



Supported by NIH grant R01DC007124 and NSF grants 1514544, 1908865

# Simulating anticipatory coarticulation in VCV utterances with a gestural articulatory synthesizer

Asterios Toutios and Shrikanth Narayanan

#### Introduction

**Context:** Development of a framework for articulatory speech synthesis from gestural score specifications directly informed by real-time MRI data (Alexander et al., 2019)

### A Simple Loop

Input: w[0], w[1],  $\omega_o[n], z_o[n], n = 2...N$ for n=2...N do Find cluster where  $\mathbf{w}[n-1]$  lies;

- Proposing a modified/enhanced gestural score
- Describing the mechanism by which such gestural score animates (a midsagittal slice of) the vocal tract
- Addressing gestural overlap ullet
- Test case: coarticulation in /adu/, as observed by Öhman (1966)  $\bullet$

## **Articulatory Model and Forward Map**



Retrieve  $\mathbf{z}_{\mathbf{c}}, F$  and J's for that cluster ; Calculate  $K(\boldsymbol{\omega}_o), B(\boldsymbol{\omega}_o);$ Solve system for  $\mathbf{w}[n]$ end

## **Simulating Anticipatory Coarticulation**

Replicating Alexander et al. (2019), /adu/ without gestural overlap (Note: In that work, targets for vowels were defined in parameter space, rather than constriction space)



JUIL SPACE IS INDUCIED AS A UNION OF CLUSTERS WHELE THE forward *map* from parameters **w** to constrictions **z** is approximately linear, i.e.  $\mathbf{z} = \mathbf{G}(\mathbf{w}) = F * \mathbf{w} + \mathbf{z}_{\mathbf{c}}$ 

#### **Discretization of Dynamical Systems**

From Saltzman and Munhall (1998):

$$\ddot{\mathbf{w}} = J^* (-BJ\dot{\mathbf{w}} - K(\mathbf{G}(\mathbf{w}) - \mathbf{z}_0)) - J^* \dot{J} \dot{\mathbf{w}}$$
$$- (I_8 - J^* J) B_N \mathbf{w} - G_N (-B_N \mathbf{w} - K_N \mathbf{w})$$

Consider the sequence of arrays w[n] at a rate h (e.g. 1msec), replace derivatives by finite differences, and after some algebra: Introducing anticipatory coarticulation, inspired by Ohman, leads to discernible difference in vocal-tract shaping dynamics (and also

acoustics):	pharyngeal	(18.42, 5.76)		(15.3	5, 13.92)	(1	(18.42, 13.44)		
	velar	(18.42, 11.76)		(15.35, 5.04)		(1	(18.42, 13.68)		
	palatal	(18.42, 1	8.96)	(15.35, 2.88)		(1	(18.42, 14.16)		
	alveolar	(18.42, 1	7.04)	(30.70, -4.80)	(30.70, 4.80)	(	(18.42, 7.44)		
	velopharyngeal	(18.42, 0.48) (18.42, 12.00)		(15.35, 1.20)		(	(18.42, 2.88)		
	bilabial			(15.3	35, 0.48)	(	(18.42, 6.24)		
	C	0 100	200	300	400 500	600	700	800	



## **Another Example: /span/**



 $(I_8 + hA_1 + h^2A_2)\mathbf{w}[n] =$ 

 $\mathbf{w}[n-2] + 2\mathbf{w}[n-1] + hA_1\mathbf{w}[n-1] + h^2J^*K(\mathbf{z}_o - \mathbf{z}_c)$ with:  $A_1 = -J^*BJ - J^*\dot{J} - B_N + J^*JB_N - G_NB_N$ 

 $A_2 = -G_N K_N - J^* KF$ 

- Given a cluster we know F, J and its derivative, and we can invert (because the map in the cluster is linear) to get J\*
- **z**<sub>o</sub> is a 6-dimensional array of **targets** and B, K are simple functions of a 6-dimensional array of natural frequencies  $\omega_{0}$
- $G_N$  and  $K_N$  are constants (neutral attractor)

#### **Future Work**

- Design gestural scores for more utterances
- Optimize scores to exactly fit recorded real-time MRI data