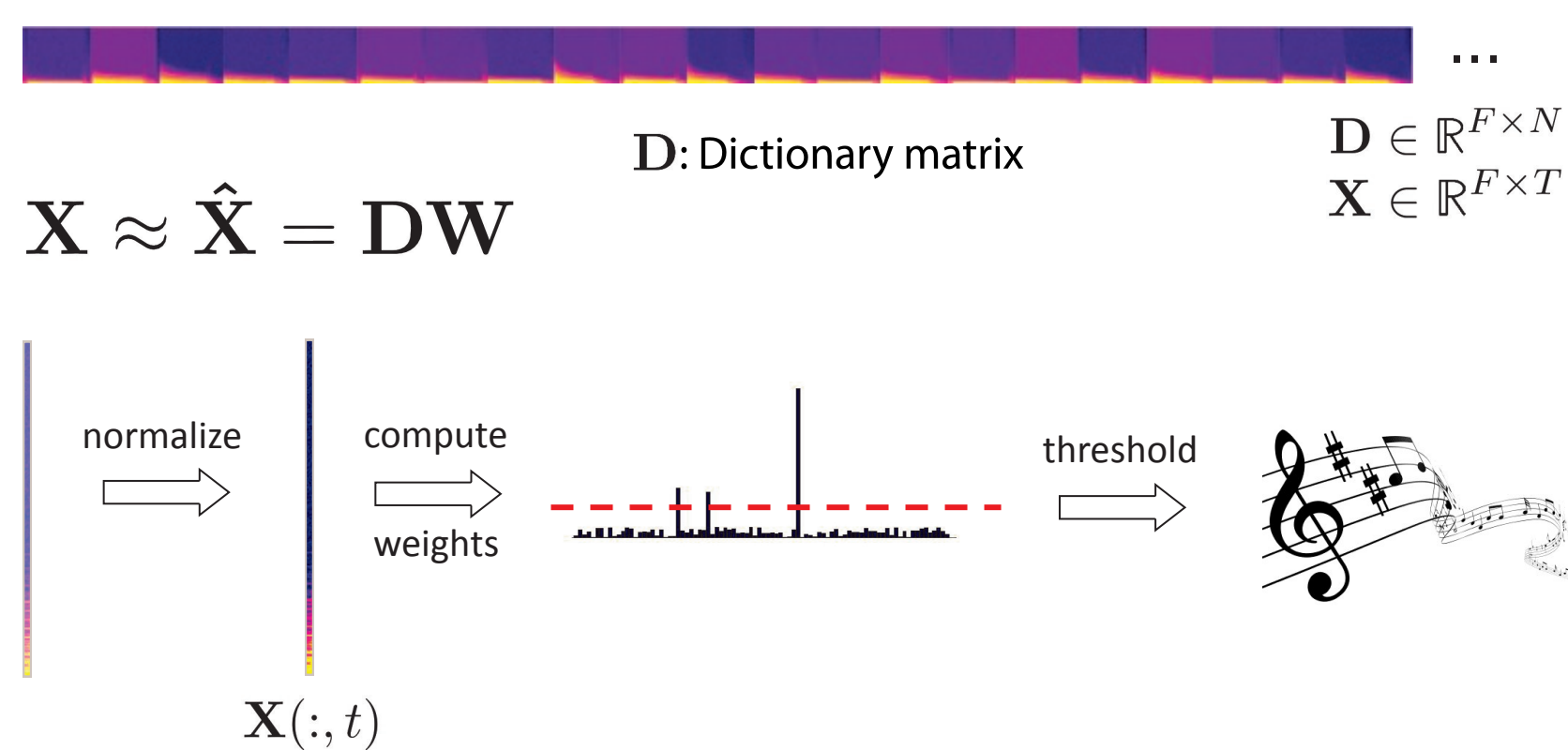
la₂ mi₃ do#₅

Requirements

- Fully automatic
- Deal with large data
- Efficiency (time & space complexities)
- Accuracy, simplicity and robustness

Materials & Methods

Linear model approach



Find \mathbf{W} that minimizes the cost $\mathcal{D}[\mathbf{X}||\mathbf{D}\mathbf{W}]$.

We choose KL-divergence as cost function:

- randomly initialize \mathbf{W}
- continue with the following step until convergence

$$\mathbf{W} \leftarrow \mathbf{W} \odot \left(\frac{\mathbf{D}^T \mathbf{X}}{\mathbf{D}^T \mathbf{1}} \right)$$

Division is done element-wise

⊙ implies element-wise (Hadamard) multiplication

$\mathbf{1}$ is a vector of ones of appropriate size

\mathbf{D}^T is the transpose of \mathbf{D}

- 😊 Simple model; easy to understand and implement
- ☹ Dictionary is huge in real apps; it may not fit into the RAM
- ☹ Matrix-products in the equation are costly

Improving efficiency via SVD

Use reduced SVD to compute rank- k approximation of \mathbf{D} .

Employ **randomized SVD** to deal with large data.

$$\arg \min_{\tilde{\mathbf{D}}, \text{rank}(\tilde{\mathbf{D}}) \leq k} \|\mathbf{D} - \tilde{\mathbf{D}}\|_F = \mathbf{A}_k \Sigma_k \mathbf{B}_k^T$$

Store \mathbf{A}_k and store $\Sigma_k \mathbf{B}_k^T$ as $\tilde{\mathbf{B}}_k^T$.

Modify the update rule:

$$\mathbf{W} \leftarrow \mathbf{W} \odot \left(\frac{\tilde{\mathbf{B}}_k \left(\mathbf{A}_k^T \frac{\mathbf{X}}{\mathbf{A}_k (\tilde{\mathbf{B}}_k^T \mathbf{W})} \right)}{\tilde{\mathbf{B}}_k (\mathbf{A}_k^T \mathbf{1})} \right)$$

- 😊 Best factorization in terms of reconstruction error
- ☹ Singular vectors may not represent *physical reality*

Improving efficiency via CUR

Use $\mathbf{D} \approx \mathbf{C}\mathbf{U}\mathbf{R}$ to represent dictionary via its columns/rows.

$k_{\text{col}}/k_{\text{row}}$ are the number of selected columns/rows.

First, compute singular vectors, then

⇒ Probability of selecting i^{th} row is

$$\rho_i = \frac{1}{k} \sum_{j=1}^k A_k(i, j)^2, \quad i = 1, \dots, F$$

⇒ Probability of selecting j^{th} column is

$$\pi_j = \frac{1}{k} \sum_{i=1}^k B_k(j, i)^2, \quad j = 1, \dots, N$$

Select $k_{\text{col}}/k_{\text{row}}$ columns/rows randomly using a multinomial distribution with the computed probabilities.

Store \mathbf{R} and store $\mathbf{C}\mathbf{U}$ as $\tilde{\mathbf{C}}$.

Modify the update rule:

$$\mathbf{W} \leftarrow \mathbf{W} \odot \left(\frac{\mathbf{R}^T \left(\tilde{\mathbf{C}}^T \frac{\mathbf{X}}{\tilde{\mathbf{C}} (\mathbf{R}\mathbf{W})} \right)}{\mathbf{R}^T (\tilde{\mathbf{C}}^T \mathbf{1})} \right)$$

- 😊 More interpretable decomposition
- ☹ Still trying to approximate the large dictionary matrix

Improving efficiency via selection of important data

Instead of using low rank approximation of the full matrix

⇒ Use only the selected columns

$$\mathbf{W} \leftarrow \mathbf{W} \odot \left(\frac{\mathbf{C}^T \frac{\mathbf{X}}{\mathbf{C}\mathbf{W}}}{\mathbf{C}^T \mathbf{1}} \right)$$

⇒ Let $\tilde{\mathbf{D}}$ be the truncated matrix involving only the selected indices and $\tilde{\mathbf{X}}$ involve the selected rows of the input

Then, use only the selected indices

$$\mathbf{W} \leftarrow \mathbf{W} \odot \left(\frac{\tilde{\mathbf{D}}^T \frac{\tilde{\mathbf{X}}}{\tilde{\mathbf{D}}\mathbf{W}}}{\tilde{\mathbf{D}}^T \mathbf{1}} \right)$$

- 😊 Use only relevant information
- 😊 Fast & low space requirement

Time & Space Complexities (in Big-O notation)

| | Time | Space |
|------------|--|--|
| Full model | FNT | FN |
| SVD-based | $(F+N)kT$ | $(F+N)k$ |
| CUR-based | $(F+N) \min\{k_{\text{row}}, k_{\text{col}}\} T$ | $(F+N) \min\{k_{\text{row}}, k_{\text{col}}\}$ |
| C-based | $F k_{\text{col}} T$ | $F k_{\text{col}}$ |
| Skeleton | $k_{\text{row}} k_{\text{col}} T$ | $k_{\text{row}} k_{\text{col}}$ |

Experimental Setup

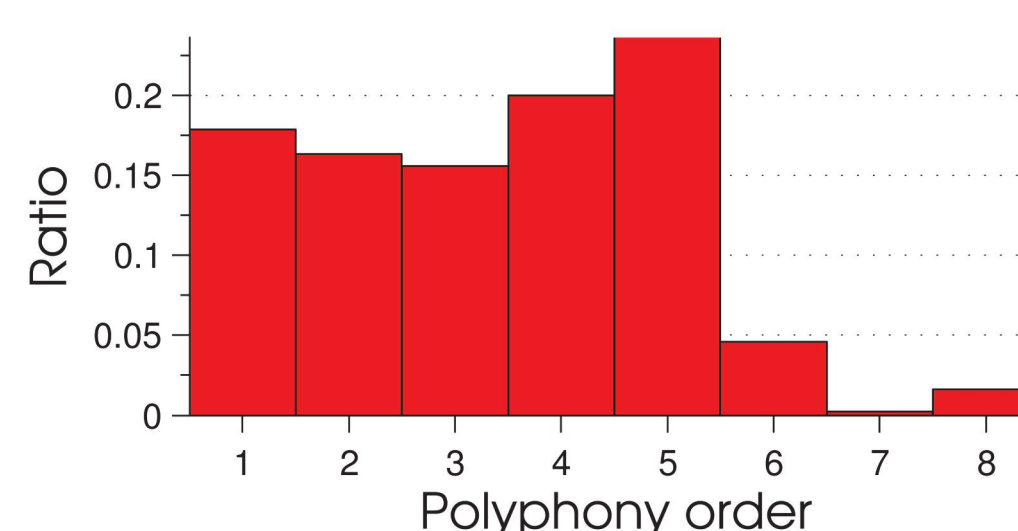
MAPS (MIDI Aligned Piano Sounds) Dataset

⇒ Training set

- 440 monophonic piano recordings
- Dictionary of size 1025×115600 (~ 860 MB)

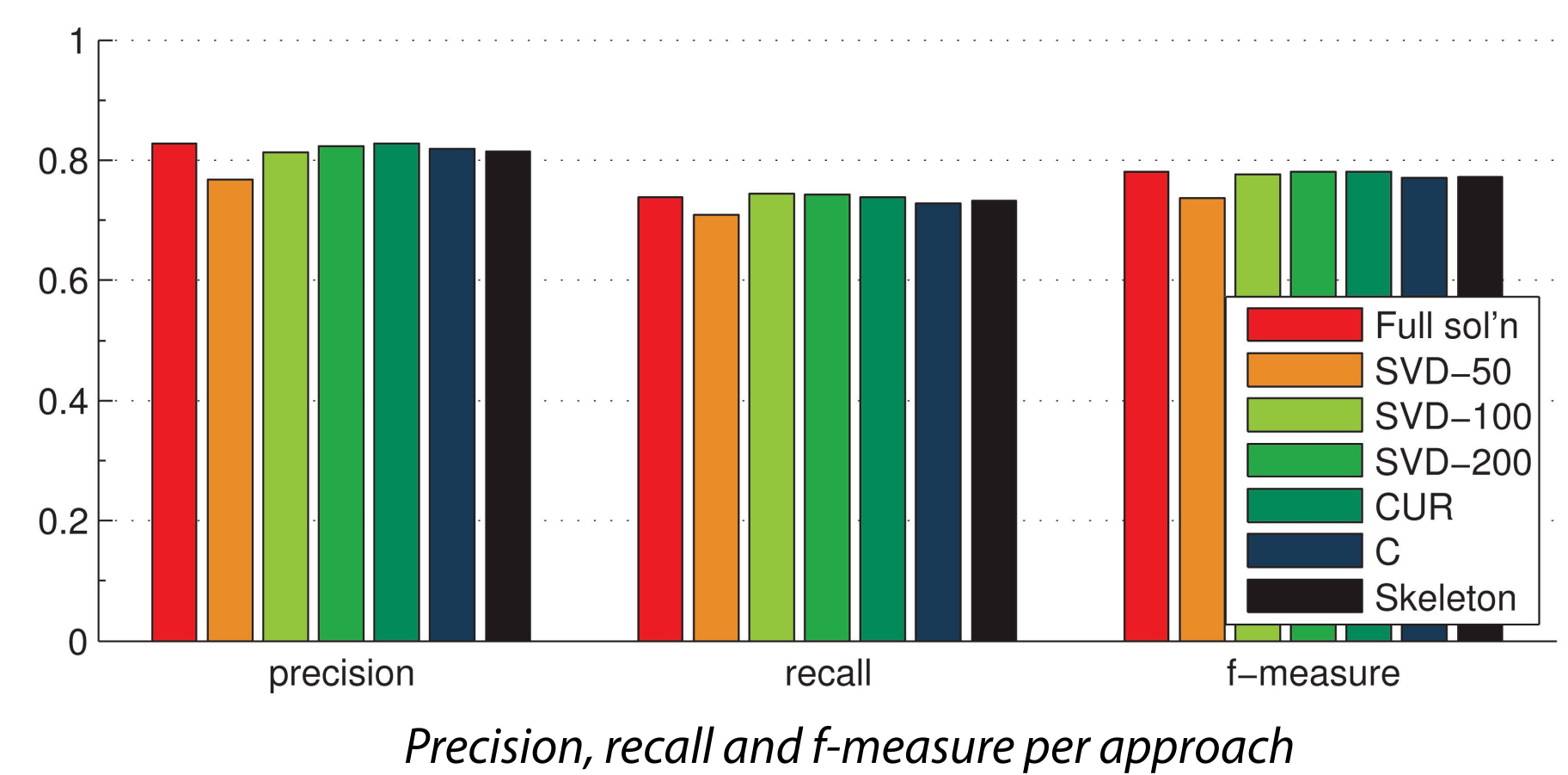
⇒ Test set

- Random sections from 5 different polyphonic pieces
- 3000 samples



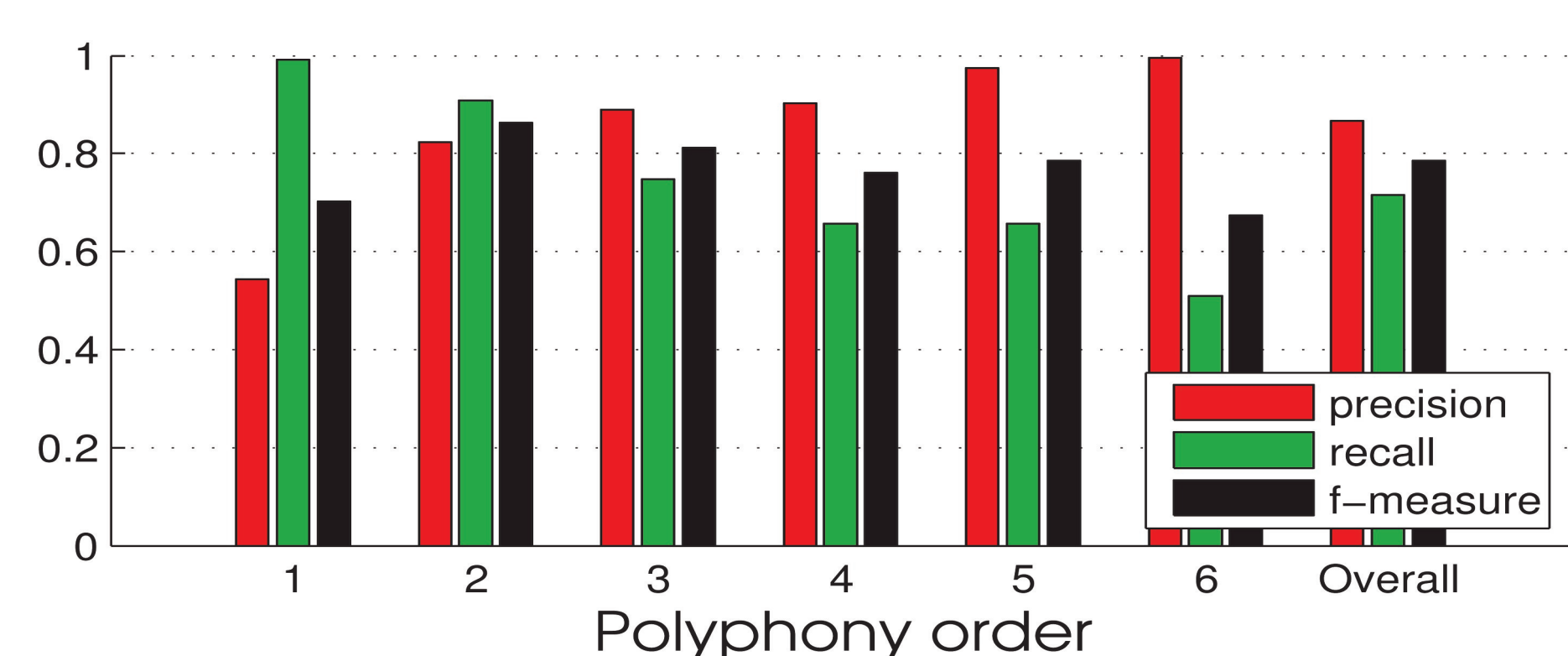
Results

Results per approach



- 😊 For SVD-based approach, 200 dimensions (with 78.14% f-measure) is enough to maintain the original f-measure score (78.07%).
- 😊 We select 400 frequency bands and 4000 samples; CUR is a good alternative to SVD with similar results.
- 😊 C-based approach and skeletonizing (reducing both the #cols and the #rows) give promising results.
- 😊 We only use less than 2% of the data after skeletonization, and get an f-measure of 77.17%.

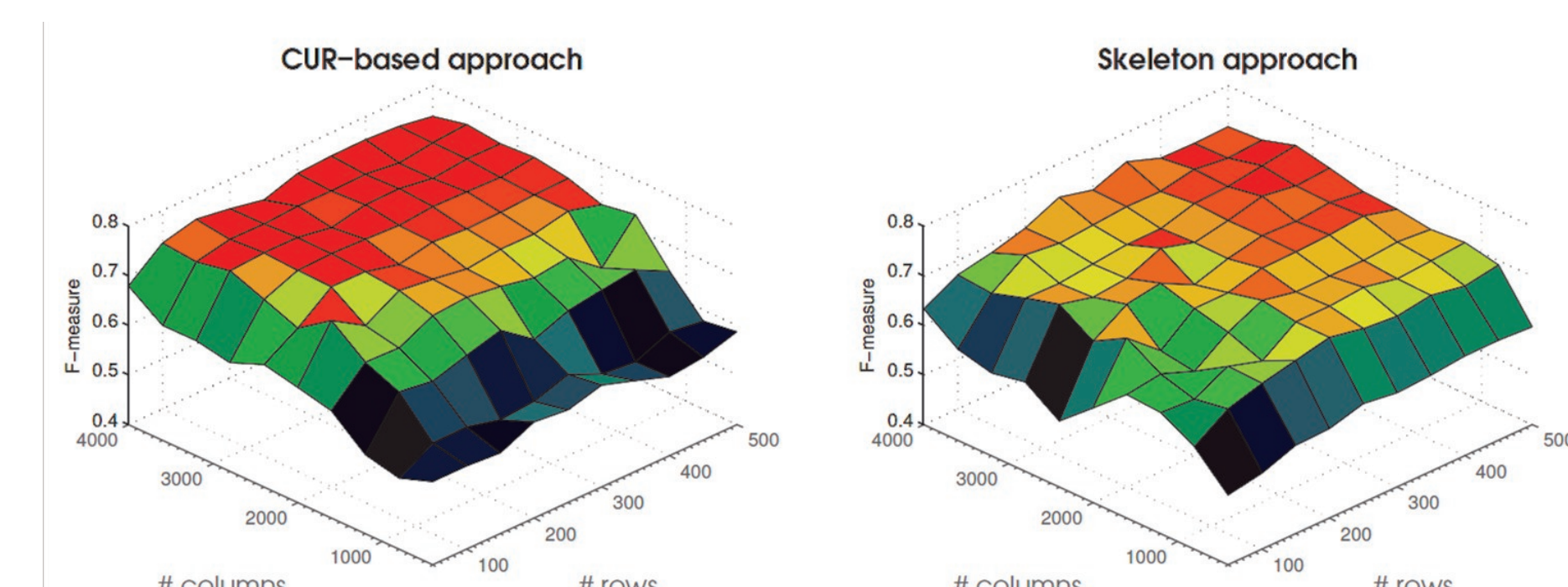
Results per polyphony order



Precision, recall and f-measure per polyphony order using skeletonization

- 😊 More than 65% f-measure for each polyphony order ⇒ Robust

Effect of #cols/#rows on the result



f-measures versus the #cols (samples) and the #rows (frequency bins)

- 😊 Only a few hundred frequency bins and a few thousand samples are sufficient to keep success ratio of the algorithm

Conclusions

- We show that even a standard matrix factorization model is prohibitive in real applications where a huge amount of data is used.
- A high f-measure value (~78%) is obtained on polyphonic recordings by using only a few hundred frequency bins and a few thousand sample columns out of a huge dictionary.
- The technique is simple, yet powerful and robust.
- Randomized matrix decompositions are crucial for practical issues.
- With abundance of data in different applications, randomized matrix decompositions are likely to get more attention in the future.

