

Use of Particle Filtering and MCMC for Inference in Probabilistic Acoustic Tube Model

Ruobai Wang¹

wangruobai11@tsinghua.org.cn

Yang Zhang²

yzhan143@illinois.edu

Zhijian Ou¹

ozj@tsinghua.edu.cn

Mark Hasegawa-Johnson²

jhasegaw@illinois.edu

¹ Tsinghua University, Department of Electronic Engineering

² University of Illinois, Urbana-Champaign, Department of Electrical and Computer Engineering

The Signal Model of PAT3

$$x(t) = \left[\left(v(t) + \tilde{b} u(t) \right) * h(t) + \tilde{e} \varepsilon(t) \right] w(t)$$

voiced source unvoiced amplitude vocal tract response noise window

Gaussian unvoiced signal

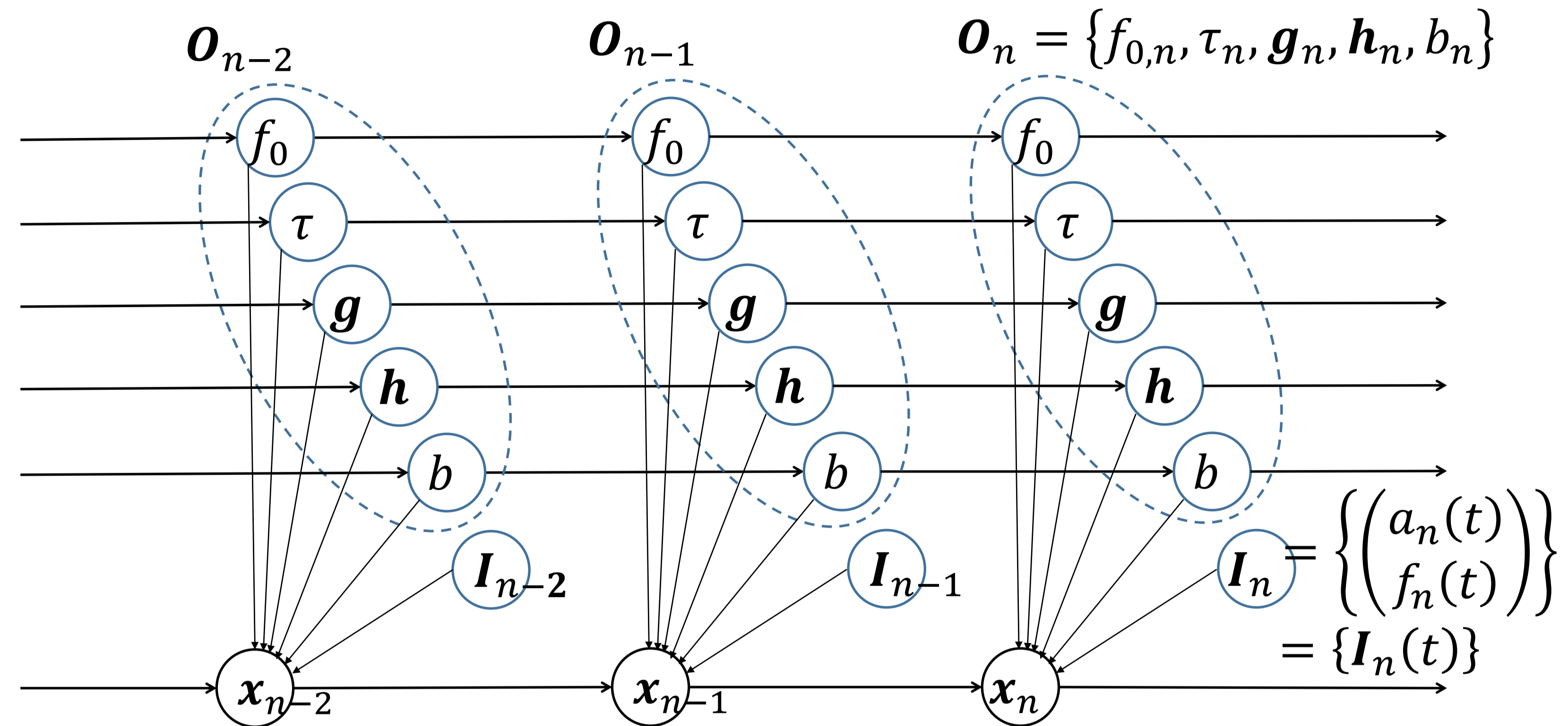
- $X_n(\omega) = \left[\left(V_n(\omega) + \tilde{b}_n U_n(\omega) \right) H_n(\omega) + \tilde{e}_n E_n(\omega) \right] \otimes W_n(\omega)$
- Stack the real and imaginary parts of $X_n \left(\frac{2\pi k}{L} \right)$, $k = 0$ to $L/2$, as \mathbf{x}_n .
- Voiced part $S_n(\omega) = [V_n(\omega)H_n(\omega)] \otimes W_n(\omega)$, vectorized as \mathbf{s}_n .
- Unvoiced part implies variance: $\sigma_{n,k}^2 = \tilde{b}_n^2 \left| H_n \left(\frac{2\pi k}{L} \right) \right|^2 + \tilde{e}_n^2$

$$v(t) = \tilde{a}(t) \sum_{d=1}^D \text{real} \left[G(d\omega_0) \exp(-d\omega_0\tau + d\phi(t)) \right]$$

amplitude envelope harmonic order glottal transfer function pitch frequency group delay instant phase for d=1

The Probabilistic Model of PAT3

Speech frames are modeled as a non-linear state-space model.



Observation likelihood :

$$\log p(\mathbf{x}_n | \mathbf{s}_n) = \sum_{k=0}^{L-1} \left[-\frac{1}{2} \log(2\pi\sigma_{n,k}) - \frac{\left| X_n \left(\frac{2\pi k}{L} \right) - S_n \left(\frac{2\pi k}{L} \right) \right|^2}{2\sigma_{n,k}^2} \right]$$

State transition :

$$p(\mathbf{O}_n | \mathbf{O}_{n-1}) = \mathcal{N}(\mathbf{O}_{n-1}, \Sigma_O)$$

$$p(\mathbf{I}_n(t) | \mathbf{I}_n(t-2), \mathbf{I}_n(t-1)) = \mathcal{N}(\lambda \mathbf{I}_n(t-1) + (1-\lambda)\mathbf{I}_n(t-2), \Sigma_I)$$

Outer Loop using APF

Frame-level MAP inference : $\hat{\mathbf{O}}_{1:N}^{MAP} = \underset{\mathbf{O}_{1:N}}{\text{argmax}} p(\mathbf{O}_{1:N} | \mathbf{x}_{1:N})$.

Suppose that $p(\mathbf{O}_{n-1} | \mathbf{x}_{1:n-1})$ is approximated by :

$$\left\{ \mathbf{O}_{n-1}^{(i)}, w_{n-1}^{(i)} \right\}_{i=1:M}$$

Then, APF aims to sample from

$$p(\mathbf{O}_n | \mathbf{x}_{1:n}) \propto p(\mathbf{x}_n | \mathbf{O}_n) \sum_{i=1}^M p(\mathbf{O}_n | \mathbf{O}_{n-1}^{(i)}) w_{n-1}^{(i)}$$

with proposal

$$q(\mathbf{O}_n | \mathbf{O}_{n-1}^{(i)}, \mathbf{x}_n) \beta_n^{(i)}.$$

For $n = 1, \dots, N$:

- Resample** j_1, \dots, j_M from $\{1, \dots, M\}$ according to $\left\{ \beta_n^{(i)} \right\}_{i=1, \dots, M}$
- Propagate**: Sample $\mathbf{O}_n^{(i)}$ from $q(\mathbf{O}_n | \mathbf{O}_{n-1}^{(j_i)}, \mathbf{x}_n)$
- Reweight** each particle $\mathbf{O}_n^{(i)}$ as $w_n^{(i)} \propto \frac{p(\mathbf{x}_n | \mathbf{O}_n^{(i)}) p(\mathbf{O}_n | \mathbf{O}_{n-1}^{(j_i)}) w_{n-1}^{(j_i)}}{q(\mathbf{O}_n | \mathbf{O}_{n-1}^{(j_i)}, \mathbf{x}_n) \beta_n^{(j_i)}}$

Inner Loop using Taylor Expansion Assisted MCMC

Viterbi approximation :

$$p(\mathbf{x}_n | \mathbf{O}_n) \approx \max_{\mathbf{I}_n} p(\mathbf{x}_n | \mathbf{O}_n, \mathbf{I}_n) = p(\mathbf{x}_n | \mathbf{O}_n, \hat{\mathbf{I}}_n) p(\hat{\mathbf{I}}_n)$$

$\hat{\mathbf{I}}_n$ is picked from the samples $\left\{ \hat{\mathbf{I}}_n^{(r)} \right\}_{r=1:R}$ drawn from $p(\mathbf{I}_n | \mathbf{O}_n, \mathbf{x}_n)$.

For each frame n , we run Metropolis-Hasting algorithm :

$$\min \left\{ 1, \frac{p(\mathbf{I}_n | \mathbf{O}_n, \mathbf{x}_n) q(\hat{\mathbf{I}}_n^{(r)} | \mathbf{I}_n, \mathbf{O}_n, \mathbf{x}_n)}{p(\hat{\mathbf{I}}_n^{(r)} | \mathbf{O}_n, \mathbf{x}_n) q(\mathbf{I}_n | \hat{\mathbf{I}}_n^{(r)}, \mathbf{O}_n, \mathbf{x}_n)} \right\}$$

with proposal

$$q(\mathbf{I}_n | \hat{\mathbf{I}}_n^{(r)}, \mathbf{O}_n, \mathbf{x}_n) \triangleq \prod_t q(\mathbf{I}_n(t) | \hat{\mathbf{I}}_n^{(r)}(1:t-1), \hat{\mathbf{I}}_n^{(r)}(t+1:T), \mathbf{O}_n, \mathbf{x}_n)$$

We exploit **2nd-order Taylor expansion** of the single-site conditional distribution

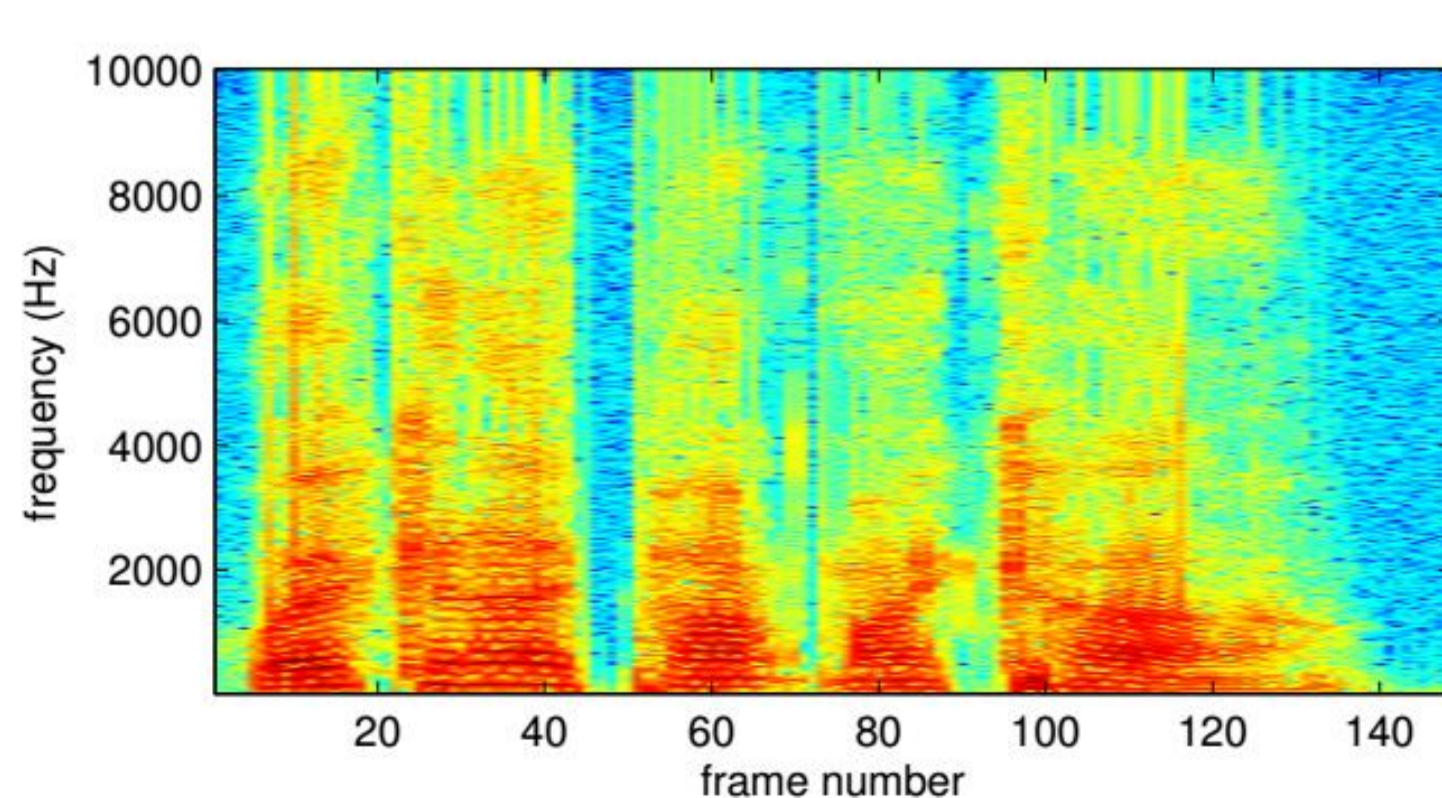
$$\ln p(\mathbf{I}_n(t) | \hat{\mathbf{I}}_n^{(r)}(1:t-1), \hat{\mathbf{I}}_n^{(r)}(t+1:T), \mathbf{O}_n, \mathbf{x}_n)$$

at the a prior conditional mean of $\mathbf{I}_n(t)$

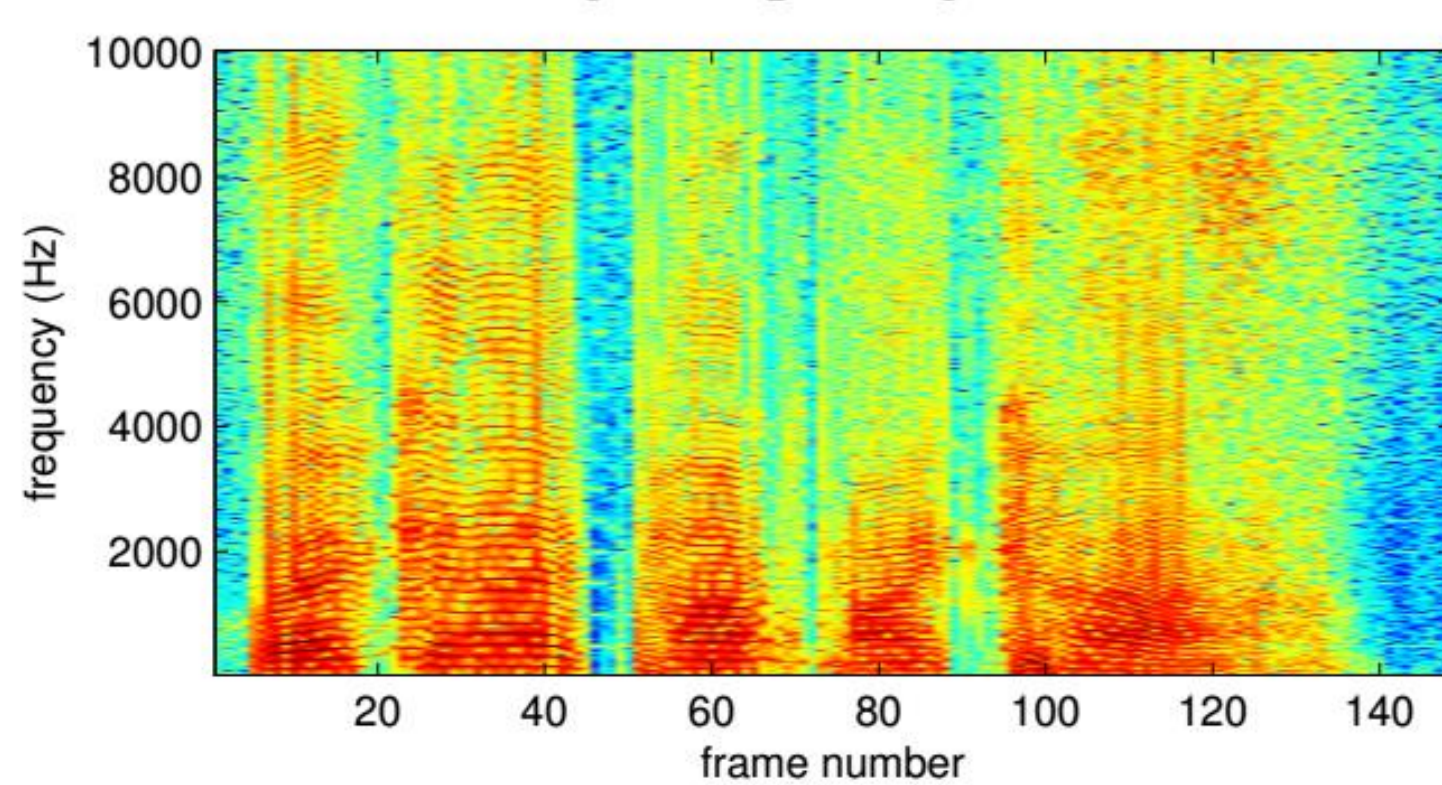
$$\mathbb{E}(\mathbf{I}_n(t) | \hat{\mathbf{I}}_n^{(r)}(t-2:t-1), \hat{\mathbf{I}}_n^{(r)}(t+1:t+2))$$

to define the single-site proposal distribution.

Experiments

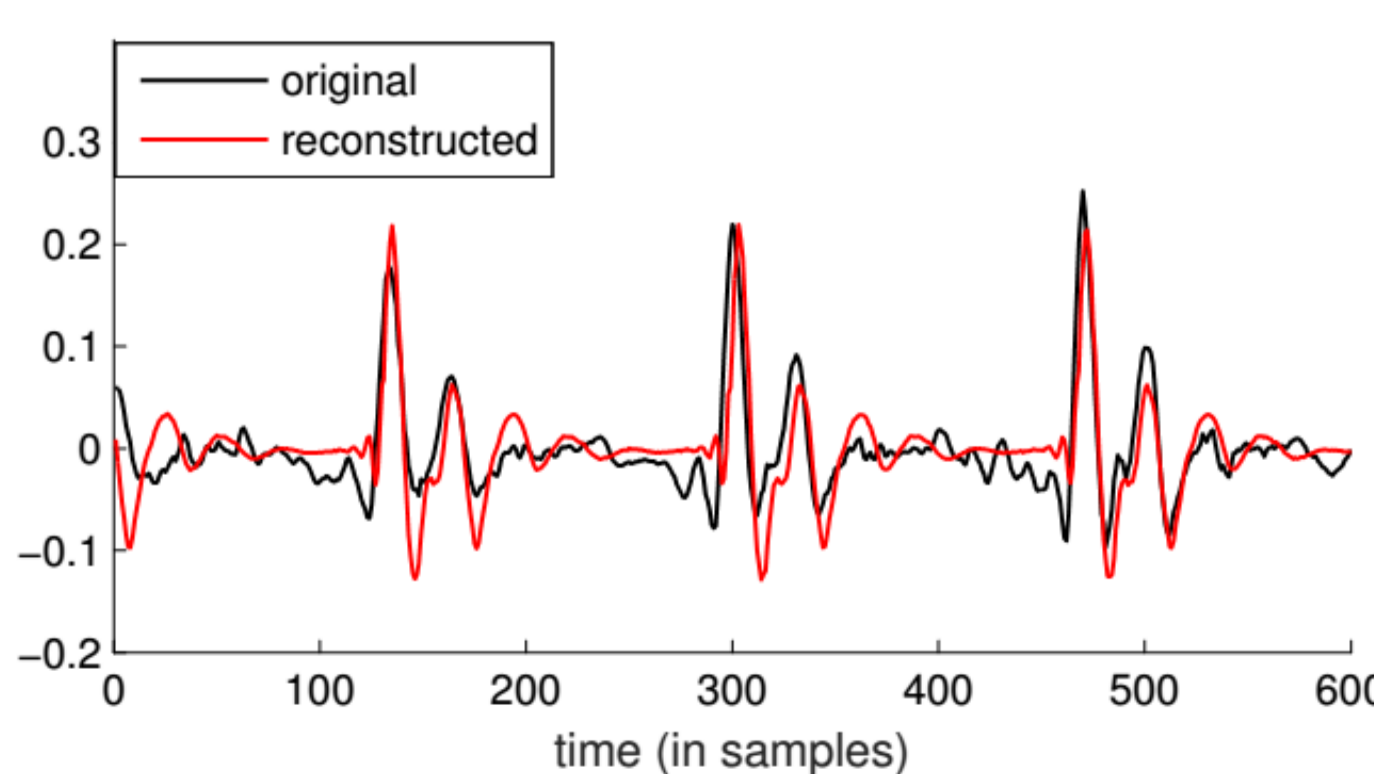


(a) Original Spectrogram

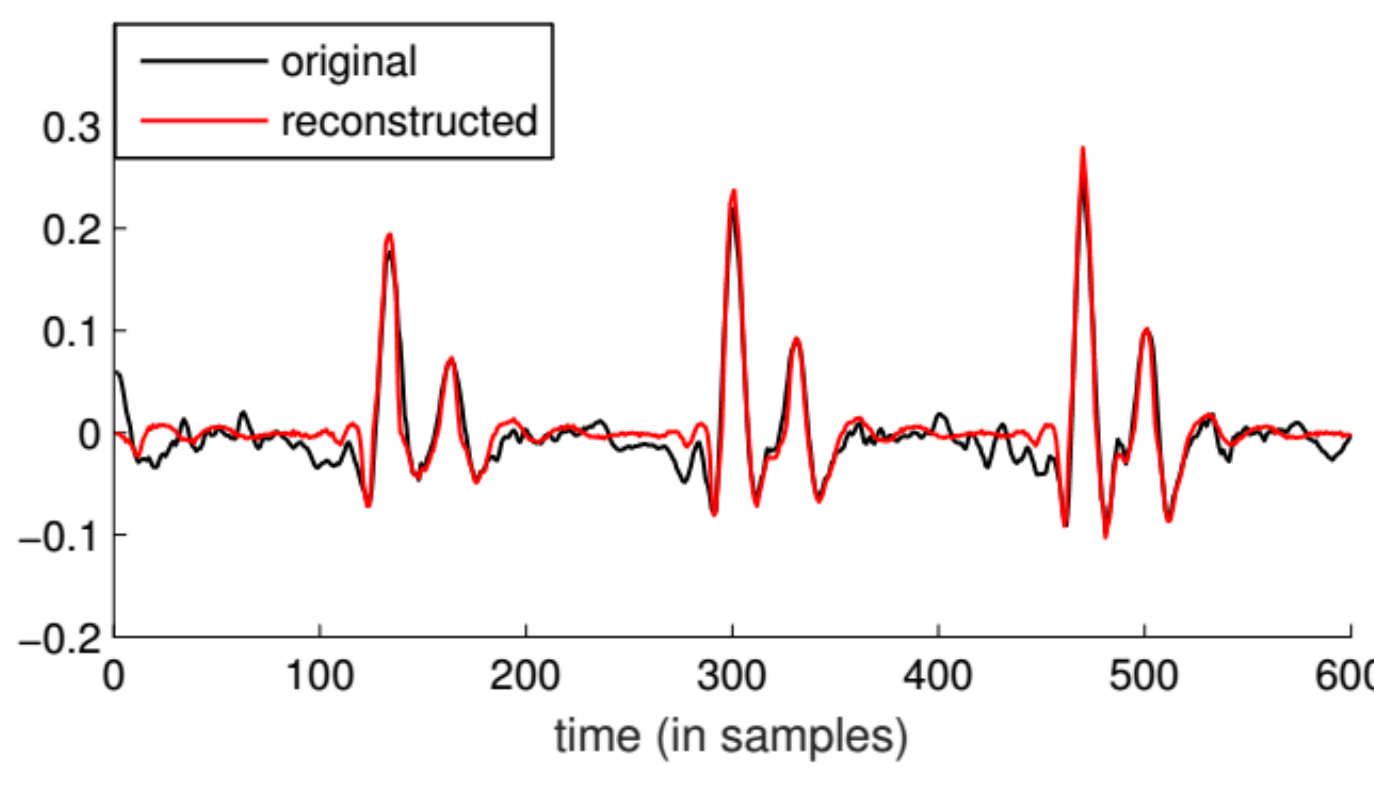


(b) Reconstructed Spectrogram

Fig. 1: Comparison of the reconstructed and original spectrogram.



(a) Without AM-FM Tracking



(b) With AM-FM Tracking

Fig. 2: The reconstructed voiced waveform (red line) with and without AM-FM tracking, compared with the original waveform (black).

PAT3	STRAIGHT
8.79	-2.57

Table 1: SNR of speech reconstruction

	PAT3	GetF0
GPE (%)	2.10	2.07
RMS (Hz)	5.052	5.780

Table 2: Pitch tracking results on Edinburgh dataset

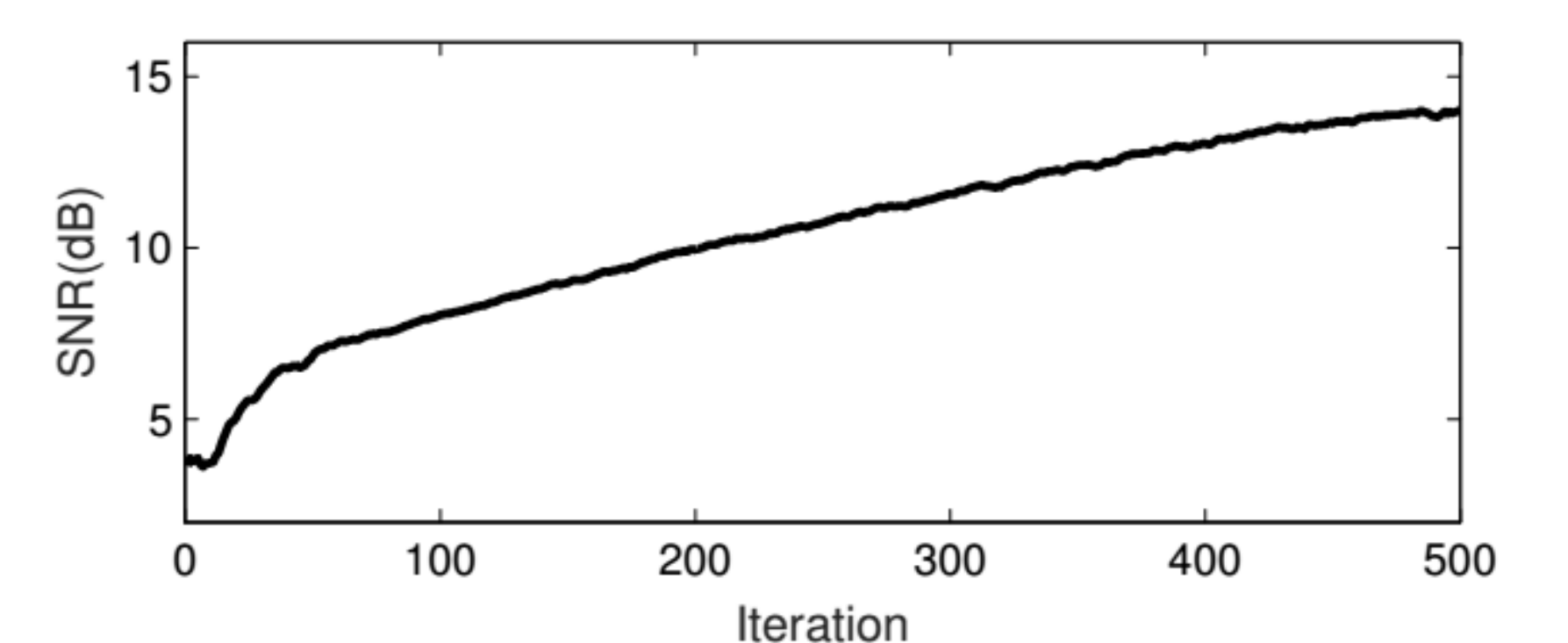


Fig. 3: SNR of reconstruction as a function of MCMC iterations.