

Compositional models for sound source separation and processing

Tuomas Virtanen

Tampere University of Technology
Laboratory of Signal Processing
Finland

www.cs.tut.fi/~tuomasv/



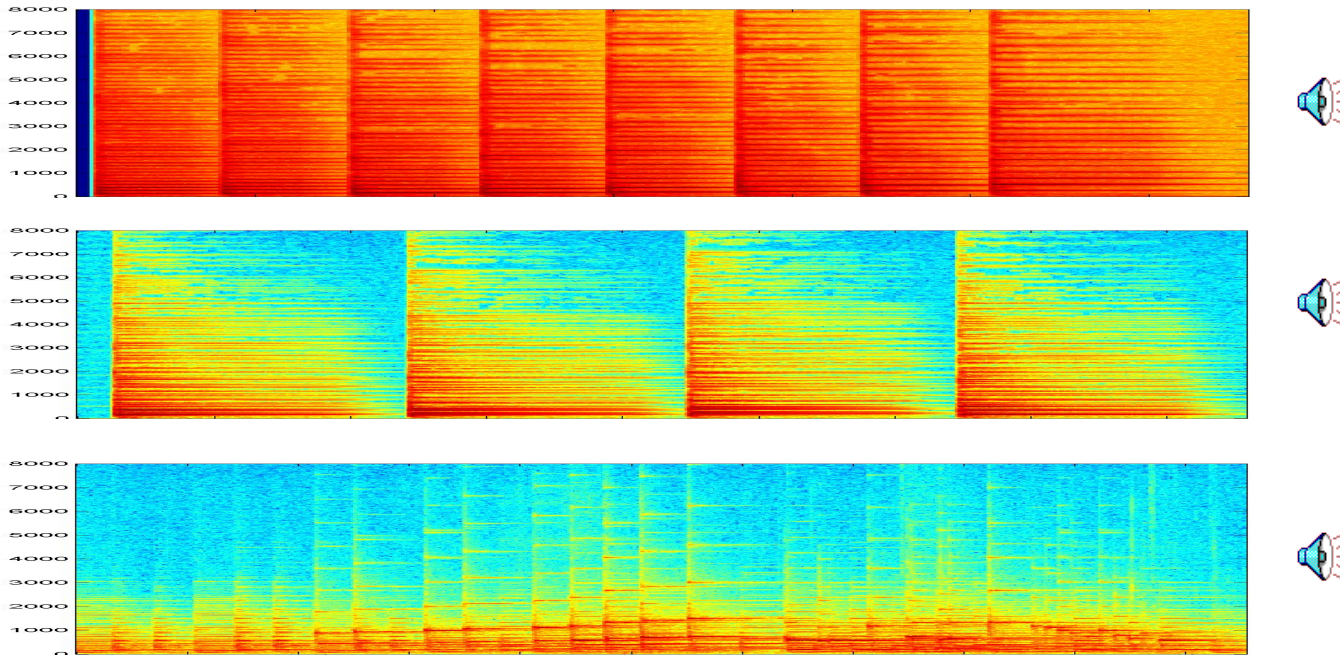
TAMPERE UNIVERSITY OF TECHNOLOGY

Outline

- Introduction
 - What are compositional models
 - Signal representation
- Application: source separation
- Algorithm: non-negative matrix factorization (NMF)
- Analyzing the semantics of sound
- Model alternatives
- Comparison to DNNs
- Missing data techniques

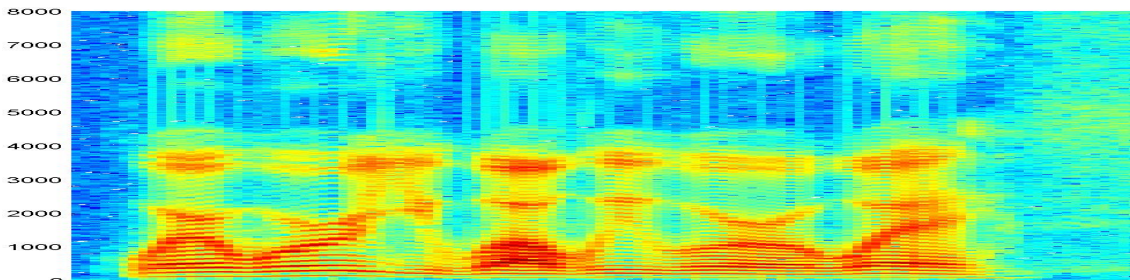
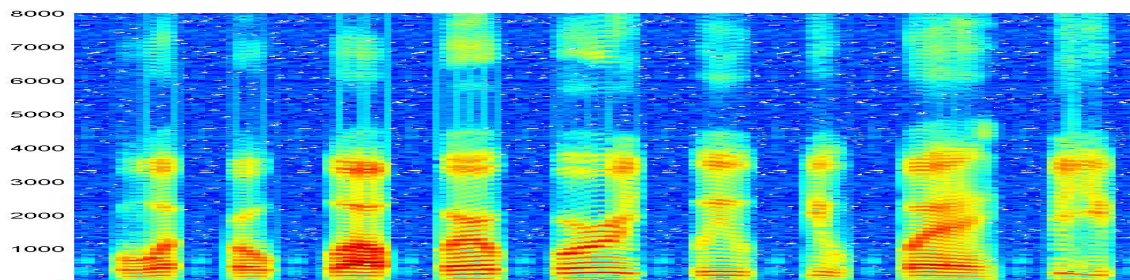
Compositional model

- In general, sounds do not cancel each other
- Typically, individual components combine to form the sounds we hear
 - Notes, but also multiple instruments, form music



Compositional model

- In general, sounds do not cancel each other
- Typically, individual components combine to form the sounds we hear
 - Notes, but also multiple instruments, form music
 - Phoneme-like sounds combine to form speech



Compositional model

- In general, sounds do not cancel each other
- Typically, individual components combine to form the sounds we hear
 - Notes, but also multiple instruments, form music
 - Phoneme-like sounds combine to form speech
- The *compositional model* is a linear, additive combination of components that do not result in subtraction or diminishment of any of the constituents

Compositional model

- Feature vector \mathbf{y}_t is decomposed into weighted sum of basis vectors \mathbf{a}_n

$$\mathbf{y}_t \approx \sum_n \mathbf{a}_n x_{nt}$$

- x_{nt} are gains of the components in observation t
- Compositional model: both the basis vectors and weights are constrained to be **non-negative**

Compositional model

- Model in a vector-matrix form

$$\begin{bmatrix} y_{1t} \\ \vdots \\ y_{Ft} \end{bmatrix} \approx \begin{bmatrix} a_{11} & \cdots & a_{1N} \\ \vdots & \ddots & \vdots \\ a_{F1} & \cdots & a_{FN} \end{bmatrix} \cdot \begin{bmatrix} x_{1t} \\ \vdots \\ x_{Nt} \end{bmatrix} \quad \rightarrow \quad \mathbf{y}_t \approx \mathbf{A}\mathbf{x}_t$$

Model for multiple observations

- We can efficiently write the compositional model

$$\mathbf{y}_t \approx \mathbf{A}\mathbf{x}_t, \quad t = 1 \dots T$$

- for all T observations (e.g. a spectrogram), as:

$$[\mathbf{y}_1 \dots \mathbf{y}_T] \approx \mathbf{A} \cdot [\mathbf{x}_1 \dots \mathbf{x}_T]$$

- Or even:

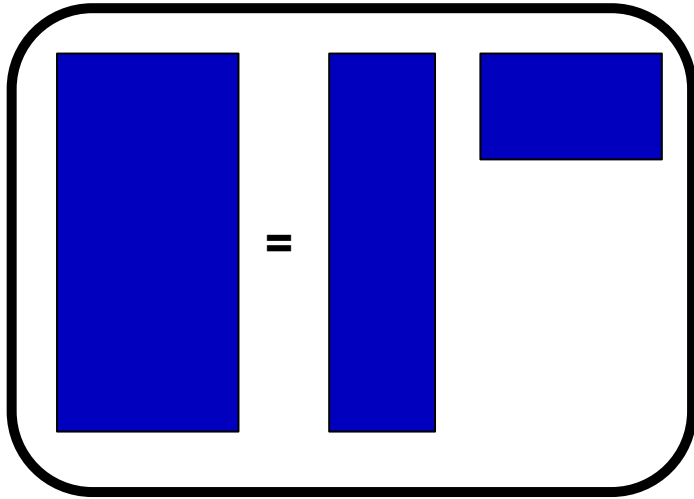
$$\mathbf{Y} \approx \mathbf{A}\mathbf{X}$$

Visual representation

$$Y \approx AX$$

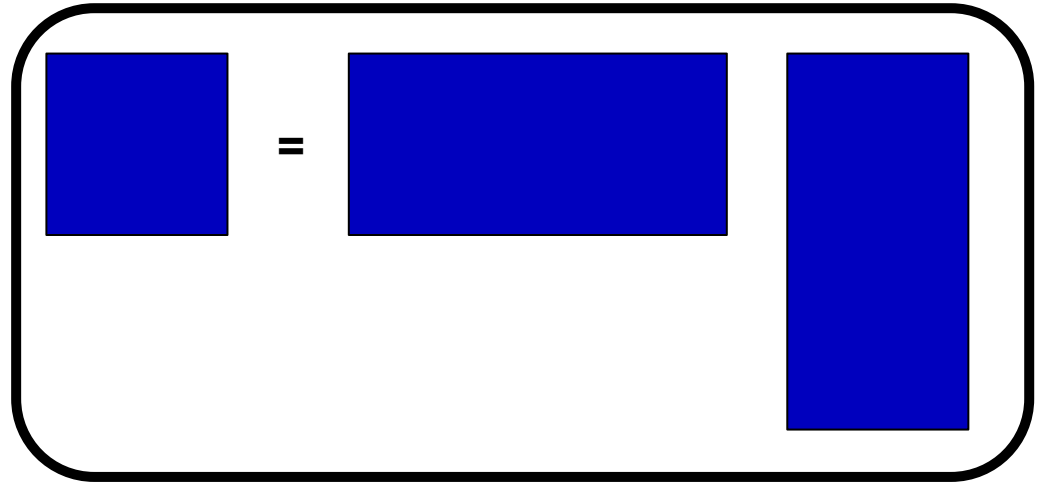
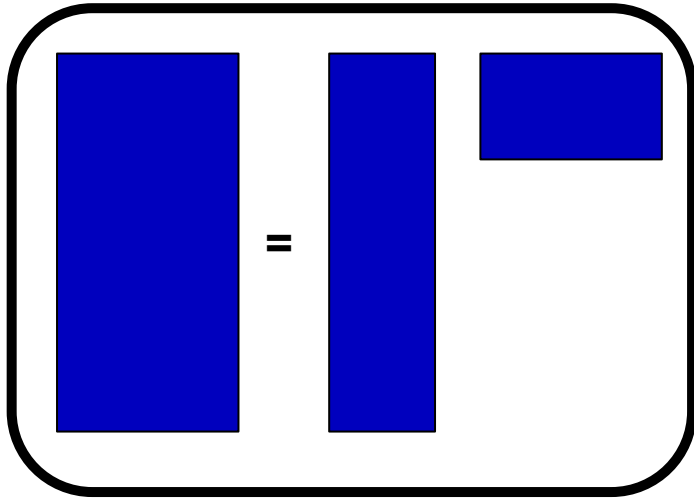
Visual representation

$$Y \approx AX$$



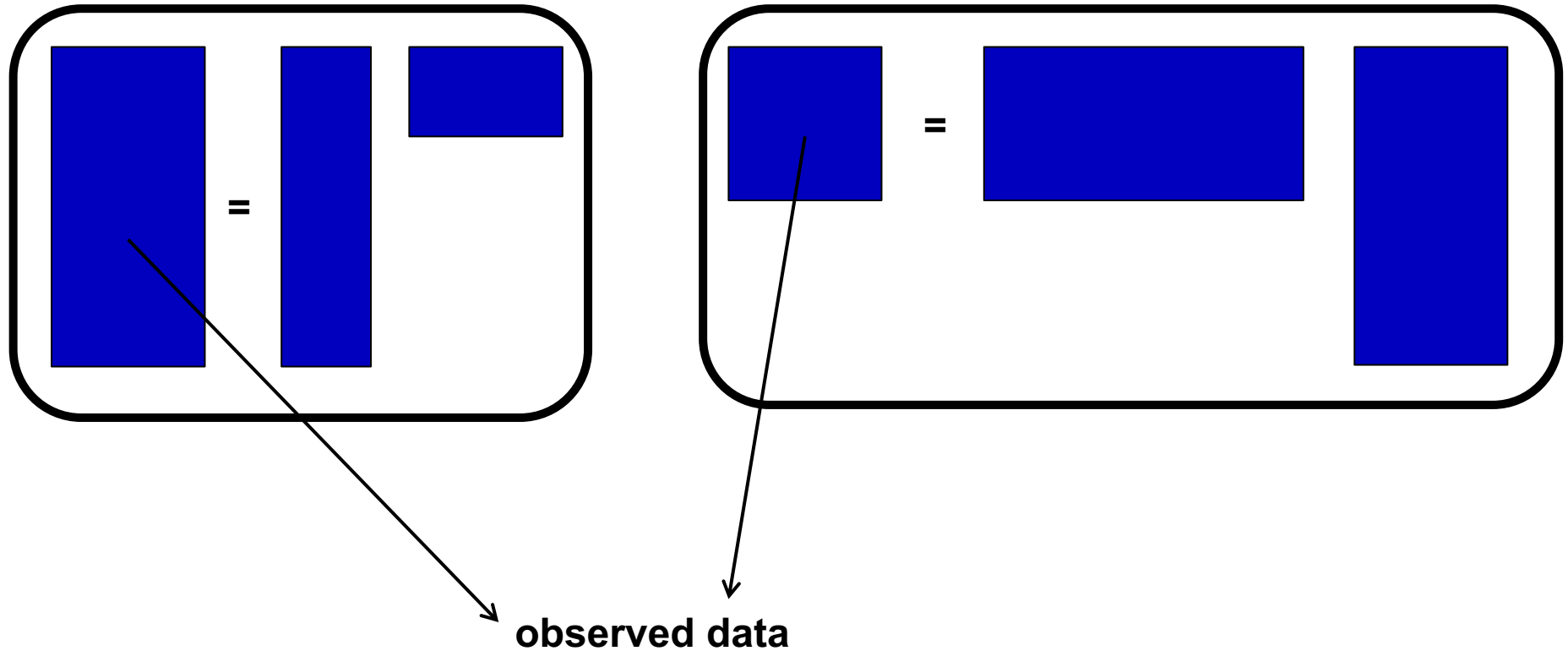
Visual representation

$$Y \approx AX$$



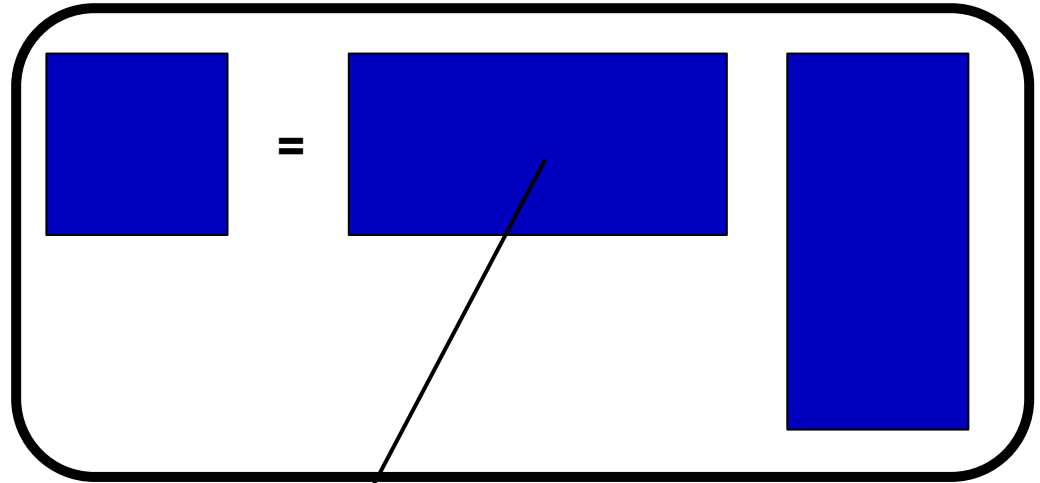
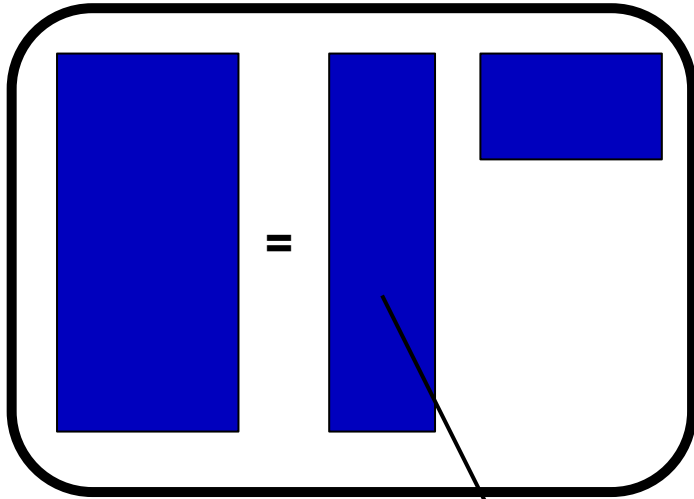
Visual representation

$$Y \approx AX$$



Visual representation

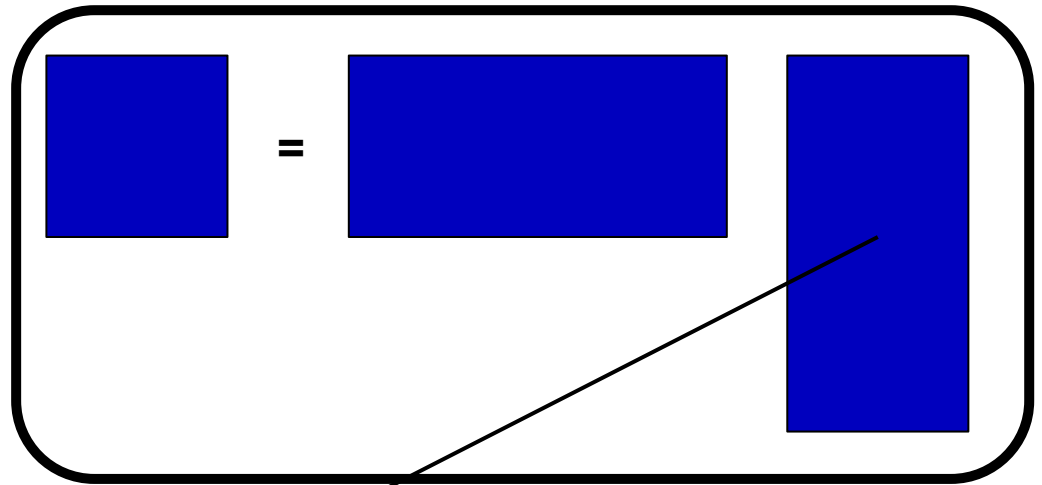
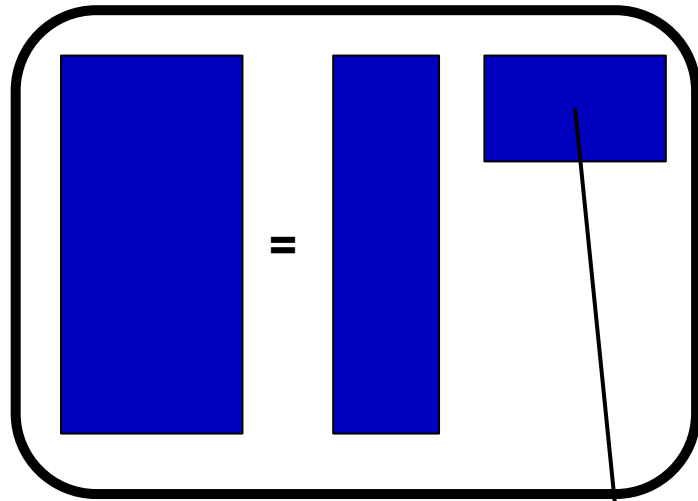
$$Y \approx AX$$



dictionary or basis (learned or constructed, updated or kept fixed)

Visual representation

$$Y \approx AX$$



mixture weights, sparse representations

Compositional generative model

- The compositional model explains how the observations are generated, given the model parameters

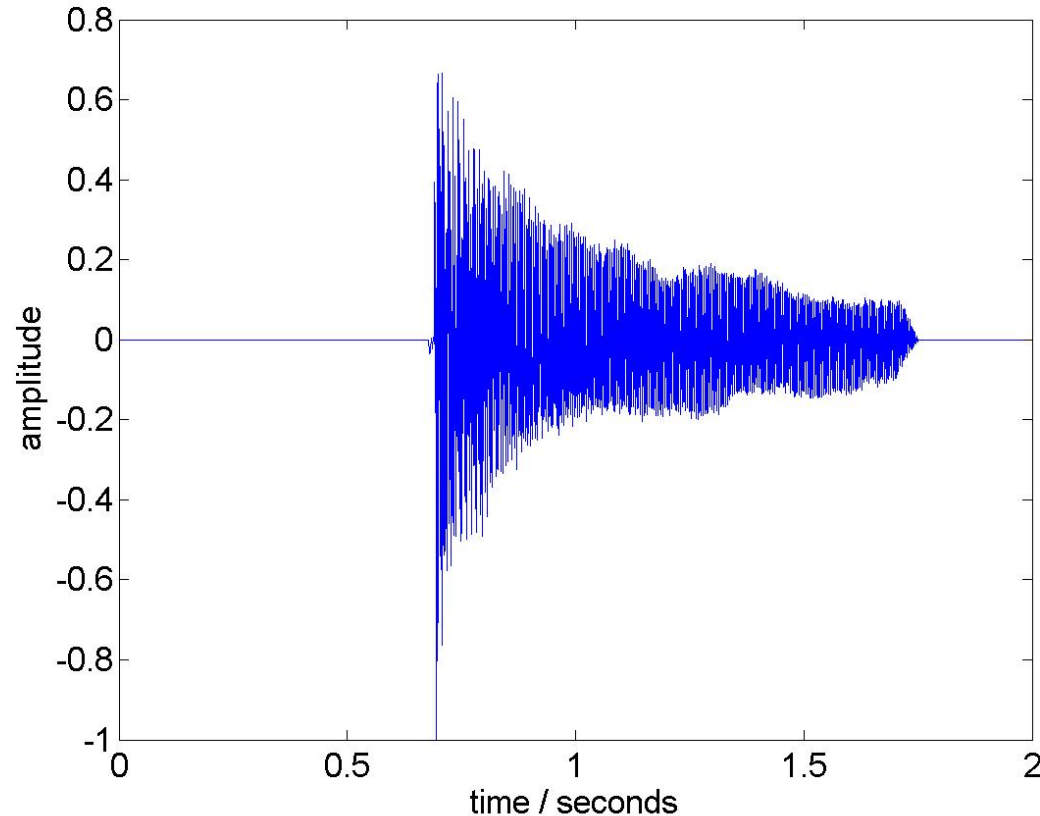
$$p(Y|A, X)$$

Audio representation

- Compositional models require a non-negative representation
- Audio *signals* have both negative and positive values
- Need for a mid-level representation that is used for processing

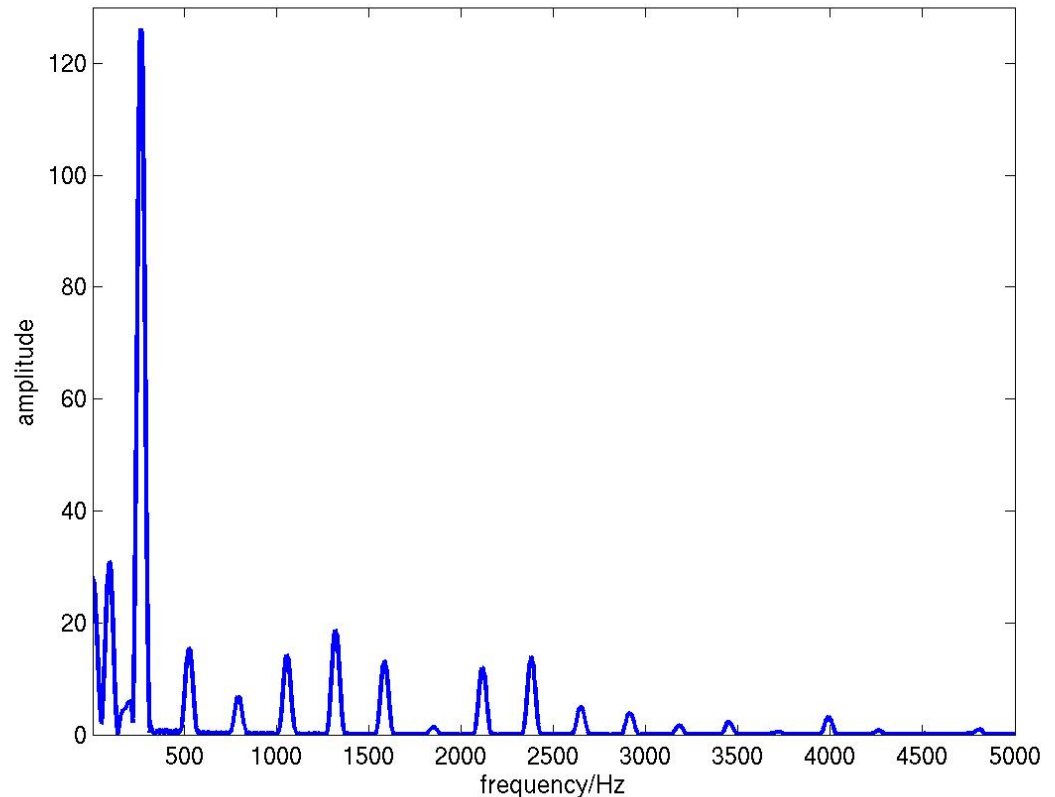
Audio representation

- Audio signal – amplitude as a function of time



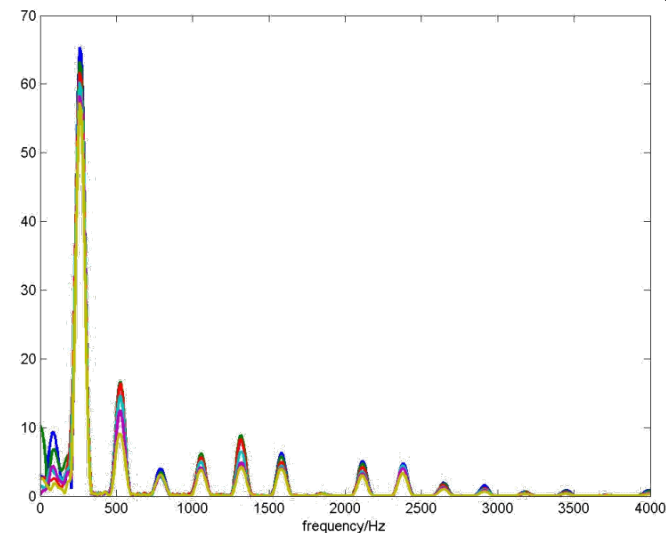
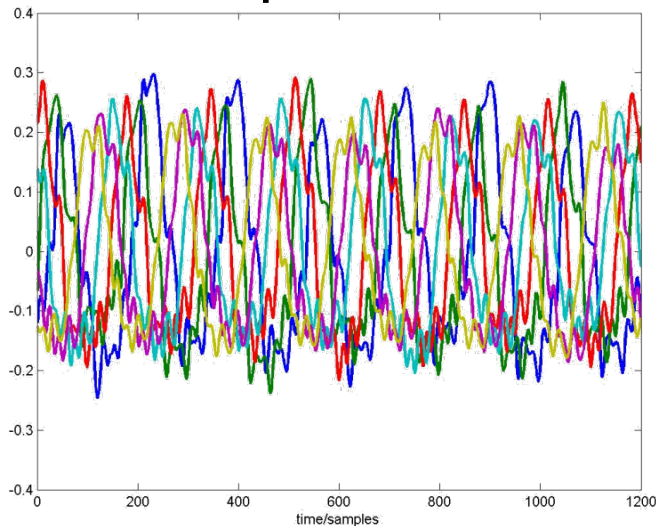
Magnitude spectrum

- Phases are discarded and only the magnitudes are used
→ non-negative representation
- Can use any spectral resolution (linear, logarithmic, perceptual...)



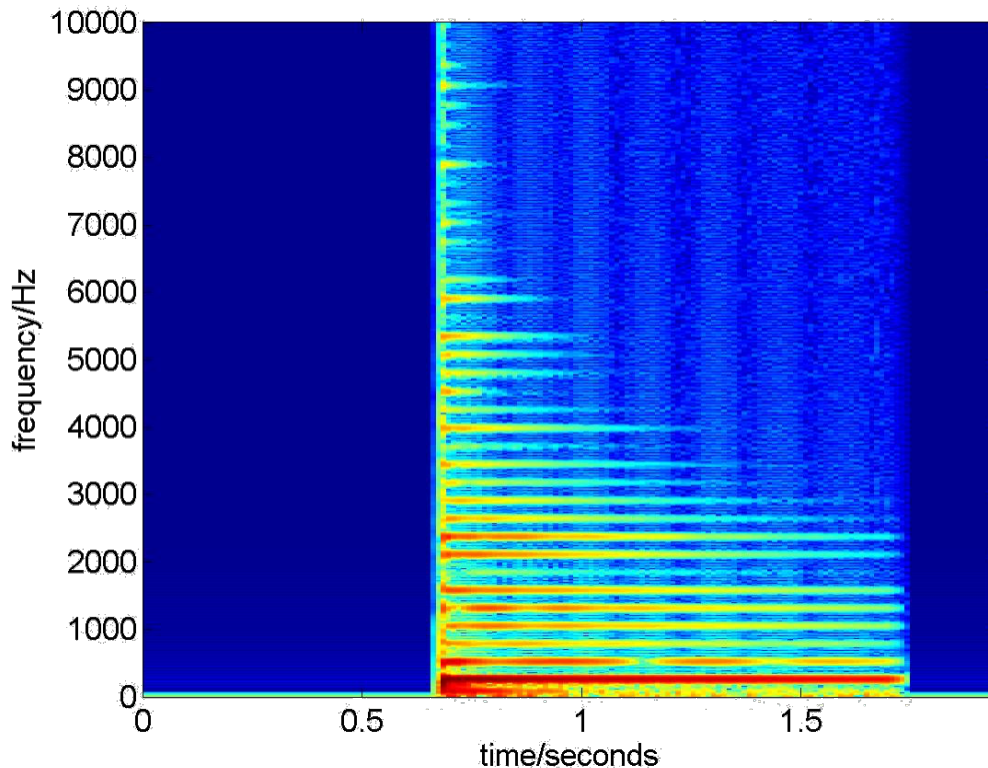
Magnitude spectrum

- Natural sounds have a clear structure in the magnitude spectrum domain
- Discarding the phases makes the representation invariant to many factors
 - Relative window position
 - Phase of the acoustic impulse response from source to microphone
- Example: five consecutive frames from the earlier signal



Spectrogram

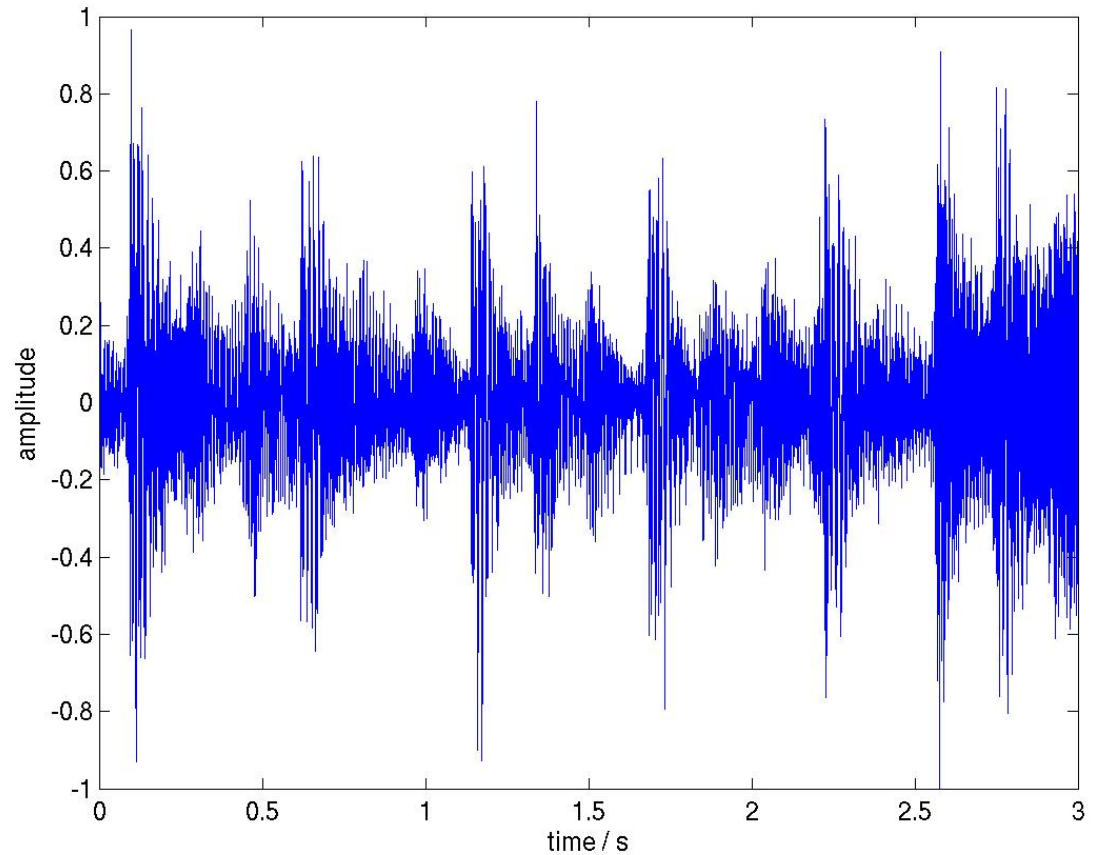
- Spectra in each frame grouped to a matrix
- Represents the intensity of a sound as a function of time and frequency



Linear superposition

- When multiple sound sources are present, the time-domain signals add linearly

$$y(t) = s_1(t) + s_2(t) \dots$$



Additivity of magnitude spectra

- In the magnitude spectral domain, sounds are *approximately* additive

$$y(t) = s_1(t) + s_2(t) \dots$$

$$|Y(f)| \approx |S_1(f)| + |S_2(f)| \dots$$

- Exactly additive only when the phases are coherent
- For independent source, power spectra are additive in the expectation sense:

$$E \left\{ |Y(f)|^2 \right\} = E \left\{ |S_1(f)|^2 \right\} + E \left\{ |S_2(f)|^2 \right\} \dots$$

- ... sounds are also approximately additive in the power spectral domain;

$$|Y(f)|^2 \approx |S_1(f)|^2 + |S_2(f)|^2 \dots$$

Additivity of magnitude spectra

- Magnitude vs. power spectrum representation?
 - I.e., $|Y(f)|$ vs. $|Y(f)|^2$
- Determines the dynamic scale of the representation
- Affects the relative importance of low vs. high-intensity observations
- Related to the compositional model estimation criterion (see later)
- Empirically observed that additivity of magnitudes works better

Additivity of magnitude spectra

- How valid is the approximation?
- Natural sounds are *sparse* and therefore *disjoint* in the time-frequency domain

$$|S_1(t, f)| |S_2(t, f)| \approx 0$$

- Additivity of magnitude or power spectrum works well enough in practice
- Lower frequency and time resolutions lead to lower sparseness and disjointness

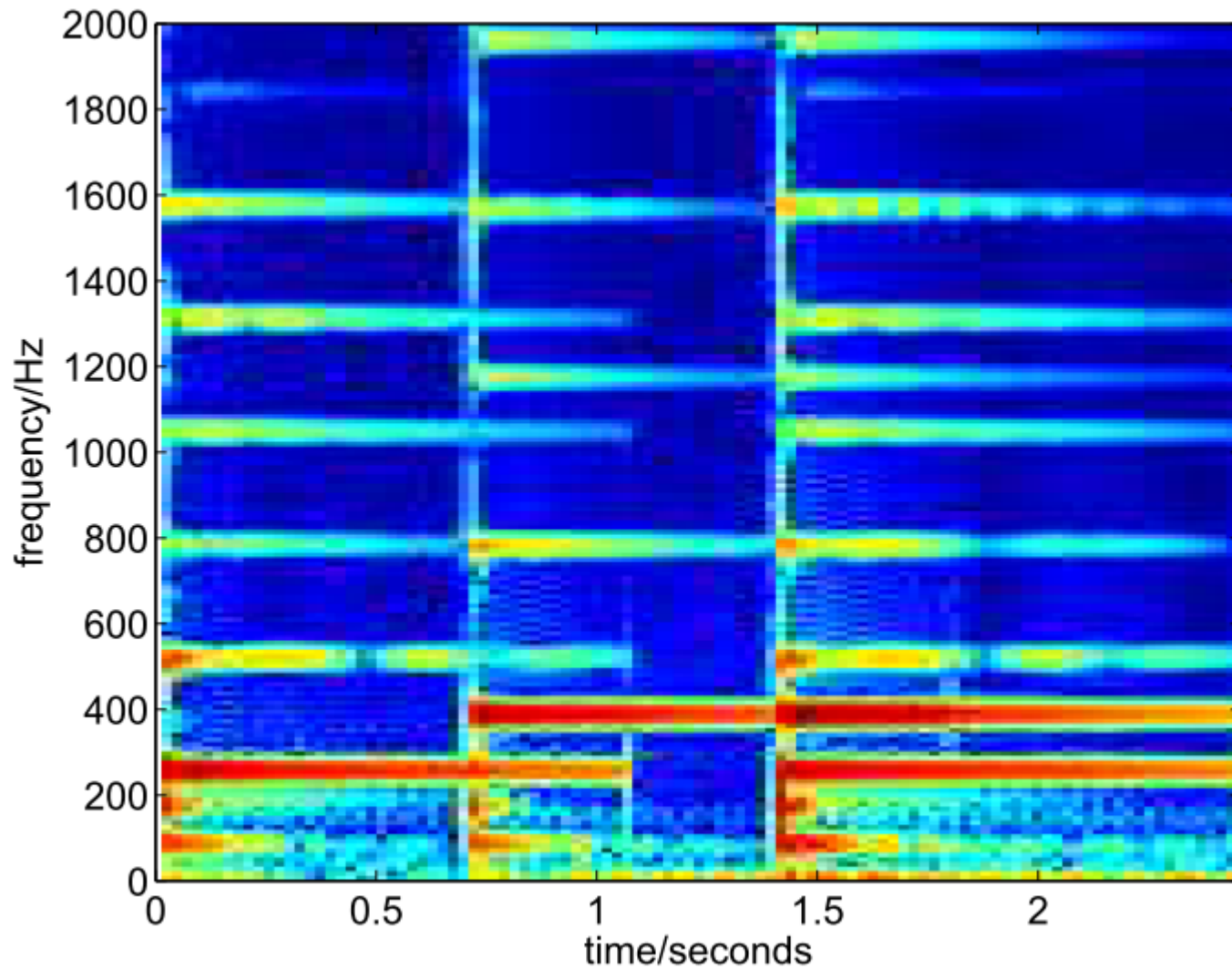
Compositional spectral model

- Clear interpretation:

- Signal modeled as a sum of components
- Each components has a fixed spectrum (basis vectors \mathbf{a}_n) and time-varying gain x_{nt}

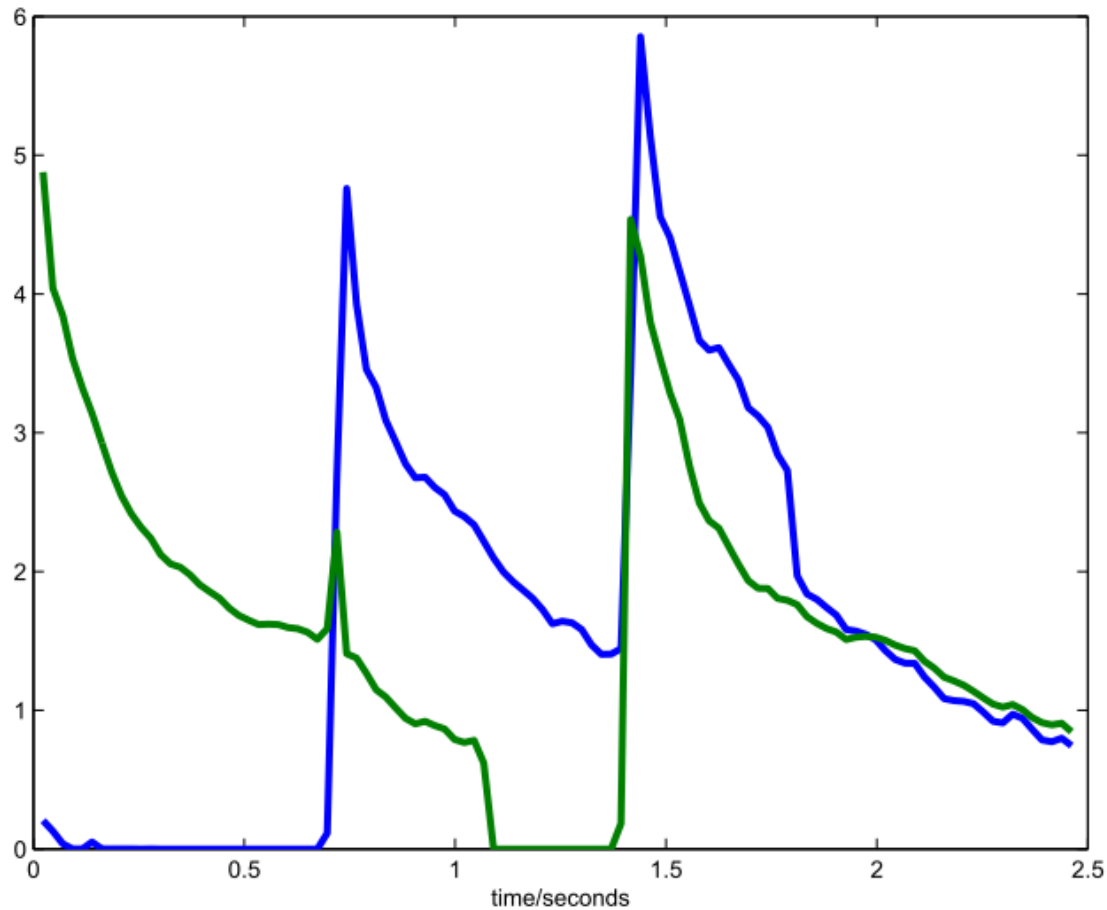
$$\mathbf{y}_t \approx \sum_n \mathbf{a}_n x_{nt}$$

Mixture spectrogram



Results with NMF

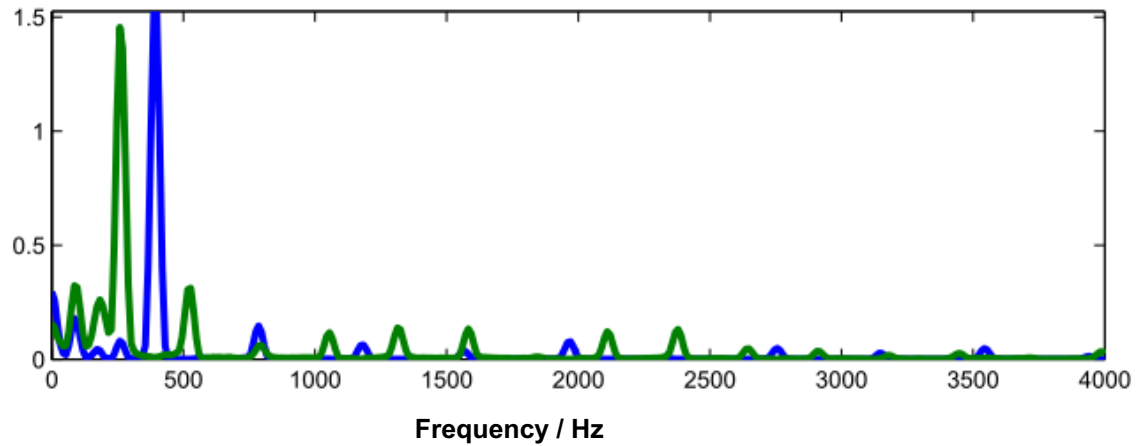
- Weights over time: separation of notes



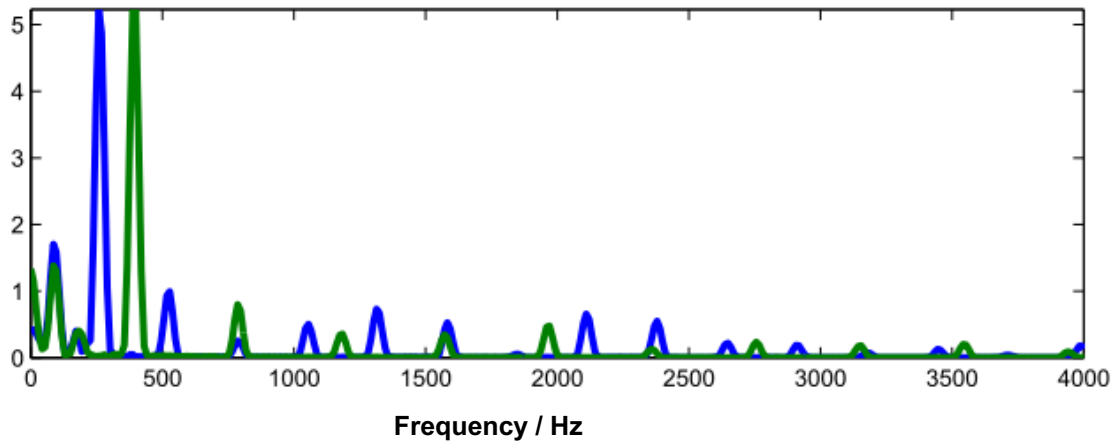
Results with NMF

- Bases correspond to individual notes

NMF

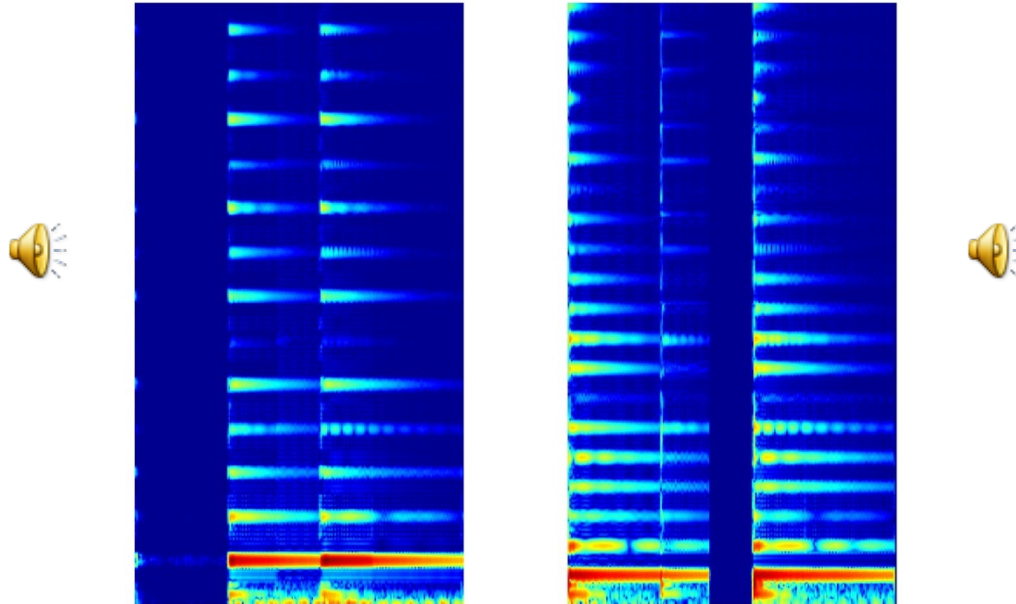


Original



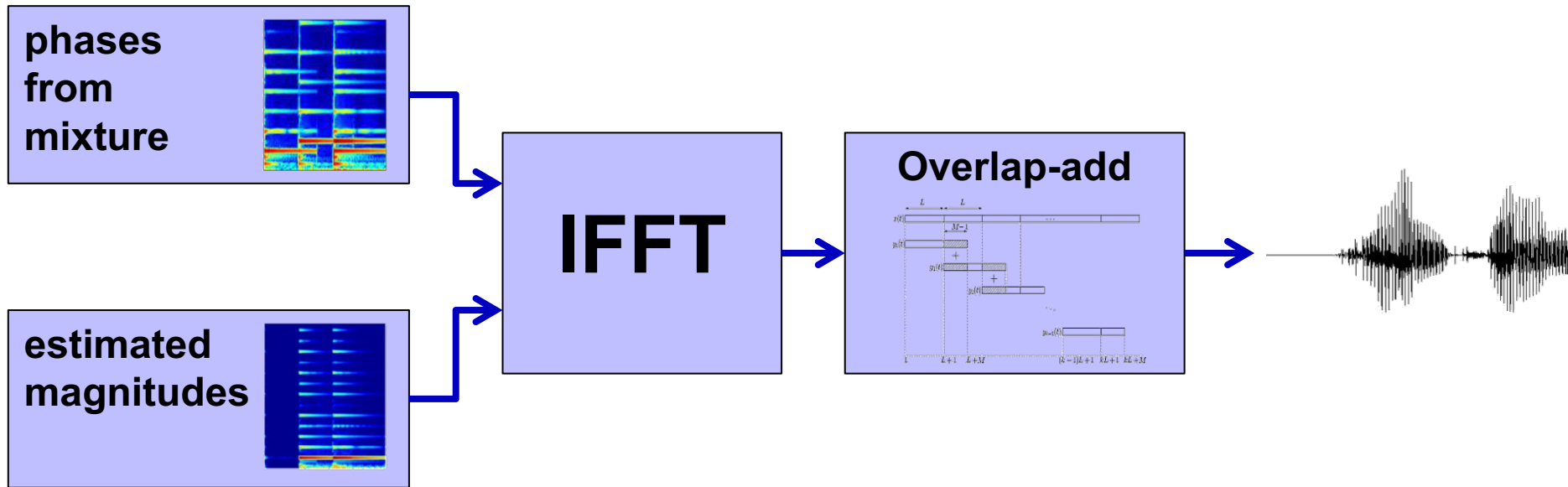
Wiener-style reconstruction

- k -th component reconstructed as: $Y \otimes \frac{A(:,k) \bullet X(k,:)}{AX}$



Signal reconstruction

- For each component:
 1. Use the phases of the mixture signal
 2. IFFT
 3. Overlap-add



Non-negative matrix factorization

- Minimize error $d(\mathbf{Y}, \mathbf{AX})$ between \mathbf{Y} and \mathbf{AX} while restricting \mathbf{A} and \mathbf{X} to be entry-wise non-negative

$$\mathbf{A}^*, \mathbf{X}^* = \arg \min_{\mathbf{A}, \mathbf{X}} d(\mathbf{Y}, \mathbf{AX})$$

- Supervised NMF – estimate only the weights, the bases are given:

$$\mathbf{X}^* = \arg \min_{\mathbf{X}} d(\mathbf{Y}, \mathbf{AX})$$

NMF criteria

- Different distance measures
- Squared error (L2 norm):

$$d_{\text{SQ}}(\mathbf{Y}, \mathbf{AX}) = \sum_{f,t} (\mathbf{Y}_{ft} - [\mathbf{AX}]_{ft})^2 = \|\mathbf{Y} - \mathbf{AX}\|_F^2$$

- Generalized Kullback-Leibler divergence:

$$d_{\text{KL}}(\mathbf{Y}, \mathbf{AX}) = \sum_{f,t} \mathbf{Y}_{ft} \log(\mathbf{Y}_{ft} / [\mathbf{AX}]_{ft}) - \mathbf{Y}_{ft} + [\mathbf{AX}]_{ft}$$

- Itakura-Saito divergence

$$d_{\text{IS}}(\mathbf{Y}, \mathbf{AX}) = \sum_{f,t} \mathbf{Y}_{ft} / [\mathbf{AX}]_{ft} - \log(\mathbf{Y}_{ft} / [\mathbf{AX}]_{ft})$$

- Each of these correspond to different generative model $p(\mathbf{Y}|\mathbf{A}, \mathbf{X})$

NMF algorithms

- The objective function is biconvex

$$\mathbf{A}^*, \mathbf{X}^* = \arg \min_{\mathbf{A}, \mathbf{X}} d(\mathbf{Y}, \mathbf{AX})$$

- Global optimum cannot be found
- Iterative algorithms which repeatedly update \mathbf{A} and \mathbf{X} so that the cost decreases at each iteration

Multiplicative update rules

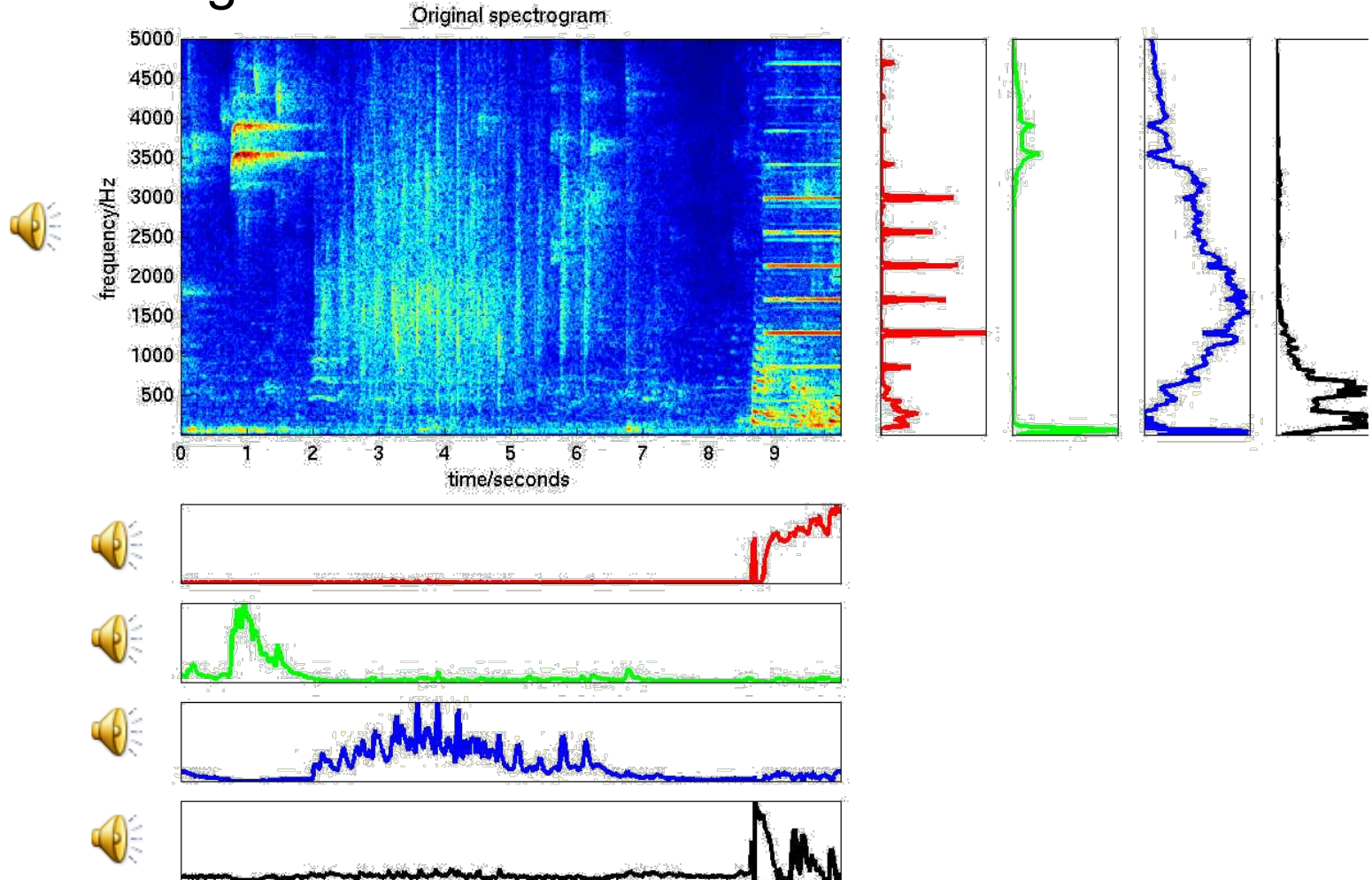
- Update rules under which the cost are guaranteed to be non-increasing
- Guarantees non-negativity of the parameters
- Easy to implement and to extend
- Updates for the KL divergence

$$\mathbf{X} \leftarrow \mathbf{X} \otimes \frac{\mathbf{A}^T (\mathbf{Y} / \mathbf{A}\mathbf{X})}{\mathbf{A}^T \mathbf{1}} \quad \mathbf{A} \leftarrow \mathbf{A} \otimes \frac{(\mathbf{Y} / \mathbf{A}\mathbf{X}) \mathbf{X}^T}{\mathbf{1}\mathbf{X}^T}$$

where $\mathbf{1}$ is all-one matrix of size \mathbf{Y}

Real-world examples

■ Basketball game



Real-world examples

- High-quality separation of complex auditory scenes in blind manner not achievable
- Multiple components required to represent an individual source
- Each component still corresponds to semantically meaningful entity

Supervised source separation

- Prior information easy to include by training the spectral basis **A** vectors in advance
- Optimization problem is convex and therefore finding the global optimum is guaranteed
- More efficient algorithms

$$\mathbf{X}^{|*} = \arg \min_{\mathbf{X}} d(\mathbf{Y}, \mathbf{A}\mathbf{X})$$

- Yields impressive results in matched conditions

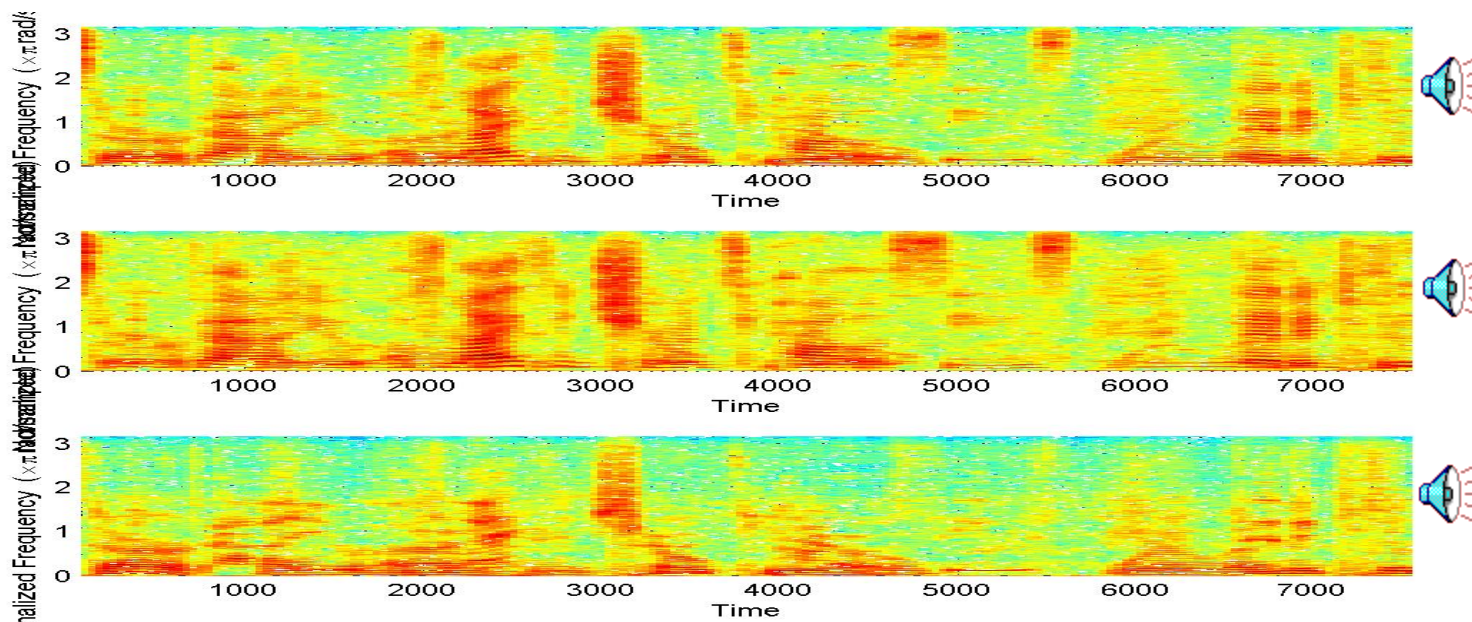
Supervised source separation

- Source separation (SS) scenario:
 - Isolated training material of speech (s) and noise (n)
 - Obtain basis spectra for each source separately
 - Concatenate the dictionaries:

$$\mathbf{Y} = \mathbf{S} + \mathbf{N} \approx \hat{\mathbf{S}} + \hat{\mathbf{N}} = \mathbf{A}_s \mathbf{X}_s + \mathbf{A}_n \mathbf{X}_n = \mathbf{A} \mathbf{X}$$

- Use NMF with the obtained dictionary – keep the dictionary fixed while updating the mixing weights
- Synthesize each source by using only its own basis vectors

Separate overlapping speech



- Bases for both speakers learnt from 5 second recordings of individual speakers

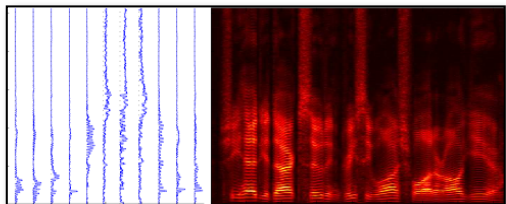
Semi-supervised separation

- We may not have training data for all sources
 - But we usually know some

Semi-supervised separation

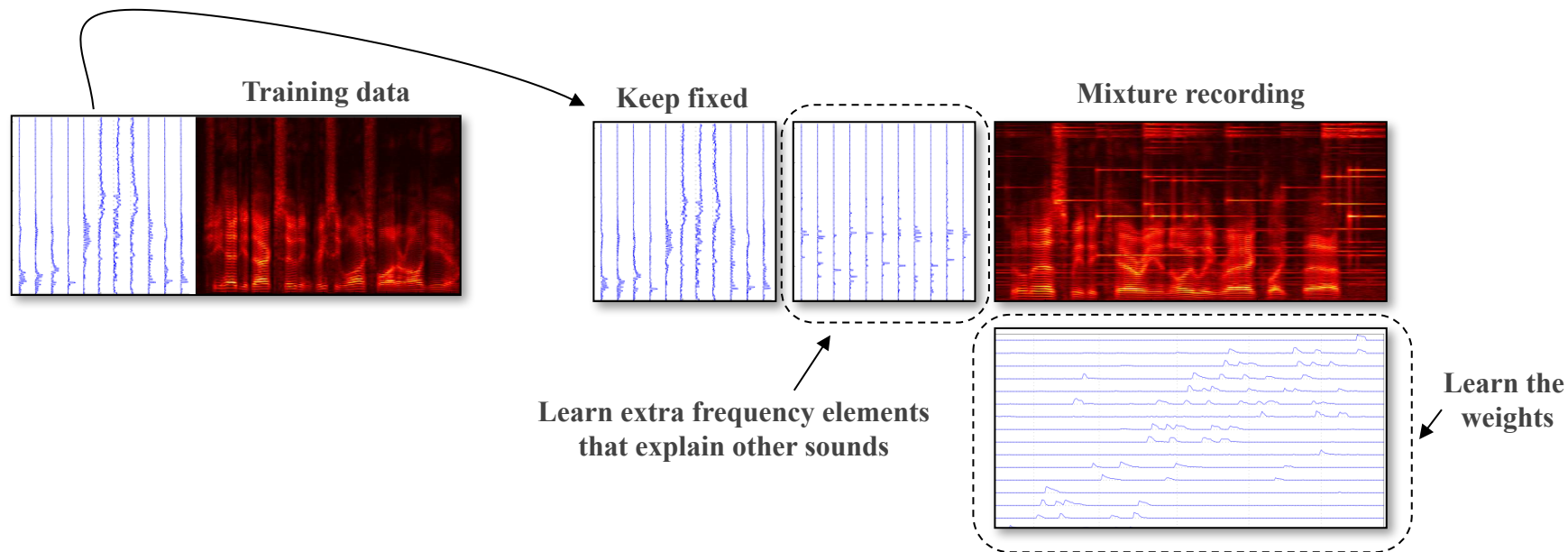
- We may not have training data for all sources
 - But we usually know some
- Two steps:
 - Supervised: Train dictionary for known sources

Training data

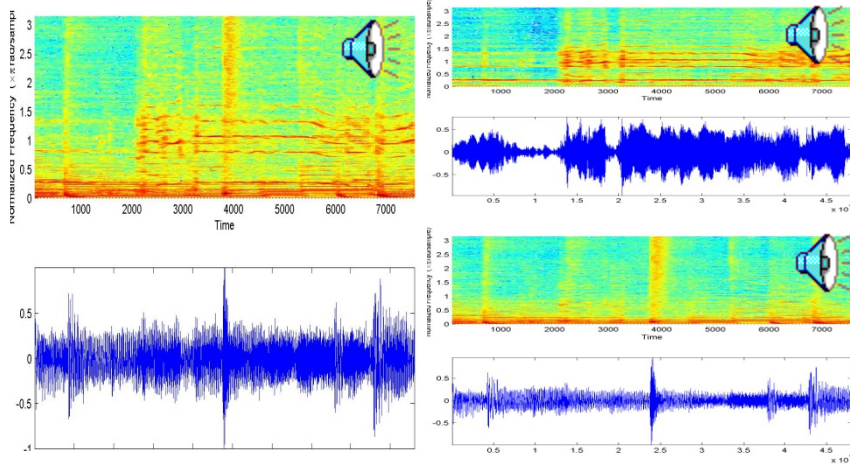


Semi-supervised separation

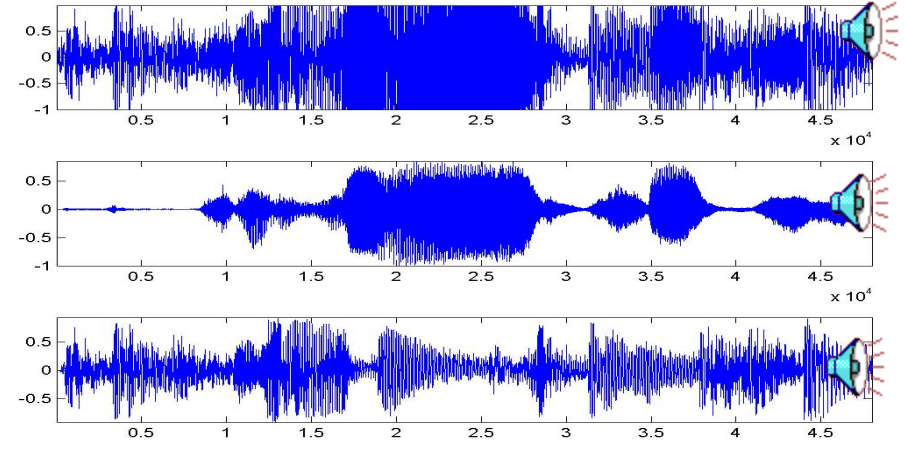
- We may not have training data for all sources
 - But we usually know some
- Two steps:
 - Supervised: Train dictionary for known sources
 - Unsupervised: Train *part* of the dictionary unsupervised on target mixture



Semi-supervised separation



“Raise my rent” by David Gilmour

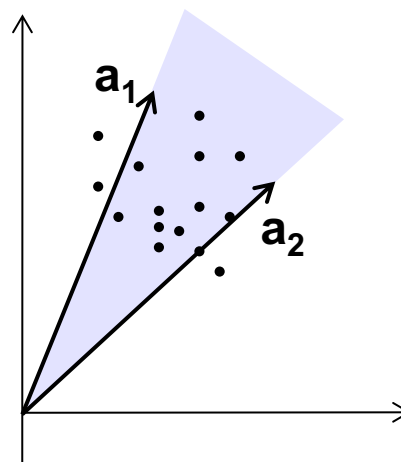
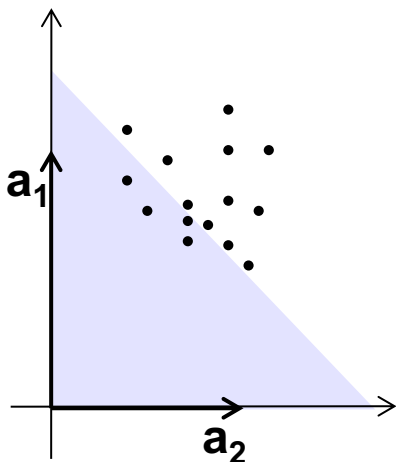


Norah Jones singing “Sunrise”

- Background music “bases” learnt from 5-seconds of music-only segments within the song
- Guitar bases/voice learnt from the rest of the song

What is a good dictionary?

- It should be kept relatively small
 - Overcomplete dictionaries have non-unique solutions without regularization
 - To reduce computational complexity
- It should be capable of accurately describing the source, and generalize well to unseen data
- It should be discriminative: sources cannot be well represented using a dictionary of another source






What is a good dictionary?

- It should be kept relatively small
 - Overcomplete dictionaries have non-unique solutions without regularization
 - To reduce computational complexity
- It should be capable of accurately describing the source, and generalize well to unseen data
- It should be discriminative: sources cannot be well represented using a dictionary of another source
 - In the field of *sparse representations*, this is stated as: a source should be sparse in one dictionary and dense in the other

Dictionary learning

- Non-negative matrix factorization (NMF)
- Clustering
 - k-means clustering
 - Hierarchical clustering
- Approaches used in the field *sparse representations* and *compressed sensing* (CS)
 - Attempt to find dictionaries which *sparse* represent sources
 - Generally no non-negativity constraints
 - Some exceptions, e.g. non-negative K-SVD

Dictionary learning

- Non-negative matrix factorization (NMF)
- Clustering
 - k-means clustering
 - Hierarchical clustering
- Approaches used in the field *sparse representations* and *compressed sensing* (CS)
 - Attempt to find dictionaries which *sparse* represent sources
 - Generally no non-negativity constraints
 - Some exceptions, e.g. non-negative K-SVD
- Pros and cons
 -  – Dictionaries generalize well to unseen data
 -  – NMF and CS approaches consider additivity: smaller, parts-based dictionaries
 -  – Parts-based representations are often less discriminative between sources






Dictionary sampling

- Approach: directly use samples from the training data
 - Often called “exemplars”
 - This may lead to very large dictionaries!

Dictionary sampling

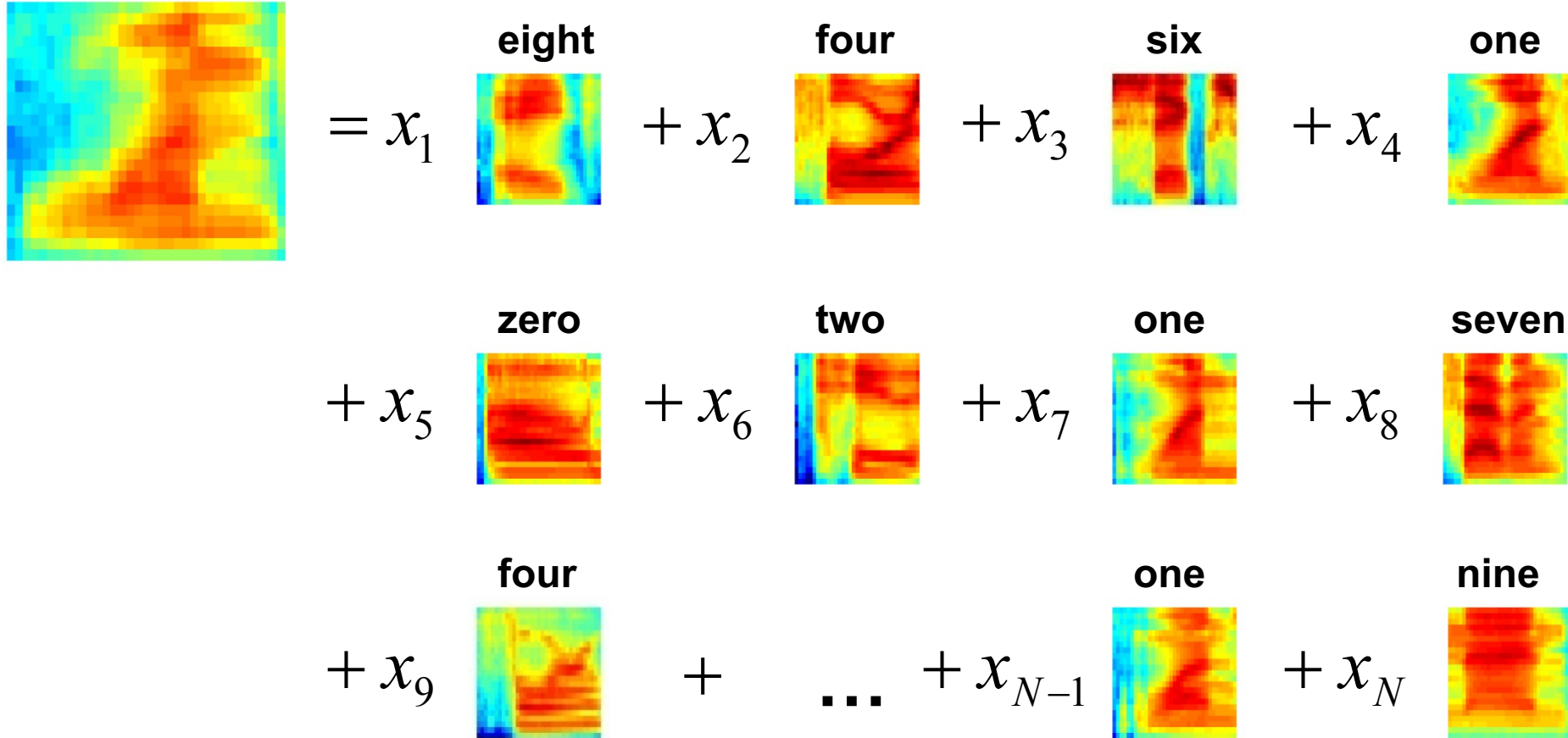
- Approach: directly use samples from the training data
 - Often called “exemplars”
 - This may lead to very large dictionaries!
- Common techniques:
 - random sampling: a random subset of the exemplars
 - pruning: select a subset using some criterion
 - Correlation between exemplars
 - How often an exemplar is activated on development data

Dictionary sampling

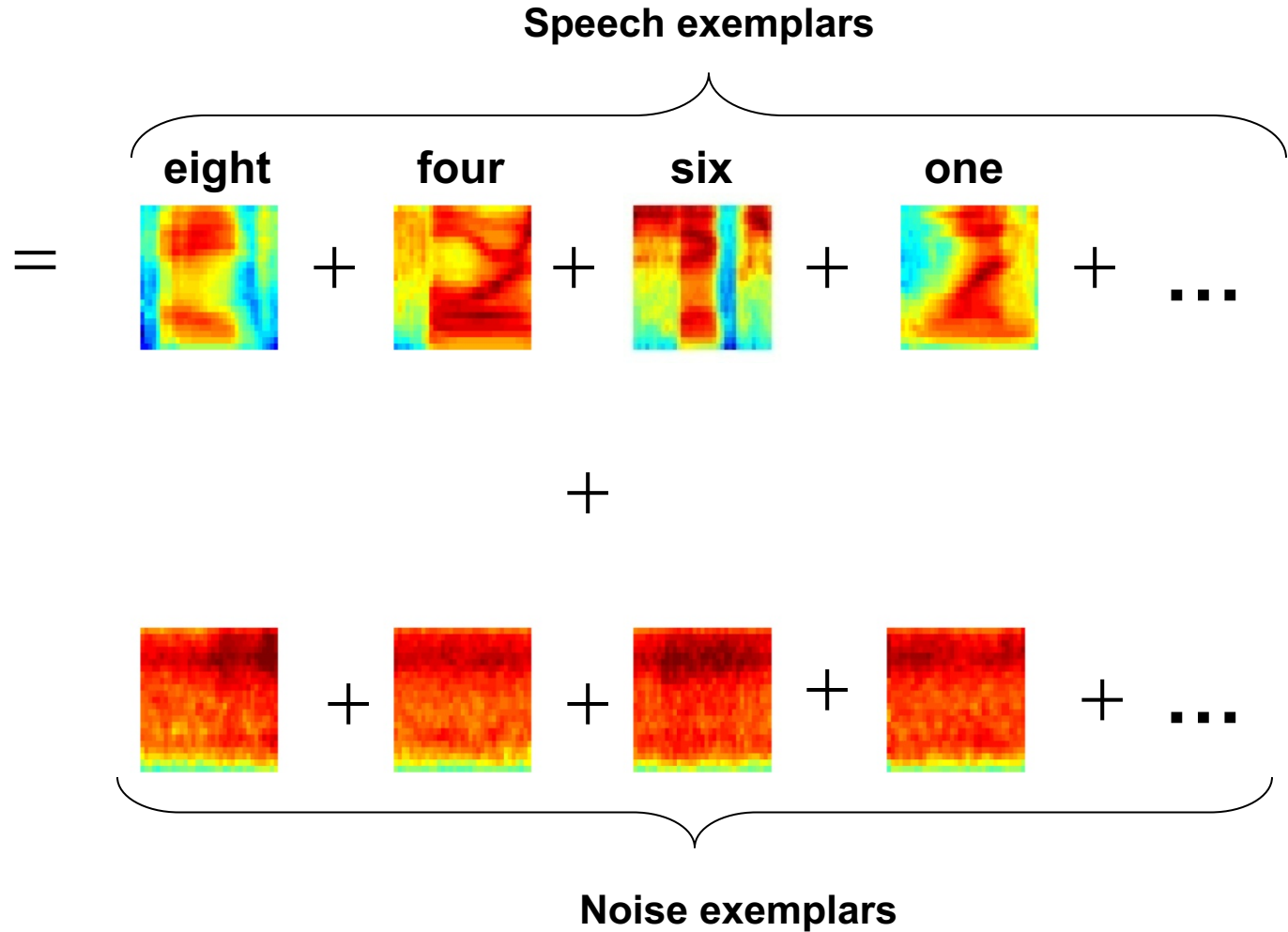
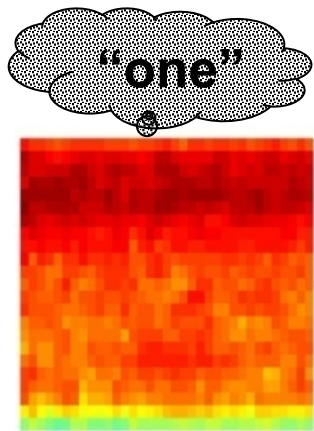
- Approach: directly use samples from the training data
 - Often called “exemplars”
 - This may lead to very large dictionaries!
- Common techniques:
 - random sampling: a random subset of the exemplars
 - pruning: select a subset using some criterion
 - Correlation between exemplars
 - How often an exemplar is activated on development data
- Pros and cons
 -  – Dictionaries do not always generalize well to unseen data
 -  – Does not consider additivity: large dictionaries (but activations are *sparse*)
 -  – Dictionaries are discriminative between sources
 -  – Simpler to use more time-context (many features)
 -  – Requires little tuning

Exemplar-based dictionary

“one”

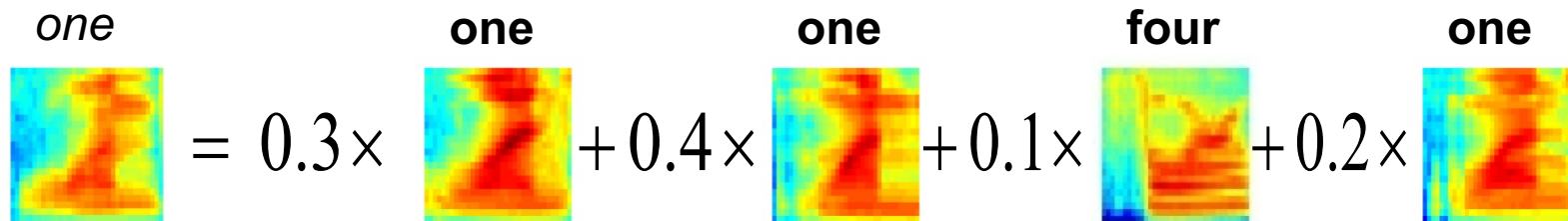


Exemplar-based source separation

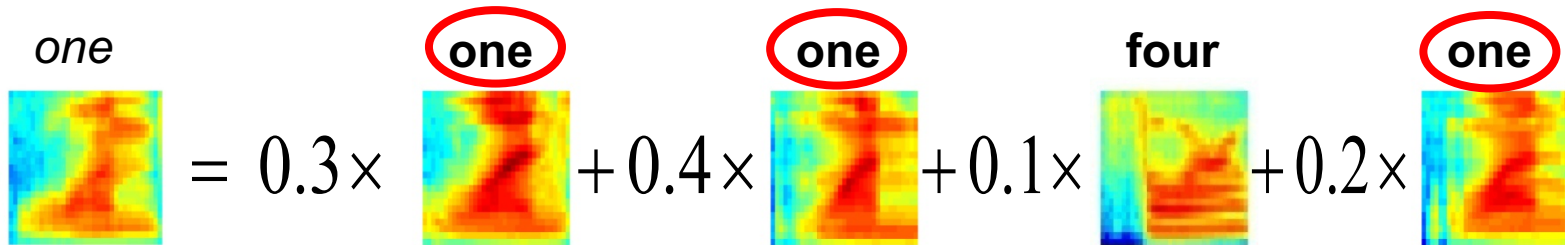


Analyzing the semantics of audio

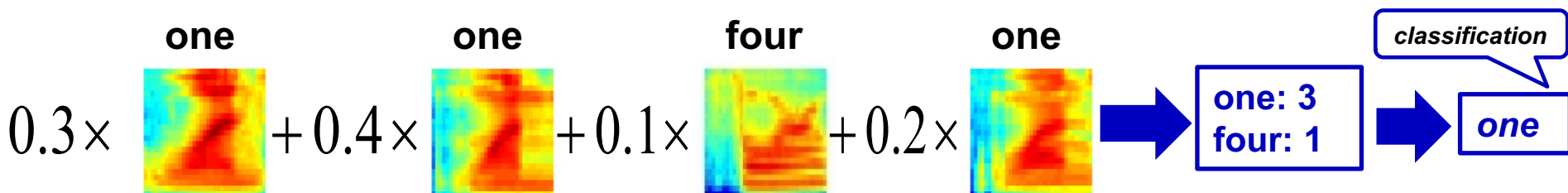
- Supervised dictionary allows using meta information about each atom



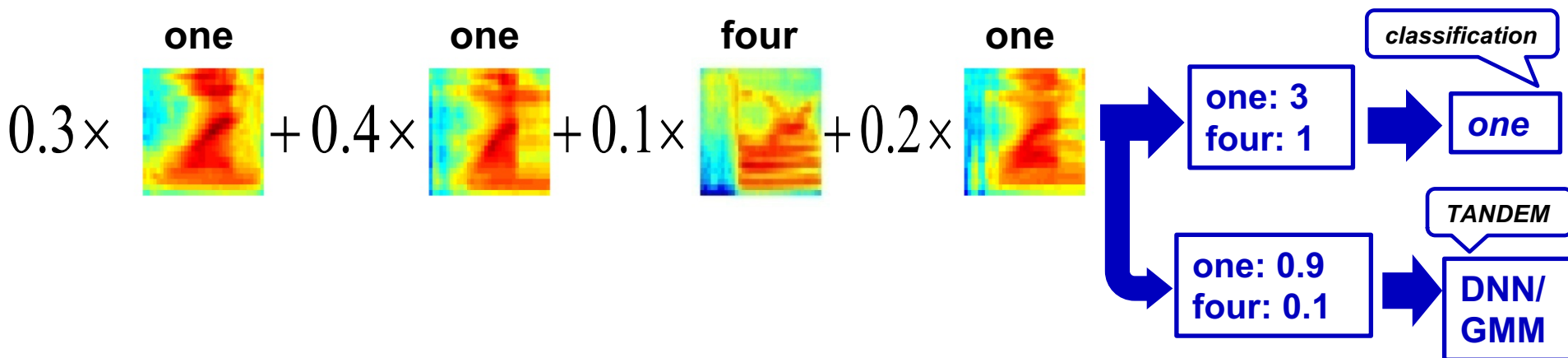
Speech recognition



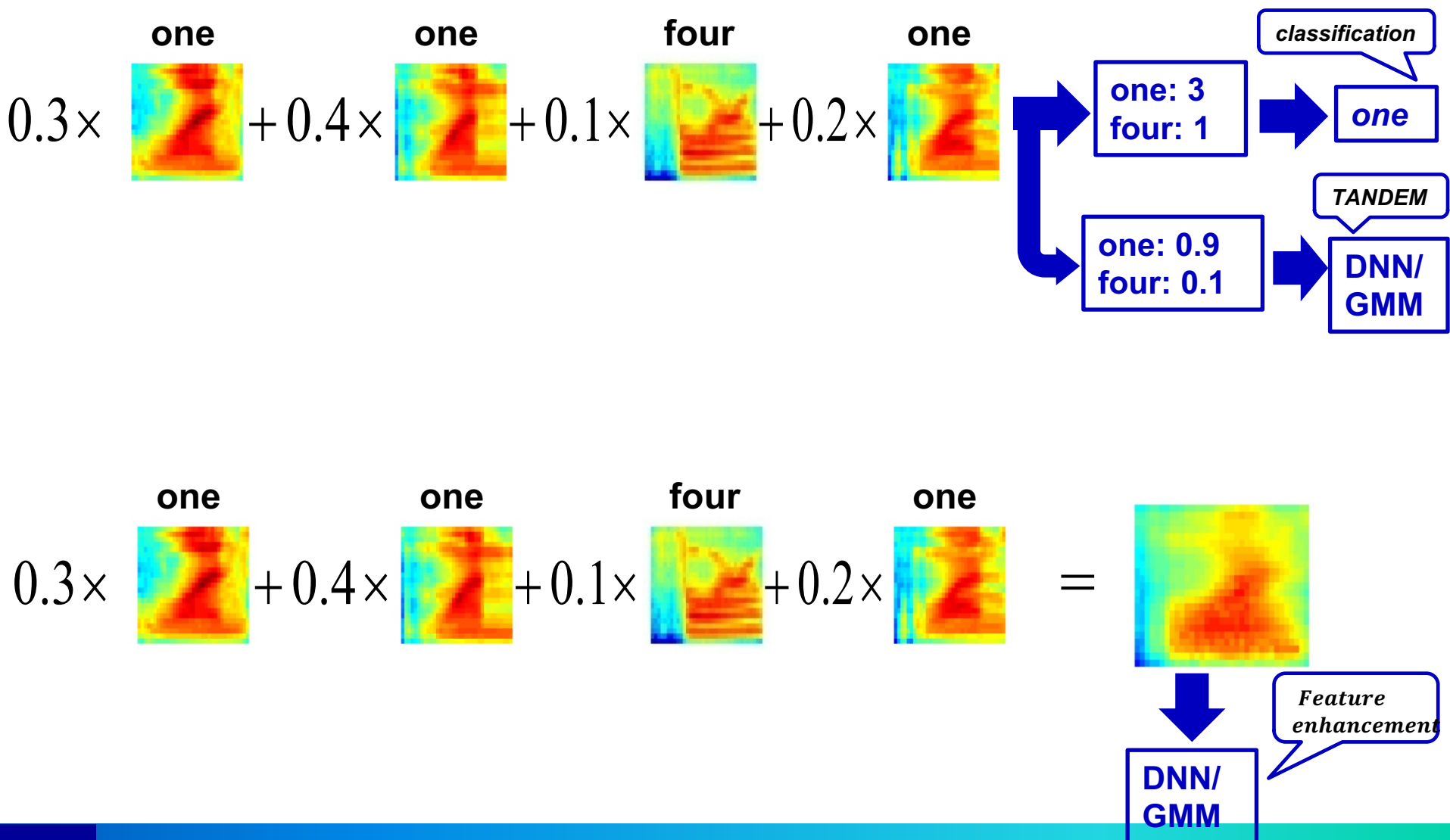
Speech recognition



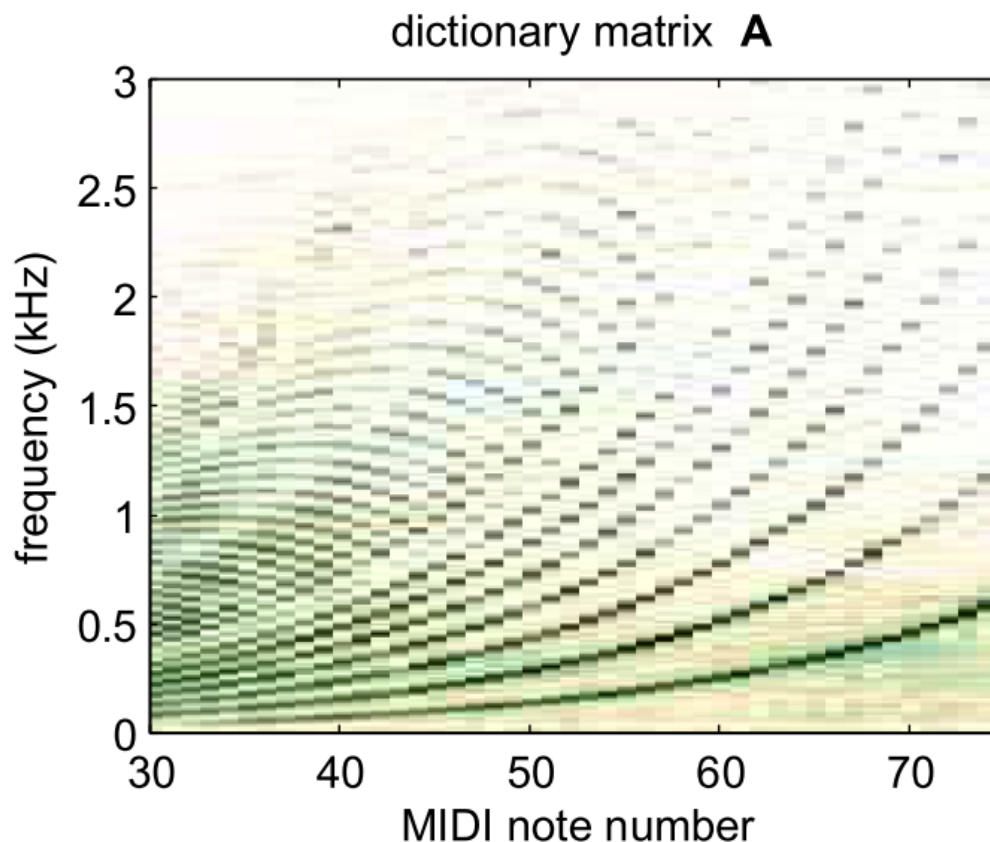
Speech recognition



Speech recognition



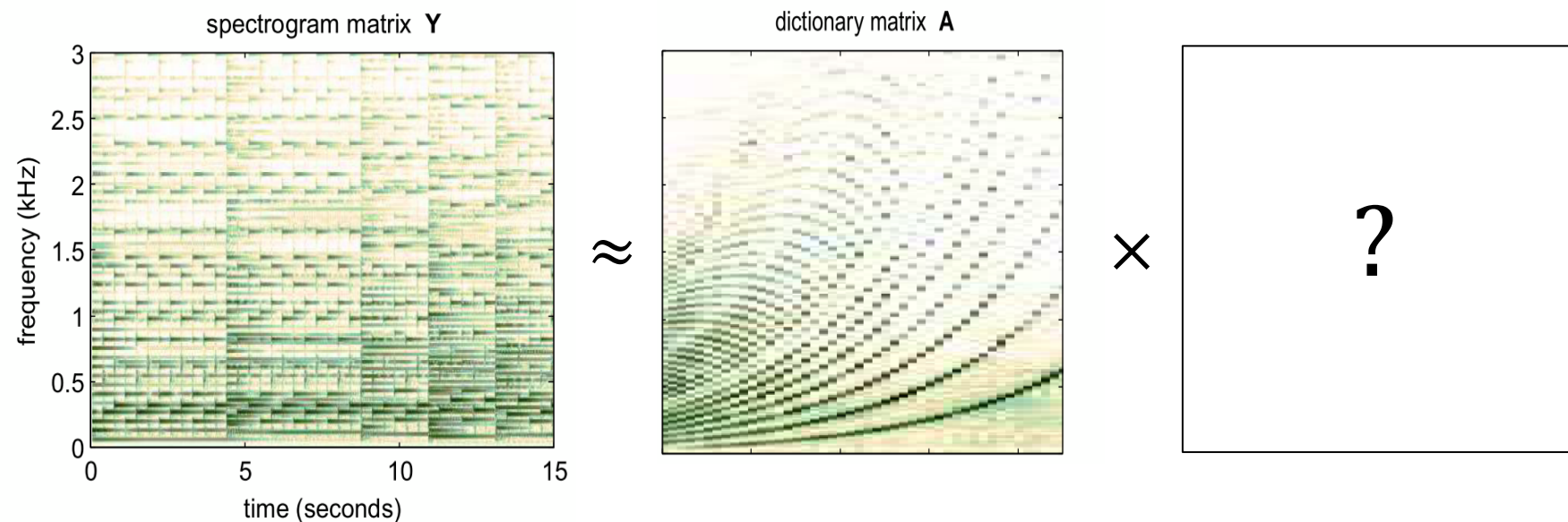
Music example



each atom corresponds to the spectrum of a piano note

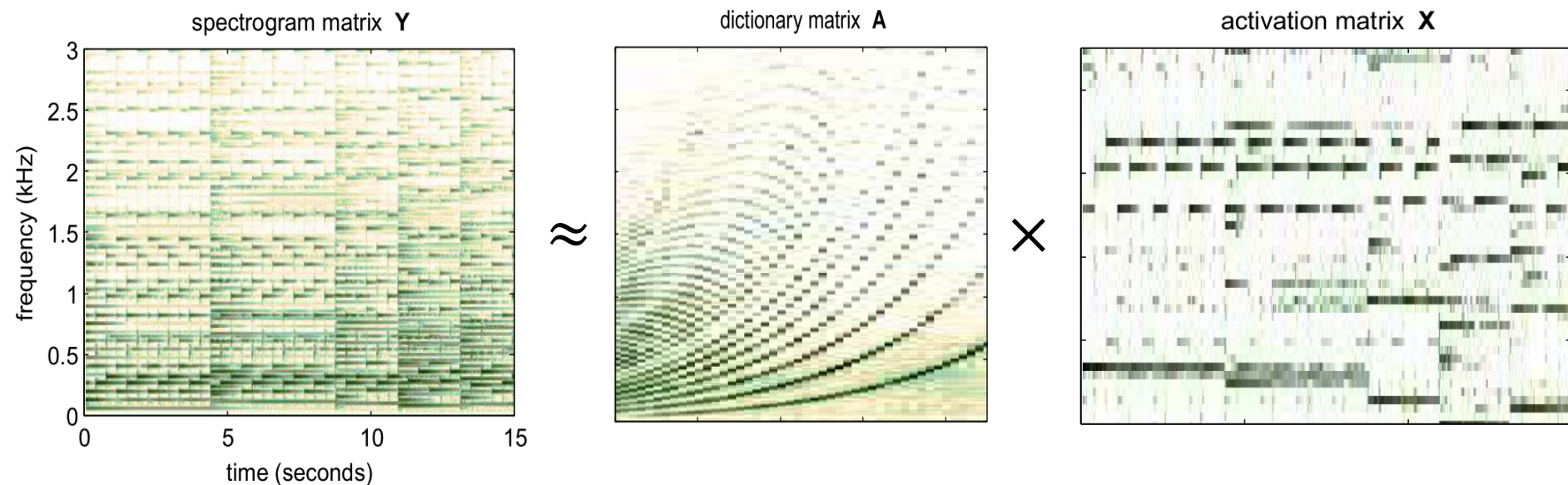
- N. Bertin, R. Badeau, and E. Vincent, "Enforcing harmonicity and smoothness in Bayesian non-negative matrix factorization applied to polyphonic music transcription," IEEE Transactions on Audio, Speech, and Language Processing, vol. 18, no. 3, pp. 538 - 549, 2010.
- T. Heittola, A. Klapuri, and T. Virtanen, "Musical instrument recognition in polyphonic audio using source-filter model for sound separation," in Proceedings of International Conference on Music Information Retrieval, Kobe, Japan, 2009.

Music analysis problem



- N. Bertin, R. Badeau, and E. Vincent, "Enforcing harmonicity and smoothness in Bayesian non-negative matrix factorization applied to polyphonic music transcription," IEEE Transactions on Audio, Speech, and Language Processing, vol. 18, no. 3, pp. 538 - 549, 2010.
- T. Heittola, A. Klapuri, and T. Virtanen, "Musical instrument recognition in polyphonic audio using source-filter model for sound separation," in Proceedings of International Conference on Music Information Retrieval, Kobe, Japan, 2009.

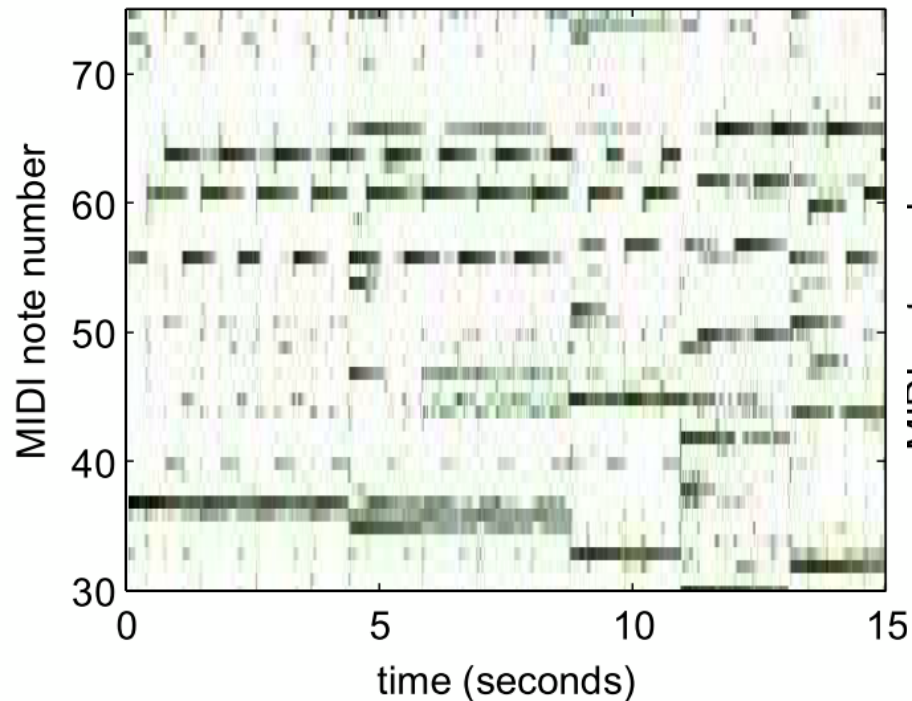
Music analysis problem



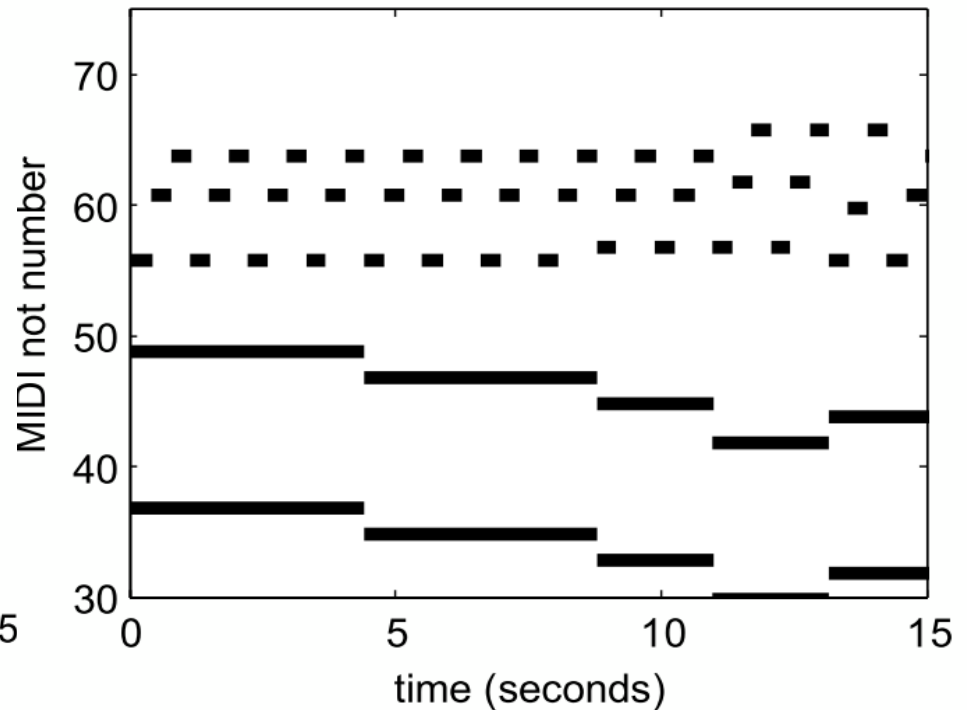
- N. Bertin, R. Badeau, and E. Vincent, "Enforcing harmonicity and smoothness in Bayesian non-negative matrix factorization applied to polyphonic music transcription," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 3, pp. 538 - 549, 2010.
- T. Heittola, A. Klapuri, and T. Virtanen, "Musical instrument recognition in polyphonic audio using source-filter model for sound separation," in *Proceedings of International Conference on Music Information Retrieval*, Kobe, Japan, 2009.

Music analysis problem

activation matrix X



reference activations



Regularization in NMF

- Place additional constraints on the NMF formulation:

$$\mathbf{A}^*, \mathbf{X}^* = \underset{\mathbf{A}, \mathbf{X}}{\operatorname{argmin}} D(\mathbf{Y} || \mathbf{A}\mathbf{X}) + \lambda \Phi(\mathbf{X})$$

- Leading to modified multiplicative updates:

$$\mathbf{X} \leftarrow \mathbf{X} \otimes \frac{\mathbf{A}^\top \frac{\mathbf{Y}}{\mathbf{A}\mathbf{X}}}{\mathbf{A}^\top \mathbf{1} + \lambda \Phi'(\mathbf{X})}$$

- With $\Phi'(\mathbf{X})$ the matrix derivative of $\Phi(\mathbf{X})$ with respect to \mathbf{X}
- This necessitates an l_2 normalization of the columns of \mathbf{A}

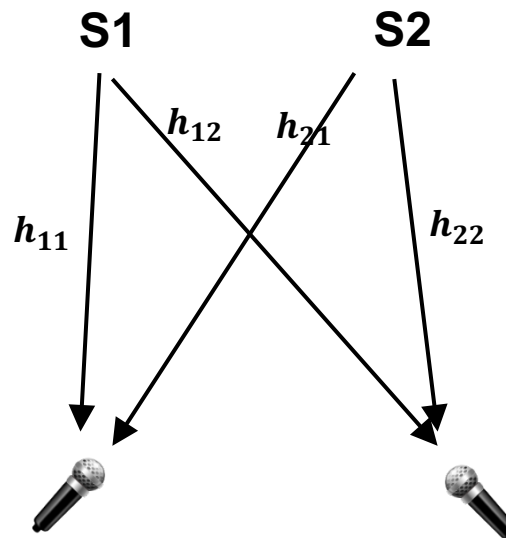
Model alternatives

Regularization in NMF

- A very popular regularizer is *sparsity*: $\Phi(\mathbf{X}) = \|\mathbf{X}\|_1$
- Sparsity regularisation allows decomposition with overcomplete dictionaries
- Other commonly used regularizers:
 - Temporal continuity (Virtanen 2007)
 - Correlation of weights (Wilson et al. 2008)
 - Correlation of spectra (Virtanen & Cemgil 2009)
 - Correlation of components (Wilson & Raj 2010)
 - Hidden Markov Models (Gemmeke et. al. 2013)

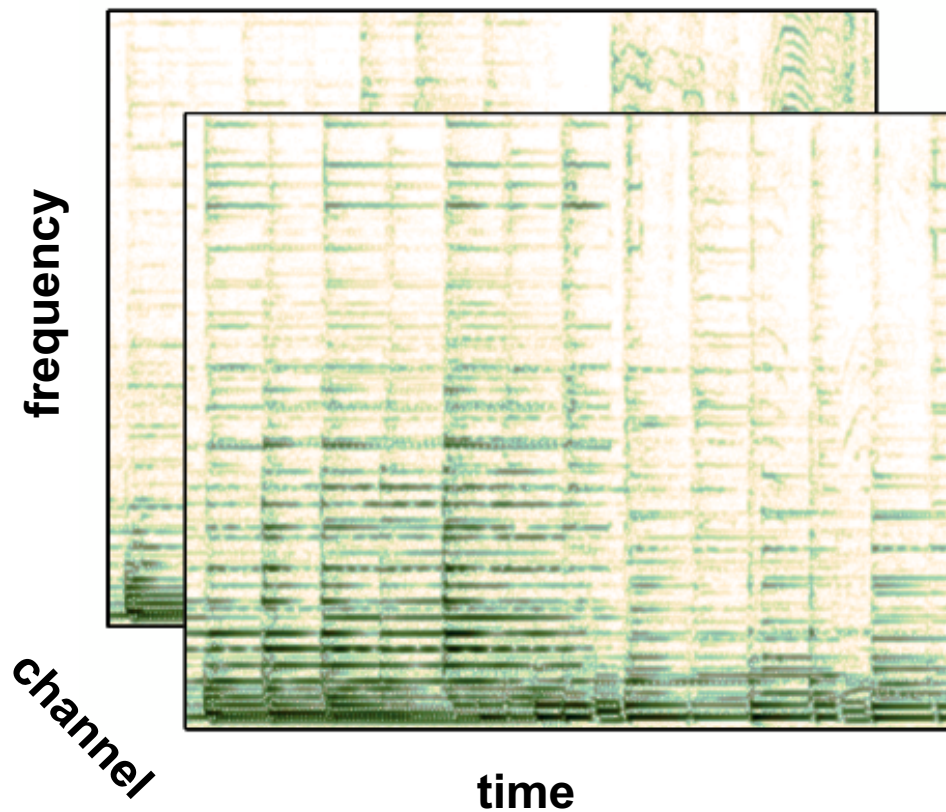
Multichannel audio

- How to use information about differences between channels?
- Phase differences: cannot be modeled with compositional models, require additional modeling components
- Possible to model amplitude differences using compositional models



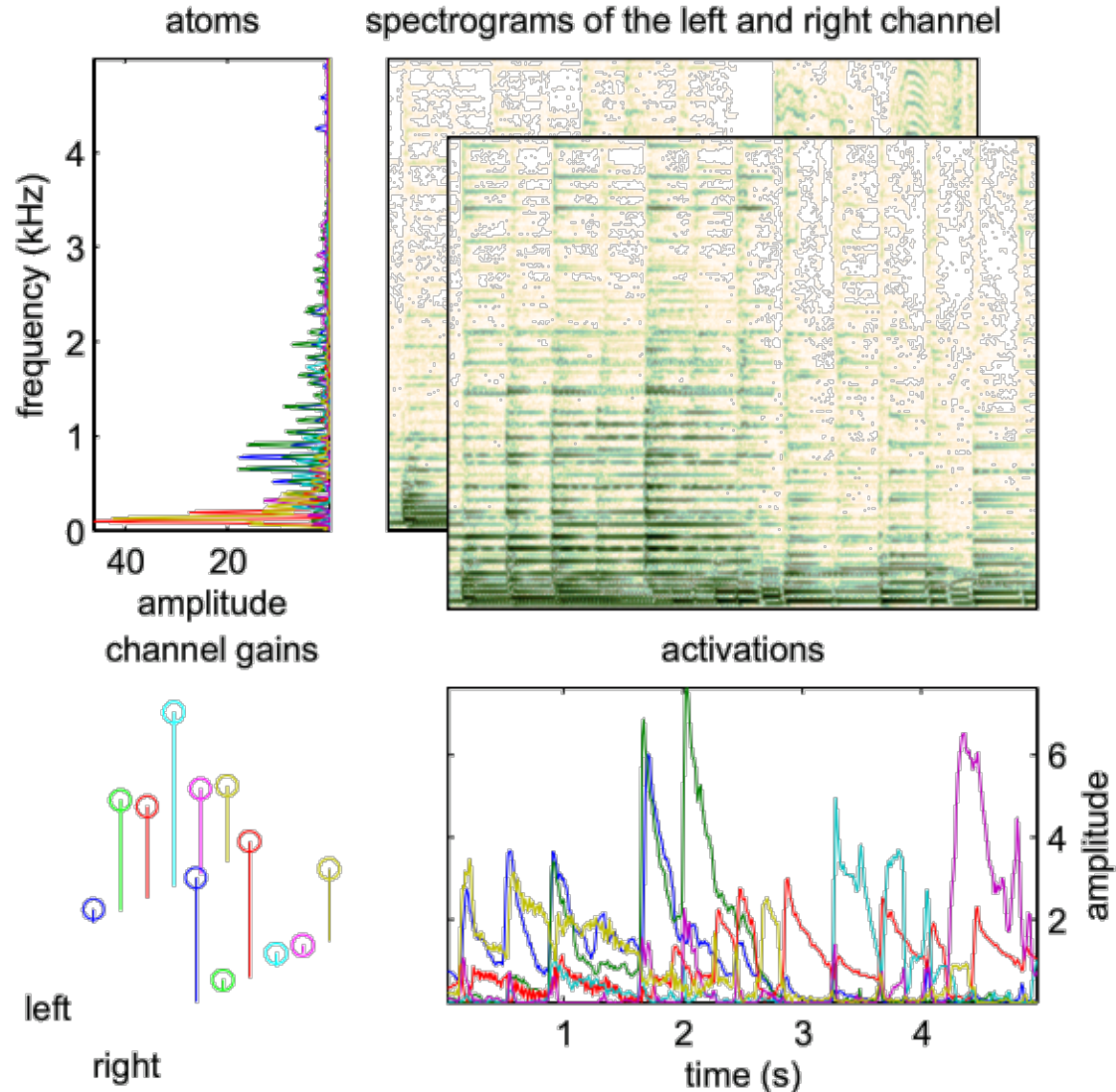
Multichannel audio

- Signal of each channel represented using spectrogram
- Combined into a 3-D *tensor*

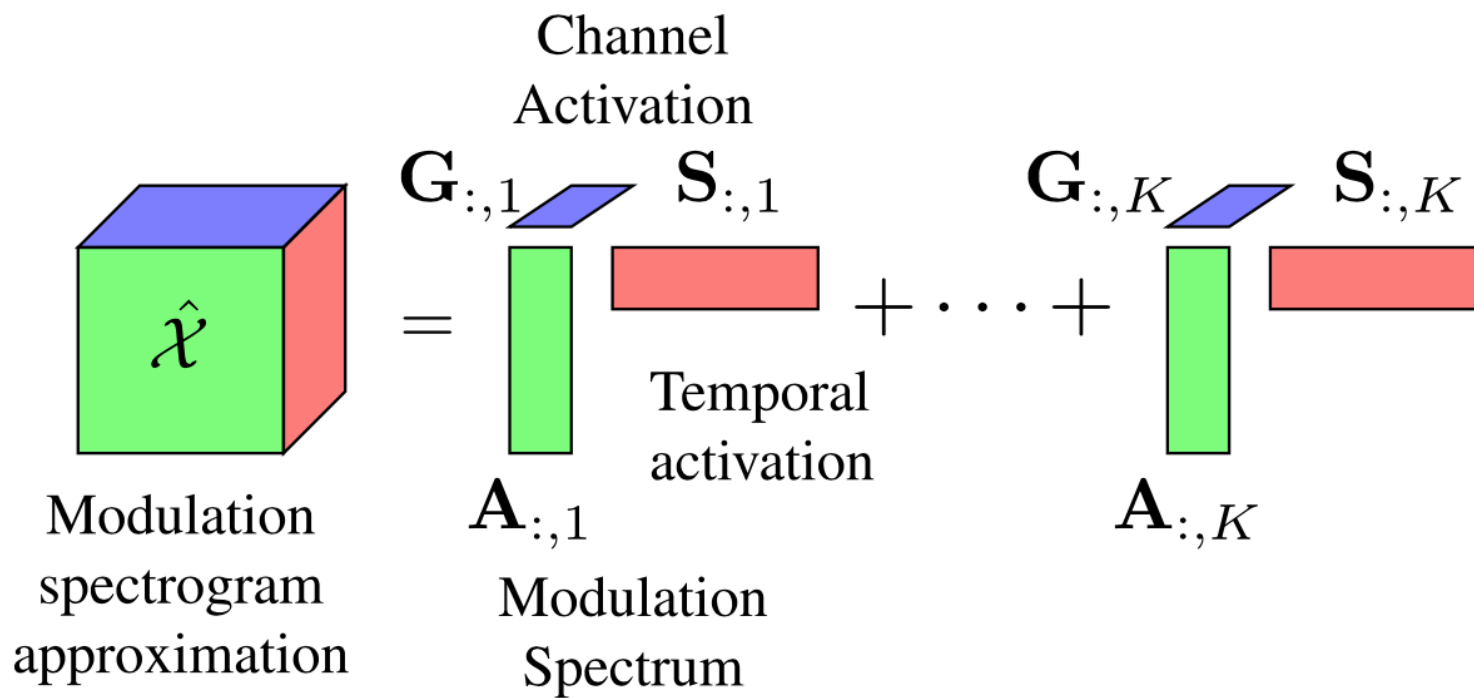


Multichannel audio

- Tensor decomposed into a sum of components
- Each component presented as an outer product of
 1. Spectral atoms
 2. Temporal activations
 3. Channel gains



Tensor factorisation of modulation spectrograms



- T. Barker, T. Virtanen, "Blind separation of audio mixtures through nonnegative tensor factorisation of modulation spectrograms", in IEEE/ACM Transactions on Audio, Speech and Language Processing, Volume 24, Issue 12, December 2016, pp. 2377-2389.

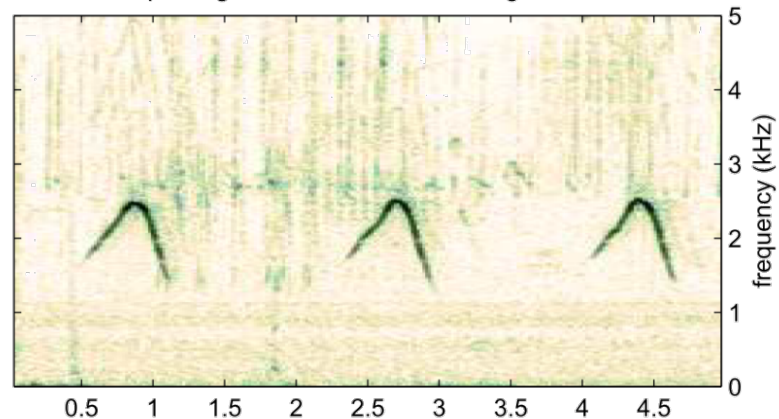
Temporal context

- Sound typically has much spectral and temporal structure
- The basic NMF model treats each frequency and frame as independent from each other
- Modelling contextual information (time-frequency patches) useful

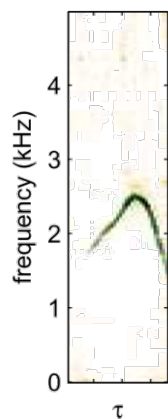
Temporal context: NMF deconvolution

$$\hat{\mathbf{y}}_t = \sum_k \sum_{\tau} \mathbf{a}_{k,\tau} \mathbf{x}_{k,t-\tau}$$

spectrogram matrix \mathbf{Y} with missing frames



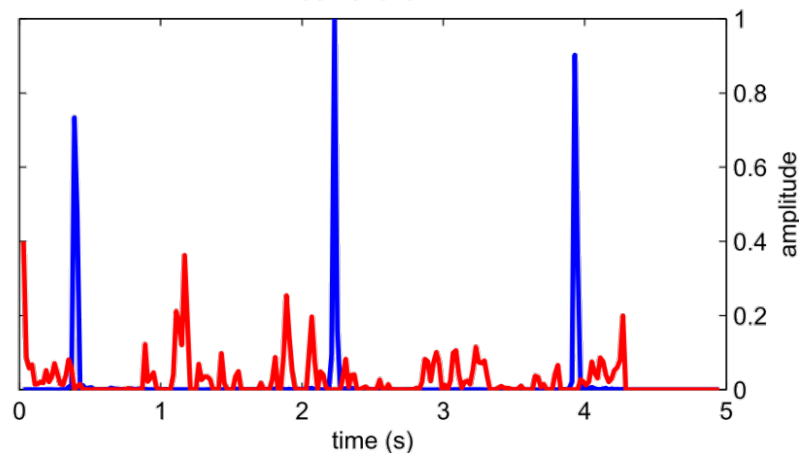
component $\mathbf{a}_{1,\tau}$



component $\mathbf{a}_{2,\tau}$



activations \mathbf{X}

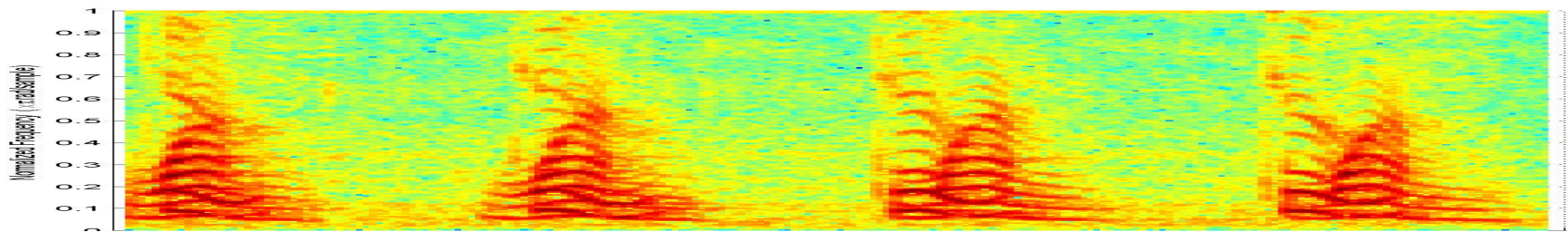


- P. Smaragdis, "Convolutional speech bases and their application to supervised speech separation," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 1, pp. 1 – 12, 2007.
- P. D. O. Grady, "Sparse separation of underdetermined speech mixtures," Ph.D. dissertation, National University of Ireland, Maynooth, 2007.

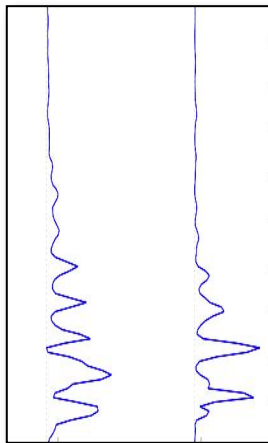
Example

- Two distinct sounds occurring with different repetition rates within a signal

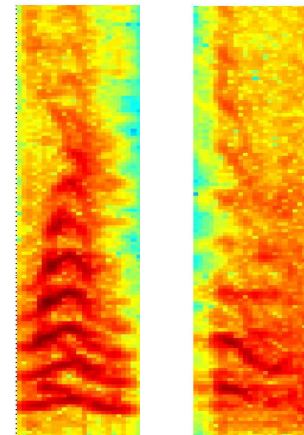
INPUT SPECTROGRAM



NMF dictionary atoms

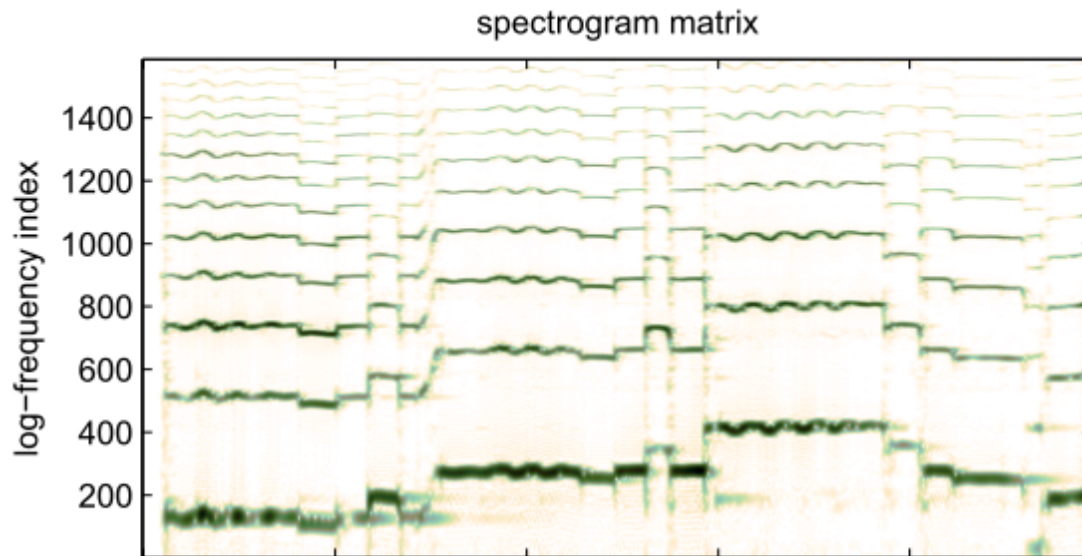


Discovered time-frequency "patch" atoms



Invariant models

- The basic NMF model requires a separate atom to model sounds with different pitch
- Modeling multiple pitches with a single atom?
- On a log-frequency scale, pitch shifting corresponds to translation

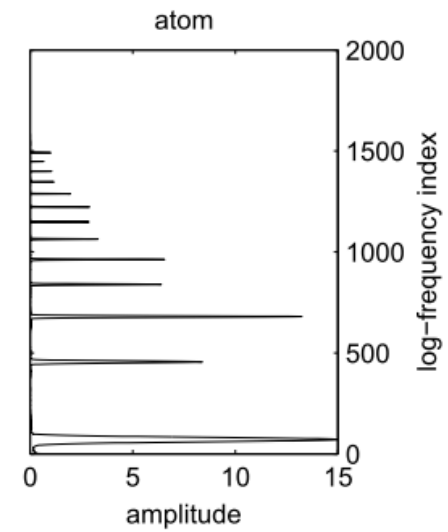
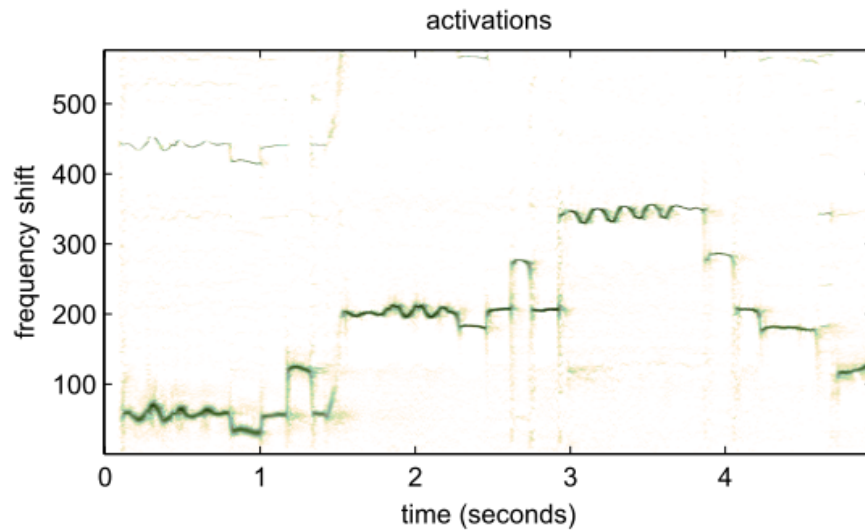
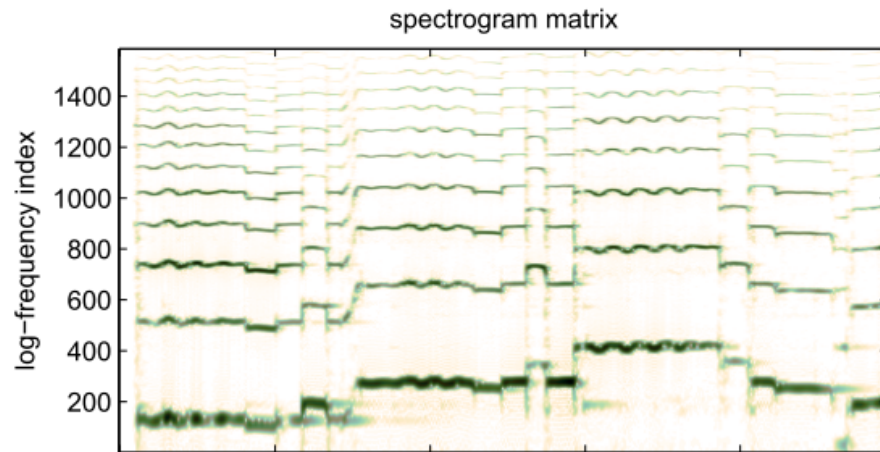


NMF deconvolution in frequency

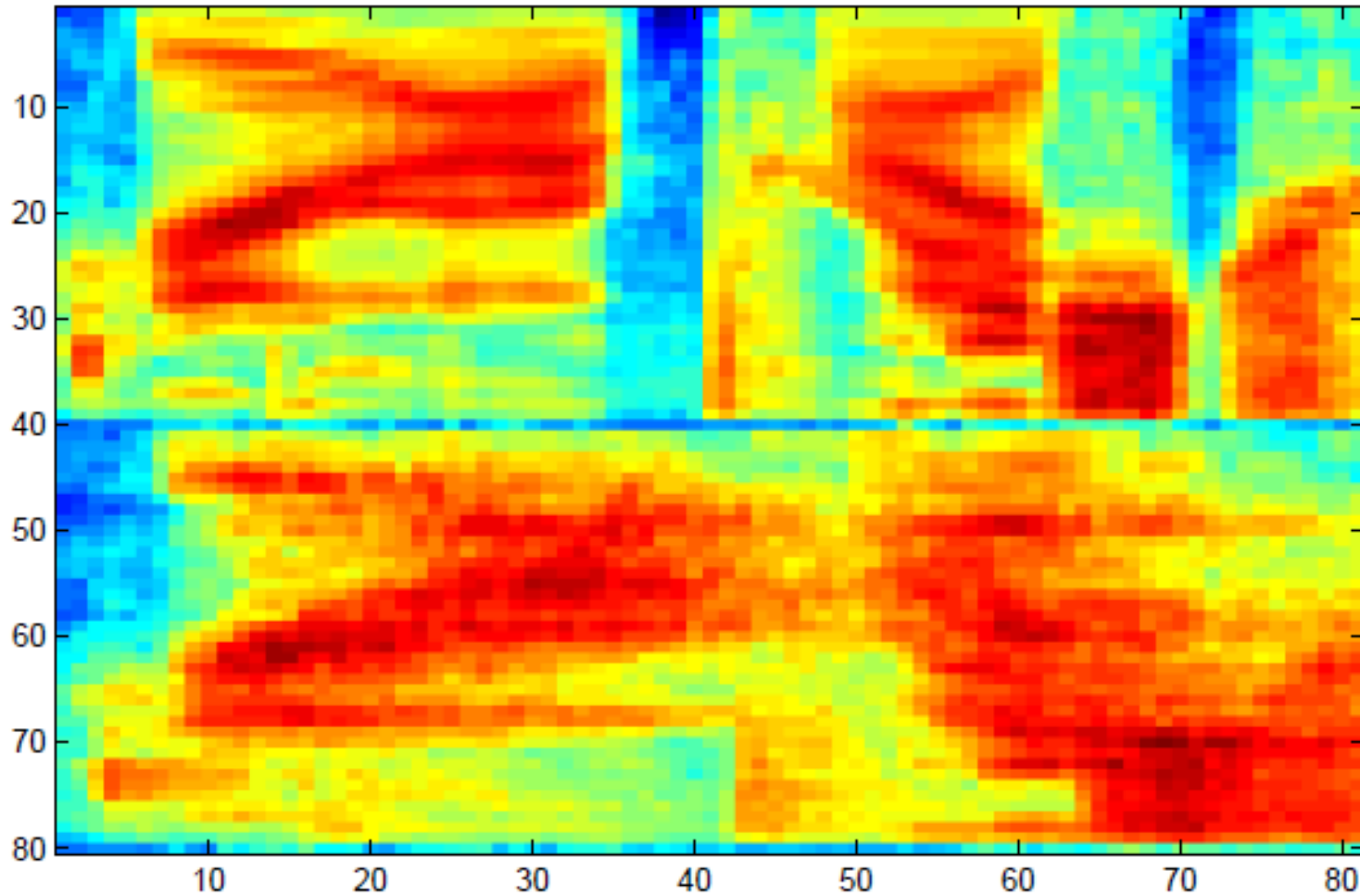
- Basic idea:
 - Use logarithmic frequency axis
 - Allow translating atoms a in frequency (convolution)
 - Estimate activation x_τ for each amount τ of translation

$$\hat{y}_{f,t} = \sum_{k=1}^K \sum_{\tau \in \mathcal{L}} a_{f+\tau,k} x_{k,\tau}$$

NMF deconvolution

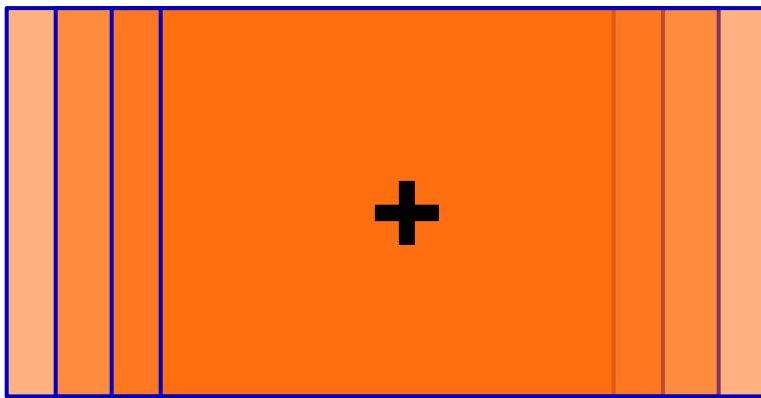


Dereverberation

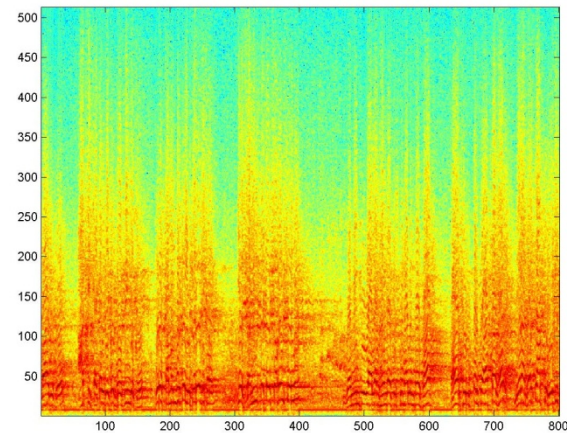


Dereverberation

- A convolutional model of reverberation:
 - The spectrogram of the reverberated signal is a sum of the spectrogram of the clean signal and several shifted and scaled versions of itself

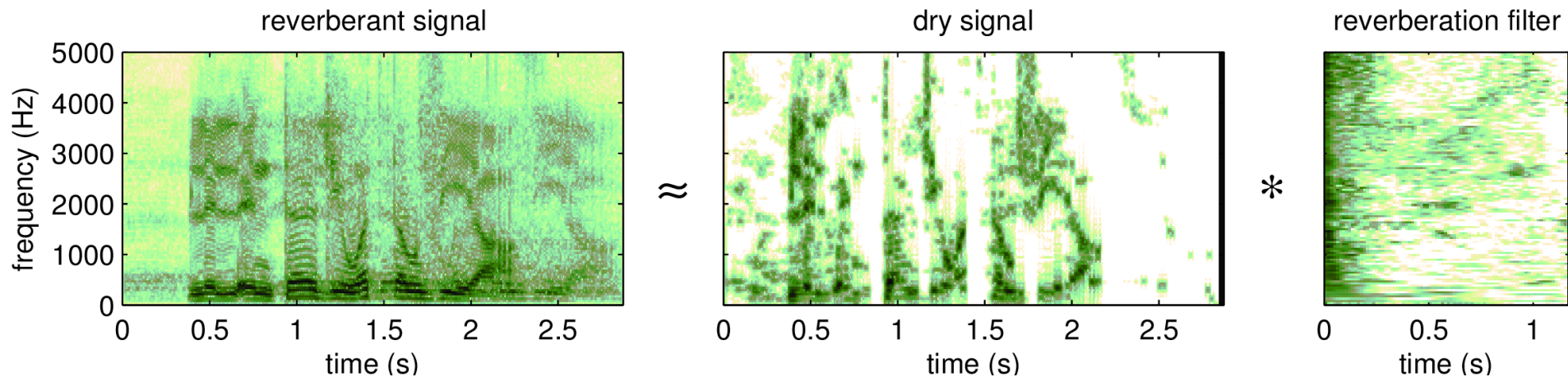


=



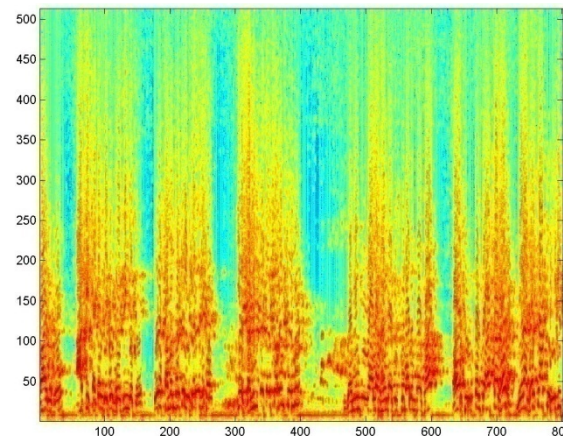
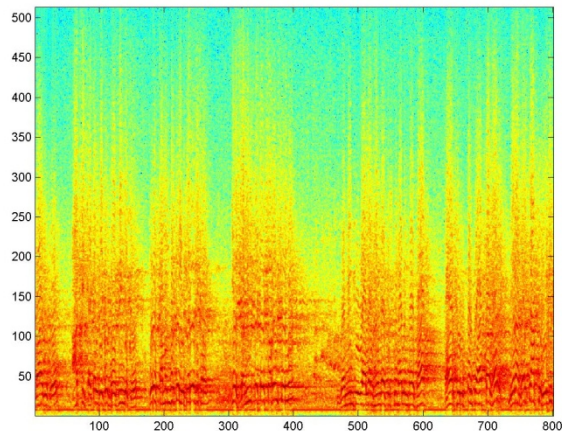
Dereverberation

- A convolutional model of reverberation:
 - The spectrogram of the reverberated signal is a sum of the spectrogram of the clean signal and several shifted and scaled versions of itself
 - A convolution of the spectrogram and a room response



Dereverberation

- A convolutional model of reverberation:
 - The spectrogram of the reverberated signal is a sum of the spectrogram of the clean signal and several shifted and scaled versions of itself
 - A convolution of the spectrogram and a room response
 - Factorial model: $\mathbf{Y}=\mathbf{S}\mathbf{H}$, with \mathbf{Y} the reverberated spectrum, \mathbf{S} the dry spectrum, and \mathbf{H} the reverberation filter
 - Sparsity must be enforced on the filter



Comparison to DNNs

- DNNs are discriminative models
 - Ideal for classification
- Compositional models are generative
 - Can explain the properties of the data better

Comparison to DNNs

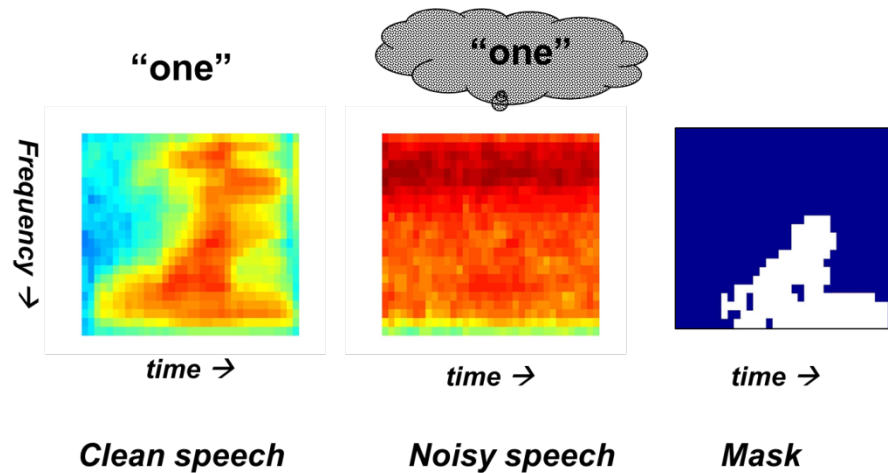
- DNNs give state of the art results in many application areas
 - Applicable also on many mentioned problems (source separation, robust recognition)
- Given large amounts of training data, DNNs typically outperform compositional models
- Compositional models can be used with small amounts of training data
 - Exemplar-based dictionaries obtained from few examples
- Compositional models enable unsupervised processing that does not require training data

Missing data

- Missing data occurs in many applications
 - Packet or frame drops
 - Signal clipping
 - Audio corrupted at specific frequencies

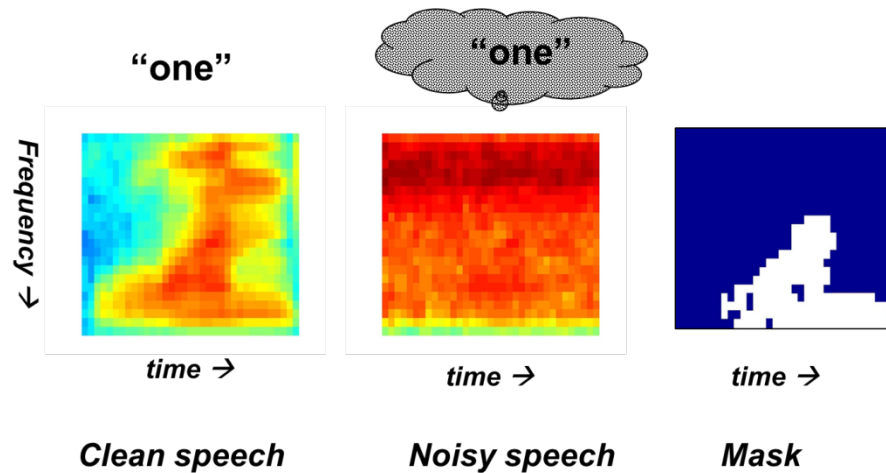
Missing data

- Missing data occurs in many applications
 - Packet or frame drops
 - Signal clipping
 - Audio corrupted at specified frequencies



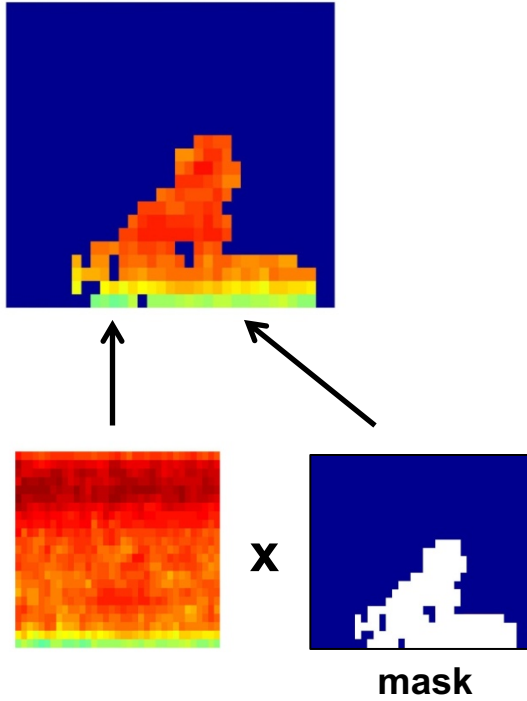
Missing data

- Missing data occurs in many applications
 - Packet or frame drops
 - Signal clipping
 - Audio corrupted at specified frequencies

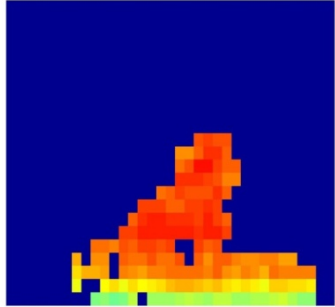


- Original audio: $\mathbf{s} \approx \hat{\mathbf{s}} = \mathbf{A}\mathbf{x}$
- Missing data: $\mathbf{M}\mathbf{y} \approx \mathbf{M}\hat{\mathbf{s}} = \mathbf{M}\mathbf{A}\mathbf{x}$

Clean speech reconstruction

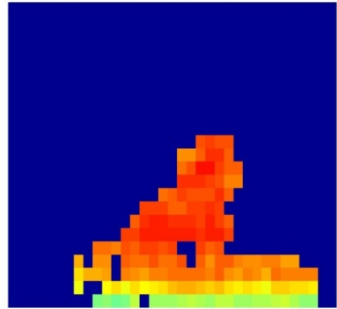


Clean speech reconstruction



$$= 0.3 \times \begin{matrix} \text{one} \\ \text{spectrogram} \end{matrix} + 0.4 \times \begin{matrix} \text{one} \\ \text{spectrogram} \end{matrix} + 0.1 \times \begin{matrix} \text{four} \\ \text{spectrogram} \end{matrix} + 0.2 \times \begin{matrix} \text{one} \\ \text{spectrogram} \end{matrix}$$

Clean speech reconstruction

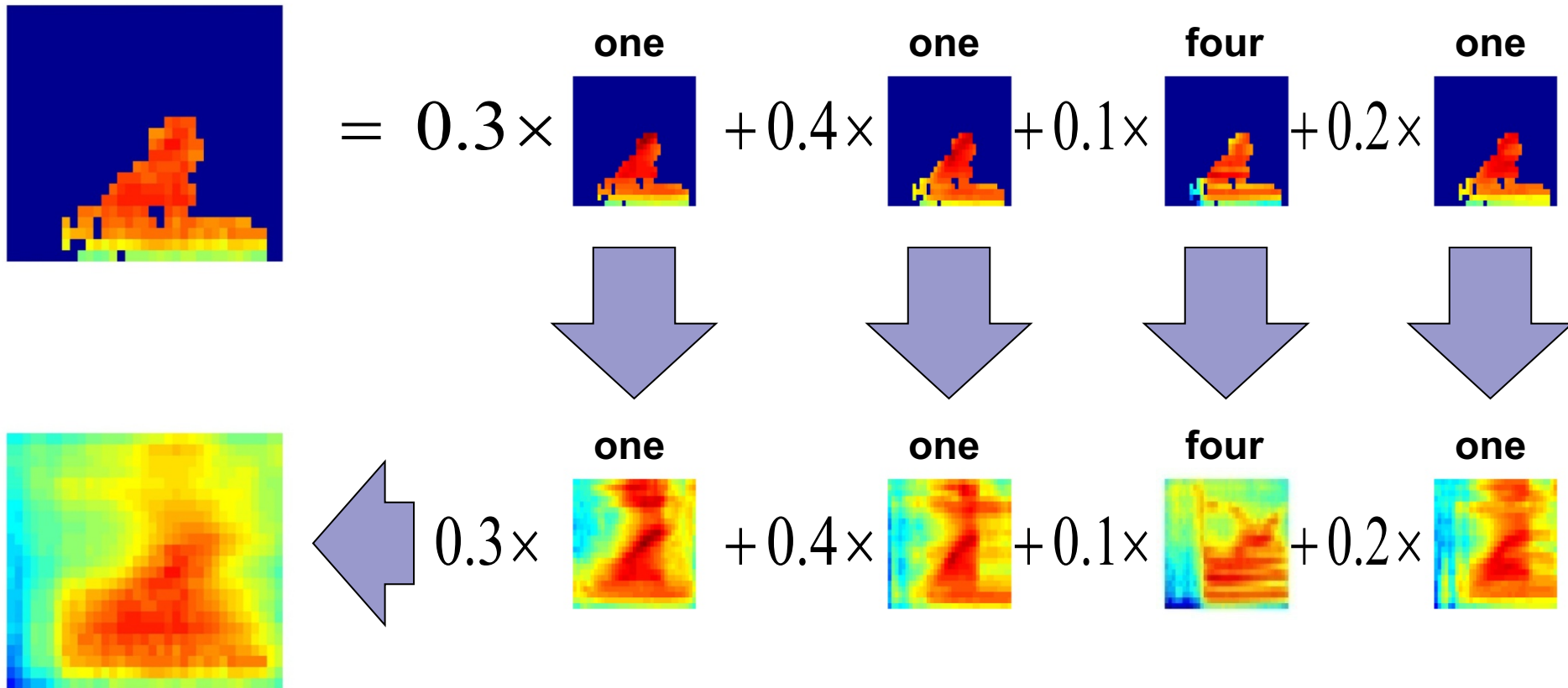


$$= 0.3 \times \begin{matrix} \text{one} \\ \text{[Spectrogram of 'one']} \end{matrix} + 0.4 \times \begin{matrix} \text{one} \\ \text{[Spectrogram of 'one']} \end{matrix} + 0.1 \times \begin{matrix} \text{four} \\ \text{[Spectrogram of 'four']} \end{matrix} + 0.2 \times \begin{matrix} \text{one} \\ \text{[Spectrogram of 'one']} \end{matrix}$$


↓ ↓ ↓ ↓

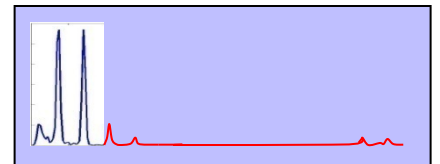
$$0.3 \times \begin{matrix} \text{one} \\ \text{[Spectrogram of 'one']} \end{matrix} + 0.4 \times \begin{matrix} \text{one} \\ \text{[Spectrogram of 'one']} \end{matrix} + 0.1 \times \begin{matrix} \text{four} \\ \text{[Spectrogram of 'four']} \end{matrix} + 0.2 \times \begin{matrix} \text{one} \\ \text{[Spectrogram of 'one']} \end{matrix}$$

Clean speech reconstruction



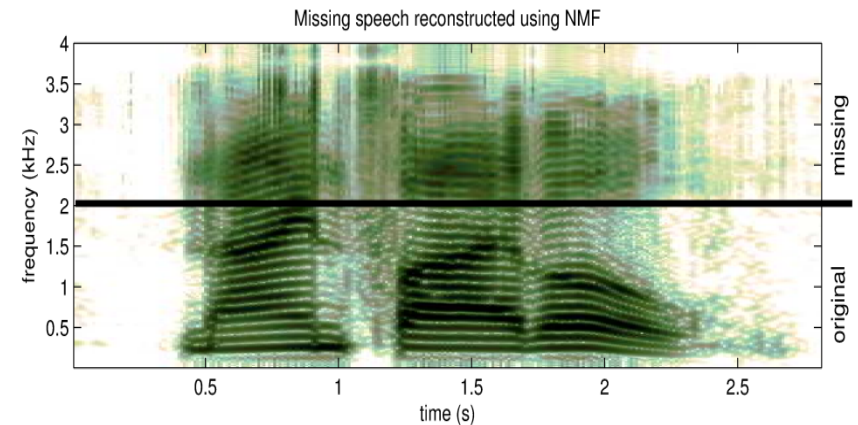
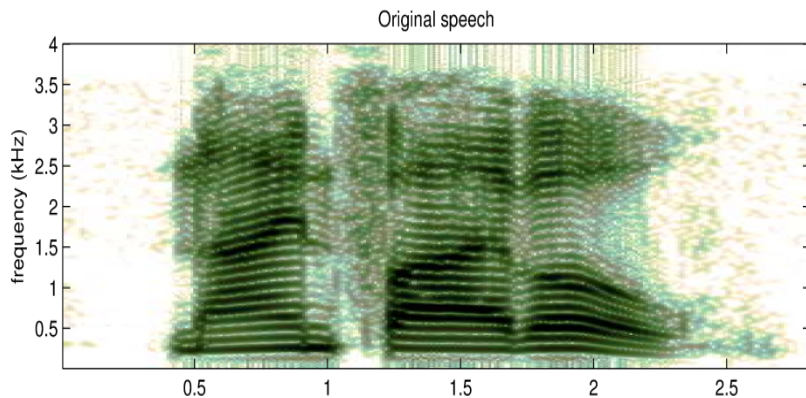
Bandwidth extension

- Problem: A given speech signal only has frequencies in the 300Hz-3.5Khz range
 - Telephone quality speech 
- Goal: restore the missing frequencies
- Assumptions:
 - We have full-bandwidth training data
 - Training data is representative
 - We know which frequencies are missing



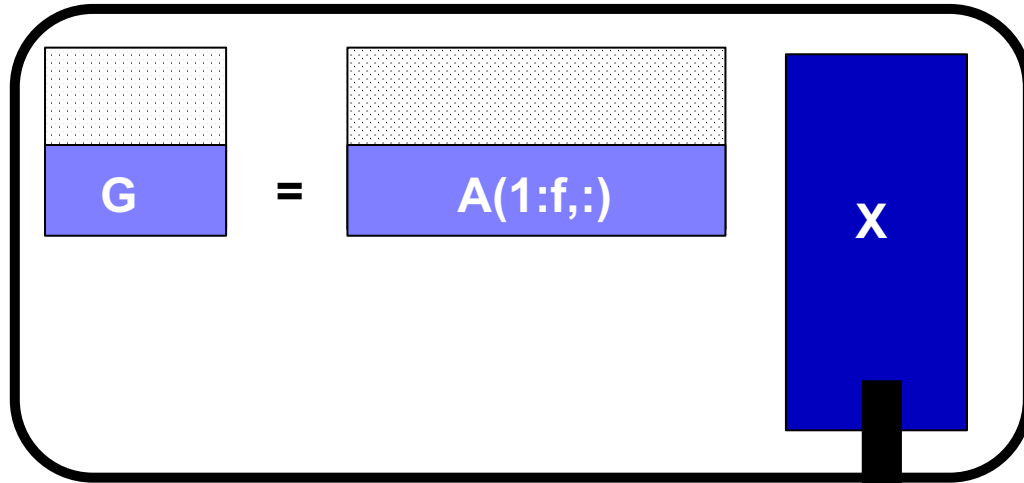
Bandwidth extension

- Almost trivial use of compositional models (using Matlab notation):
 - Step 1: Create a full-bandwidth dictionary \mathbf{A}
 - Step 2: Given the limited bandwidth observation \mathbf{G} , in which only the frequency bands $1 \dots f$ are retained, use a bandwidth-limited $\mathbf{A}(1:f,:)$ dictionary to obtain activations \mathbf{X}
 - Step 3: Reconstruct the full-bandwidth estimate \mathbf{Y} using the full bandwidth dictionary \mathbf{A}

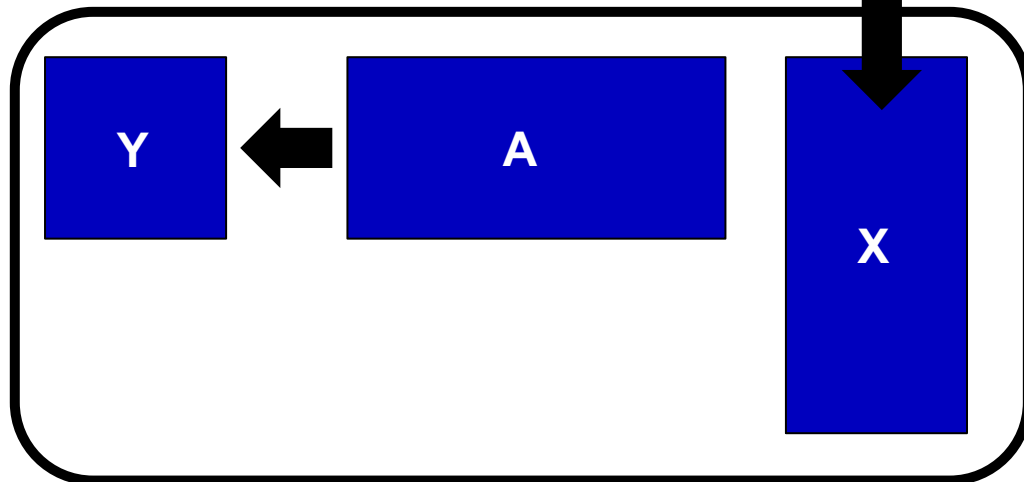


Visual representation

Step 1

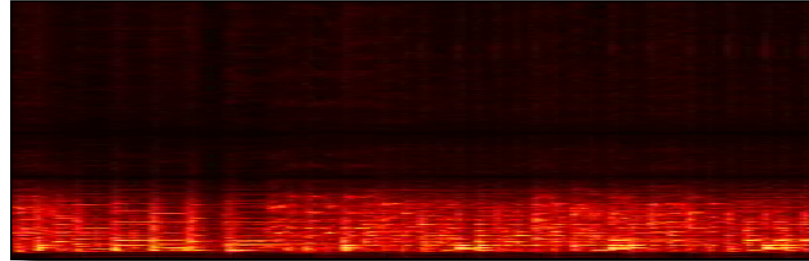


Step 2

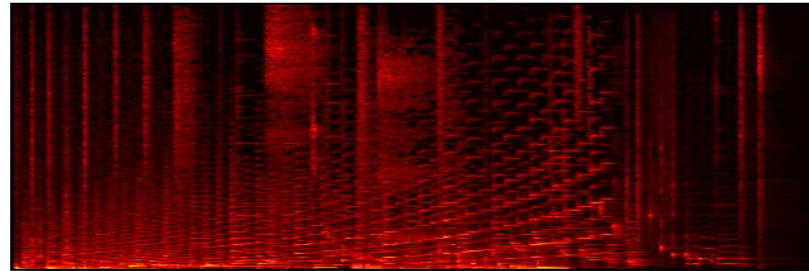


Audio example

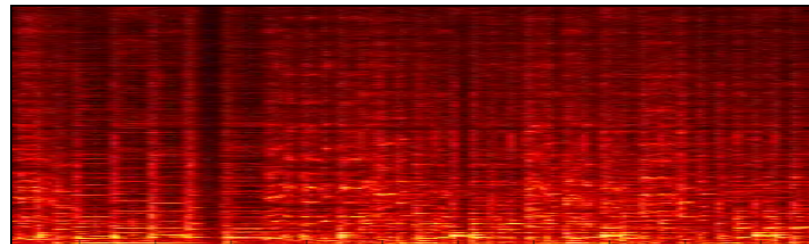
Reduced BW data



Training material



Bandwidth expanded version



Getting started

■ Literature

- Check the references 😊
- Tutorial article: T. Virtanen, J. F. Gemmeke, B. Raj, and P. Smaragdis. [Compositional Models for Audio Processing](#). IEEE Signal Processing Magazine, March 2015

■ Matlab code

- Supervised NMF-based SS: <http://www.cs.tut.fi/~tuomasv/software.html>
- SS, recognition and imputation: <http://www.amadana.nl/software>
- FASST, evalation, etc: <http://www.loria.fr/~evincent/soft.html>
- NMFlab: <http://www.bsp.brain.riken.jp/ICALAB/nmflab.html>
- PLCA: <http://www.cs.illinois.edu/~paris/pubs/>

■ C++ code

- openBliSSART: <http://openblissart.github.io/openBliSSART/>

Summary

- Realistic sounds consist of components that combine purely additively
- Composition models are purely additive models that are powerful in modeling sound mixtures
- Good models for spectral representations of sound
- Applications in source separation and audio processing

The end...

