# Probabilistic Acoustic Tube:
# a probabilistic generative model of speech for speech analysis/synthesis

**Zhijian Ou**

ozj@tsinghua.edu.cn

**Yang Zhang**

zhangyangbill@gmail.com

Department of Electronic Engineering, Tsinghua University, Beijing 100084, China
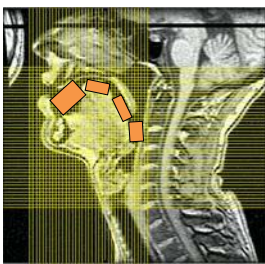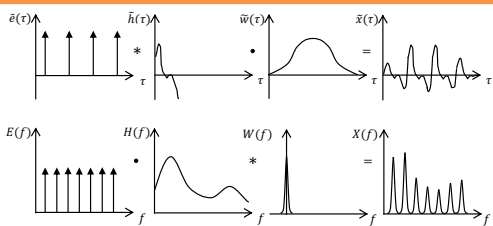
## Motivation

| Current Speech Analysis Methods | Probabilistic Acoustic Tube (PAT) |
|---|---|
| Signal processing front-end + probabilistic back-end. Most features are nonlinear operators of speech (autocorrelation, cepstrum). | Directly model the spectrogram. Preserve additivity. |
| Different tasks of speech analysis are carried out separately. Ignore the chicken-and-egg relationship. | A unified probabilistic model to integrate the pitch, energy and spectral envelope. |

## Model Formulation

### Physical Acoustic Tube



### Signal Processing Modeling



**Voiced Frame:**
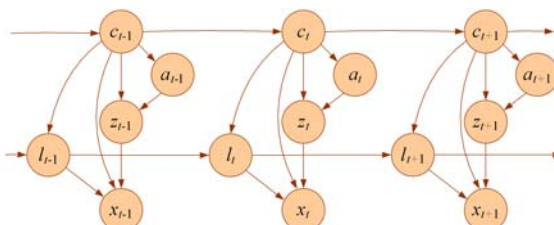
$$|X(\omega)| = |[E_l(\omega)H(\omega)] * W(\omega)|$$
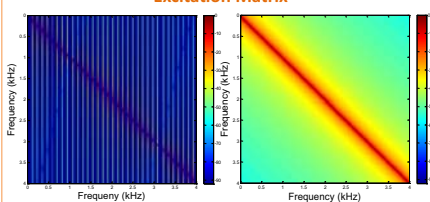$$\approx \sum_n |H(n\eta_l)||W(\omega - n\eta_l)|$$

**Unvoiced Frame:**

$$|X(\omega)| = \left|\sum_\xi H(\xi)W(\omega - \xi)\right|$$
$$\approx \sum_\xi |H(\xi)||W(\omega - \xi)|$$

### Probabilistic Modeling



**Excitation Matrix**



| $c_t$ | $l_t$ | $a_t$ | $z_t$ | $x_t$ |
|---|---|---|---|---|
| Phoneme Cluster | Discretized Frequency | Excitation Gain | DCT Coef. of Spectral Envelope | Observed Spectrogram |
| $p(c_t|c_{t-1})$ Constrain UV transition | $p(l_t|l_{t-1}, c_t)$ Constrain abrupt changes in pitch | $p(a_t|c_t = c)$ $= \mathcal{N}(a_t; m_c, \sigma_c^2)$ | $p(z_t|c_t = c, a_t)$ $= \mathcal{N}(z_t; a_t\mu_c, \Phi_c)$ $\mu_c'\mu_c = 1$ | $x_t = E_{l_t}Cz_t + n_t$ $p(n_t|c_t = c)$ $= \mathcal{N}(n_t; 0, m_c^2\Psi)$ |

**The parameters of PAT - $\Theta \triangleq \{\mu_c, \Phi_c, \Psi, m_c, \sigma_c^2\}$ can be solved using the EM algorithm**

## Experimental Results

**The capability of PAT is demonstrated for a number of speech analysis/synthesis tasks.**

### Pitch Tracking

**Pitch tracking result**

| | PAT | Get_f0 |
|---|---|---|
| UE (%) | 5.38 | 8.84 |
| VE (%) | 4.83 | 4.29 |
| GPE (%) | 0.91 | 2.86 |
| RMS (Hz) | 5.46 | 5.83 |

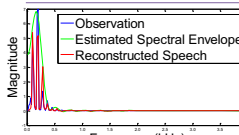**Pitch tracking result with UV labeling**

| | PAT | Get_f0 |
|---|---|---|
| GPE (%) | 1.51 | 2.07 |
| RMS (Hz) | 5.4556 | 5.7792 |

### Speech Synthesis

**MOS grading result**

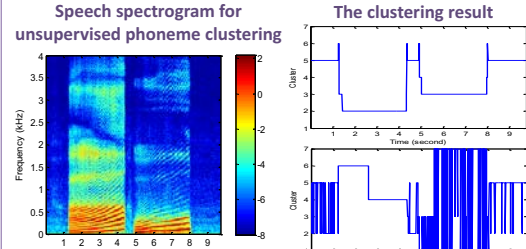| | PAT | others |
|---|---|---|
| Z_SYNTHESIS vs. LPC | 4.33 | 2.21 |
| Z_SYNTHESIS vs. original | 4.37 | 4.69 |
| MU_SYNTHESIS vs. LPC | 3.24 | 2.31 |
| MU_SYNTHESIS vs. original | 3.34 | 4.98 |



$$x_t^{Z\_SYNTHESIS} = \Gamma_{\hat{l}_t} \cdot E[z_t|x_{1:T}]$$
$$x_t^{MU\_SYNTHESIS} = E[a_t|x_{1:T}] \cdot \Gamma_{\hat{l}_t}\mu_{\hat{c}_t}$$

### Phoneme Clustering

**Speech spectrogram for unsupervised phoneme clustering**

**The clustering result**



### Speech Enhancement

**Speech enhancement result for the vowel /ɔː/**



Original signal    Noisy signal (0dB)    Enhanced signal