# Tongue Tracking in Ultrasound Images with Active Appearance Models

## Anastasios Roussos, Athanassios Katsamanis and Petros Maragos
### School of ECE, National Technical University of Athens 15773, Greece

{troussos,nkatsam,maragos}@cs.ntua.gr          http://cvsp.cs.ntua.gr

## 1. Outline

- Ultrasound (US) imaging of speakers' tongues
  - Widely used for human speech production analysis and modeling
  - Captures the shape & dynamics of tongue during speech
  - Simple to use, no radiation, high frame rates
- Automatic tongue tracking
  - Extremely helpful for large datasets of acquired US videos
  - Difficulties: high speckle noise, non-visible tongue parts

### Contributions

- We propose a novel tracking method
  - Built on a variant of Active Appearance Models
  - Incorporates prior about tongue shape variation
  - Bayesian formulation of the tracking
- Properties of the method
  - Robust even in cases of bad tongue visibility
  - Also extrapolates the contour in the non-visible parts
  - Improved performance compared to other, previously proposed techniques

## 2. Preliminaries

- Acquired speech articulation data of the same speaker:
  - Ultrasound imaging @ 66 Hz
  - Electro-Magnetic sensors on US probe & head @ 40 Hz
  - Magnetic Resonance Imaging (static)
  - X-ray videos @ 25 Hz
- Exploitation of the X-rays to model the tongue shape
  - The entire tongue contour is visible, in contrast to US images
  - Usage of a Vocal Tract (VT) grid to represent the tongue shape
    - fixed pose w.r.t. the speaker's palate
    - it bypasses the point correspondence problem
- Estimation of the VT grid's pose at every US frame, using EM sensors data & the head's MRI
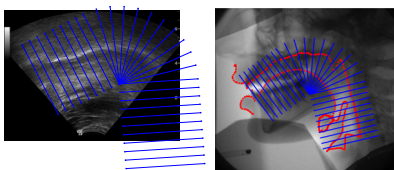


**Fig. 1.** Ultrasound and X-ray images of the speaker, with the registered Vocal Tract grid superimposed.
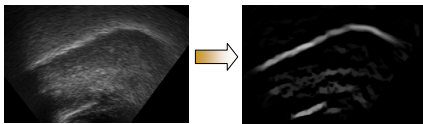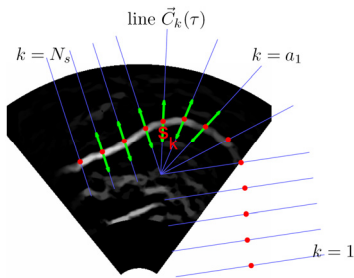
- Preprocessing of the US frames using our method of [2]



**Fig. 2**. Original and filtered US frame, using the method of [2]

## 3. Tongue Appearance Representation



- Tongue appearance
  - Shape
    $$s = [s_1, .., s_{N_s}]^T$$
    - $s_k$ : scalar that determines the intersection of the tongue contour with the grid line $k$
  - Texture
    $$g(s) = \left[ [u_{a_1}(s_{a_1}+t)]_{t \in W}^T \cdots [u_{a_{N_a}}(s_{a_{N_a}}+t)]_{t \in W}^T \right]^T$$
    - Only the texture-active grid lines $G_{act}$ are used for texture, since some parts of the tongue contour are never or rarely visible
    - $W = \{-d, -d+1, .., d\} \cdot \delta\ell$ : sampling window
    - $u_k(\tau) = u(\vec{C}_k(\tau))$ : restriction of the image to grid line $k$
- Differences from classic AAMs
  - Various modifications to exploit application-specific properties
  - Reduced complexity of the appearance representation & model
  - Lighter optimization problem for the model fitting

## 4. Modeling Appearance Variation

- Shape model
  $$s \approx s_0 + Q_s b$$
  - $b$: normalized shape parameters vector with $p(b) = \mathcal{N}(b|0, I_{N_b})$
  - Principal Component Analysis (PCA) to learn $s_0$ and $Q_s$
    - Training vectors from manually annotated tongue contours on 700 X-ray frames
- Texture model
  $$g = g_0 + Q_g \lambda + \varepsilon$$
  - $\lambda$ : texture parameters with $p(\lambda) = \mathcal{N}(\lambda|0, I_{N_\lambda})$
  - $\varepsilon$ : texture reconstruction error with:
    $$p(\varepsilon) = \mathcal{N}(\varepsilon|0, \Sigma_\varepsilon), \quad \Sigma_\varepsilon = \widetilde{Q}_g \mathrm{diag}(\rho_1, .., \rho_{N_g}) \widetilde{Q}_g^T$$
  - Training of the model
    - Manual annotations at 400 US frames. This training set is divided into 2 subsets $T_1$ and $T_2$
    - Subset $T_1$ is used to learn $g_0$ and $Q_g$ using PCA
    - Subset $T_2$ is used to learn the optimum parameters $\rho_1, .., \rho_{N_g}$

## 5. Tongue Tracking

- Tracking via fitting of the appearance model in every US frame
- MAP estimation of shape & texture parameters $b$ and $\lambda$ by maximizing:
  $$p(b, \lambda | u(x,y)) \propto p(u|b, \lambda) p(b, \lambda) = p(\varepsilon) p(b) p(\lambda)$$
  *filtered US frame*          $\varepsilon = g(s(b)) - g_0 - Q_g\lambda$
- Equivalently: minimization of the energy:
  $$E(b, \lambda) = -\ln p(b, \lambda | u)$$
  $$= C + \tfrac{1}{2}\left\{ \|b\|^2 + \|\lambda\|^2 + \varepsilon^T \Sigma_\varepsilon^{-1} \varepsilon \right\}$$
- Gradients of the energy:
  $$\nabla_b E = b + Q_s^T (\partial g / \partial s)^T \Sigma_\varepsilon^{-1} \varepsilon$$
  $$\nabla_\lambda E = \lambda - Q_g^T \Sigma_\varepsilon^{-1} \varepsilon$$
  where:
  $$\frac{\partial g}{\partial s_k} = \begin{cases} [0 \cdots\cdots 0]^T, & \text{if } k \notin G_{act} \\ \left[ \underbrace{0 \cdots 0}_{(k-1)N_W} [u'_k(s_k+t)]_{t \in W}^T \underbrace{0 \cdots 0}_{(N_s-k)N_W} \right]^T, & \text{if } k \in G_{act} \end{cases}$$
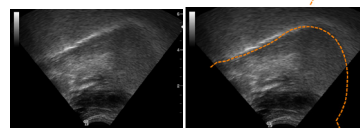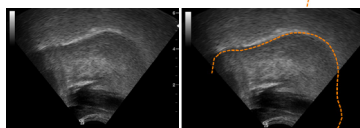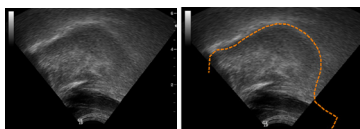- Optimization algorithm
  - Gradient descent
  - Initialization of parameters:
    - $b_0$: from previous frame result
    - $\lambda_0$: maximization of the posterior $p(\lambda|g(s(b_0)))$

## 6. Experimental Results

### I. Results of the proposed method





*Original US frames*          *Same frames + extracted contour*

**Fig. 3.** Tongue tracking & extrapolation in a US image sequence, using the proposed method.

- Parameters of the method:

|  | Dimensionality of original vector | Number of model parameters | Variance explained (% of the total) |
|---|---|---|---|
| Shape | 30 | 6 | 96% |
| Texture | 1215 | 35 | 93% |

### II. Comparisons with other tracking methods
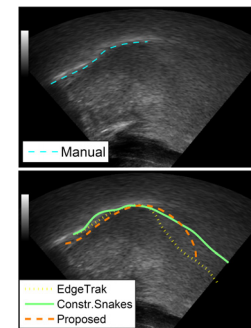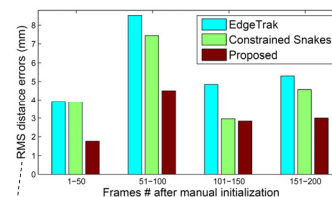


- - - Manual

EdgeTrak
Constr.Snakes
Proposed

**Fig. 4.** Frame from a sequence where the tongue tracking methods have been applied. *Top*: manually annotated contour. *Bottom*: comparison of the methods' results.

- Quantitative evaluation



$$e_d = \sqrt{(d_{om}^2 + d_{mo}^2)/2}$$

where $d_{om}$ ($d_{mo}$) is the RMS distance of the points of the output (manual) contour from the manual (output) contour.

## References

[1] M. Li, X. Khambhamettu, and M. Stone, "Automatic contour tracking in ultrasound images," *Clinical Linguistics and Phonetics*, vol. 6, no. 19, pp. 545–554, 2005.

[2] M. Aron, A. Roussos, M.O. Berger, E. Kerrien, and P. Maragos, "Multimodality Acquisition of Articulatory Data and Processing," in *Proc. EUSIPCO*, 2008.

[3] T. F. Cootes, G. J. Edwards, and C. J. Taylor, "Active appearance models," *IEEE Trans. PAMI*, vol. 23, no. 6, pp. 681–685, 2001.

[4] G. Papandreou and P. Maragos, "Adaptive and constrained algorithms for inverse compositional active appearance model fitting," in *Proc. CVPR*, 2008.

[5] M. Aron, A. Toutios, M.-O. Berger, E. Kerrien, B. Wrobel-Dautcourt, and Y. Laprie, "Registration of multimodal data for estimating the parameters of an articulatory model," in *Proc. ICASSP*, 2009.

[6] S. Maeda, *Compensatory articulation during speech: evidence from the analysis and synthesis of vocal tract shapes using an articulatory model*, chapter in Speech Production and Speech Modeling, pp. 131–149, Kluwer, 1990.