

# Comparison of a Joint Iterative Method for Multiple Speaker Identification with Sequential Blind Source Separation and Speaker Identification

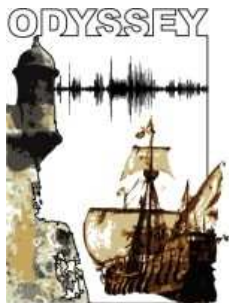
**Youngmoo E. Kim, John M. Walsh, & Travis M. Doll**

Department of Electrical and Computer Engineering

Drexel University

Philadelphia, PA

[jwalsh@ece.drexel.edu](mailto:jwalsh@ece.drexel.edu) & {[ykim](mailto:ykim@drexel.edu),[tmd47](mailto:tmd47@drexel.edu)} @drexel.edu



# Overview

1. Motivation: Multiple Speaker Identification
2. “Off the Shelf” Approach: BSS then SID
  - (a) System Diagram
  - (b) Possible Drawbacks
3. Joint SS and SID: Exact Bayesian Approach
  - (a) Bayesian Model Formulation
  - (b) Infeasibility of Exact Bayesian Estimation / Detection / Inference
4. Approximate Bayesian Inference via Expectation Propagation
  - (a) Overall Idea and Approximating Family
  - (b) System Architecture & Efficient Computation
5. Bounds on Performance Increase
  - (a) Simulations Description
  - (b) Results
6. Extensions

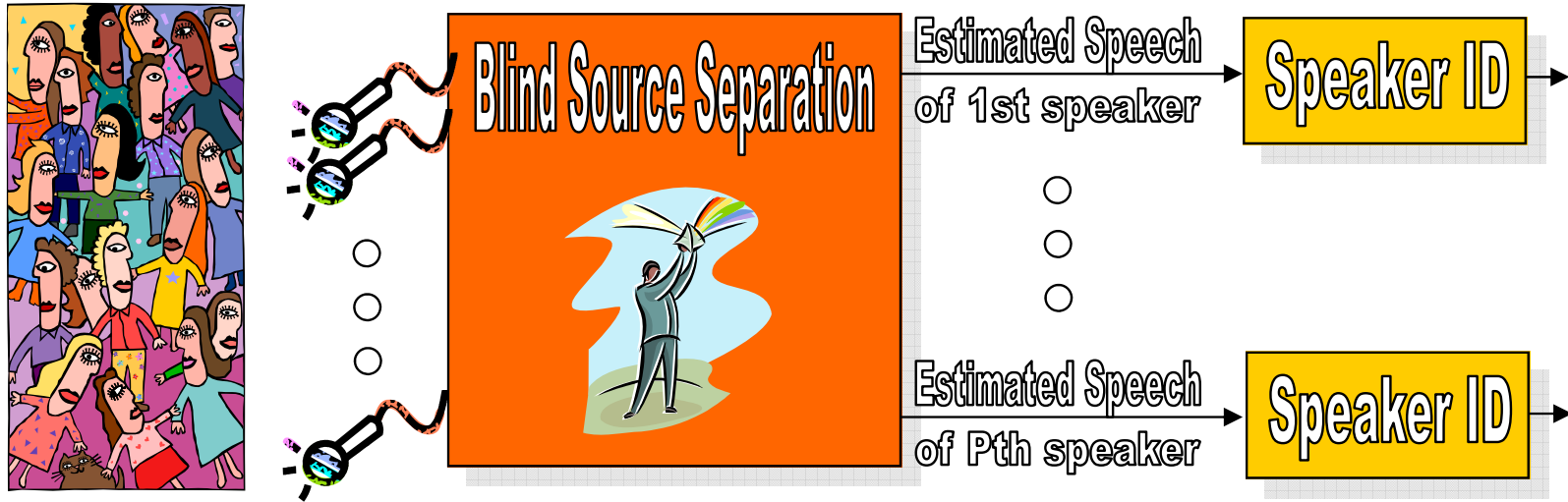
# Multiple Speaker Identification



- Augmented “cocktail party” problem where in addition to separating the sources, we are ultimately interested in determining their identities among a library of known speakers. [1, 2, 3]
- Applications in signal intelligence, meeting minutes transcription, hearing aids.
- Involves both SS and SID.

Figure 1: The multiple speaker identification problem.

## “Off the Shelf” Approach: BSS then SID



- A basic approach to the joint SS and SID problem would capitalize on decades of research in SS [4, 15] and SID [5, 6] by concatenating solutions to the two problems.
- SID is usually considered for a single speaker, with the noise model not following characteristics of background speech.
- BSS is usually considered for ideal independent in time sources.

## BSS then SID : Possible Drawbacks

- *inconsistent models:*
  - plain BSS doesn't incorporate special statistics of speech, let alone even more special statistics of speech from speakers in the library!
  - speaker ID unit does not know incoming signal comes from demixing several speakers (and thus may contain "cross talk")
  - statistical models used for inputs / outputs in the two units are mismatched
- *data reliability not encoded*
  - the reliability of the separated speech symbols is not included among the data passed to the SID unit, yet some may be significantly more reliable than others
- *feedforward structure*
  - feed-forward nature of the system precludes use of SID information in BSS

## Joint SS and SID: Exact Bayesian Approach

$$p_{\theta, \mathbf{r}} = \underbrace{p_{\ell^{(1)}, \dots, \ell^{(P)}}}_{\text{Identity Prior}} \underbrace{p_{\mathbf{r}|\xi}}_{\text{Source Separation}} \prod_{p=1}^P \underbrace{p_{\xi^{(p)}|\mathbf{u}^{(p)}}}_{\text{features}} \underbrace{p_{\mathbf{u}^{(p)}|\ell^{(p)}}}_{\text{speaker ID}}$$

- Bayesian inference/estimation/detection provides a complete framework for the joint approach
  - affords use of existing SS and SID models (capitalizing on decades of disjoint research)
  - Bayesian Speaker Detection:

$$\arg \min_{\ell} \sum_{\ell} C(\hat{\ell}(\mathbf{r}), \ell) \int p_{\theta|\mathbf{r}}(\theta|\mathbf{r}) d\xi d\mathbf{u}, \quad \theta := [\ell, \xi, \mathbf{u}]$$

- Downside: Necessary integral with  $p_{\theta|\mathbf{r}}(\cdot|\mathbf{r})$  nowhere near computationally feasible

# Approximate Bayesian Inference via Expectation Propagation

- try to approximate  $p_{\theta|r}$  with

$$g_{\ell^{(1)}, \dots, \ell^{(P)}} g_{\xi} \prod_{p=1}^P g_{\xi^{(p)} | \mathbf{u}^{(p)}} g_{\mathbf{u}^{(p)} | \ell^{(p)}}$$

among  $g_{\theta^{(i)}}$ s in particular exponential families for which the marginal calc. integral is easy.

- refine  $g_{\theta^{(i)}}$  with local Bayes computation, e.g. refined  $g_{\xi}$  solves

$$\begin{aligned} & \beta \int \mathbf{t}(\xi) g_{\ell^{(1)}, \dots, \ell^{(P)}} g_{\mathbf{r} | \xi} \prod_{p=1}^P g_{\xi^{(p)} | \mathbf{u}^{(p)}} g_{\mathbf{u}^{(p)} | \ell^{(p)}} d\theta \\ &= \alpha \int \mathbf{t}(\xi) g_{\ell^{(1)}, \dots, \ell^{(P)}} p_{\mathbf{r} | \xi} \prod_{p=1}^P g_{\xi^{(p)} | \mathbf{u}^{(p)}} g_{\mathbf{u}^{(p)} | \ell^{(p)}} d\theta \end{aligned}$$

- $\mathbf{t}$  the sufficient statistic for the app. fam. for  $g_{\xi}$

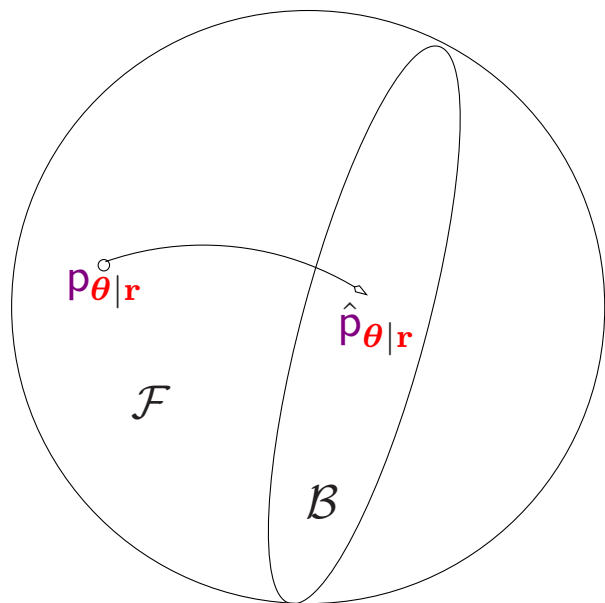


Figure 2: EP [7, 8, 9] refines an approximation  $\hat{p}_{\theta|r}$  to the true a posteriori distribution  $p_{\theta|r}$  among PDFs in an approximating exponential family [10, 11]  $\mathcal{B}$ .

# Joint Iterative System Architecture

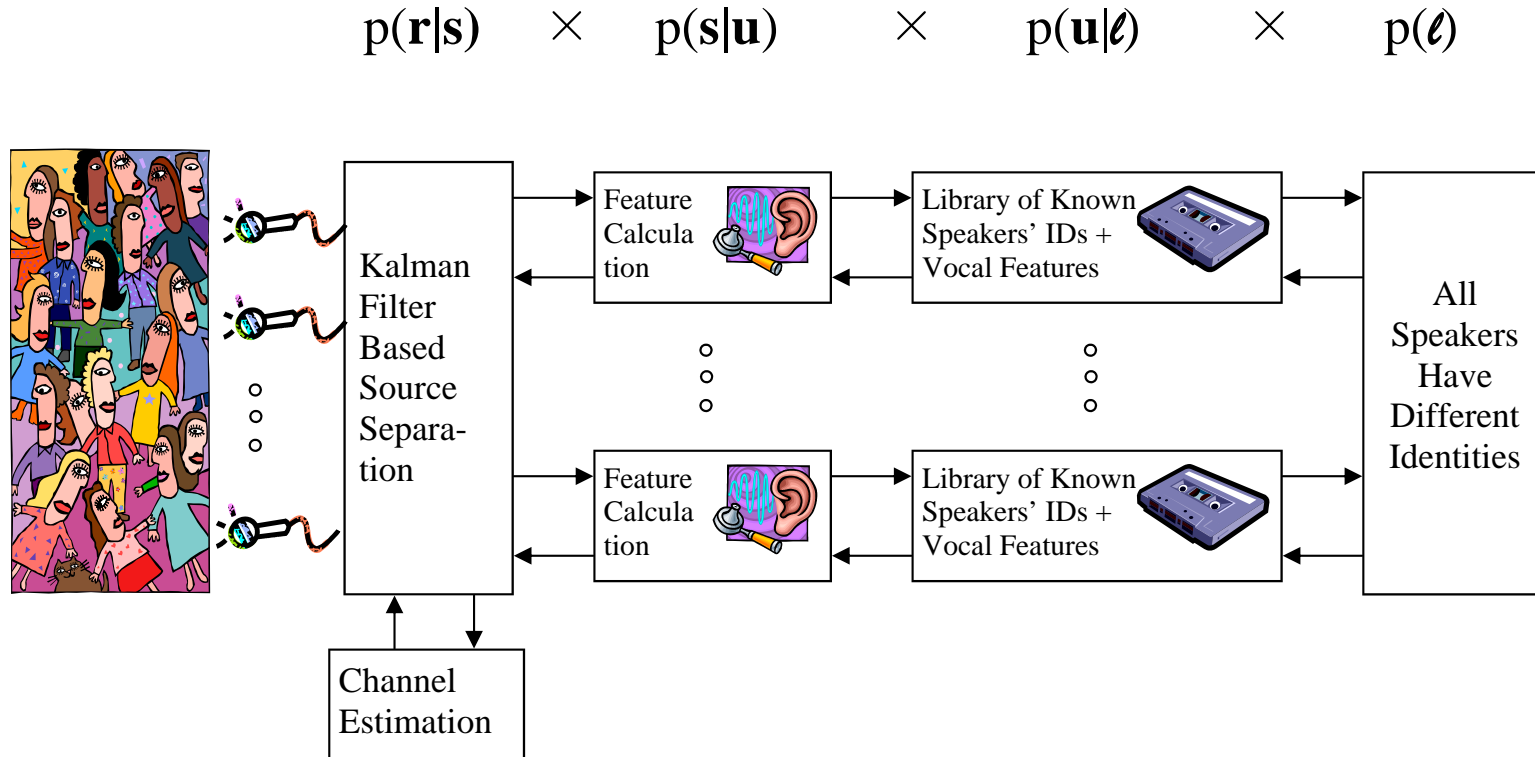


Figure 3: Studied joint iterative SS and SID system.



## Saving Further Computational Complexity

- LTI AWGN acoustic channel model, MFCC features, i.i.d. m.v. Gaussian MFCC model with identity dependent mean and covariance matrix.
- indep. model separated audio samples  $\xi$  as indep. Gaussian r.v.s, MFCC vectors  $\mathbf{u}_k^{(p)}$  as indep. Gaussian r.vec.s, and identity indices  $\ell^{(p)}$  as indep. d.r.v.s
- Gives RTS smoother (forward backward algorithm) [12, 13] as EP source separator.
- Linearize nonlinear MFCC Calculation about mean to propagate covariances.
- Use Kalman filter instead of RTS (forward backward) smoother source separator. [14]
- Use union bound in all speakers have different identities to cut down on number of possibilities to be enumerated in sum.

# Genie Aided Bounds on Performance Increase Obtained

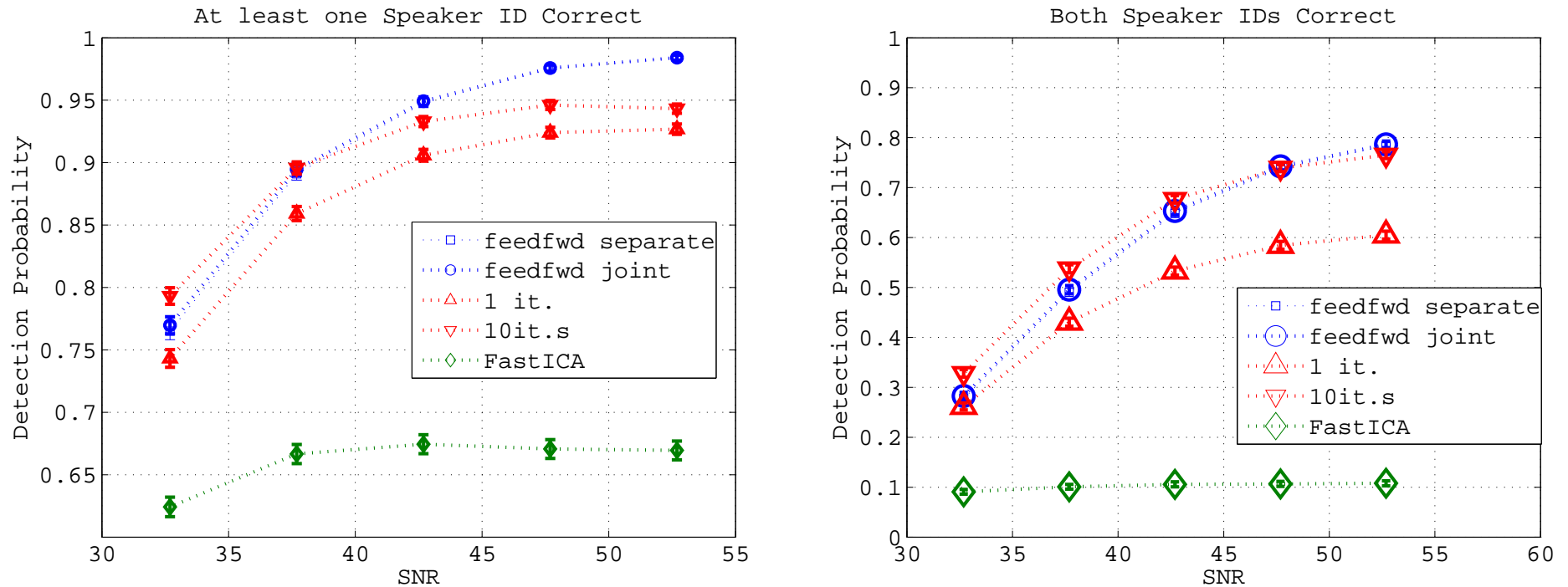


Figure 4: Performance of joint separation and identification system in correctly identifying at least one speaker (left) and both speakers (right) compared to that of reference system (FastICA [15] source separation, followed by classification via MFCC KL distance). (40 speakers from TIMIT database [16], 2.5 seconds, random Rayleigh fading channel w/ exponentially decreasing power profile.)

# Genie Aided Bounds on Performance Increase Obtained

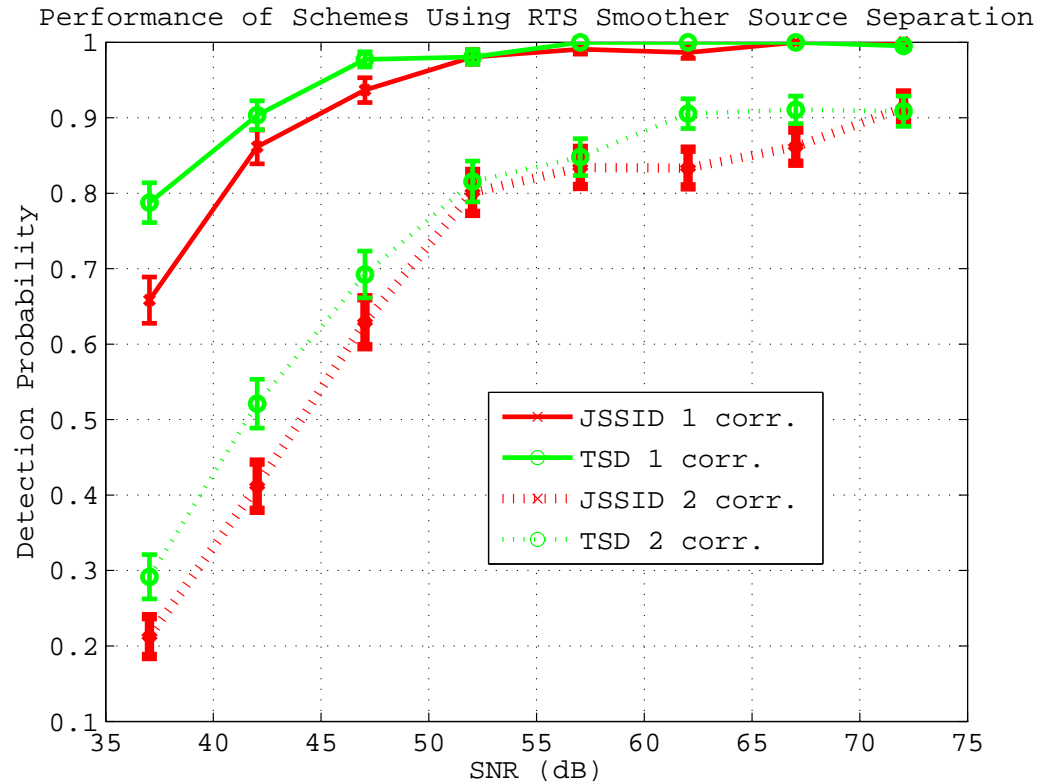


Figure 5: Performance of two joint separation and identification systems using a RTS smoother for the source separation. The systems labeled with JSSID use a speaker feature MFCC feature model that is i.i.d. multivariate Gaussian. The systems labeled with TSD use a speaker feature MFCC model that is a GLSSM. (160 TEST speakers from TIMIT database [16], 2.5 seconds, random Rayleigh fading channel w/ exponentially decreasing power profile.)

## Extensions, Current & Future Work

- incorporating EM algorithm based channel estimation
- alternative features (PCA like eigenbasis for multispeaker ACF)
- frequency domain source separation
- use GLSSM MFCC model: detector becomes turbo decoder like
- how many iterations? 1 or 2 currently probably enough, can have convergence problems (due to loops) with more.
- non-simulations (i.e. analysis) based performance / convergence guarantees?
- music source separation
- Better approximating families including a temporal dependence structure in separated speech model.

# References

- [1] M. A. Przybocki, A. F. Martin, and A. N. Le, "NIST speaker recognition evaluations utilizing the mixer corpora—2004, 2005, 2006," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 15, no. 7, pp. 1951–1959, 2007.
- [2] A. Stolcke *et al*, "Progress in meeting recognition: The ICSI-SRI-UW spring 2004 evaluation system," in *Proc. NIST ICASSP 2004 Meeting Recognition Workshop*, Montreal, May 2004.
- [3] M. Przybocki, A. Martin, and A. Le, "NIST speaker recognition evaluation chronicles - Part 2," in *Proc. IEEE Speaker and Language Recognition Workshop (Odyssey)*. IEEE, 2006, pp. 1–6.
- [4] S. Haykin and Z. Chen, "The Cocktail Party Problem," *Neural Computation*, vol. 17, no. 9, pp. 1875–1902, Sept. 2005.
- [5] D. A. Reynolds, "An overview of automatic speaker recognition technology," in *Proc. IEEE Int'l Conf. on Acoustics Speech and Signal Processing (ICASSP)*, 2002, pp. 4072–4075.
- [6] W. M. Campbell, J. P. Campbell, D. A. Reynolds, E. Singer, and P. A. Torres-Carrasquillo, "Support vector machines for speaker and language recognition," *Computer Speech & Language*, vol. 20, no. 2-3, pp. 210–229, 2006.
- [7] T. P. Minka, "A family of algorithms for approximate bayesian inference," Ph.D. dissertation, Massachusetts Institute of Technology, 2001.
- [8] —, "Expectation propagation for approximate Bayesian inference," in *Uncertainty in AI'01*, 2001.
- [9] J. M. Walsh, "Distributed Iterative Decoding and Estimation via Expectation Propagation: Performance and Convergence," Ph.D. dissertation, Cornell University, 2006.
- [10] S. Amari and H. Nagaoka, *Methods of Information Geometry*. AMS Translations of Mathematical Monographs, 2004, vol. 191.
- [11] L. D. Brown, *Fundamentals of Statistical Exponential Families with Applications in Statistical Decision Theory*. Institute of Mathematical Statistics, 1986.
- [12] B. D. O. Anderson and J. B. Moore, *Optimal Filtering*. Prentice Hall, Inc., 1979.
- [13] O. Cappé, E. Moulines, and T. Rydén, *Inference in Hidden Markov Models*. Springer Science and Business Media, 2005.
- [14] J. M. Walsh, Y. E. Kim, and T. M. Doll, "Joint iterative multi-speaker identification and source separation using expectation propagation," in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, 2007, submitted.
- [15] A. Hyvärinen and E. Oja, "Independent component analysis: Algorithms and applications," *Neural Networks*, vol. 13, no. 4-5, pp. 411–430, 2000.
- [16] National Institute of Standards and Technology (NIST), "The DARPA TIMIT acoustic-phonetic continuous speech corpus," NIST, 1990.