# A Coordinated Muscle Model for Speech Animation

## K. Waters and J. Frisbie[1]

Digital Equipment Corporation, Cambridge Research Lab, One Kendall Square, Cambridge MA 02139

## Abstract

*During the production of speech, coordinated muscle contractions generate a wide variety of lip and mouth postures. These postures are orchestrated by complex muscle interactions that are difficult to capture using traditional computer animation algorithms and techniques. We present a framework for producing natural-looking speech on a facial image focusing on the muscles around the mouth.*

**Keywords:** *Facial Animation, Speech Animation, Muscle Model, Muscle-Based, Physics-Based*

## 1 Introduction

During the past twenty years, efforts directed toward animating facial expressions have been more successful than those directed toward animating speech. The hardest challenge in facial animation has been to develop computational models of the mouth capable of synthesizing realistic-looking speech. Accurate speech production would allow not only believable facial animation, but also aids for the hearing impaired, an additional communication channel in noisy environments, and perhaps even surrogate videophone images for automated services. This paper develops a different approach to the state-of-the-art by creating a framework for investigating speech animation using muscle-based models.

The complexity of the oral anatomy coupled with the enormous range of motion makes lip animation difficult. The following factors make speech harder to animate than facial expression:

- *Separability.* For the most part, facial expressions can be created by decoupled facial actions. For example, the elevation of the mouth corners during smiling do not influence the upper regions of the face [EF77]. Conversely, when the upper brow

is raised in surprise, the lower face is unaffected. This is in sharp contrast to the interaction between the lips, jaw and cheeks during speech, where every muscular contraction has some influence on the mouth shape [WWDB89].

- *Cataloging.* Facial expressions, and their meanings, have been categorized and cataloged by psychologists, thereby providing a valuable resource for facial expression modeling [EF77]. In contrast, identifying the important characteristics of visual speech is an open research question [Sum92]. It is difficult because the physical realization of an utterance can vary widely depending on many factors, including the overall speaking rate, how clearly the speaker is speaking, stress and emphasis, and loudness. In addition, the discrete phonemes of an utterance are translated into sequences of overlapping articulations. The rules underlying this coarticulation are not well understood [KM77, Per95] and are not simply inertial effects. For example, the lips are already quite rounded when 't' is pronounced in the word 'too' anticipating the following vowel as compared with the 't' in 'tee'.

- *Time Sequencing.* Whereas facial expressions denote emotional states, which change slowly and persist for long durations, speech requires continuous and rapid movements in a precise time sequence. The timing specification for different phonemes may be different: for vowels, it may be sufficient to move toward a prototypical lip configuration without actually achieving it as long as a certain width and height ratio is reached. However, for labial consonants such as 'p', 'b', and 'm', it is critical that that lips close in order for the motion to be perceived correctly [MM86].

In this paper we present a novel approach to animated speech by focusing on the muscles around the mouth. We use a muscle-based approach to implement a two-dimensional model of the mouth that emulates motion dynamics using a physics-based system. This straightforward approach can be incorporated into larger-scale

[1]Current address: Sensory Communication Group Research Lab of Electronics, MIT 36-759, Cambridge MA 02139

systems such as DECface, a real-time lip-synchronized synthetic face [WL94]. Ultimately, we believe that a muscle-based approach for facial articulation will produce results superior to the ad-hoc approaches developed to date. Furthermore, it provides a cohesive and consistent interface to other muscle-based animation systems [WT92].

## 1.1 Overview

The remainder of the paper is organized as follows: in Section 2, we describe previous work on computer-generated faces. This is followed by a brief description of the musculature of the lower face in Section 3. Next, Section 4 describes our muscle-based framework and Section 5 describes the system we implemented. In Section 6, we present an example utterance, and finally in Section 7, we discusses our results and suggestions for future work. Appendix I describes the mathematics necessary to implement a physics-based system and Appendix II describes a bilinear warping scheme used for texture mapping.

## 2 Previous Work

Facial animation systems of the mouth can be divided into two categories: parametric and physics-based systems.

### 2.1 Parameteric Systems

Parametric systems manipulate geometric primitives, such as polygons [Par74, CM93] or spline surfaces [NHS88], to create facial images. Expressions, including utterances, are created by displacing a small number of vertices or control points to create the desired facial posture. In principle, displacements must be chosen for all vertices or control points for each different image. In practice, vertices are grouped into sets which are controlled by a single parameter [Par72] or procedure [MPT88]. Animations are generated by interpolating between geometric configurations specified for key-frames.

Parametric systems have some advantages: 1. they are not computationally demanding; 2. any particular facial configuration can be reproduced explicitly; 3. the technology is mature – public domain and commercial programs exist, and libraries of expressions and utterances have been tuned over a period of years. A central weakness of parametric systems is their inability to blend expressions. For example, at the lowest level, if a vertex position is controlled by two parameters, an *ad-hoc* mechanism for resolving the conflict must be devised. This rarely results in a natural motion or configuration for all possible expression sequences. Consequently, parameters typically operate over very local-

ized regions that are controlled independently. Unfortunately, this introduces artificial boundaries where deformations abruptly end.

## 2.2 Physics-based Systems

Originally, physics-based facial animation systems were developed to address the interacting parameters problem. Generally, mass-spring or finite element networks are used to model the elastic properties of skin [Pla80, TW90, Pie89]. In these models, parameters, instead of specifying displacements for geometric primitives, specify forces that are applied to the network, thereby deforming it. Consequently, conflicts between parameters are resolved because forces are summed together and deformations blend naturally between and across regions as the forces are distributed by the network. Unfortunately, physics-based systems tend to be computationally expensive, especially if an elaborate elasticity model is used, or if the network resolution is increased in an effort to improve graphical quality.

## 3 Musculature of the Lower Face

The musculature of the lower face is perhaps the most complex in the human body. Unlike most other muscles, which originate and insert in bone, facial muscles insert into skin or into other muscles. The *orbicularis oris* encircles the mouth, lying beneath and around the lips. Other muscles merge into the *orbicularis oris*. The intertwining of individual muscle fibers makes it difficult to determine the exact structure and precise function of these muscles [WWDB89]. The principle muscles of the mouth are shown schematically in Figure 1.

The *orbicularis oris* determines the shape of the lips and mouth opening directly by independently contracting its various sections – inner and outer, left and right, upper and lower – and indirectly by varying its stiffness, thus affecting the propagation of forces applied by other muscles. The other muscles either attach or merge with the *orbicularis oris* or lips radially, or they attach radially to the *modiolus*, a mobile, fibro-muscular mass found at the corner of the mouth. Although the radial muscles overlap, in general they occur at different depths from the skin and are free to slide over one another; their actions are only coupled where they merge at the *orbicularis oris* or *modiolus*.

## 4 A Muscle-Based Framework

Our approach follows from the observation that muscles are the active elements driving expressions and articulations, whereas the surface is just passively carried along. The skin's elastic properties can be incorporated into the muscle's elastic behavior. Skin does not oscillate

Figure 1: The facial muscles of the lower face and a table that describes muscle origins and insertions.

| Muscle | Origin | Insertion |
|---|---|---|
| Obicularis oris | (from other facial muscles) | lips |
| Levator labii superioris alaeque nasi | Frontal process of maxilla | Upper lip |
| Levator labii superioris | Margin maxilla and zygomatic bone | Upper lip |
| Zygomatic major | Zygomatic bone | Angle of mouth |
| Zygomatic minor | Zygomatic bone | Angle of mouth |
| Risorius | Masseteric facia | Angle of mouth |
| Depressor anguli oris | Oblique line of mandible | Angle of lower lip |
| Depressor labii inferioris | Oblique line of mandible | Lower lip |
| Buccinator | Mandible | Angle of mouth |
| Mentalis | Incisive fossa of mandible | Skin of chin |



Figure 2: Framework.

around the muscles, nor does it under- or overshoot the configuration specified by the muscles. Therefore, we explicitly model the interactions between facial muscles and let the skin surface follow the muscle behavior.

The framework, illustrated in Figure 2, is broken down into two independent parts: the model and the display. The model, in our instantiation, is a physical model coupled via a separate module to the display function that renders the mouth using texture mapping.

# 5 Implementation

To validate our ideas, we implemented a relatively simple mass-spring model with the goal of producing realistic animated speech. The model's mathematical details are developed in Appendix I. We implemented a mass-spring model because of its simplicity and because we had previously written code for similar systems. The framework, however, supports any type of muscle model, for example finite element or spline based.

Although the teeth and tongue are important for visually perceiving speech gestures, we were most interested in capturing the behavior of lips; consequently, we did not include them in our current system.

## 5.1 A Muscle Model

The muscle model is a simplified version of the real anatomy. Geometrically, the upper and lower halves of the *orbicularis oris* are represented by two cross-braced bands. The muscle geometry is specified independently of the model resolution and image geometry. We used this feature to explore different tradeoffs between model complexity and image realism. We found that a relatively crude resolution (and hence computationally speedy model) was sufficient. We achieved acceptable results using the schematized muscles shown in Figure 1. The number of segments around the *orbicularis oris* can be varied, and the overall dimensions are chosen to match the dimensions of the facial image to be animated. A radial muscle attached to the *orbicularis oris* is specified by three parameters: its attachment arc on the *orbicularis oris*, its length, and its direction.

There is one control parameter in addition to the muscle contractions: vertical jaw displacement. One displacement is sufficient because horizontal motions typically occur only during chewing, and protrusion isn't relevant for a two-dimensional model. Opening the jaw displaces the origins of the muscles attached to the lower half of the *orbicularis oris*, thus causing forces to be applied to the lower portion of the *orbicularis oris*.

## 5.2 Image Coupling and Display

We create an image from the model by texture mapping a polygonal surface geometry that is linked to the
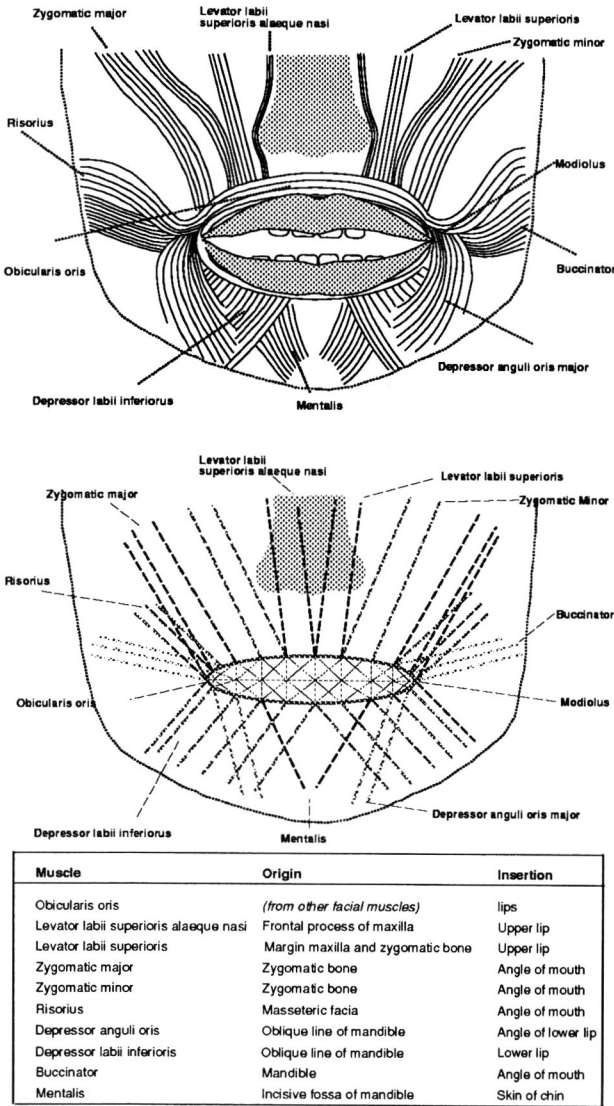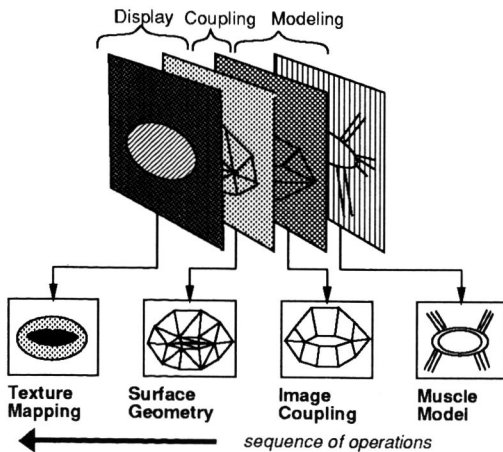
muscle model via a bilinear warping function [Wol91]. As shown in Table 1(a), two sets of quadrilaterals are generated around the upper and lower halves of the *orbicularis oris*. Surface polygon vertices (shown in Table 1(b)) are assigned $(u, v)$ coordinates in the quadrilateral that encloses them. During the simulation, the quadrilateral corners lying in the *orbicularis oris* move, and new surface vertex positions are then calculated based on the deformed quadrilaterals new shape. Appendix II describes the mathematical details of the bilinear warping.

## 5.3 Parameter Determination

To generate the control parameters for our model, we began by video taping a real person saying:

1. *"What's up Doc?"*

2. *"When do the boats come into town?"*

3. *"How did the butcher cut the red beef?"*

In each frame of the sequences, seven $x - y$ pairs, the left and right lip corners, the upper and lower vermillion margins, the upper and lower lip edges and the chin base were digitized (see Table 2). The Nelder-Mead optimization procedure [PFTV86] then determined the input vector which, when applied to the model, best matched the digitized data. Best was defined as the minimum total distance between the reference points on the digitized frames and the corresponding points on the model.

## 6 Example

Table 3 contains ten images generated by our system while reproducing the utterance *"What's up doc?"*. The first and last frames show the model in its at-rest position. The second frame occurs during the initial inhalation just prior to speaking. The third frame shows the lips puckering for the *w* and the next three show the transition from the vowel *uh* through the consonants *ts* to the second vowel *uh*. Frame seven shows the closure for *p* and frame eight the large oral opening for *ah*. The ninth frame coincides with *k* being produced while the mouth is closing.

In figures 3, 4, and 5, displacements are plotted against time in order to compare the model's performance to the actual data measured from the video recording. The first reference point is centered horizontally on the lower lip at the margin between the vermillion and flesh colored regions. The second reference point is the right corner of the lips. The displacements shown were chosen because they exhibit the largest movements and are representative. In each figure the solid line is the measured data and the two sets of
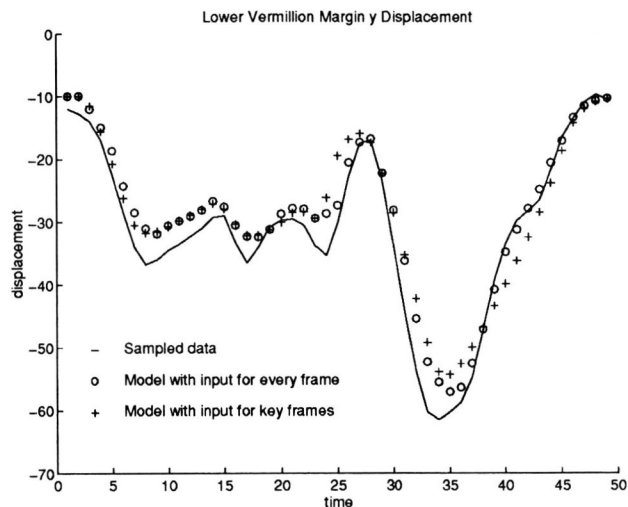
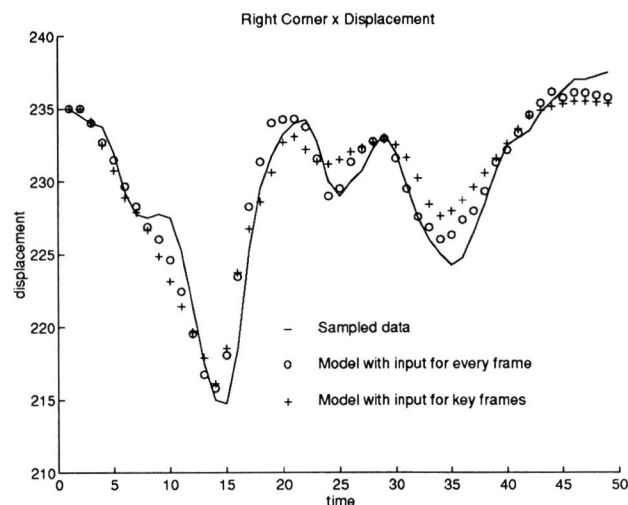

Figure 3: Horizontal displacement.
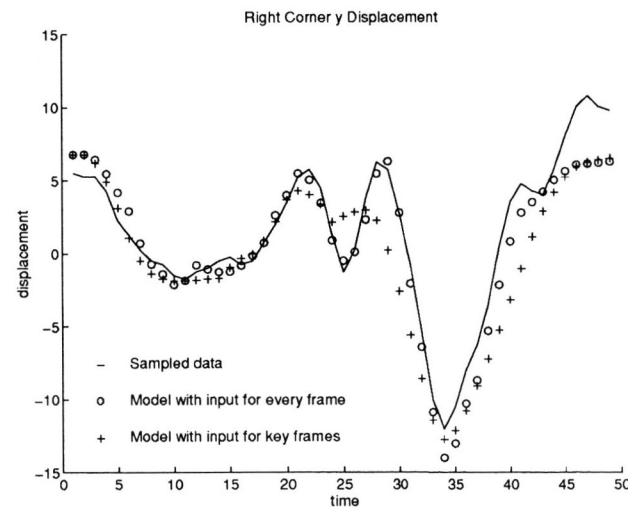


Figure 4: Vertical Displacement.



Figure 5: Phase portrait.

symbols are data from the model. In the first case, the model was driven with a new input vector at the start of each frame. In the figures, this data is plotted with circles. Since the input vectors are optimal in that they achieve the closest realizable fit between the measured data and the model, the degree to which the circles follow the line indicates the degree to which the model can reproduce the motions of natural lips. The second set of data, indicated with crosses, was generated by specifying only nine input vectors; the nine that generated the frames in table 7. The key input vectors were chosen because they represented characteristic articulations. To the extent this data tracks the measured data, we can say our model captures the dynamic behavior of the lips. This data also suggests it may be possible to produce realistic speech using input vectors specified only once per phoneme.

## 7   Discussion and Future Work

Many extensions can be accommodated in the framework. The mass-spring muscle model could be more sophisticated, or even replaced with a spline model where contracting muscles might correspond to moving control points. The current two-dimensional model could be extended to three dimensions by operating in a cylindrical coordinate space [WT91]. Alternatively, a full three-dimensional model taking into account sliding over bones and other tissues could be developed.

In speech modeling, mass-spring systems have been used to model the phonetic structure of speech [KVBSK85]. These models are strikingly successful at fitting real articulatory data [BG85]. The appeal of this approach lies both in its simplistic description of articulatory movements, and in its physical generality. For example, in [BG85], a vocal-tract simulation was controlled sucessfully by only two mass-spring systems – one for lip aperture, the other for protrusion. Likewise [KVBSK85] a similar mass-spring system was developed and compared to a real person's face in reiterant speech production. Although a vocal-tract model is not the focus of this paper, we follow the same philosophy of simplicity and generality characteristic of mass-spring systems.

Presently, we are developing an experimental setup where we can evaluate different speech animation systems and quantify their intelligibility. It is intriguing to observe temporal motions that appear to match patterns found when people really speak [BG85]. This is espeically true when the dynamics are observed in real-time, on high-performance graphics workstations.

## 8   Conclusions

We have described an approach to animating speech; namely, modeling muscle behavior rather than surface or skin behavior. This is appropriate because muscles animate faces; their movement and properties determine the other facial tissues' configuration. We believe realistic muscle behavior can be modeled with fewer elements than a comparable surface model because it is not necessary to increase the number of model elements just to improve the image quality. This is a significant savings; increasing the model resolution by a factor of $n$ increases the number of modeled elements by a factor of $n^2$. More significantly, using standard integration techniques, the minimum number of integration steps required to propagate a change on one side of the model to the other also increases by a factor of $n$.

The research presented in this paper demonstrates that a coordinated muscle model is capable of capturing the significant static and dynamic characteristics evident in real speech. Furthermore, realistic articulations can be achieved with small computational effort and sparse input control data.

## Appendix I: Discrete Mass-Spring Systems

Networks of masses and springs can model non-rigid curves, surfaces, and solids [TF88]. The fundamental elements are springs connected to point mass nodes. The beauty of this discrete model is that diverse topologies can be constructed allowing springs to share nodes by chaining them together to form curves; they maybe assembled into more complex composite units that can become surfaces and solids. In this paper we adopt a deformable topology that reflects the structure of the lower facial muscles.
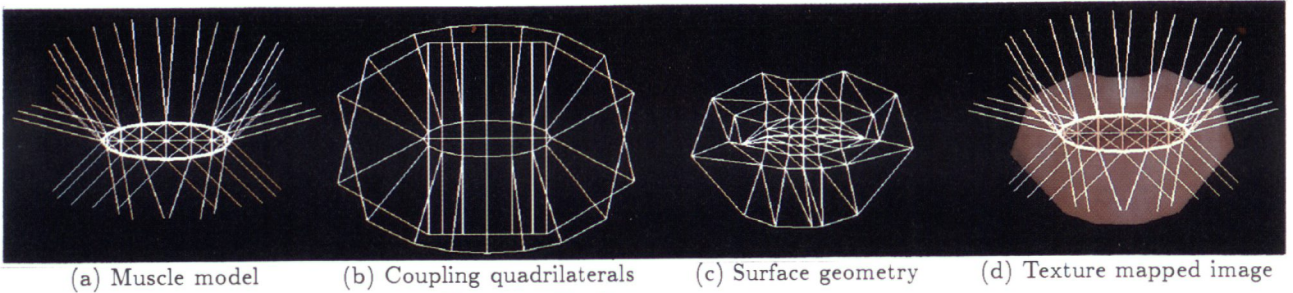
In the simplest case, a node $i$ has a mass $m_i$, and a two-space position $\vec{x}_i(t) = [x(t)y(t)]^T$. The velocity of the node can be described by $\vec{v}_i = d\vec{x}_i/dt$ and its acceleration by $\vec{a}_i = d^2\vec{x}_i/dt^2$. The nodes are subtended by springs with rest lengths $l_0$, and spring constants $k$. The springs are Hookean in nature, exerting a force on the nodes $i$ and $j$ that it connects:

$$\vec{f}_{ij} = k(\|\vec{x}_i - \vec{x}_j\| - l_0)\frac{\vec{x}_i - \vec{x}_j}{\|\vec{x}_i - \vec{x}_j\|} \qquad (1)$$

so that the spring if extended, tries to contract and if contracted, tries to expand. With this model it is possible to compute the net force acting on node from all connected springs:
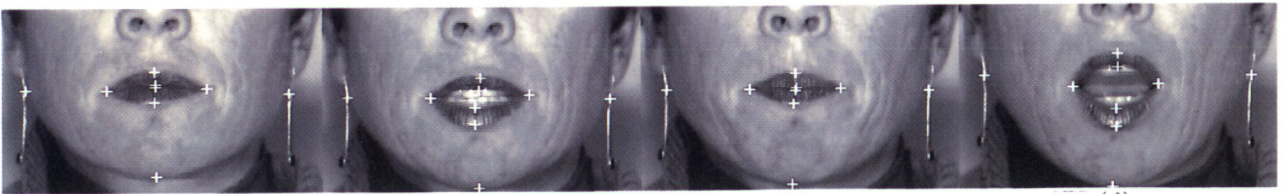
$$\vec{f}_i^{net} = \sum_{j \in \mathcal{N}_i} \vec{f}_{ij} \qquad (2)$$

(a) Muscle model    (b) Coupling quadrilaterals    (c) Surface geometry    (d) Texture mapped image

Table 1: The Framework components.



(a)    (b)    (c)    (d)

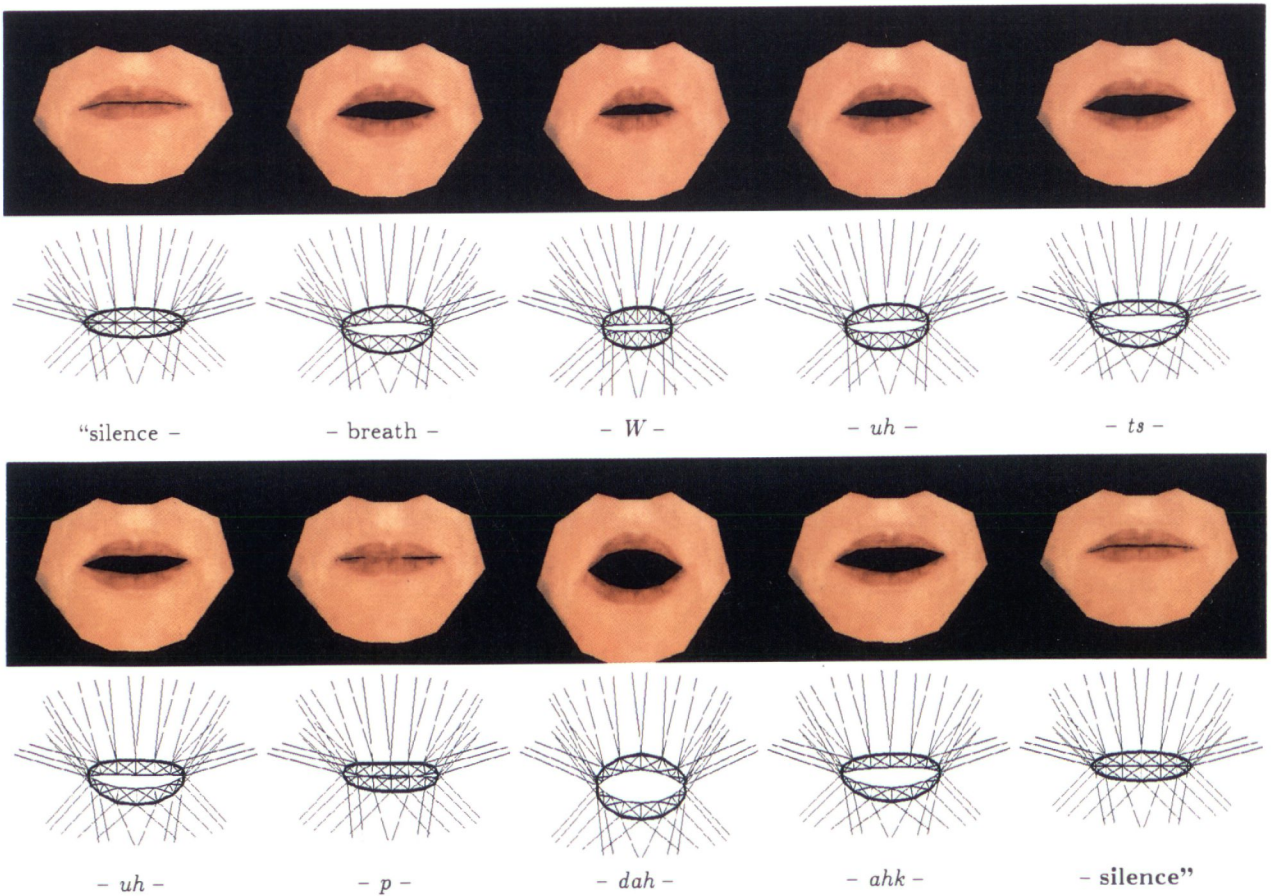Table 2: Digitized optimization nodes.



"silence –    – breath –    – W –    – uh –    – ts –



– uh –    – p –    – dah –    – ahk –    – silence"

Table 3: Ten key postures generated from the optimization procedure.

$$\begin{bmatrix} \beta_x \\ \beta_y \end{bmatrix} = \begin{bmatrix} D_x \\ D_y \end{bmatrix} + u(\begin{bmatrix} C_x \\ C_y \end{bmatrix} - \begin{bmatrix} D_x \\ D_y \end{bmatrix}) \quad (8)$$

$$\begin{bmatrix} S_x \\ S_y \end{bmatrix} = \begin{bmatrix} \alpha_x \\ \alpha_y \end{bmatrix} + v(\begin{bmatrix} \beta_x \\ \beta_y \end{bmatrix} - \begin{bmatrix} \alpha_x \\ \alpha_y \end{bmatrix}) \quad (9)$$

Substituting 7 and 8 into 9:

$$\begin{bmatrix} S_x \\ S_y \end{bmatrix} = u \begin{bmatrix} B_x \\ B_y \end{bmatrix} + v \begin{bmatrix} D_x \\ D_y \end{bmatrix}$$

$$+ uv(\begin{bmatrix} C_x \\ C_y \end{bmatrix} - \begin{bmatrix} D_x \\ D_y \end{bmatrix} - \begin{bmatrix} B_x \\ B_y \end{bmatrix})$$

These two equations, one for each coordinate component, can each be solved for $u$ and then equated:

$$\frac{S_x - vD_x}{B_x + v(C_x - D_x - B_x)} = \frac{S_y - vD_y}{B_y + v(C_y - D_y - B_y)}$$

Simplifying yields the quadratic:

$$\begin{aligned}
& (B_x D_y - C_x D_y - D_x B_y + D_x C_y)\, v^2 \\
+\ & (C_x S_y - B_x D_y - D_x S_y - B_x S_y - S_x C_y + D_x B_y \\
+\ & \quad S_x D_y + S_x B_y)\, v \\
+\ & B_x S_y - S_x B_y \\
=\ & 0
\end{aligned}$$

whose roots, although tedious to write down in closed form, are trivial to compute. For $v$, we select the root that falls between 0.0 and 1.0. Once $v$ is obtained, $u$ follows:

$$u = \frac{S_x - vD_x}{B_x + v(C_x - D_x - B_x)}$$

# References

[BG85]   C. Browman and L. Goldstein. Dynamic modeling of phonetic structure. In V. Fromkin, editor, *Phonetic Linguistics*, pages 35–53. Academic Press, New York, 1985.

[CM93]   M. Cohen and D. Massaro. Modeling coarticulation in synthetic visual speech. In Thalmann N.M. and Thalmann D., editors, *Model and Techniques in Computer Animation*, pages 139–156. Springer-Verlag, Tokyo, 1993.

[EF77]   P. Ekman and W.V. Friesen. *Manual for the Facial Action Coding System*. Consulting Psychologists Press, Palo Alto CA, 1977.

[KM77]   R.D. Kent and F.D. Minifie. Coarticulation in recent speech production. *Journal of Phonetics*, 5:115–133, 1977.

Figure 6: Bilinear coupling.

where $\mathcal{N}_i$ is the set of nodes connected by springs to node $i$.

To achieve a dynamic behavior, the system's collection of second-order ordinary differential equations:

$$m_i \frac{d^2 \vec{x}_i}{dt^2} + \gamma_i \frac{d\vec{x}_i}{dt} = \vec{f}_i^{net} \quad (3)$$

where $\gamma_i$ is a velocity-dependent damping coefficient that dissipates kinetic energy through friction, can be numerically integrated by the explicit forward Euler method:

$$\vec{a}_i(t) = \frac{1}{m_i}(\vec{f}_i^{net}(t) - \gamma_i \vec{v}_i(t)) \quad (4)$$

$$\vec{v}_i(t + \Delta t) = \vec{v}_i(t) + \Delta t \vec{a}_i(t) \quad (5)$$

$$\vec{x}_i(t + \Delta t) = \vec{x}_i(t) + \Delta t \vec{v}_i(t) \quad (6)$$

# Appendix II: Bilinear Mapping

First, consider the forward problem. Given a surface vertex $S$ with bilinear coordinates $(u, v)$ within quadrilateral $ABCD$ whose corner coordinates, $(A_x, A_y)$, $(B_x, B_y)$, etc., are known, we must calculate $(S_x, S_y)$, the coordinates of $S$. Refer to Fig. 6.

$$\begin{bmatrix} \alpha_x \\ \alpha_y \end{bmatrix} = \begin{bmatrix} A_x \\ A_y \end{bmatrix} + u(\begin{bmatrix} B_x \\ B_y \end{bmatrix} - \begin{bmatrix} A_x \\ A_y \end{bmatrix})$$

$$\begin{bmatrix} \beta_x \\ \beta_y \end{bmatrix} = \begin{bmatrix} D_x \\ D_y \end{bmatrix} + u(\begin{bmatrix} C_x \\ C_y \end{bmatrix} - \begin{bmatrix} D_x \\ D_y \end{bmatrix})$$

$$\begin{bmatrix} S_x \\ S_y \end{bmatrix} = \begin{bmatrix} \alpha_x \\ \alpha_y \end{bmatrix} + v(\begin{bmatrix} \beta_x \\ \beta_y \end{bmatrix} - \begin{bmatrix} \alpha_x \\ \alpha_y \end{bmatrix})$$

The inverse problem, finding $(u, v)$ given $(S_x, S_y)$ and the corner coordinates, is more complicated but is performed only once, at the beginning of the simulation. Consider the surface vertex $S$ contained in the quadrilateral $ABCD$ as shown in Fig. 6. $S$ must be assigned bilinear coordinates $u$ and $v$, each between 0.0 and 1.0. To reduce the number of terms, without loss of generality, we can assume $A$ is the origin. Starting with the forward equations:

$$\begin{bmatrix} \alpha_x \\ \alpha_y \end{bmatrix} = u \begin{bmatrix} B_x \\ B_y \end{bmatrix} \quad (7)$$

170

[KVBSK85] J. Keslo, E. Vatikiotis-Bateson, E. Saltz-
man, and B. Kay. A qualitive dynamic
analysis of reiterant speech production:
Phase portraits, kinematics, and dynamic
modeling. *J. Acoust. Soc. Am*, 1(77):266–
288, Jan 1985.

[MM86] H. McGurk and J. MacDonald. Hearing
lips and seeing voices. *Nature*, 264:126–
130, 1986.

[MPT88] N. Magnenat-Thalmann, N.E. Primeau,
and D. Thalmann. Abstract muscle actions
procedures for human face animation. *Vi-
sual Computer*, 3(5):290–297, 1988.

[NHS88] M. Nahas, H. Huitric, and M. Saintourens.
Animation of a B-spline figure. *The Visual
Computer*, 3:272–276, 1988.

[Par72] F.I. Parke. Computer generated anima-
tion of faces. Master's thesis, University of
Utah, Salt Lake City, June 1972. UTEC-
CSc-72-120.

[Par74] F.I. Parke. *A Parameteric Model for Hu-
man Faces*. PhD thesis, University of
Utah, Salt Lake City, Utah, December
1974. UTEC-CSc-75-047.

[Per95] J.S. Perkell. Articulatory processes. In
W.J. Hardcastle and J. Laver, editors, *A
Handbook of Phonetic Science*. (In Press),
1995.

[PFTV86] W. Press, B. Flanney, S. Teukolsky, and
W. Vettering. *Numerical Recipes: The Art
of Scientific Computing*. Cambridge Uni-
versity Press, Cambridge, 1986.

[Pie89] S.D. Pieper. More than skin deep: Physi-
cal modeling of facial tissue. Master's the-
sis, Massachusetts Institute of Technology,
Media Arts and Sciences, 1989.

[Pla80] S.M. Platt. A system for computer simu-
lation of the human face. Master's thesis,
The Moore School, Pennsylvania, 1980.

[Sum92] Q. Summerfield. Lipreading and audio-
visual speech perception. *Phil. Trans. R.
Soc Lond. B*, 355(1273):71–78, Jan 1992.

[TF88] D. Terzopoulos and K. Fleischer. De-
formable models. *The Visual Computer*,
4(6):306–331, 1988.

[TW90] D. Terzopoulos and K. Waters. Physically-
based facial modeling, analysis, and ani-
mation. *Journal of Visualization and Com-
puter Animation*, 1(4):73–80, 1990.

[WL94] K. Waters and T. Levergood. An au-
tomatic lip-synchronization algorithm for
synthetic faces. In *Proceedings of the Mul-
timedia Conference*, pages 149–156, San
Francisco, California, Sept 1994. ACM.

[Wol91] G. Wolberg. *Digital Image Warping*. IEEE
Computer Society Press, Los Alamitos,
CA, 1991.

[WT91] K. Waters and D. Terzopoulos. Modeling
and animating faces using scanned data.
*Journal of Visualization and Animation*,
2(4):123–128, December 1991.

[WT92] K. Waters and D. Terzopoulos. The com-
puter synthesis of expressive faces. *Phil.
Trans. R. Soc. Lond. B*, 355(1273):87–93,
Jan 1992.

[WWDB89] P.L. Williams, R. Warwick, M. Dyson, and
L.H. Bannister. *Grey's Anatomy 37th Edi-
tion*. Churchill Livingstone, London, 1989.