# Speech-Driven Realtime Lip-Synch Animation with Viseme-Dependent Filters

*Supplemental Material for SIGGRAPH 2013 Poster*　　　　Shin-ichi KAWAMOTO (JAIST)

JAIST
JAPAN ADVANCED INSTITUTE OF SCIENCE AND TECHNOLOGY
1990

## Introduction

**Motivation**

Lip-synching is one of fundamental components for creating facial animation. Mouth movement is synchronized along with the speech, when a character utters a word or phrase. Especially, speech-driven realtime lip-synching animation is useful for helping speech communication.

**Aim**

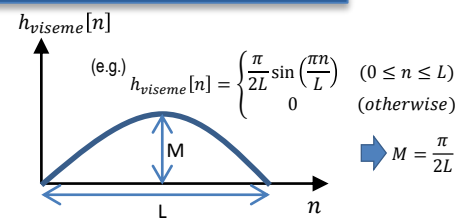Realizing speech-driven realtime lip-synching based on blendshapes, linear shape interpolation model.

**Problem**

Simple solution is to construct a mapping between speech and mouth-shape directly. These direct mapping approaches can realize lip-synching with small delay. However it is sometimes unnatural since mouth movement is mismatched between the speaker and the pre-designed characters.

**Our solution**

we consider customization of mouth movement by viseme-dependent filters designed for each mouth shape of given characters.
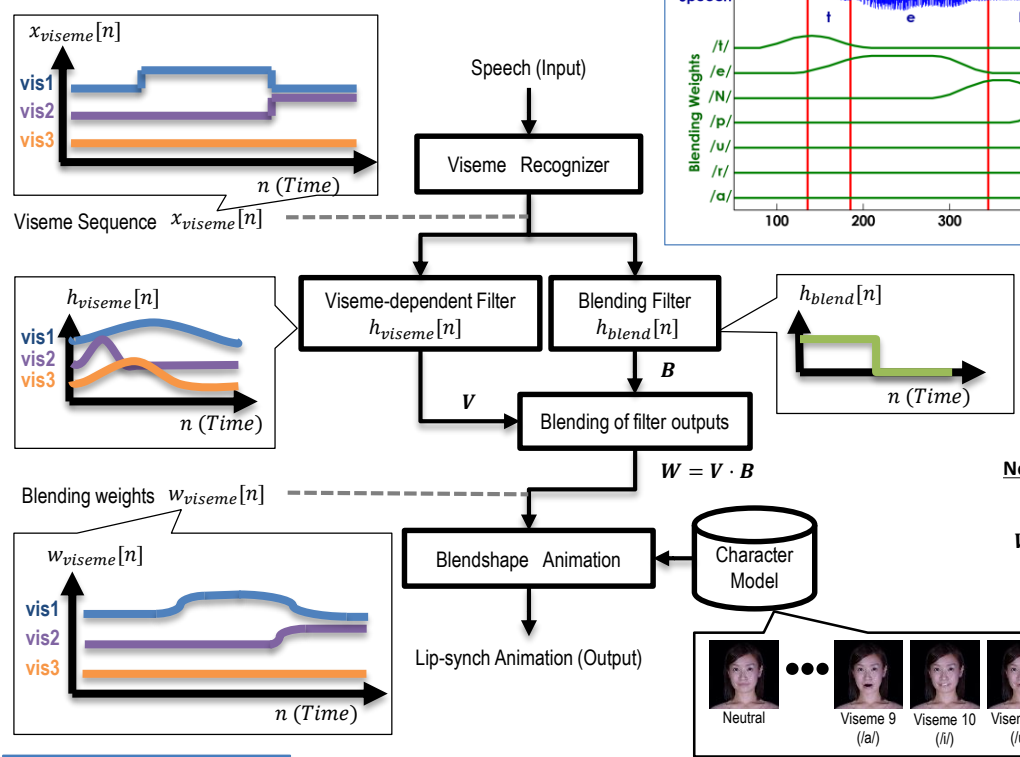
## Viseme-dependent Filter

$$h_{viseme}[n]$$

(e.g.) $h_{viseme}[n] = \begin{cases} \dfrac{\pi}{2L}\sin\left(\dfrac{\pi n}{L}\right) & (0 \le n \le L) \\ 0 & (otherwise) \end{cases}$

$M = \dfrac{\pi}{2L}$

L : filter length (#tap)
$\Rightarrow$ Transition Time between Mouth Shapes

$M = \max_n (h_{viseme}[n]) = \dfrac{\pi}{2L}$

$\propto$ Maximum Mouth Movement Speed

**How to decide this parameters?**

**[default]** depends on Euclid distance between a target-shape and a neutral-shape (*normalize maximum mouth movement speed*)

**[customize]** depends on the desirable transition time (*eg. Consonants /p/, /b/, and /m/*)

## Processing Flow

$x_{viseme}[n]$

vis1　vis2　vis3

$n$ (Time)

Viseme Sequence $x_{viseme}[n]$

Speech (Input)

Viseme Recognizer

$h_{viseme}[n]$

vis1　vis2　vis3

$n$ (Time)

Viseme-dependent Filter $h_{viseme}[n]$　　Blending Filter $h_{blend}[n]$

$h_{blend}[n]$

$V$　　　　$B$

Blending of filter outputs

$W = V \cdot B$

Blending weights $w_{viseme}[n]$

$w_{viseme}[n]$

vis1　vis2　vis3

$n$ (Time)

Blendshape Animation　　←　Character Model

Lip-synch Animation (Output)

Neutral　Viseme 9 (/a/)　Viseme 10 (/i/)　Viseme 11 (/u/)

**Constraints of Blending Filter**

$$\sum_n h_{blend}[n] = 1, \quad h_{blend}[n] > 0$$

**Notations**　Note: '*' means a convolution operator

$$V = \begin{bmatrix} h_{vis_1} * x_{vis_1} & h_{vis_2} * x_{vis_1} & \cdots & h_{vis_N} * x_{vis_1} \\ h_{vis_1} * x_{vis_2} & h_{vis_2} * x_{vis_2} & \cdots & h_{vis_N} * x_{vis_2} \\ \vdots & \vdots & \ddots & \cdots \\ h_{vis_1} * x_{vis_N} & h_{vis_2} * x_{vis_N} & \cdots & h_{vis_N} * x_{vis_N} \end{bmatrix}$$

$$B = \begin{bmatrix} h_{blend} * x_{vis_1} \\ h_{blend} * x_{vis_2} \\ \vdots \\ h_{blend} * x_{vis_N} \end{bmatrix} \quad W = \begin{bmatrix} w_{vis_1} \\ w_{vis_2} \\ \vdots \\ w_{vis_N} \end{bmatrix}$$

## Example of Our Result (Speech & Blending Weights)

Speech

t　e　N　p　u　r　a

Blending Weights: /t/ /e/ /N/ /p/ /u/ /r/ /a/

Time [msec]: 100　200　300　400　500　600　700　800　900
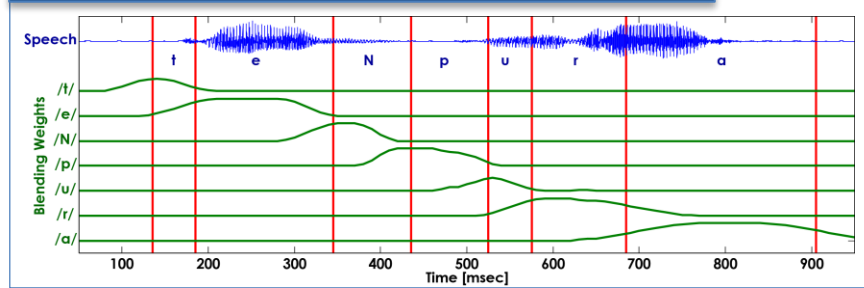
## Viseme Recognizer

Viseme recognizer iterates the following processes at a constant frequency:

**[Step.1]** Estimates phonemes and its duration by using an HMM-based speech recognizer from observing input speech at this time.

**[Step.2]** Convert phoneme to viseme based on table lookup.

Note: Viseme is a basic unit of mouth shapes that are classified visually

Phoneme-Viseme Mapping Table (Japanese)

| Viseme ID | Phonemes |
|---|---|
| 1 | r, ry |
| 2 | b, by, m. my, p, py |
| 3 | t |
| 4 | d, n, ny |
| 5 | g, gy, hy, k, ky, N |
| 6 | f |
| 7 | ch, dy, j, s, sh, ts, z |
| 8 | w |
| 9 | a |
| 10 | I |
| 11 | u |
| 12 | e |
| 13 | o |

## Result & Discussion

Our lip-synch system worked well with about 0.3 sec of delay from the input speech.

Factors of its delay:

(1) Viseme recognition also needs some processing time.

Delay　←tradeoff→　Accuracy of viseme recognition results

(2) Length of the filters for generating blending weights is also related to the delay of the animation

Delay　←tradeoff→　Smoothness of mouth shape transition

(3) Minimum delay of speech output depends on the hardware specification (*eg. delay of sound device*)