# A VARIATIONAL EM ALGORITHM FOR LEARNING EIGENVOICE PARAMETERS IN MIXED SIGNALS

## Ron J. Weiss and Daniel P. W. Ellis

LabROSA · Dept of Electrical Engineering · Columbia University, New York, USA

{ronw,dpwe}@ee.columbia.edu

## 1. Summary

- Model-based monaural speech separation where the precise source characteristics are not known a priori
- Extend original adaptation algorithm from Weiss and Ellis (2008) to adapt Gaussian covariances as well as means
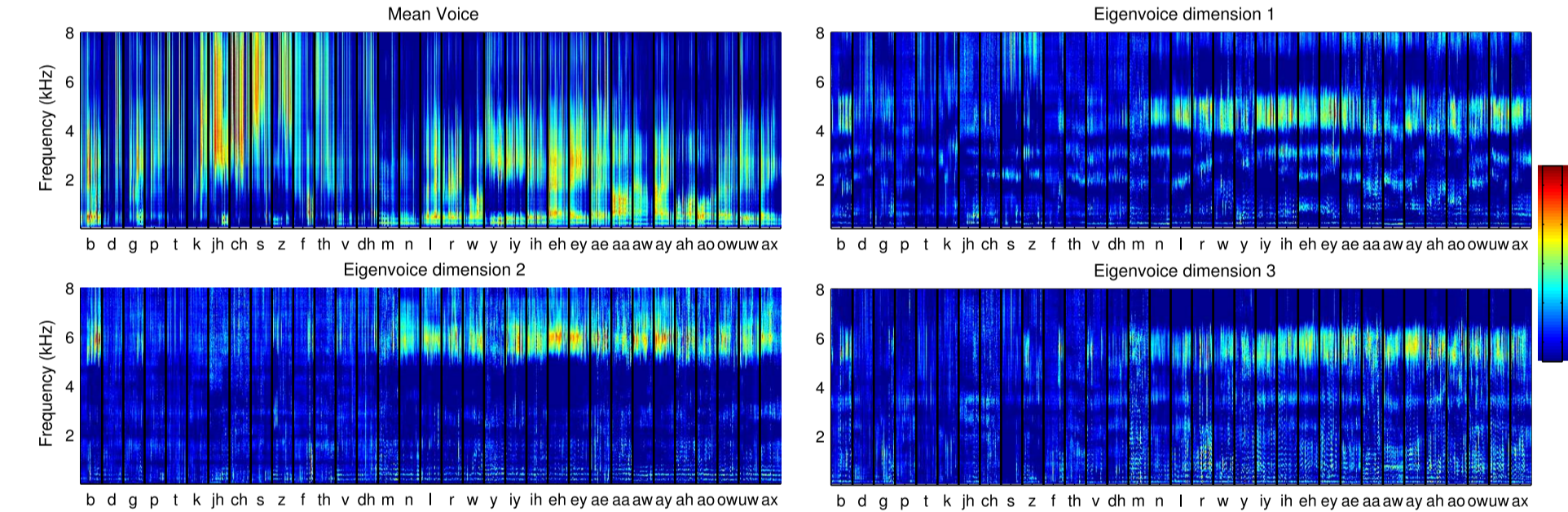- Derive a variational EM algorithm to speed up adaptation

## 2. Mixed signal model

- Model log power spectra of source signals using hidden Markov model (HMM):

$$P(x_i(1..T), s_i(1..T)) = \prod_t P(s_i(t) \mid s_i(t-1)) P(x_i(t) \mid s_i(t))$$

- Represent speaker-dependent model as linear combination of eigenvoice bases (Kuhn et al., 2000):

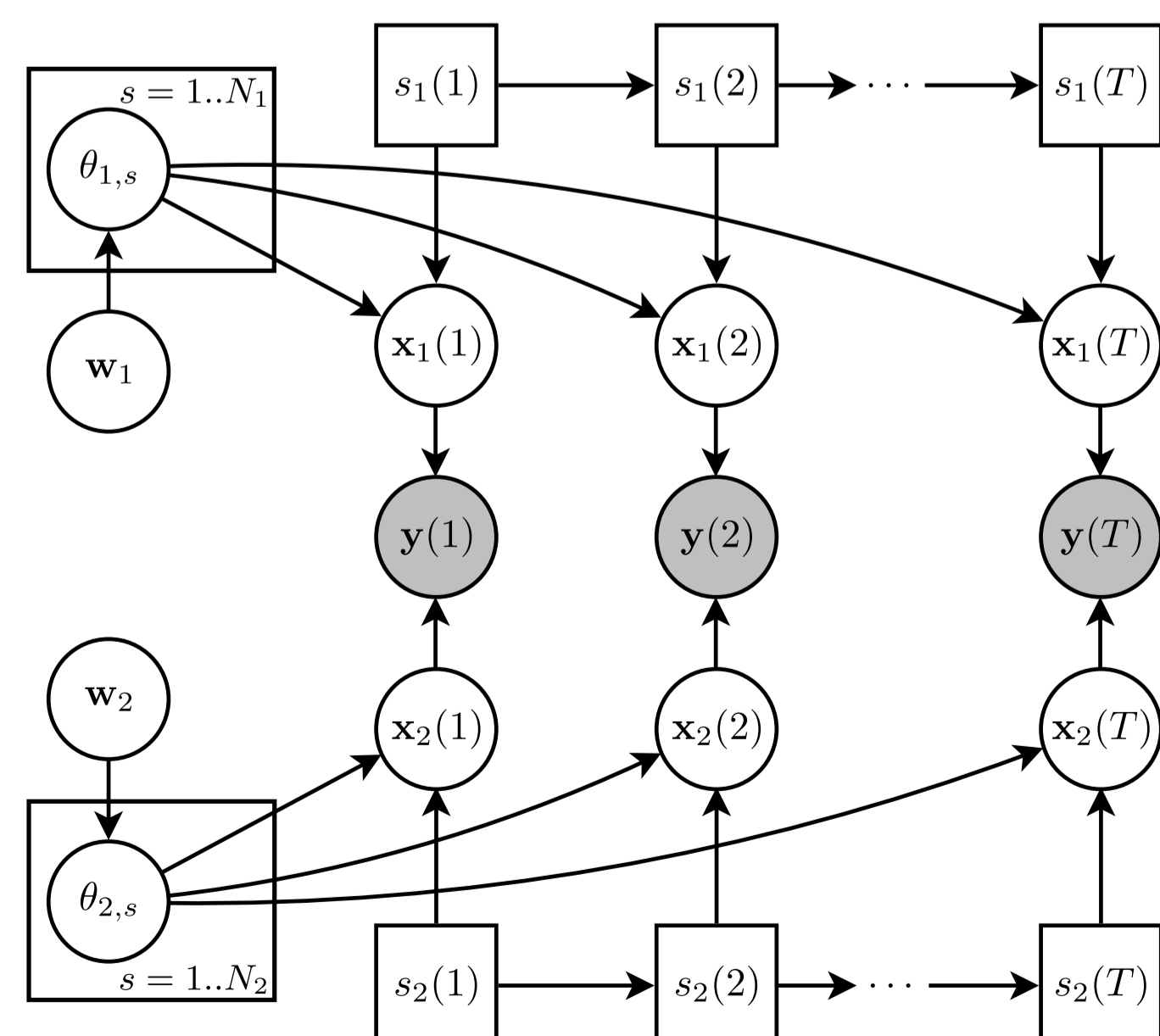$$P(x_i(t) \mid s) = \mathcal{N}(x_i(t); \bar{\mu}_s + U_s \mathbf{w}_i, \bar{\bar{\Sigma}}_s)$$



- Can incorporate covariance parameters into eigenvoice bases to adapt them as well:

$$\log \Sigma_s(w_i) = \log(S_s) \mathbf{w}_i + \log \bar{\bar{\Sigma}}_s$$

- Combine adapted source models into factorial HMM to model mixture:

$$P(y(1..T), s_1(1..T), s_2(1..T))$$
$$= \prod_t P(s_1(t) \mid s_1(t-1)) P(s_2(t) \mid s_2(t-1)) P(y(t) \mid s_1(t), s_2(t))$$
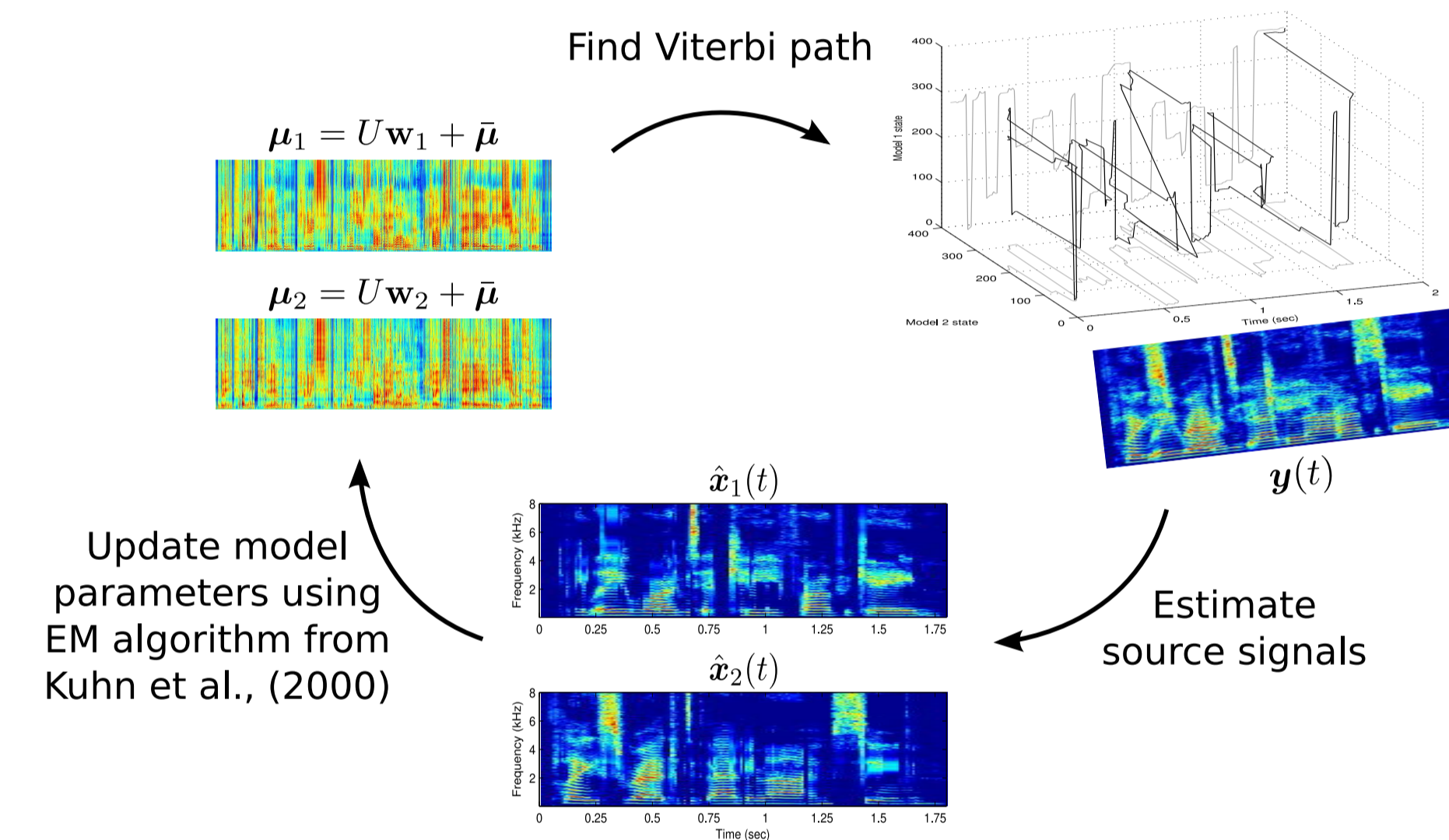


## 3. Adaptation algorithms

- Need to learn eigenvoice adaptation parameters $\mathbf{w}_i$ from mixture
- Exact inference in factorial HMM is intractable – $O(TN^3)$
- Propose two approximate adaptation algorithms:

1. Hierarchical algorithm (Weiss and Ellis, 2008)
   - Iteratively separate sources and learn adaptation parameters from each reconstructed source signal
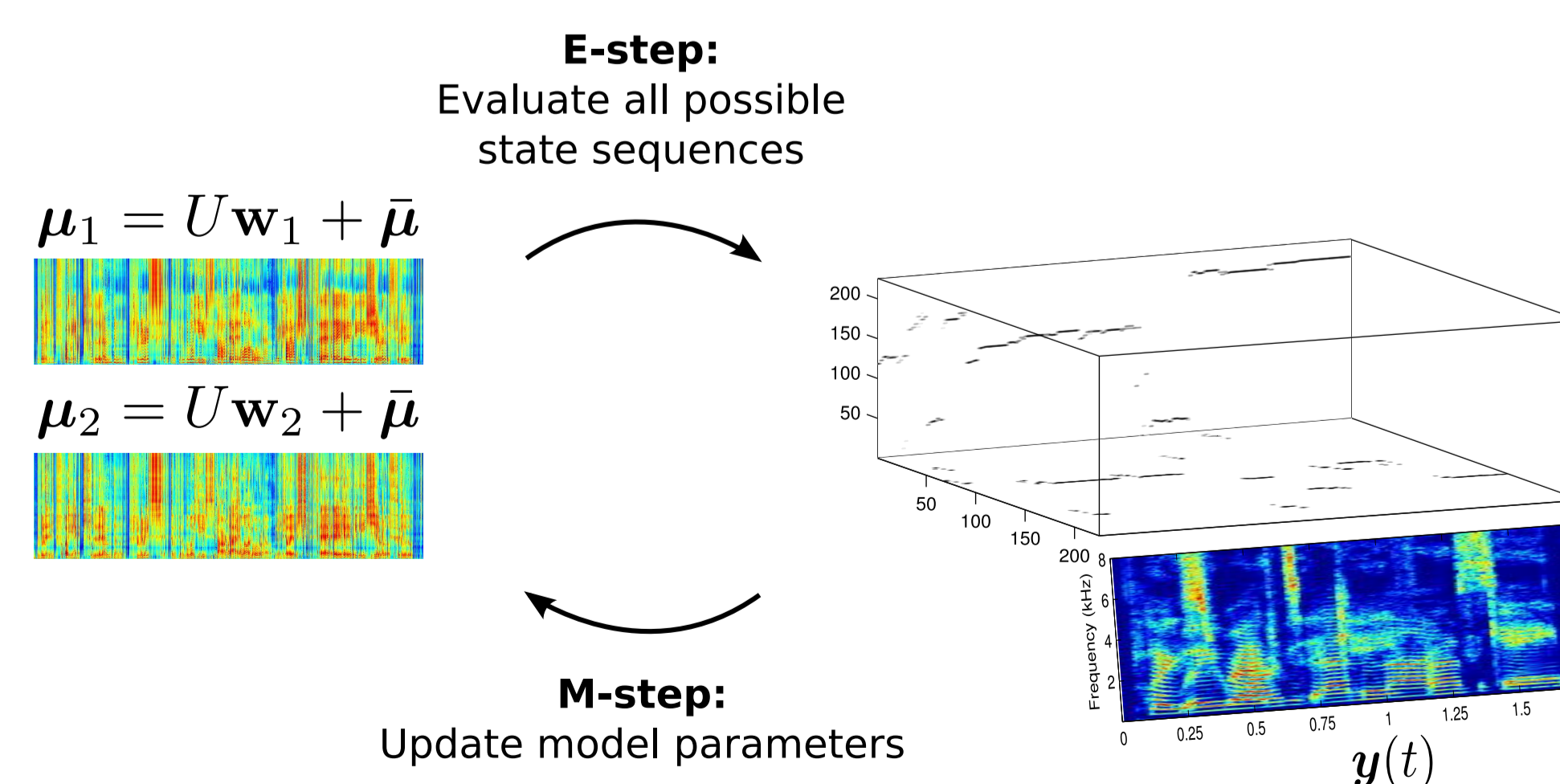   - Use aggressive pruning in factorial HMM Viterbi algorithm to make separation feasible



2. Variational EM algorithm
   - EM algorithm based on structured variational approximation to mixed signal model (Ghahramani and Jordan, 1997)
   - Treat each source HMM independently:

$$P(y(1..T), s_1(1..T), s_2(1..T)) \approx \prod_i Q_i(y(1..T), s_i(1..T))$$

   - Introduce variational parameters to couple them:

$$Q_i(y(1..T), s_i(1..T)) = \prod_t P(s_i(t) \mid s_i(t-1)) h_{i,s_i}(t)$$



**E-step:**
Evaluate all possible state sequences

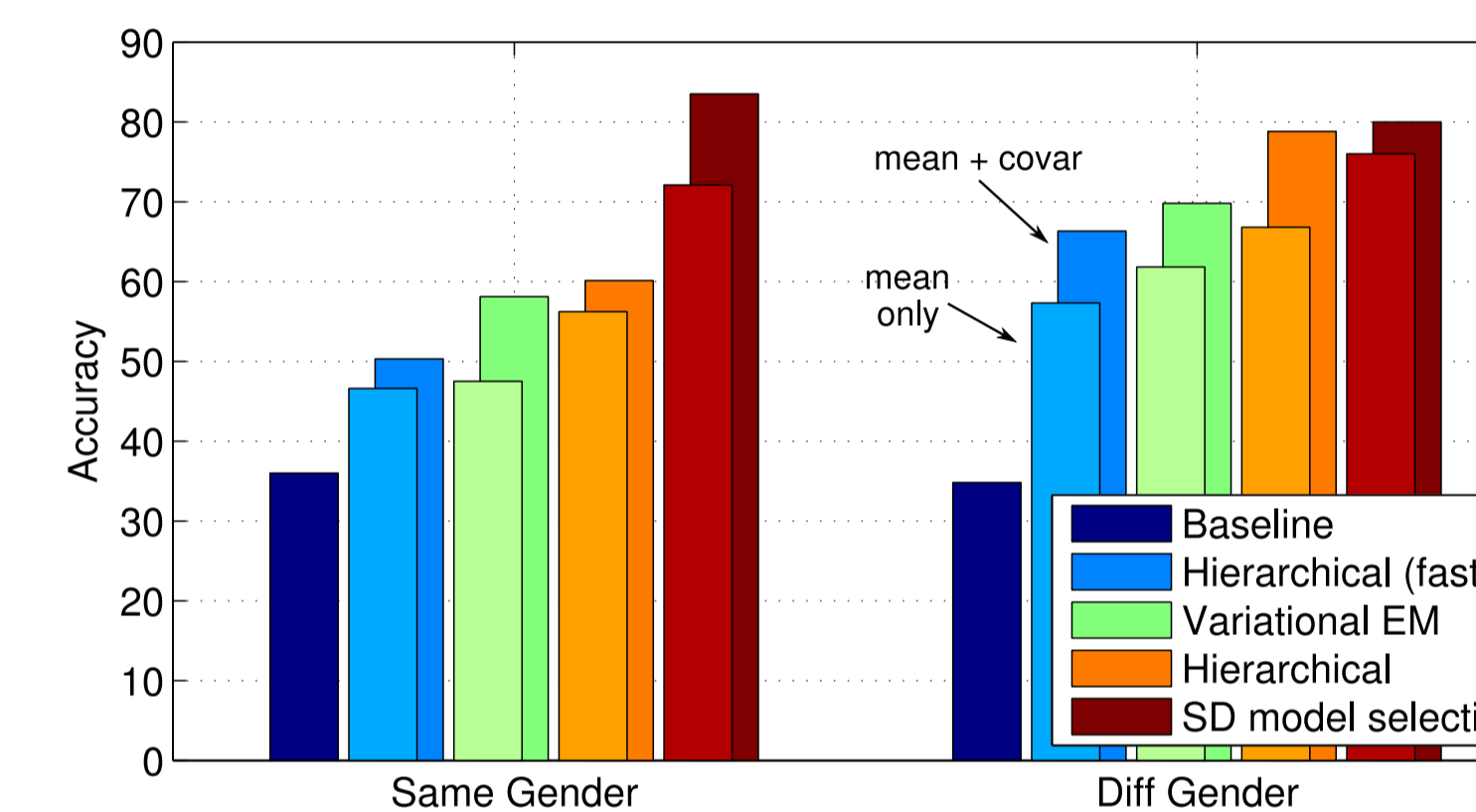**M-step:**
Update model parameters

## 4. Experiments

- Compare two adaptation algorithms with separations based on speaker-dependent (SD) models using speaker identification algorithm from Rennie et al. (2006)
- 0 dB SNR subset of 2006 Speech Separation Challenge data set (Cooke and Lee, 2006)
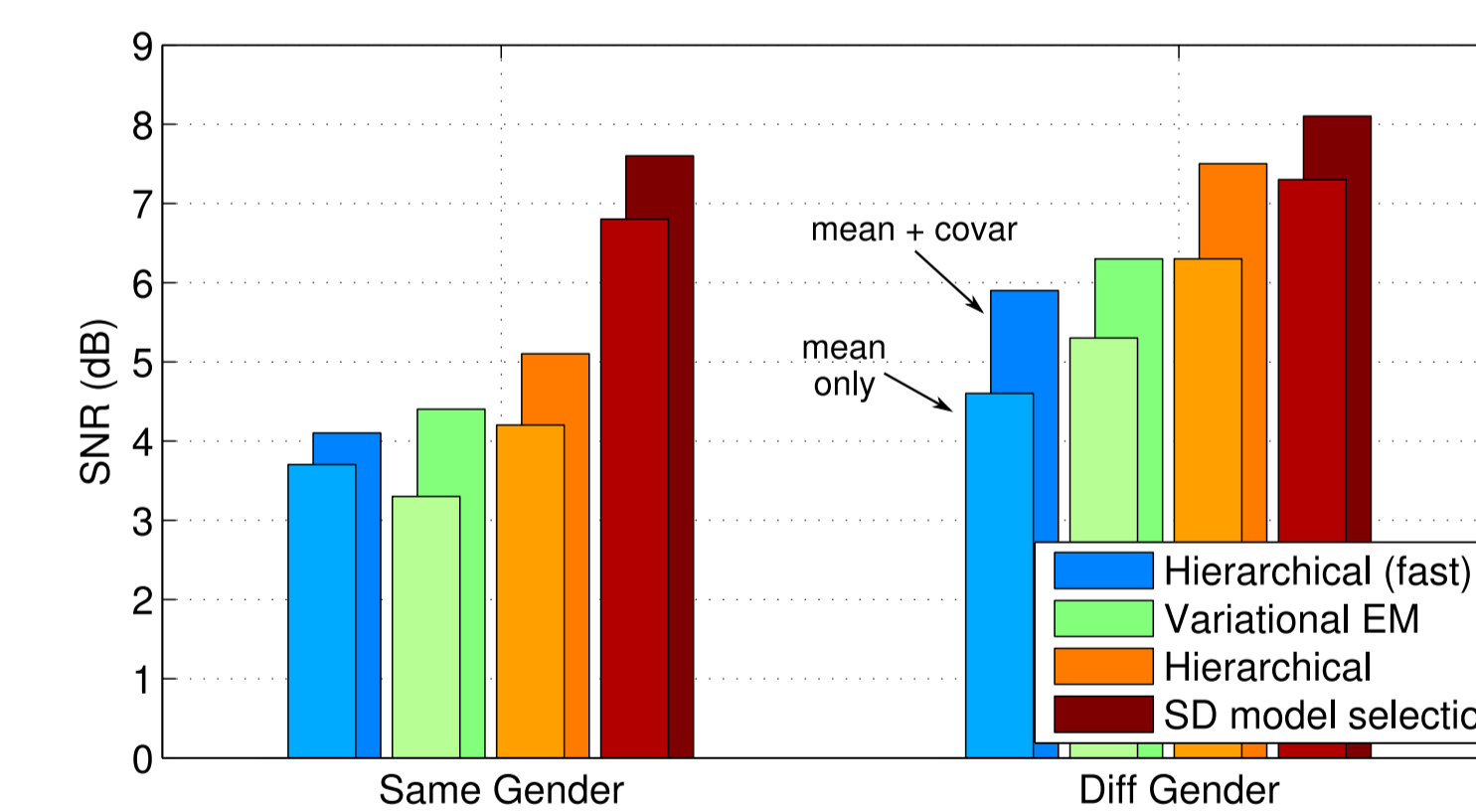- Mixtures of utterances derived from simple grammar:

| command | color | preposition | letter | digit | adverb |
|---------|-------|-------------|--------|-------|--------|
| bin | blue | at | a-v | | again |
| lay | green | by | x-z | 0-9 | now |
| place | red | in | | | please |
| set | white | with | | | soon |

- Task: determine letter and digit spoken by source whose color is "white"
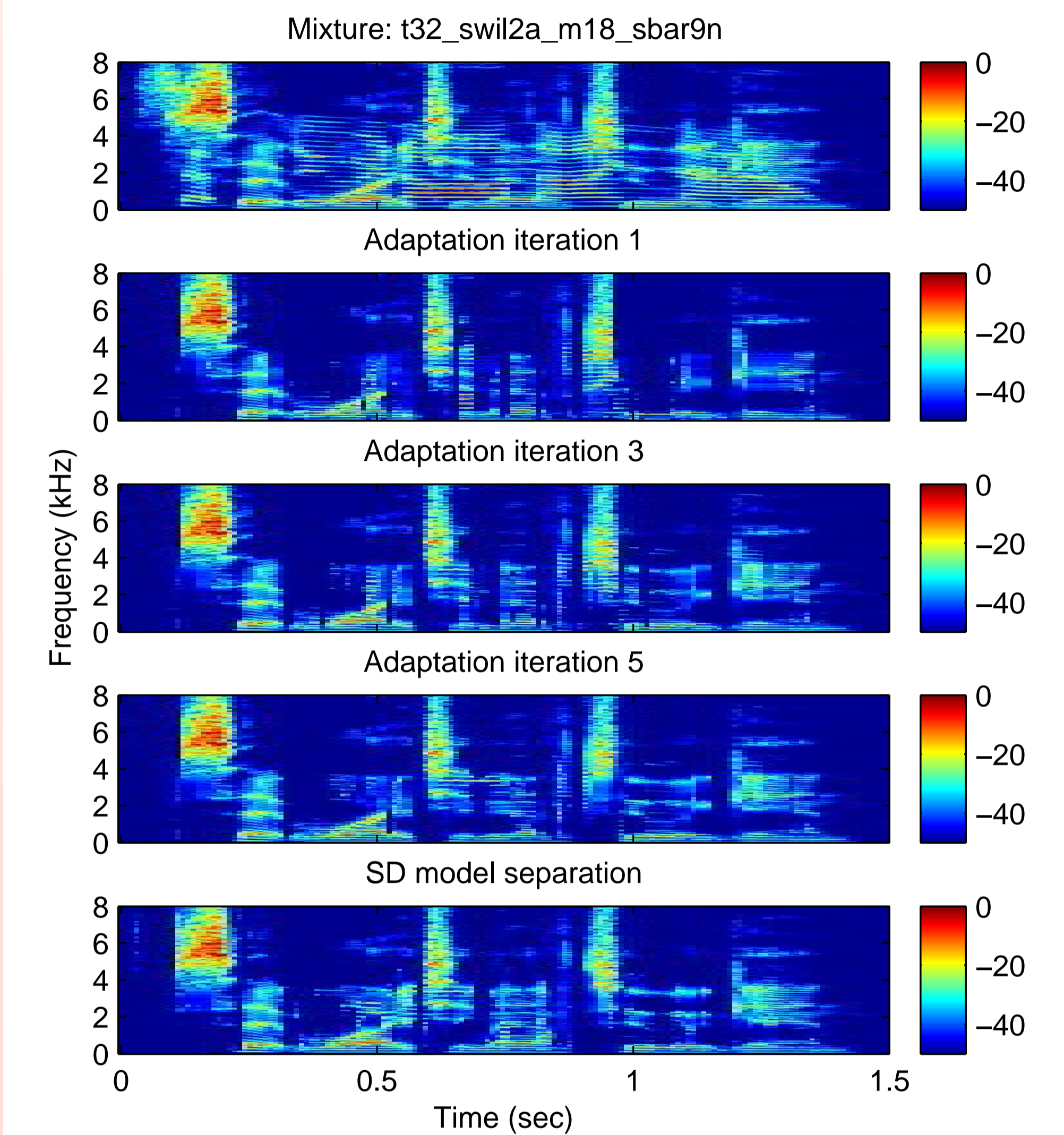  Digit-letter recognition accuracy:



SNR of target source reconstruction:



## 5. Discussion

- Adapting Gaussian covariances as well as means significantly improves performance of all systems
- Adaptation comes to within 23% to 1.2% of best-case SD model performance
- Hierarchical algorithm outperforms variational EM
- But variational algorithm is significantly ($\sim$ 4x) faster
- Performance of the hierarchical algorithm suffers when it is sped up to be as fast as the variational algorithm by pruning even more aggressively ("Hierarchical (fast)" in figures above)

## 6. Example



## 7. References

M. Cooke and T.-W. Lee. The speech separation challenge, 2006. URL http://www.dcs.shef.ac.uk/~martin/SpeechSeparationChallenge.htm.

Z. Ghahramani and M.I. Jordan. Factorial hidden markov models. *Machine Learning*, 29(2-3):245–273, 1997.

R. Kuhn, J. Junqua, P. Nguyen, and N. Niedzielski. Rapid speaker adaptation in eigenvoice space. *IEEE Transations on Speech and Audio Processing*, 8(6):695–707, November 2000.

S. Rennie, P. Olsen, J. Hershey, and T. Kristjansson. The Iroquois model: Using temporal dynamics to separate speakers. In *Workshop on Statistical and Perceptual Audio Processing (SAPA)*, Pittsburgh, PA, September 2006.

R. J. Weiss and D. P. W. Ellis. Speech separation using speaker-adapted eigenvoice speech models. *Computer Speech and Language*, 2008. In press.