

## 1. Introduction

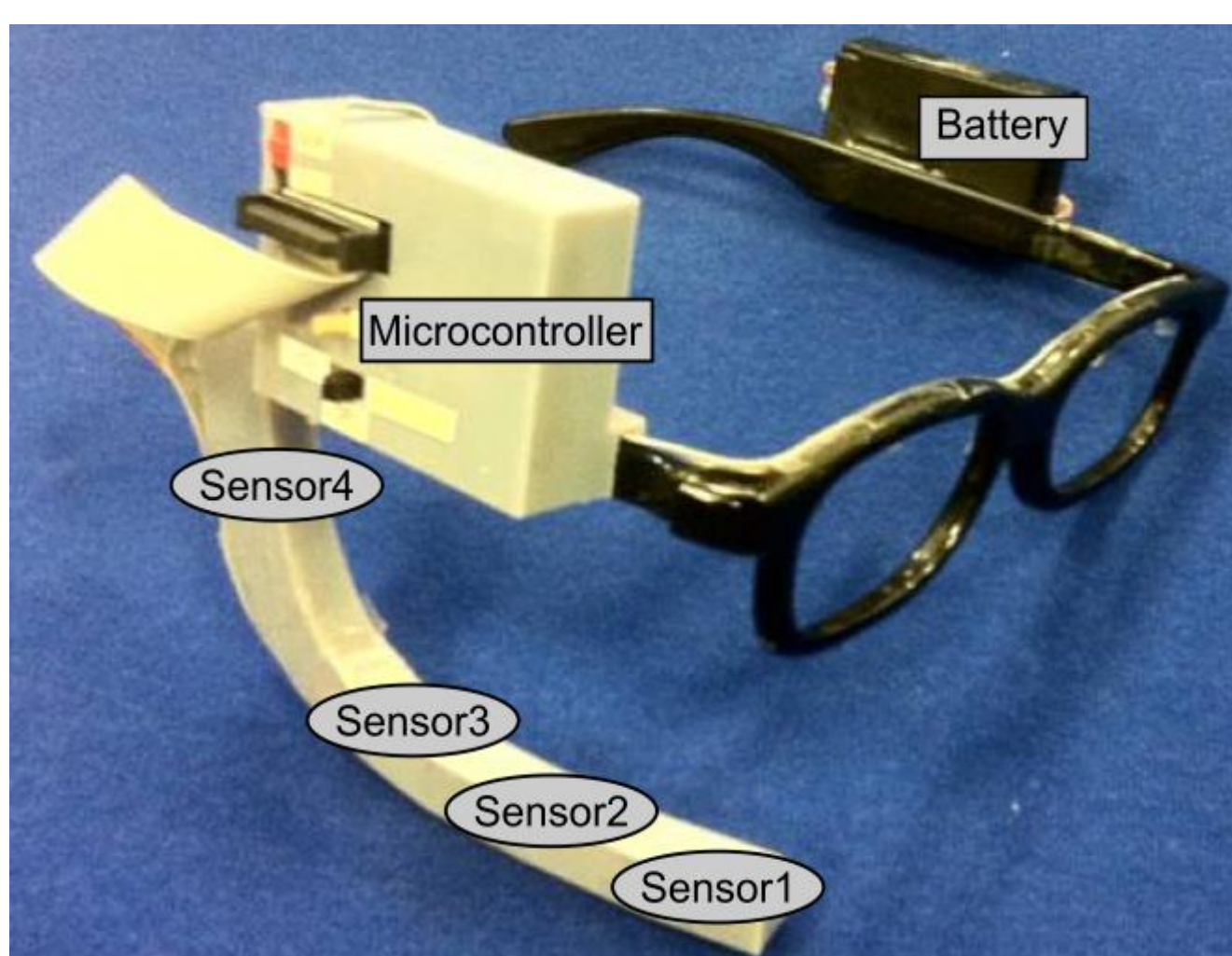
### Silent Speech Interfaces (SSIs)

- **Motivation**
  - Patients with larynx cancer often lose their voice after *laryngectomy*.
  - Existing methods for voice restoration are unsatisfactory.
  - SSIs enable speech communication when the audible acoustic signal is unavailable by exploiting other *speech-related biosignals*.
  - Devices for capturing articulator motion data: cameras, ultrasound, surface electrodes or **PMA**.
- **SSI approaches**
  - a) ASR from articulator motion data + TTS synthesis.
  - b) **Direct transformation of the articulator data to audible speech.**

### About this Work

- **Summary**
  - In previous work we have shown that it is possible to recognise speech from PMA data.
  - **Here, we investigate the use of shared Gaussian process dynamical models (SGPDMs) for articulatory-to-acoustic conversion.**
  - Results are reported in which audible speech is synthesised from PMA data for two speakers with no speech impairment.
  - Preliminary results are very promising, outperforming state-of-the-art GMM-based conversion, but further investigation is needed.
  - The ultimate goal is to **restore the ability to communicate to laryngectomees.**

## 2. Permanent Magnet Articulography (PMA)



### How PMA works

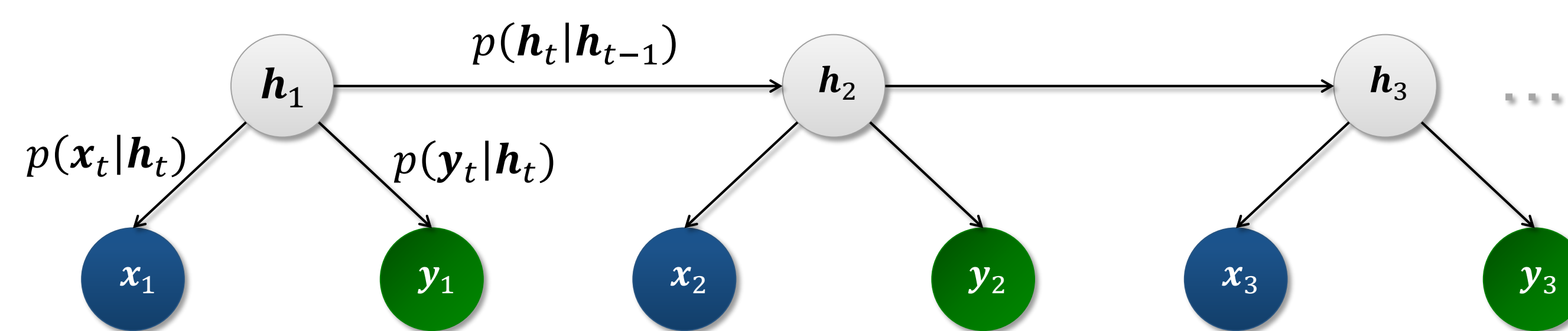
- Small magnets are attached to the lips and tongue of the patient.
- The magnetic field generated when the patient 'speaks' is captured by the magnetic sensors.
- PMA does not provide the exact position of the magnets.

## 3. Articulatory-to-Acoustic Mapping

### Problem Formulation

- **Generative model**
  - Mapping between PMA vectors  $x_t$  and speech parameter ones  $y_t$ :
$$y_t = f(x_t)$$
  - We assume that  $x_t$  and  $y_t$  are the outputs of an underlying stochastic process with hidden state  $h_t$ :
$$x_t = f_x(h_t) + \epsilon_x$$

$$y_t = f_y(h_t) + \epsilon_y$$



- Two problems
  - **Training:** estimation of  $p(x_t|h_t)$ ,  $p(y_t|h_t)$  and  $p(h_t|h_{t-1})$ .
  - **Conversion:** estimate the most likely sequence of speech parameter vectors for the source sequence  $\mathbf{X} = (x_1, x_2, \dots)$ .

### Shared Gaussian Process Dynamical Models

#### Statistical modelling

##### Gaussian processes

$$p(\mathbf{z}|\mathbf{h}) = N(m(\mathbf{h}), k(\mathbf{h}, \mathbf{h}'))$$

$m(\mathbf{h})$  and  $k(\mathbf{h}, \mathbf{h}')$  are the mean and covariance (kernel) functions.

- In a SGPDM we have several GPs sharing the same latent space + a dynamical model in the latent space.

##### Data modelling

- $m(\mathbf{h}) = \mathbf{0}$ .

- An RBF kernel is used for the observation models:

$$k_x(\mathbf{h}, \mathbf{h}') = \beta_1 \exp\left(-\frac{\beta_2}{2} \|\mathbf{h} - \mathbf{h}'\|^2\right) + \frac{\delta_{\mathbf{h}, \mathbf{h}'}}{\beta_3}$$

- For the dynamical model, we use an RBF+linear kernel:

$$k_H(\mathbf{h}, \mathbf{h}') = \alpha_1 \exp\left(-\frac{\alpha_2}{2} \|\mathbf{h} - \mathbf{h}'\|^2\right) + \alpha_3 \mathbf{h}^T \mathbf{h}' + \frac{\delta_{\mathbf{h}, \mathbf{h}'}}{\alpha_4}$$

#### Training and conversion phases

Training	Conversion
<ul style="list-style-type: none"> <li>• <b>Model parameters:</b> kernels hyperparameters <math>\{\alpha, \beta, \gamma\}</math> and shared latent coordinates <math>\mathbf{H}</math>.</li> <li>• <b>Loss function:</b> <math display="block">\mathcal{L} \equiv p(\mathbf{H}, \alpha, \beta, \gamma   \mathbf{X}, \mathbf{Y}) \propto p(\mathbf{X}   \mathbf{H}, \alpha) p(\mathbf{Y}   \mathbf{H}, \beta) p(\mathbf{H}   \gamma) p(\alpha) p(\beta) p(\gamma)</math> </li> <li>• Uninformative priors are chosen for <math>\{\alpha, \beta, \gamma\}</math>.</li> <li>• <math>\mathcal{L}</math> is optimised using the SCG algorithm.</li> <li>• <math>\mathbf{H}</math> is initialised using canonical correlation analysis (CCA).</li> </ul>	<ul style="list-style-type: none"> <li>• The latent sequence <math>\mathbf{H}^*</math> for <math>\mathbf{X}</math> is initialised using the Viterbi algorithm. <ul style="list-style-type: none"> <li>• Transition probabilities given by <math>p(\mathbf{h}_t   \mathbf{h}_{t-1})</math>.</li> <li>• Observation probabilities given by <math>p(x_t   \mathbf{h}_t)</math>.</li> </ul> </li> <li>• Next, <math>\mathbf{H}^*</math> is refined using the SCG algorithm.</li> <li>• Finally, <math>\hat{\mathbf{Y}}</math> is just the mean of <math>p(\mathbf{Y}   \mathbf{H}^*)</math>.</li> </ul>

## 4. Experiments

### Conditions

<b>Database</b>	<ul style="list-style-type: none"> <li>– Isolated digits.</li> <li>– PMA and speech data were recorded simultaneously.</li> <li>– Two native English speakers (with no speech impairment): male &amp; female.</li> <li>– Amount of data: 7.2 minutes (male) &amp; 8.46 minutes (female speaker).</li> </ul>
<b>Feature extraction</b>	<ul style="list-style-type: none"> <li>– Speech signal: 25 MFCCs computed every 10ms [Fs:16kHz, window length:20ms].</li> <li>– PMA signal: features extracted by Partial Least Squares [9 channels @ 100Hz].</li> <li>– Speech is synthesised with no voicing (i.e. as whispered speech).</li> </ul>
<b>Objective evaluation</b>	<ul style="list-style-type: none"> <li>– The Mel-Cepstral distortion measure is used to evaluate reconstruction accuracy.</li> <li>– 10-fold cross-validation scheme is used.</li> <li>– SGPDM mapping is compared with GMM-based mapping proposed by Toda'2007. <ul style="list-style-type: none"> <li>– 32-component GMM is employed.</li> <li>– Both MMSE and MLE estimation algorithms are evaluated.</li> </ul> </li> </ul>

### Results

#### Experiment 1

- Conversion is performed using a model trained on the same digit.

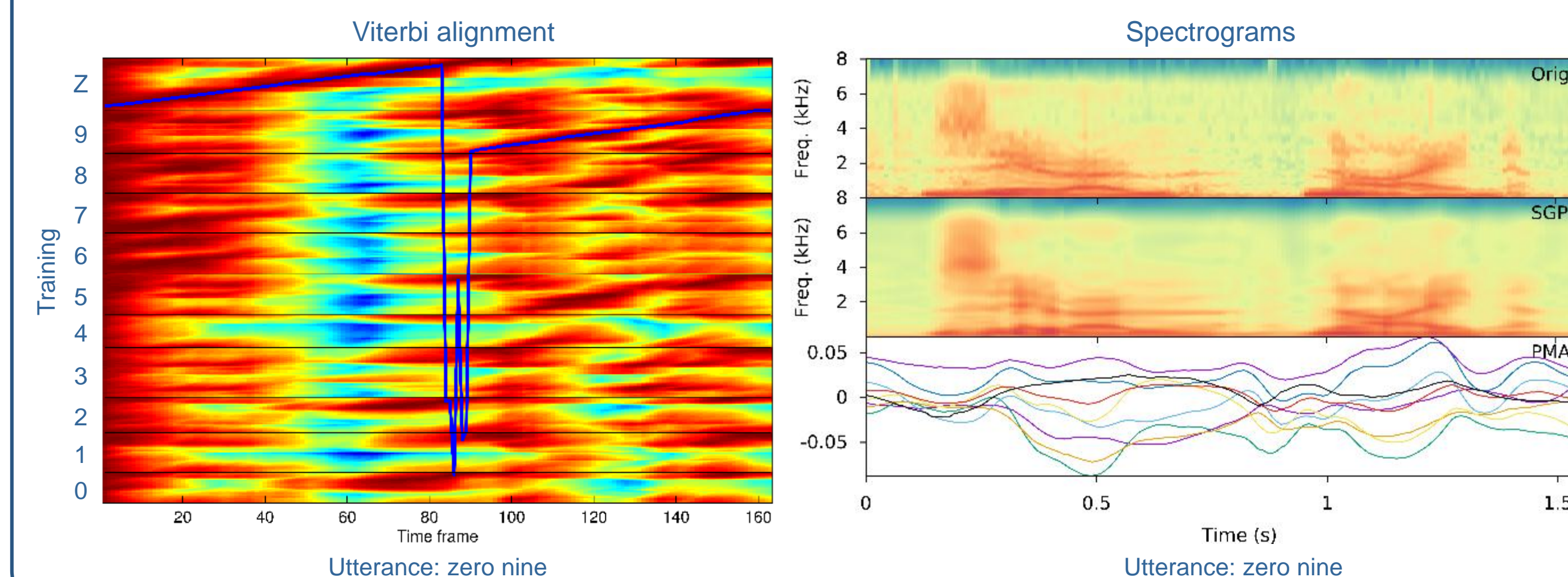
Speaker	GMM		SGPDM		
	MMSE	MLE	$D_h = 3$	$D_h = 5$	$D_h = 7$
Male	5.71	5.04	4.37	4.64	4.72
Female	5.99	5.92	4.70	4.89	5.01
<b>Average</b>	<b>5.85</b>	<b>5.48</b>	<b>4.54</b>	<b>4.77</b>	<b>4.87</b>

#### Experiment 2

- The transformation is now estimated from sequences of isolated digits.

Speaker	GMM		SGPDM		
	MMSE	MLE	$D_h = 3$	$D_h = 5$	$D_h = 7$
Male	5.04	5.05	4.74	5.22	5.05
Female	5.57	5.64	4.82	5.71	5.97
<b>Average</b>	<b>5.31</b>	<b>5.35</b>	<b>4.78</b>	<b>5.47</b>	<b>5.51</b>

#### Example: Digit sequence reconstruction



## 5. Conclusions

- We have presented a non-parametric approach for articulatory-to-acoustic conversion using shared Gaussian process dynamical models.
- Results demonstrate that the approach outperforms state-of-the-art mapping based on GMMs.
- Future research
  - Evaluation: more complex task & more speakers.
  - Model: introduce switching states & evaluate other kernel functions.