# Current Models of Speech Motor Control: A Control-Theoretic Overview of Architectures & Properties

Benjamin Parrell [†,1] Adam C. Lammert [†,2, a)] Gregory Ciccarelli,[2] and Thomas F. Quatieri[2]

[1]*Department of Communication Sciences & Disorders, University of Wisconsin-Madison, Madison, Wisconsin, USA*

[2]*Bioengineering Systems & Technologies, MIT Lincoln Laboratory, Lexington, Massachusetts, USA*

(Dated: 15 January 2019)

1 This paper reviews the current state of several formal models of speech motor con-

2 trol with particular focus on the low level control of the speech articulators. Fur-

3 ther development of speech motor control models may be aided by a comparison

4 of model attributes. The review builds an understanding of existing models from

5 first principles, before moving into a discussion of several models, showing how each

6 is constructed out of the same basic domain-general ideas and components – e.g.,

7 generalized feedforward, feedback, and model predictive components. This approach

8 allows for direct comparisons to be made in terms of where the models differ, and

9 their points of agreement. Substantial differences among models can be observed

10 in their use of feedforward control, process of estimating system state, and method

11 of incorporating feedback signals into control. However, many commonalities exist

12 among the models in terms of their reliance on higher-level motor planning, use of

13 feedback signals, lack of time-variant adaptation, and focus on kinematic aspects of

14 control and biomechanics. Ongoing research bridging hybrid feedforward/feedback

15 pathways with forward dynamic control, as well as feedback/internal model-based

16 state estimation is discussed.

---

a)Corresponding Author: Adam.Lammert@LL.mit.edu; †Both authors contributed equally to this work.

## I.  INTRODUCTION

Several formal models of speech motor control have been formulated and presented in the speech production literature. Based on decades of observation, it seems clear that the mechanisms of speech motor control are complex, and consequently benefit from the detailed and rigorous description that formal, mathematical models can provide. Speech motor control is, indeed, one of the most intricate sensorimotor activities in which humans engage. Producing speech requires fine timing and coordination of muscles that are interwoven, redundant and have complex mechanical properties, in order to move the diverse articulatory structures of the tongue, lips, jaw, velum and larynx into a wide range of configurations, all of which have a nonlinear relationship with the vocal tract's acoustic output. Control mechanisms are additionally modulated by higher-level processes that determine motor planning, and also mediate semantic, syntactic, prosodic and phonological organization. The various aspects of speech motor control can be conceptualized as layered modules (see Figure 1). In such a layered description, the bridge between higher-level planning processes and the movements of the biomechanical speech-producing structures is a layer which produces motor commands that drive kinematics given some motor plan and potentially in light of some monitoring or prediction of action. The central role filled by this layer – hereafter, simply referred to as the *control layer* – has ensured that all formal models of sensorimotor control for speech have defined architectures that govern its functionality. The field of models that have provided a formal description of the control layer comprises: DIVA (Guenther, 1994, 2016), Task Dynamics (Saltzman and Kelso, 1987; Saltzman and Munhall, 1989), State Feedback Control

38  (Houde and Nagarajan, 2011), ACT (Kröger *et al.*, 2009), GEPPETO (Perrier *et al.*, 2005),

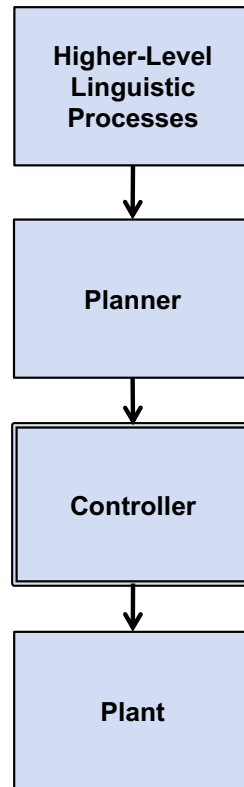39  FACTS (Ramanarayanan *et al.*, 2016).



FIG. 1. Representation of the distinct levels of speech production modeling. This paper focuses on modeling the speech controller, the system that takes in a speech plan and potentially feedback from the plant and issues motor commands to the plant. Other components of the speech production hierarchy include higher level linguistic processes (prosody, semantics, syntax), the planner (low level sequencing of motor actions), and the plant itself (e.g. speech synthesizers including but not limited to articulatory synthesizers such as CASY, Birkholz or Maeda).

40  An impediment to progress in developing rigorous speech motor control models appears

41  to be the variety of distinct approaches, taken in the published literature, to explaining

42  the attributes of the more prominent models of speech motor control. There is very lit-

tle agreement, for instance, even concerning the terminology used to describe the models.

Nevertheless, there is reason to believe that a direct comparison of speech control models

is possible, based on the important, high-level observation that the models presented in

the literature are all closely related to engineering approaches to motor control, and bear

a strong resemblance to classical control-theoretic architectures. Given that the theory be-

hind current understanding of biological motor control largely grew out of early advances in

engineering fields (Bellman, 1957; Wiener, 1948), it is perhaps unsurprising that the same is

true specifically in the area of speech motor control. Indeed, engineering approaches are a

sensible place to begin investigations into the nature of speech motor control, in part because

our current understanding of the functional interpretation of motor control neuroanatomy

follows the engineering architectures closely (consider, e.g., Brainard and Doupe (2002);

Shadmehr and Krakauer (2008); Takakusaki (2017); Wolpert *et al.* (1998)).

Progress in the development of speech motor control models may be facilitated by a direct

comparison of the various models, using a common framework of domain-general (i.e., not

speech-specific) motor control principles and unified terminology to describe their attributes.

The purpose of the present paper is to provide such a direct comparison for models of the

control layer that utilize mechanisms to move the plant in support of accomplishing speech

tasks in accordance with higher-level speech goals. These models have been developed

to attempt meaningful reproduction of speech behavior, including potentially acoustics,

articulatory and neural signals. Demonstrations of the ability of these models to capture

aspects of human speech production kinematics have been presented in the literature, and

the extent and quality of these efforts may differ by model. No systematic review will be
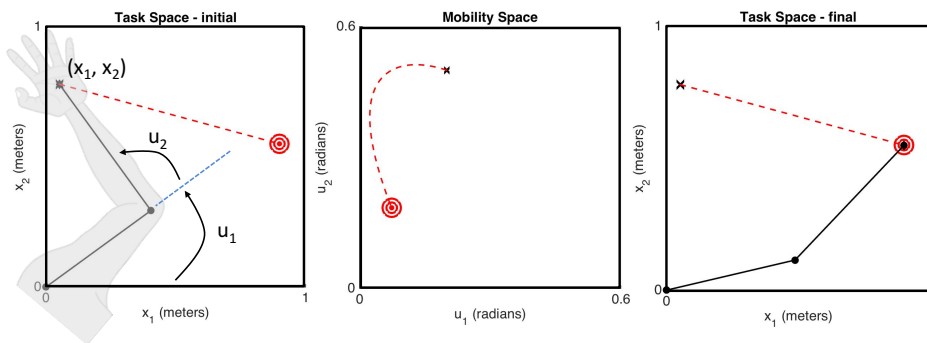
65 offered here of experimental data, either behavioral or neurological, that has been or could

66 be used to support the expressivity or biological plausibility of any model. However, a brief

67 summary of the demonstrated capabilities of each model is included. This choice reflects an

68 intention to focus on the model architectures themselves.

69 Our review begins with general motor control principles and approaches, before moving

70 into basic, domain-general models of motor control. The paper then proceeds to provide

71 detailed discussions of currently proposed models of speech motor control, showing how

72 each model is constructed out of these basic domain-general ideas and components. By

73 showing how each model is built up on these basic elements, this approach allows for a clear

74 comparison between the proposed models, showing where they differ as well as points of

75 agreement. The present review focuses specifically on control of the speech articulators in

76 fully developed, adult speech. Control that is adaptive (i.e., time variant), which may be

77 relevant for speech acquisition and development, will only be considered in the discussion,

78 and not in the primary overview framework. Formal explanations, including an appendix

79 with full equations for each model, is provided where possible. Other important aspects of

80 speech production, including learning and optimization, higher-level linguistic processing,

81 motor program generation (i.e. the "planner"), the neurological basis of hypothesized model

82 components, and biomechanical details of the speech articulators (i.e. the "plant") will only

83 be discussed to the extent necessary to clarify the nature and operation of the proposed

84 control mechanisms.

## II.  BACKGROUND

### A.  Motor control principles and terminology

The first step in discussing speech motor control models is to define certain key concepts

and terminology. To illustrate these ideas, a simple example is borrowed from the control

of upper extremity reaching control, as shown in Figure 2, which is based on the descrip-

tion of a simple two-link robotic arm moving on a planar surface. This commonly-used

example, though taken from a completely different domain of motor control, shares many

of the same concepts and terminology with speech motor control, and has the benefit of

being low-dimensional, which makes it possible to represent the relevant spaces in a two-

dimensional plot. Fundamental similarities and distinctions between this simple example

and the (considerably more complex) speech production system, in terms of their assump-

tions and structure, will be drawn where appropriate throughout the present section.



(left panel) Robot arm in its initial configuration at $(x_1, x_2)$ in task space, and the final goal (red circle). The arm's state variables $(u_1, u_2)$ are defined as the angles of the shoulder and elbow. $(u_1, u_2)$ are the parameters directly changed by the controller and therefore exist in mobility space. (middle panel) The trajectory in mobility space. The evolution of the mobility space variables $(u_1, u_2)$ over time may be a non-linear trajectory despite a linear trajectory in task space. (right panel) The final orientation of the arm in task space at the goal.
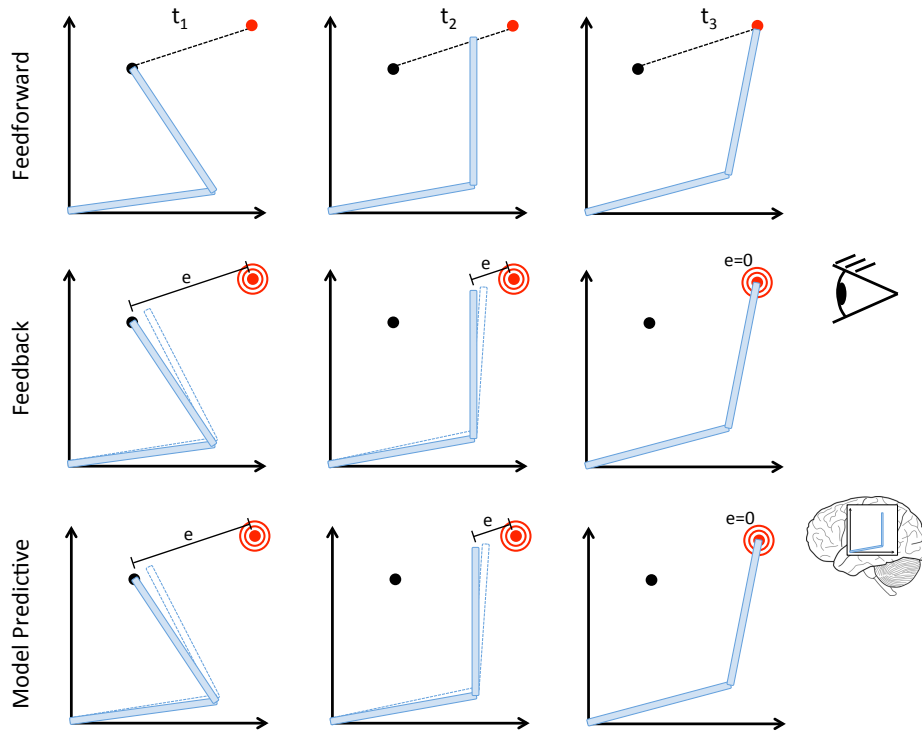
FIG. 2.

Illustration of the difference between feedforward (top row), feedback (middle row) and model predictive (bottom row) control using a simple reaching example. In feedforward control, the arm traces out a fully preplanned trajectory with no feedback about the position of the arm at any point in time. In feedback control, an error is computed between an observed state of the system (observation represented by the eye) and the target. The arm progressively works to minimize this error which drives the end effector towards the target. In model predictive control, an error is computed internally as opposed to being derived from feedback of the state of the system (represented by the brain with an internal model of the robot arm). The arm's position is updated to minimize the predicted error of the system.

FIG. 3.

97    The robotic arm, as a physical structure, is the apparatus to be controlled, and can

98    be referred to as the plant $(G)$. Note that the term *plant* is not specific to this example,

99    and could be used in the domain of speech production to specify the vocal tract and its

100    component articulators, as well as possibly the larynx and the respiratory system. The

101    plant's two links are connected to each other at a revolute joint that changes the angle

102    between the links, $u_2$. The proximal end of the robot's first link is fixed at the origin of

8

103 the planar surface, defined as $(x_1, x_2) = (0, 0)$, but is free to rotate about this point which

104 changes the angle $u_1$. These two variables, $u_1$ and $u_2$ describe the configuration of the plant,

105 and also define the set of possible configurations of the plant, known as *mobility space*[1]. The

106 variables $u_1$ and $u_2$ can be considered as elements of a single $1 - by - 2$ vector, $\mathbf{u}$, which can

107 be said to specify the state of the plant in mobility space (sometimes, the *mobility state*).

108 The distal end of the second link (i.e., the "hand") is considered the end-effector of the

109 robot, the precise positioning of which is typically the focus of controlling the plant in the

110 context of reaching tasks. The variables $x_1$ and $x_2$, already used to define locations on the

111 planar surface, can also be used to describe the location of the end-effector on that surface.

112 The space of possible locations for the end-effector is known as *task space*, and the desired

113 outcome of a controlled movement is known as a *task*. The variables $x_1$ and $x_2$ can be

114 considered as elements of a single $1 - by - 2$ vector, $\mathbf{x}$, specifying the state of the plant

115 in task space (sometimes, the *task state*). Tasks with respect to the robotic arm might be

116 putting the end-effector as a specified location in task space (i.e., achieving a state where

117 $\mathbf{x}$ takes on a particular value), or alternatively achieving a specific trajectory through task

118 space (i.e., tracking some sequence of values for $\mathbf{x}$). In speech production, task spaces might

119 include, for instance, formant space or vocal tract constriction degree/location space.

120 Task and mobility spaces can be viewed as "high" and "low" level spaces, respectively,

121 with the variables comprising each space having a hierarchical arrangement where the task

122 variables are composed of, but distinct from, mobility variables. Often this arrangement is

123 many-to-one, such that many different (or, potentially infinite) locations in mobility space

124 will map to the same location in task space. Task variables consequently describe the state

125 of the plant in a way that is directly relevant to the task, and which abstracts away from

126 a certain amount of detail as to how that task state was achieved via some mobility state.

127 Mobility variables describe the state of the plant in a way that is more relevant to control, in

128 the sense that motor commands are typically defined so as to affect some change in mobility

129 state. Using the robotic arm example, motor commands would typically be given in terms

130 of the joint angles, and not in terms of the end-effector position. In a speech context, a

131 model might assert that motor commands are issued in terms of the positions of the speech

132 articulators (e.g. upper lip, lower lip, tongue tip, etc.), and not in terms of some desired

133 formant values (e.g., F1 = 500 Hz) or vocal tract constrictions (e.g., lip aperture = 2 mm).

134 The details of the task are specified in the *reference*, $\mathbf{r}$, a vector representing a desired

135 state. The reference vector typically resides in task space ($\mathbf{r_x}$), but may also be given in

136 mobility space ($\mathbf{r_u}$) for specific applications. Reference vectors originate in the planner ($P$),

137 and may be part of a larger motor program maintained by the planner, toward achieving

138 some higher-level sensorimotor or cognitive goal (e.g. reach to a series of targets in space,

139 utter the word "dad"). As implied above, however, reference vectors will typically be insuf-

140 ficient for use directly as motor commands to the plant because they reside in task space.

141 The reference will need to be transformed into a motor command in mobility space. This is

142 the function of the controller.

143 The controller ($C$) is the bridge between the planner and any feedback, on the one

144 hand, and movements of the plant, on the other. The ultimate purpose of the controller

145 is to issue motor commands that produce movement (or lack thereof) in the plant. Note

146 that the present paper assumes that motor commands take the form of vectors in mobility

147  space, $\mathbf{u}$, and that those vectors can be used directly as commands to the plant. In a real

148  biological system, several transformations may be required for encoding motor commands

149  as neural signals, and to elicit muscle activations that bring about changes in mobility state.

150  This assumption is made to promote consistency with the speech motor control modeling

151  literature, and for the sake of simplicity. In any case, the motor command issued by the

152  controller will depend either upon the reference directly, or upon the *state error*, $\mathbf{e}$, a vector

153  representing the difference between the reference and the plant's state (or an estimate of

154  that state, see below).

155  In biological systems, the plant's actual state may not always be directly accessible to the

156  controller. It can be therefore important to develop the notion of a *state estimate* ($\hat{\mathbf{x}}$ or $\hat{\mathbf{u}}$),

157  which is an internal estimate of the plant's state, either in task space or in mobility space.

158  The state estimate may be informed by sensory measurements of the plant's actual state –

159  represented by the *sensory state* vector $\mathbf{y}$ – and by predictions generated from an internal

160  model of the plant – represented by the *predicted state* vector $\tilde{\mathbf{x}}$ or $\tilde{\mathbf{u}}$. The sensory state

161  vector, an approximation to either $\mathbf{x}$ or $\mathbf{u}$, may be corrupted by some combination of noise

162  (e.g., neuronal noise), delays (e.g., slowed synaptic/axonal propagation) or transformations

163  (e.g., warping). The predicted state vector may also be imperfect, since the internal model

164  may be inaccurate or biased. For the robotic arm example, the sensory state vector would

165  represent measured joint angles ($\mathbf{y_u}$). This contrasts with the sensory output for speech

166  production, which is typically considered to be a combination of auditory ($\mathbf{y}_{aud}$) and so-

167  matosensory signals ($\mathbf{y}_{somat}$), where the somatosensory signal may include proprioceptive

168  and/or tactile sensation.

11

169   In general, motor control can be viewed as a collection of transformations between vectors

170   and spaces of different types, and the planner, the controller, and the plant can all be

171   described – using the conventions developed above – as functional transformations from

172   specific inputs to specific outputs. The planner generates the reference vector, $\mathbf{r} = P(\alpha)$ as a

173   function of some high-level motor program-related information $\alpha$, and possibly as a function

174   of time: $\mathbf{r} = P(\alpha, t)$. The controller takes a reference vector or an error vector as input

175   and generates a motor command in mobility space: $(\mathbf{u}, \dot{\mathbf{u}}) = C(\mathbf{r})$ or $(\mathbf{u}, \dot{\mathbf{u}}) = C(\mathbf{e})$. The

176   plant, which can also be viewed as a transformation, converts motor commands, through

177   movement, into different plant states which can be measured in both mobility and task

178   space: $(\mathbf{u}, \dot{\mathbf{u}}, \mathbf{x}, \dot{\mathbf{x}}) = G(\mathbf{u}, \dot{\mathbf{u}})$. These variables are used in Figure 4, and in related diagrams

179   throughout the paper. The state of the plant can then be measured by some sensory system:

180   $(\mathbf{y}, \dot{\mathbf{y}}) = S(\mathbf{u}, \dot{\mathbf{u}}, \mathbf{x}, \dot{\mathbf{x}})$, the details of which are often not explicitly treated in the literature.

181   Therefore, the present review will often lump $G$ and $S$ together into a single component.

## B.   Types of motor control models

183   The purpose of this section is to lay out, in a general way, some common control architec-

184   tures that are employed in various control applications, including both controlling robotic

185   systems as well as describing the functional aspects of physiological control. These general

186   architectures are presented as a scaffold for understanding the specific architectures used in

187   various speech motor control models, and also for the purpose of clarifying the terms used in

188   the present paper to refer to those architectures. To illustrate these various architectures in

189   an intuitive way, the example of the planar robotic arm will continue to be employed as in

12

(a) feedforward

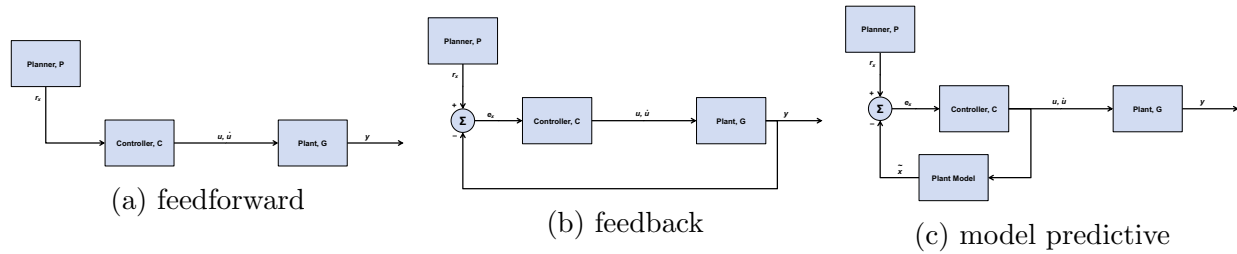(b) feedback

(c) model predictive

FIG. 4. Control architecture of a generic (a) feedforward, (b) feedback, and (c) model predictive controller. The feedforward control architecture is distinguished from the other two because the controller only receives information from the planner, not information from the plant or predicted information from the plant. The feedback control and model predictive control architectures differ in the nature of the feedback received by the controller. In feedback control, the state of the plant (different than the output) is sent back to the controller. By contrast, in model predictive control, the state of the plant is sent back to the controller using an estimate of the plant based on a copy of the issued control signal.

190 the previous section. However, these same architectures can be used to control much more

191 complex systems, such as the speech production system.

### *1.  Feedback control*

193    Figure 4b shows an example of a *feedback* system architecture that, by definition of the

194 term, makes use of outputs from the plant for maintaining control. These feedback signals,

195 which convey the sensory state vector, are compared with the reference vector from the

196 planner in order to generate an error vector. The error vector, in its most basic form, simply

197 represents the difference between the current state and the reference. The error vector is

198 passed to the controller for determining the motor command. This type of controller is

13

199  often referred to as a *closed-loop controller* in the control theory literature, since the flow

200  of signals through the system form a loop from the motor command to the error signal and

201  back again. Many types of controllers exist which match this general description, only a few

202  of which will be discussed here. What all feedback controllers share is the basic idea that

203  the error between the state of the plant (or an estimate thereof) and the reference forms the

204  basis for the motor command issued to the plant. The simplest feedback controller design is

205  the proportional controller, in which the motor command is simply proportional to the error

206  signal – e.g., $C(\mathbf{e_x}) = \mathbf{K}_p\mathbf{e_x}$, where the term $\mathbf{K}_p$ is a matrix of weights known as the *gains*.

207  Larger gains lead to larger motor commands (i.e. the error has more of an effect on the

208  system) while smaller gains result in smaller commands. Smaller gains are often preferable

209  as large gains can lead to instability and oscillatory behavior.

210  The second row in Figure 3 shows, across times $t_1$, $t_2$ and $t_3$, the progress of the robotic

211  arm as controlled by a feedback controller. At the beginning, the task is defined as a

212  desired point in task space $\mathbf{x} = (x_1, x_2)$. This type of task is sometimes referred to as a

213  point-attractor, or a target, since the system should evolve to approach this point in task

214  space regardless of its initial position, given sensible motor commands that reduce the error

215  signal over time. The motor commands issued at each time step are a function of the error,

216  $\mathbf{e}_x$, between the current position of the end-effector and the point target. The error is

217  determined by sensory feedback, which provides monitoring of the current state of the arm

218  with respect to the position of the target.[2] Although the error signal is in task space, the

219  motor command issued by the controller must be given in mobility space since the only way

220  to change the position of the end effector is to change the joint angles $\mathbf{u} = (u_1, u_2)$. The

14

process of determining those commands requires some kind of transformation (i.e., kinematic inversion) from the desired coordinates in task space to corresponding coordinates in mobility space. Alternatively, it is also possible for the target to be a pre-specified trajectory rather than a point in task space. In this case, the error would be computed between the current position of the end-effector and the current desired position along the trajectory (typically time-locked).

Feedback control architectures have wide applicability in engineered and biological systems. Even simple designs typically lead to systems that accurately produce desired behaviors, and which can naturally handle unstable or unpredictable environments, including external perturbations to the plant. However, feedback systems can require careful calibration to ensure stability of control. Incorrectly tuned feedback systems can result in movements that grow uncontrollably or oscillate indefinitely. Feedback architectures are also heavily dependent on the quality of feedback signals. If those signals are slow to propagate, or if they require extensive processing once received, this can lead to motor commands being issued based on outdated state information, resulting in poor and/or slow performance. Additionally, if feedback signals are corrupted or otherwise inaccurate, this can lead to inaccurate movements. These final considerations are particularly important for biological systems, as there are substantial delays and noise inherent to neural processing of sensory feedback.

### 2. *Feedforward control*

One way to avoid the problems of delayed and noisy sensory information is to cut out the use of feedback entirely. Figure 4a shows an example of a system architecture that makes

²⁴² no use of any outputs from the plant for maintaining control. Rather, the signals issued

²⁴³ in the system are entirely *feedforward*, with the motor commands depending only on the

²⁴⁴ reference signal. This architecture is commonly referred to as an *open-loop* control system,

²⁴⁵ although the terms *feedforward* and *open-loop* will be used interchangeably in the present

²⁴⁶ paper. The term feedforward control is sometimes used more specifically to refer to control

²⁴⁷ architectures that can monitor perturbations to the plant, and adjust the motor commands

²⁴⁸ to compensate without the need for explicitly monitoring outputs from the plant, usually

²⁴⁹ by employing a highly accurate mathematical model of the plant (see the section on model

²⁵⁰ predictive control, below). To date, the authors are aware of only one modeling effort in the

²⁵¹ domain of speech motor control to utilize this kind of architecture (Baraduc *et al.*, 2017),

²⁵² with preliminary results presented.

²⁵³ The first row in Figure 3 shows the progress of the robot arm as controlled by a feed-

²⁵⁴ forward controller. From the beginning, the trajectory of the end-effector is defined in

²⁵⁵ task space as a straight line originating at the end-effector's current position. The motor

²⁵⁶ commands issued to the arm at each time step are directly determined by this pre-specified

²⁵⁷ trajectory. As in a feedback controller, the reference signal is defined in task space but motor

²⁵⁸ commands must be issued in mobility space. Again, this requires some kind of transforma-

²⁵⁹ tion from the desired coordinates in task space to corresponding coordinates in mobility

²⁶⁰ space. Although the trajectory in this example is specified in task space, as is often done,

²⁶¹ an alternative feedforward controller could define the plan in mobility space (that is, for our

²⁶² robot example, in terms of joint angles) or even simultaneously in mobility and task space.

²⁶³ In any case, a key aspect of feedforward control is that no estimate of the state (that is, the

264   arm's estimated position) is used by the controller at any point throughout its movement.

265   In the absence of feedback, the simplest method of generating reasonable control signals is

266   simply to have the plan pre-specify the entire trajectory in task or mobility space, and then

267   issue motor commands that attempt to carry out that plan step-by-step from beginning to

268   end.

269   Feedforward control architectures are unsuitable for unstable or unpredictable environ-

270   ments, where the plant can be perturbed by interference external to the system. Without

271   the ability to detect and monitor errors in the system output, errors tend to persist, or even

272   compound over time. Despite this obvious disadvantage, feedforward architectures are some-

273   time attractive because they are capable of issuing motor commands quickly and without

274   the need for complex handling of feedback signals.

275   **3.   *Model predictive control***

276   An alternative to feedforward and feedback control is *model predictive* control. A model

277   predictive controller, like the feedforward controller, makes no use of outputs from the plant

278   for maintaining control. However, this architecture does make use of an *internal model* of

279   the plant, which takes motor commands as input and transforms them into a prediction of

280   the system's subsequent state, to predict the effects of the issued motor command. This

281   effectively replaces feedback from the plant with a prediction of what the controller thinks

282   that the feedback should be (Garcia *et al.*, 1989; Miall and Wolpert, 1996). An example

283   of this architecture is shown in Figure 4c. This state prediction acts as a kind of pseudo-

284 feedback which can be compared against the reference, producing an error signal that is

285 provided to the controller.

286 Note that model predictive control can be viewed as a special case of feedforward control,

287 if the plant model is considered to be part of the controller. This special case has been

288 separated out as a distinct architecture in the present framework because it is central to

289 several models of speech motor control. Therefore, feedforward architectures, as discussed

290 here, will specifically discount architectures that are model predictive.

291 The third row in Figure 3 shows the progress of the robot arm as controlled by a model

292 predictive controller. The functioning of such a controller is similar to the feedback controller

293 example, above, in that the target is defined as desired point in task space, and the motor

294 commands issued are a function of the error, $\mathbf{e_x}$, between the current position of the end-

295 effector and the point target. The difference is that the error is determined by comparing

296 the desired state to the output of an internal model.

297 In terms of performance, the primary advantage of such an architecture is speed, since

298 the delays associated with predicting the plant's state can often be much shorter than those

299 associated with feedback propagation. Additionally, a model predictive controller is one way

300 to avoid the need for having an entire trajectory formulated before movement begins, as is

301 often the case with feedforward architectures. Rather, plans can be more compact, such as a

302 single, time invariant point in task or mobility space (this is the same type of plan often used

303 in feedback controllers). The disadvantage of these systems is that accurate internal models

304 can be difficult to design or learn, especially for complex, nonlinear plants such as the vocal

305 tract. A poor internal model would mean that the predicted state may not match the true

306 state of the plant, which can result in inaccurate control. Even small errors in the prediction

307 will accumulate over time, since there is no way of correcting the prediction. Additionally,

308 model predictive controllers have similar problems as feedforward control architectures in

309 dealing with unpredictable environments and perturbations.

### 4. *Combining feedforward and feedback controllers*

311 Each basic type of control system, feedback and feedforward control, has its own strengths

312 and weaknesses. Feedback control is stable in the face of external perturbations, but becomes

313 inaccurate or slow when sensory information is noisy or delayed (respectively), as in most

314 biological systems. Feedforward control can be accomplished quickly, but is unstable when

315 the state of the system cannot be predicted due to external perturbations.

316 It is possible to combine some of the strengths of feedforward and feedback systems, and

317 mitigate the weaknesses of each, by constructing a *hybrid feedforward/feedback* controller, as

318 shown in Figure 5a. This hybrid architecture comprises separate feedforward and feedback

319 pathways that each issue their own motor commands, a (potentially weighted) combination

320 of which becomes the final motor command that is issued to the plant. Such an architecture

321 has the speed of a feedforward controller, but remains sensitive to unexpected perturbations

322 and accumulating errors. Typically, the presence of the feedforward pathway allows for

323 lower gains to be utilized in the feedback controller, leading to better stability. The primary

324 disadvantage of combining feedforward and feedback pathways into a single system is the

325 introduction of more complex designs. Complex designs may be more difficult to maintain,

326 and allow the potential for unnecessary or underutilized components. For instance, if output

(a) hybrid
feedforward/feedback

(b) hybrid model
predictive/feedback

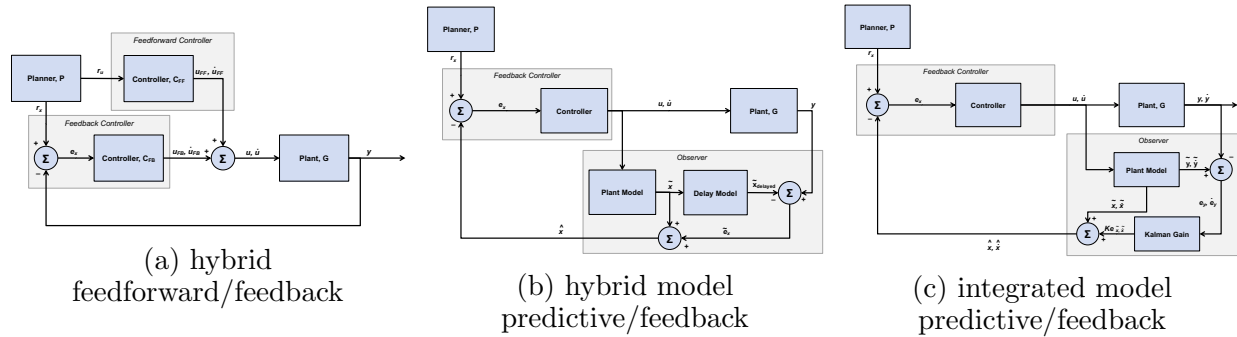(c) integrated model
predictive/feedback

FIG. 5. Control architectures of generic hybrid controllers. A hybrid controller uses two of the three simple control architectures discussed in Figure 4. Diagrammed here are (a) a feedforward-feedback hybrid, (b) a model predictive-feedback hybrid with simple summation of state predictions (i.e., a Smith predictor) and (c) a model predictive-feedback hybrid with full integration of state predictions (i.e., a Kalman filter). Architectures (b) and (c) are distinguished by the specific way in which model predictions and feedback are combined. In (b), the current state is estimated through a three-part error comparison processes. Architecture (c) also uses a three-part comparison, but also incorporates an observation model that maps the model prediction into sensory space, and a gain that allows for potentially variable weighting of model predictions and sensory measurement error.

327 from the plant always equals the reference (e.g., if the environment is entirely predictable), then the feedback pathway is not utilized and essentially unnecessary, since the feedforward pathway would be sufficient for control by itself.

330 One of the most useful applications of model predictive control is as a component of larger, hybrid architectures. For instance, internal model predictions can provide quick pseudo-feedback that can be used in conjunction with true feedback to provide fast, reliable control even in the face of long feedback propagation delays. Such methods are more stable than true

334 model predictive control, since internal predictions do not need to be perfectly accurate, and

335 small deviations between the predicted and actual states of the plant can be corrected via

336 the feedback signal. An example of an architecture that exemplifies this concept is the Smith

337 predictor (Ghosh, 2005; Smith, 1959), as shown in Figure 5b. A Smith predictor effectively

338 has three error comparison processes, generating state errors serially through comparing the

339 state with a delayed version of the internal model prediction, which in turn is compared

340 to a non-delayed internal model prediction, with this final comparison being subsequently

341 compared against the desired state from the reference signal. The integrated mechanisms

342 involved in combining model predictions with feedback signals are sometimes referred to in

343 the literature as the "observer". The present view adopts this terminology. Note that the

344 observer and speaker, in this conceptualization, are the same individual, as speakers observe

345 their own speech.

346 A Smith predictor is not the only controller that uses both state predictions from an

347 internal model and feedback signals. Prominent alternative approaches also use a three-part,

348 cascaded error comparison process, but incorporate (a) an observation model, that maps the

349 model prediction into sensory space for direct comparison with sensory measurements, and

350 (b) a gain that allows for potentially variable weighting of model predictions and sensory

351 measurement error. These additional aspects can afford more accurate estimation of the

352 plant's current state. This is the approach taken by such classic control designs as the

353 Kalman filter (Kalman *et al.*, 1960) (Figure 5c), which provides an optimal[3] state estimate

354 with noisy feedback under certain strict assumptions. Importantly, the estimated state that

355 results from combining internal predictions and feedback can be compared with the desired

356 state to generate a motor command (Todorov, 2004), just as in a pure feedback controller.

357 This type of controller is sometimes referred to as *state feedback control*.

## C.    Speech models

359 The present discussion will now move from domain-general motor control theory to models

360 of speech motor control. Among the speech production models presented in the literature,

361 perhaps the two most prominent are DIVA (Directions Into Velocities of Articulators) and

362 the Task Dynamics model.  The development of DIVA has been driven since the mid-

363 1990's (Guenther, 1994) primarily by a team of researchers at Boston University, led by

364 Frank Guenther. Task Dynamics has been developed by researchers associated with Haskins

365 Laboratories, with Elliot Saltzman playing a key role, and with the theoretical groundwork

366 being laid about five years prior to DIVA (Saltzman and Kelso, 1987; Saltzman and Munhall,

367 1989). More recent models include State Feedback Control (Houde and Nagarajan, 2011),

368 the Feedback Aware Control of Tasks in Speech (FACTS) model (Parrell *et al.*, 2006), ACT

369 (Kröger *et al.*, 2009), and GEPPETO (Perrier *et al.*, 2005).

370 Any model of speech production control must include, at a basic level, the ability to

371 generate motor commands based on some motor plan.  Those motor commands in turn

372 activate a vocal tract model, possibly resulting in the generation of an acoustic signal. While

373 complete models of speech production also need to include the formulation of motor plans,

374 these elements are beyond the scope of the present paper, which focuses more narrowly

375 on controlling the vocal tract for speech.  An important reason for limiting the scope of

376 the present paper is that the longstanding debate over acoustic vs.  articulatory targets

377 of speech production tasks is often intertwined with the critical issue of how the vocal

378 tract is controlled. For example, DIVA's tasks are formulated primarily in acoustic space,

379 whereas applications of Task Dynamics (e.g., the Articulatory Phonology of (Browman and

380 Goldstein, 1986) often assume tasks to be constrictions in the vocal tract. The choice of

381 task space, however, is almost completely independent of the control formulations that are

382 the focus of the current paper, and it is generally possible to reformulate any given control

383 architecture using different task spaces. Therefore, the present work will discuss the task

384 space used for each model, as the specific choice of task variables comprising the task spaces

385 does differ between models, but will make no attempt to discuss the relative merits of the

386 different task spaces used in different models. The concept of a task space is general enough

387 to sit over and above the specific choice of task variables, while being well-defined enough as

388 a concept to allowing meaningful comparisons of the control architectures underlying task

389 space control.

390 Control elements that are relevant to any model of speech motor control, and which will

391 be discussed in depth for each model in the following section, include: (a) the nature of

392 feedforward mechanisms of control, including the formulation of the planner, (b) the nature

393 and importance of feedback signals, (c) modeling of potentially imperfect sensory systems

394 and/or perceptual processing of feedback, (d) the consequences of delays in feedforward and

395 feedback pathways (e) the potential role of forward models in state prediction, and (f) the

396 potential integration of both feedback and state predictions for state estimation, (g) the

397 implementation of transformations between task space, mobility space, and sensory space,

398 (h) the design of the controller for generating and issuing motor commands to the plant.

₃₉₉  It is noted here that most current speech models are examples of purely kinematic con-

₄₀₀  trollers. That is, they do not account for dynamics or biomechanical considerations of the

₄₀₁  vocal tract. It is typically assumed that inertial parameters, centrifugal/coriolis forces and

₄₀₂  stationary external forces like gravity can all be ignored for the purposes of controller design

₄₀₃  and forward modeling. This may owe to the fact that several prominent models of the plant

₄₀₄  are purely kinematic: for instance, Maeda's model (Maeda, 1982) and the Haskins Config-

₄₀₅  urable Articulatory Synthesizer (CASY) (Iskarous *et al.*, 2003; Rubin *et al.*, 1981, 1996).

₄₀₆  The focus on kinematics may also reflect an implicit assumption that dynamics of the plant

₄₀₇  can be ignored in the domain of speech motor control. Such an assumption is quite common

₄₀₈  in robotics and human motor control, and amounts to conceptualizing the plant as a collec-

₄₀₉  tion of stiff articulators, akin to an industrial robotic arm. However, there is evidence that

₄₁₀  biomechanical factors play non-negligible roles in speech motor control (Buchaillard *et al.*,

₄₁₁  2009; Derrick *et al.*, 2015; Nazari *et al.*, 2011; Ostry *et al.*, 1996; Perrier *et al.*, 2003; San-

₄₁₂  guineti *et al.*, 1998; Shiller *et al.*, 2002), and more recent vocal tract models such as Artisynth

₄₁₃  (Lloyd *et al.*, 2012) incorporate dynamic and biomechanical aspects in their design.


₄₁₄  **III. PROMINENT MODELS OF SPEECH PRODUCTION**


₄₁₅  In the following section, each of the current models of speech motor control will be

₄₁₆  discussed in turn, explaining the architecture of the control system as it relates to the simple,

₄₁₇  domain-general systems discussed previously. Where necessary, additional components of

₄₁₈  each model will be touched upon, such as motor program generation. How each model

₄₁₉  addresses the key control elements listed above will also be discussed.

### A. DIVA

The Directions Into Velocities of Articulators (DIVA) model is a hybrid control system combining a model-predictive controller with separate auditory and somatosensory feedback controller loops (Golfinopoulos *et al.*, 2011; Guenther *et al.*, 2006; Tourville and Guenther, 2011). Being arguably the most complete computational model of speech motor control, DIVA has been developed to address a number of theoretical issues, primarily focused around replicating human speech production at behavioral, neurological, and developmental levels. The use of both model predictive and feedback control in DIVA is conceptually similar to a Smith Predictor. However, while a Smith Predictor uses serial error calculations to issue a single motor command, DIVA generates independent errors from each controller simultaneously. Each error is then individually transformed into a separate motor command. These three commands are then combined into a single motor command which is passed to the plant. The plant in DIVA has historically been Maeda's model (Tourville and Guenther, 2011), but this has recently been replaced with a custom plant model (Guenther, 2016).

The basic component of the planning process in DIVA is the "speech sound", which can be a phoneme, syllable, or multisyllabic chunk. Each speech sound is linked to three distinct tasks, each a function of time: an articulatory trajectory (often called "motor" trajectory in the DIVA literature) defined in mobility space $\mathbf{r}_u(t)$, an auditory sensory trajectory $\mathbf{r}_{aud}(t)$, and a somatosensory trajectory $\mathbf{r}_{somat}(t)$. The "speech sound map", which corresponds to the planner in Figure 4b, stores all three-component sets of mobility and sensory state trajectories. Each trajectory of the set serves as the reference signal to one of the controllers
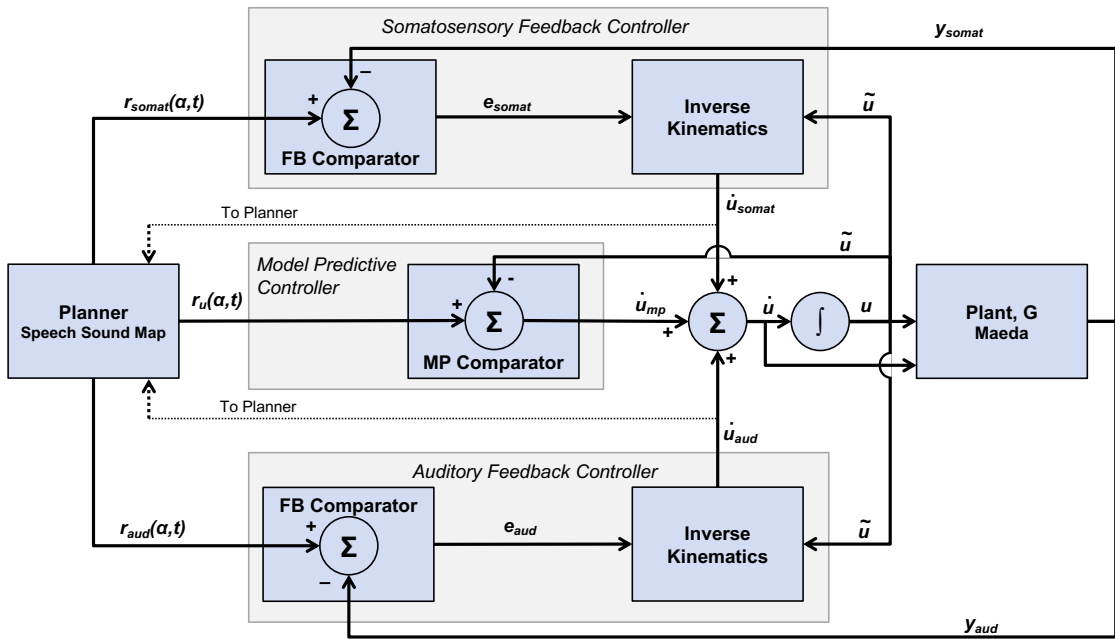
FIG. 6. Control architecture of the DIVA model. The DIVA model has two feedback paths, auditory and somatosensory, that are schematically identical, and a model-predictive pathway. The feedback pathways compute an error between the planner's signal and the output of the plant. This error is then used in conjunction with the state of the plant, $u$, to create a feedback control signal similar to the integrated model predictive-feedback control in Figure 5c. The model predictive pathway compares the desired position of the speech articulators with their current predicted position.

441 in DIVA: the articulatory trajectory serves as input to the model-predictive controller, and

442 the sensory trajectories serve as input to the respective auditory and somatosensory feedback

443 controllers. The three-component representation of speech sounds in DIVA means that each

444 speech unit has a fully-specified articulatory trajectory and time-locked sensory expectations.

445 Uniquely among models discussed in the present paper, the sensory expectations are not

446 generated online through an internal model, as in a state feedback controller.

447   The model predictive component of DIVA compares the predetermined desired position

448   of the speech articulators at each point in time, $\mathbf{r}_u(t)$, with their current predicted position,

449   $\tilde{\mathbf{u}}$, generating a control signal, $\dot{\mathbf{u}}_{mp}$. Implicitly, this assumes the existence of an internal

450   model (not explicitly shown) that is able to predict the kinematic consequences of the motor

451   commands with perfect accuracy. In order to generate the mobility state prediction, DIVA

452   integrates the control signal over time. This enables comparison of the estimated state of the

453   vocal tract articulators $\tilde{\mathbf{u}}$ with the reference signal $\mathbf{r}_u(t)$ independently of sensory feedback.

454   Although the model-predictive controller is typically referred to as the "feedforward" con-

455   troller in the DIVA literature, it is not a typical feedforward controller in the sense of "open

456   loop" control traditionally described in control systems, because it relies on a comparison be-

457   tween the predicted current model state and a reference. In its current implementation, the

458   predicted state also incorporates some auditory and somatosensory feedback information,

459   as well, since those pathways converge with the model predictive pathway. However, if the

460   auditory and somatosensory feedback controllers in DIVA are entirely removed, the model

461   predictive controller would function appropriately in the absence of sensory information.

462   In the model predictive controller, the control signal is generated from the following

463   equation: $\dot{\mathbf{u}}_{mp} = g_{mp} G[\mathbf{r}_u(\alpha, t) - \tilde{\mathbf{u}}]$, where $g_{mp}$ is a scalar amplification gain applied to the

464   motor command, and G is an additional gain that can be interpreted as a "go" signal, ranging

465   between 0 (no movement) and 1 (maximal movement speed) as in Bullock and Grossberg

466   (1988). Thus, the motor command is essentially a scaled version of an error signal, where the

467   relevant error is between the articulatory reference signal and the predicted current position

468   of the plant in mobility space. Note that the version of $u$ that is used in computing the

⁴⁶⁹ error signal is neither the true position of the articulators in mobility space, nor the one

⁴⁷⁰ measured from the plant via sensory feedback, but an internal estimate of this state, $\tilde{\mathbf{u}}$. This

⁴⁷¹ estimate is generated by integrating the summed motor commands from all three controllers,

⁴⁷² and is equal to the motor position command issued to the plant. Effectively, the quantity

⁴⁷³ $\mathbf{r}_u(\alpha, t) - \tilde{\mathbf{u}}$ is an approximation of $\dot{\mathbf{u}}$ prior to scaling. The predicted current position of

⁴⁷⁴ the plant is used purely as a way of converting the reference signal into a velocity, because

⁴⁷⁵ the reference signal (a set of articulatory positions) cannot be used directly as a motor

⁴⁷⁶ command (which must specify a change in those positions). Alternative ways of computing

⁴⁷⁷ the motor command would eliminate the need for the model-predictive component of the

⁴⁷⁸ feedforward controller, converting it into a true "open-loop" system. For example, the

⁴⁷⁹ planner could approximate the first derivative of the entire articulatory plan, and issue that

⁴⁸⁰ as the reference signal. Alternatively, the planner could issue the reference signal within a

⁴⁸¹ window surrounding the current time point, which would allow the controller to approximate

⁴⁸² the first derivative. Further details can be found in Appendix A.

⁴⁸³ The auditory and somatosensory feedback controllers closely follow the generic feedback

⁴⁸⁴ control architecture. The auditory task space in DIVA is defined as the first three formants

⁴⁸⁵ (F1-F3) and the somatosensory task space is defined as the positions of the individual ar-

⁴⁸⁶ ticulators (proprioception) as well as the degree of contact between separate articulators

⁴⁸⁷ (tactile sensation). Several publications have also envisioned the somatosensory space in-

⁴⁸⁸ cluding representations of constriction locations and degrees, as in Task Dynamics (refer

⁴⁸⁹ to sections describing Task Dynamics, below). The computations performed by the sen-

⁴⁹⁰ sory feedback controllers in DIVA begin with a comparison between the reference signal

491 and the sensory output of the plant to produce an error signal in sensory space. For the

492 sake of simplicity, only the auditory feedback computations will be presented here, but the

493 form is the same for the somatosensory pathways. The auditory error signal is defined as:

494 $\mathbf{e}_{aud} = \mathbf{r}_{aud}(\alpha, t) - \mathbf{y}_{aud}$. This auditory task-space error is then transformed into a mobility-

495 space error via the inverse kinematic equation: $\dot{\mathbf{u}}_{aud} = g_{aud}\mathbf{J}(\tilde{\mathbf{u}})^{-1}\mathbf{e}_{aud}$. The matrix $\mathbf{J}(\mathbf{u})$ is

496 known as the Jacobian, which provides a mapping between changes in mobility space and

497 changes in task space. This mapping is dependent on the current mobility state ($\mathbf{u}$) or, as in

498 DIVA, a prediction of that state ($\tilde{\mathbf{u}}$). Specifically in DIVA, the Jacobian contains the rate of

499 change for each of the dimensions of the task space for a corresponding change in mobility

500 space. The matrix $\mathbf{J}(\mathbf{u})^{-1}$, is a pseudoinverse of the Jacobian, which allows for transforming

501 task-space changes into mobility-space changes. The final motor command is then generated

502 as the transformed error signal multiplied by a fixed gain, $g_{aud}$. This represents a kind of

503 proportional control, where the motor command, ignoring transformations for the moment,

504 is simply a scaled version of the error signal. Further details can be found in Appendix A.

505 The output of the model predictive controller and sensory feedback controllers are

506 summed to generate the final control signal, $\dot{\mathbf{u}}$. Thus, the final control signal passed to

507 the plant is the velocity of the articulators (or $\dot{\mathbf{u}}$) needed to produce the desired change in

508 the position of the articulators (termed *motor movement command*). The control signal

509 additionally includes the integration of $\dot{\mathbf{u}}$ over time ($\mathbf{u}$, or *motor position command*). This

510 combined motor movement and position command is passed to the plant to drive changes

511 in the position of the articulators. The plant also produces sensory outputs based on the

512 position of articulators at each time point, $\mathbf{y}_{aud}$ and $\mathbf{y}_{somat}$. In DIVA, the output of the

513 plant is in the space of the reference signal (F1-F3 for the auditory reference, position of

514 the articulators as well as articulator contact for the somatosensory reference). This avoids

515 needing to model an auditory or somatosensory perceptual system.

516 An important detail to note is that the auditory and somatosensory reference signals

517 are specified not as unique trajectories with a single value at each time point, but as time-

518 varying regions. The error signal for each space (auditory or somatosensory) is the distance

519 from the current state to the edge of these regions. Thus, larger regions will allow greater

520 variability in production, as no corrective error signal will be generated for any production

521 that falls within the target region.

522 DIVA simulations have been able to qualitatively match human behavioral responses

523 to auditory and mechanical perturbations (Guenther *et al.*, 2006; Tourville *et al.*, 2008;

524 Villacorta *et al.*, 2007). The model has also been used to derive variable productions of /r/

525 (Nieto-Castanon *et al.*, 2005) based on a particular auditory target (low F3), a so-called

526 "trading relationship" or "motor equivalence" where multiple articulatory configurations

527 can be used for the same phoneme. Some older versions of DIVA that used time-invariant

528 targets are also able to model carry-over and anticipatory coarticulation through the use of

529 convex target regions (Guenther *et al.*, 1995).

530 Speech acquisition and learning have also received substantial consideration in the devel-

531 opment of DIVA. The primary mechanism for learning within the model involves updating

532 the motor plan based on generated auditory and somatosensory feedback motor commands.

533 Details of this adaptive modification to the motor plan fall outside the scope of the present

534 review. Nonetheless, this pathway is indicated by an open, labelled arrow in Figure 6.

535    In addition to establishing the architecture of the speech motor control system, one of the

536    primary motivations behind DIVA is establishing the neural basis of speech motor control.

537    Individual components of DIVA have been mapped onto particular brain regions based on

538    experimental neuroimaging results and model simulations (Bohland *et al.*, 2006; Ghosh *et al.*,

539    2008; Golfinopoulos *et al.*, 2011; Guenther *et al.*, 2006; Tourville *et al.*, 2008), and simulation

540    studies have provided good matches to behavioral and neural activity recorded from human

541    speakers during auditory and somatosensory perturbation experiments (Golfinopoulos *et al.*,

542    2011; Niziolek *et al.*, 2013; Tourville *et al.*, 2008; Villacorta *et al.*, 2007).

543    **B.  Task Dynamics**

544    The primary focus of the Task Dynamics model has been to model how invariant linguistic

545    targets can generate continuous and context-dependent articulatory movements. The central

546    hypothesis of this model is that articulatory movements are directed by the evolution of a

547    task-level dynamical system whose invariant parameters are determined by the linguistic

548    content of an utterance.  TD was formulated by Saltzman and Kelso (1987) in general

549    motor terms, and then by Saltzman and Munhall (1989) in the particular context of speech

550    production (see Figure 7). TD is essentially a feedback control architecture, as described in

551    Figure 7. The controller uses a feedback comparator to relate the desired state issued by the

552    planner ($\mathbf{r}_x(\alpha, t)$) to the current state of the system ($\mathbf{x}$). On the basis of this comparison ($\mathbf{e_x}$),

553    the controller computes a desired acceleration in task space ($\ddot{\mathbf{x}}$) which is then transformed

554    into a desired acceleration in mobility space ($\ddot{\mathbf{u}}$). A crucial aspect of Task Dynamics is that

555    both the desired state issued by the planner and the comparison performed by the controller

556 occur in task space, not mobility space. This necessitates a transformation of the desired

557 acceleration in task space into mobility space before it can be utilized as a motor command.

558 The plant in the Task Dynamics model is the CASY model (Iskarous *et al.*, 2003; Rubin

559 *et al.*, 1996), which is a geometric model of the vocal tract, similar in spirit to Maeda's
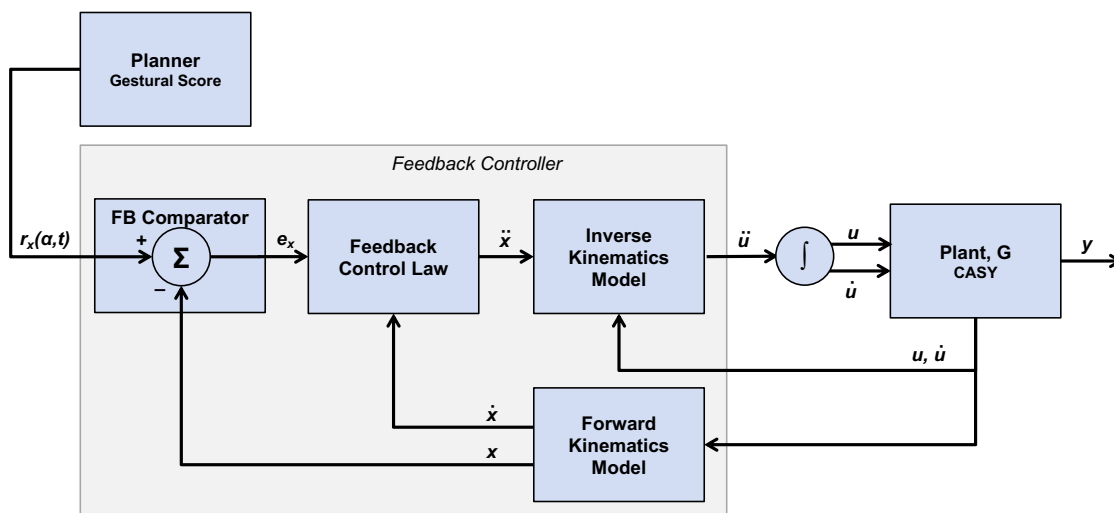
560 model.



FIG. 7. Control architecture of the Task Dynamics model. The system state, **x**, is broken out as both the state and change in state (first derivative), **ẋ**. This information is used by the controller in the rectangle. Comparing this diagram to Figure 4b, one can see TD is a feedback control architecture.

561 One view represented in the literature and in the community of Task Dynamics is that

562 it does not incorporate a feedback process. This misconception was perhaps most recently

563 mentioned in print by Kröger and Birkholz (2007) who stated that a serious problem with

564 the Task Dynamics approach has been the fact that "perception [presumably feedback] as

565 a control instance for production is not considered". Based on the discussion above, it

566 should be clear that Task Dynamics is, in fact, primarily a feedback-driven system. One

567 criticism that could be made of task dynamics, however, is that the model, as implemented,

568 treats the feedback process as noiseless and instantaneous, which is overly simplistic. Given

569 that the focus in task dynamics was on the development of the dynamic control law, this

570 simplification would seem to stem from the specific emphases and interests of the authors,

571 rather than some central conceptualization of speech motor control. Such was suggested

572 by the authors in at least one publication (Saltzman and Kelso, 1987). It is also true that

573 TD does not incorporate auditory feedback, which may, indeed, be a central property of

574 the model. Similarly, the model assumes that the current state of the plant in mobility

575 space is directly reflected via somatosensory feedback. Note that this is essentially the

576 same assumption that DIVA makes, where part of the sensory feedback signal is simply the

577 positions of the articulators.

578 The computations performed by the controller in TD begin with a comparison between

579 the (task-space) reference signal and the task-space position of the plant to produce an error

580 signal: $\mathbf{e}_x = \mathbf{r}_x(\alpha, t) - \mathbf{x}$. The error signal is then used, along with the task-space velocity

581 of the plant, $\dot{\mathbf{x}}$, to update the task-space acceleration of the plant via the feedback control

582 law (called the "forward dynamics equation" in the literature): $\ddot{\mathbf{x}} = -M^{-1}B\dot{\mathbf{x}} - M^{-1}K\mathbf{e}_x$,

583 where M is a diagonal matrix of inertial parameters, B is is a diagonal matrix of damping

584 coefficients, and K is a matrix of stiffness coefficients. Thus, the feedback control law takes

585 the form of a second-order dynamical system that transforms the error signal into the second

586 derivative of the task-space variable $\mathbf{x}$. Since the task-space acceleration cannot be used

587 directly as a motor command, it is necessary to transform this task-space acceleration into

588 a mobility-space acceleration ($\ddot{\mathbf{u}}$). This is accomplished through the use of a pseudo-inverse

589 Jacobian function: $\ddot{\mathbf{u}} = \mathbf{J}^{-1}(\mathbf{u})[\ddot{\mathbf{x}} - \dot{\mathbf{J}}(\mathbf{u}, \dot{\mathbf{u}})\dot{\mathbf{u}}]$. This mobility space acceleration can then be

590 integrated to produce mobility-space velocity and position signals, $(\mathbf{u}, \dot{\mathbf{u}})$, that can be used

591 by the plant to drive changes in the position of the speech articulators. Further details can

592 be found in Appendix A.

593    TD views speech motor control as a problem of point attractor dynamics. That is, motor

594 tasks are conceptualized as points in task space, toward which the system is drawn by means

595 of some governing control law which is a function of the system state. Task Dynamics de-

596 scribes the control law as a damped oscillator system (i.e., second-order dynamical system).

597 Damped oscillator dynamics have a number of desirable properties in terms of defining a

598 control law. In addition to the fact that damped oscillator dynamics are well-understood

599 and easily characterized, the use of such dynamics to model task-directed behavior has the

600 advantages that action patterns will be globally smooth and continuous.

601    TD is closely related to proportional-derivative control. It is common practice in engi-

602 neering control systems to take integral or derivative information of the error signal into ac-

603 count (e.g., the ubiquitous proportional-derivative, PD, and proportional-integral-derivative,

604 PID, controllers – e.g., Åström and Hägglund (1995)). Integrating the feedback error, for

605 instance, allows a controller to recognize accumulated errors, which it can then attempt

606 to nullify. Using the derivative of the feedback error, on the other hand, can minimize

607 undesirable future trends in the error signal, such as overshoot, oscillation and instabil-

608 ity. In PD control, the control signal $\mathbf{u}_{PD}$ is simply a weighted combination (given some

609 weight matrices $K_P$ and $K_D$) of the error signal and its first derivative with respect to time:

610 $\mathbf{u}_{PD} = K_P\mathbf{e}_x + K_D\dot{\mathbf{e}}_x$. This equation looks remarkably similar to the feedback control law

611 from TD: $\ddot{\mathbf{x}} = -M^{-1}B\dot{\mathbf{x}} - M^{-1}K\mathbf{e}_x$, except that weights are specified, and $\dot{\mathbf{x}}$ is substituted

612 for $\dot{\mathbf{e}}_x$. It can be easily shown that $\ddot{\mathbf{x}} = \mathbf{u}_{PD}$, given that $K_P = M^{-1}K$ and $K_D = M^{-1}B$,

613 and knowing that $\mathbf{r}_x$ has a constant value, and therefore $\dot{\mathbf{r}}_x = 0$. Thus, TD is equivalent to

614 PD control up to the generation of the task variable acceleration signal, but differs in the

615 additional transformation of the task space variables into mobility space, and integration of

616 the mobility space variables.

617     The task space in TD is defined in terms of high-level articulatory tasks (in contrast to

618 the positions of the individual articulators themselves). For speech, the tasks are suggested

619 to be constriction actions (i.e., gestures) of the vocal tract, such as achieving closure of the

620 lips, rather than the positions of the individual speech articulators (for the lip closure task,

621 these would include the upper and lower lips as well as the jaw). A point attractor task

622 is derived by the planner from a time-varying "gestural score" that issues the desired task

623 state as a function of the currently active articulatory gestures. This definition allows TD to

624 be easily put together with Articulatory Phonology (Browman and Goldstein, 1986). These

625 two components form the basis for the perspective on speech production widely associated

626 with Haskins Laboratories. Nevertheless, Task Dynamics and Articulatory Phonology are

627 separate models that address different questions. Articulatory Phonology – proposed roughly

628 in parallel with Task Dynamics – asserts that articulatory gestures are the primitive units

629 of spoken language. Gestures themselves are conceptualized with AP as discrete vocal tract

630 constriction actions, which can be composed into gestural "scores" that function as a motor

631 program for a given utterance. Therefore, in broad terms, Articulatory Phonology addresses

632 the question of how speech tasks should be defined, and how they can be composed into a

633 motor program, whereas Task Dynamics addresses the question of how those tasks can be

634 achieved and how that motor program can be realized in a physical system.

635 Use of second-order dynamics directly connects TD to research on action planning and

636 execution in biological systems. For instance, the VITE model is an influential neural-

637 inspired network model for explaining kinematic trajectory formation of directed movement

638 (Bullock and Grossberg, 1988). VITE comprises a network of three interacting hypothesized

639 neural populations, each coding a distinct quantity that is needed in the generation of the

640 motor command, given some target position. These neural populations encode quantities

641 related to the present position of the system, the desired target position, and the difference

642 between the target and the present position. These interacting populations are configured

643 in such a way that there are many structural similarities to the control architecture of TD.

644 The result of these similarities is that the present position of a population will move in a way

645 that is consistent with a $2^{nd}$-order dynamical system, much like Task Dynamics (as pointed

646 out by, e.g., (Beamish $et$ $al.$, 2006)).

647 One of the strengths of the model is accounting for coarticulatory effects. Coarticulation

648 in this model is seen as arising from temporal overlap of independent and invariant articu-

649 latory gestures – the so-called coproduction model of coarticulation (Browman $et$ $al.$, 1992,

650 1995; Fowler $et$ $al.$, 1993). Other coarticulatory effects, such as clear vs. dark /l/ alterations,

651 have been modeled at the planning level as changes in the temporal organization of gestures

652 (Browman $et$ $al.$, 1992, 1995; Zsiga $et$ $al.$, 1994).

653    Very early results from the Task Dynamics model showed that it was capable of repro-

654    ducing the compensatory behavior seen in mechanical perturbation experiments, where a

655    lowered jaw position during production of a bilabial stop is compensated for by a higher

656    lower lip and lower upper lip (Saltzman *et al.*, 1986).  However, the model is unable to

657    account for auditory perturbations, as there is no auditory feedback channel.

658    Task Dynamics can produce simple speech-rate effects by changing the dynamical pa-

659    rameters of the control law – e.g., by making the task-space motions more or less damped.

660    In addition to these linear rate effects, the Task Dynamics model is able to produce a wide

661    range of non-linear temporal effects seen in speech.  Through the $\pi$-gesture model (Byrd

662    *et al.*, 2003), the model is able to capture the non-linear slowing found adjacent to prosodic

663    boundaries as well as capture many of the spatial effects, such as larger movements (Fougeron

664    *et al.*, 1997), seen at those boundaries within a single framework.  More recent work has ex-

665    tended the model to account for syllable structure and prosodic prominence (Saltzman *et al.*,

666    2008).  While some recent work has started to explore neural mechanisms for some of the

667    components of the model (Tilsen *et al.*, 2016), and a connection to the VITE neural model

668    (Lammert *et al.*, 2018) has been established, the components of TD have not been explicitly

669    related to specific neural structures.

670    **C.   State Feedback Control**

671    The State Feedback Control for speech production (SFC) model is a speech-specific in-

672    stantiation of the general Kalman filter-type architecture in Figure 5c (Houde and Chang,

673    2015; Houde and Nagarajan, 2011).  The primary focus of SFC has been to apply the in-

674   sights gained from state feedback approaches in other motor domains to speech. This type

675   of model is used widely in current theories of motor control in non-speech domains (e.g.,

676   work from Diedrichsen *et al.* (2010); Scott (2004); Shadmehr and Krakauer (2008); Todorov

677   (2004); Todorov and Jordan (2002)), and is an evolution of a traditional feedback control

678   system (Fig 4b). Recall that a primary challenge of feedback control is that sensory feed-

679   back is typically noisy and delayed, making the instantaneous state of the plant impossible

680   to know with perfect accuracy. By adopting a Kalman filter-type architecture (Fig 5c), SFC

681   presents, in a speech motor control context, one method by which sensory feedback may be

682   integrated with internal model predictions to produce improved estimates of the state of the

683   plant.

684       In the SFC model (shown in Figure 8), estimation of the plant state is done by an *observer*

685   (refer to Fig 5c). This observer receives a copy of the outgoing motor command issued by

686   the control law (also known as the efference copy) [4]. Based on this signal, the observer

687   predicts how the plant will move at the next time step $((\tilde{\mathbf{x}}, \dot{\tilde{\mathbf{x}}}))$ as well as the auditory and

688   somatosensory feedback that will be received based on that predicted movement $((\tilde{\mathbf{y}}, \dot{\tilde{\mathbf{y}}}))$.

689   The predicted sensory feedback is then compared with actual sensory feedback to calculate

690   a sensory error $((\mathbf{e_y}, \dot{\mathbf{e}}_\mathbf{y}))$. This error is then converted to a task state error (or task gain),

691   via a gain function. Finally, the task state $((\hat{\mathbf{x}}, \dot{\hat{\mathbf{x}}}))$ is estimated using the predicted state as

692   well as the weighted sensory errors for both auditory and somatosensory predictions. As the

693   gains associated with the sensory errors are assigned to optimize the final estimation, the

694   observer in SFC functions is a Kalman filter (Todorov and Jordan, 2002), which provides

695   the optimal *a posteriori* estimate of the state, under the assumption of linear processes of
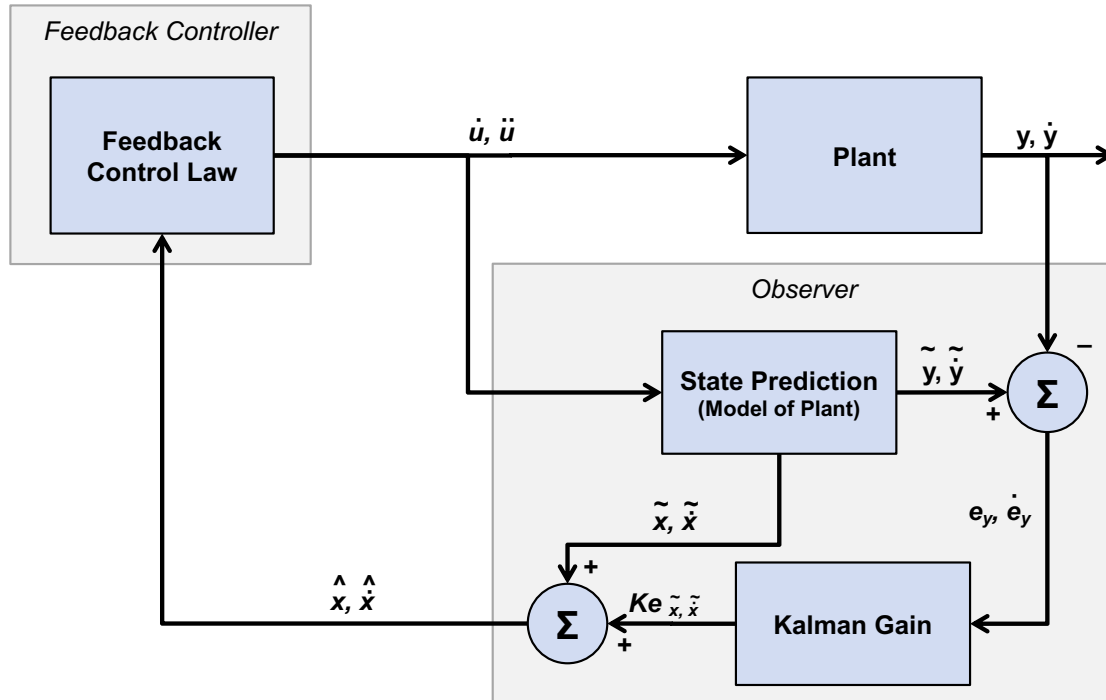
FIG. 8. Control architecture of the State Feedback Control (SFC) model. The final state estimate passed back to the controller as a feedback signal, $(\hat{\mathbf{x}}, \hat{\dot{\mathbf{x}}})$, is derived from a combination of a state prediction process and sensory processes. Comparing this diagram to Figure 5c, one can see that SFC is an integrated model predictive feedback control architecture.

prediction and sensory feedback. Note that the sensory feedback the observer receives at any time point reflects the past state of the plant, while the state prediction reflects the current state. This delay is accounted for by delaying the sensory prediction before computation of sensory errors.

The model does not make explicit mention of a reference signal or a planner, and by extension does not make explicit mention of any comparison between sensory feedback and a reference. Providing a detailed description of the controller has not been a focus in the development of SFC, and therefore the controller, as presented in the literature, is

704 represented by a generalized feedback control law which is a function $U(\hat{\mathbf{x}}, \dot{\hat{\mathbf{x}}})$ of only the state

705 estimate. This control law could take almost any form. However, the authors of this review

706 expect any feedback control law that produces reasonable speech production behavior would

707 need to be a function of some kind of reference, whether an explicitly planned trajectory or

708 a gestural score. Indeed, specifying the details of this feedback control law in SFC, and the

709 addition of a planner module, have been a primary motivation for the development of the

710 FACTS model, described below.

711 By combining a state prediction with sensory feedback to estimate the current state, the

712 SFC model is able to act quickly by operating principally on an internal prediction of the

713 plant state. This also allows the system to operate in the absence of sensory feedback, either

714 when that feedback is too delayed to be of use (as for very fast speech movements) or when

715 sensory feedback is unavailable (as when speaking in loud noise or in cases of non-congenital

716 deafness). Yet, the system is still able to respond when the internal predictions do not

717 match the incoming sensory feedback (either due to errors in the prediction process or due

718 to external perturbations of the plant). Thus, this system combines the major advantage of

719 traditional feedback control systems (robustness to perturbations) with that of feedforward

720 control (fast, accurate movement even in the absence of sensory feedback).

721 Note that, in SFC as currently implemented, there is no distinction between task space

722 and mobility space; they are effectively collapsed into a single space, such that commands

723 are issued in task space. This means that the current implementation of SFC is only able to

724 model a system where the goals of speech production are the same as the mobility space of

725  the system. SFC has been implemented to control pitch, where the fundamental frequency

726  of vocal fold vibration maps onto a one-dimensional mass-spring system.

727  This model has been shown to accurately reproduce the behavior patterns of human

728  participants in pitch-alteration studies (Houde *et al.*, 2006). The model has also been shown

729  to reproduce two neural effects seen in human speech: 1) the reduction seen in cortical

730  electroencephalography (EEG) or magnetoencephalography (MEG) signals when speaking

731  compared to listening to the one's own speech played back over headphones or speakers

732  (speech induced suppression) and 2) the enhancement of the EEG /MEG signals when seen

733  when one's speech is externally perturbed compared to when it is unperturbed (speech

734  perturbation).

735  **D.   FACTS**

736  Recently, a new model – the Feedback Aware Control of Tasks in Speech (FACTS) model

737  – has been proposed that combines aspects of both Task Dynamics and State Feedback

738  Control (Parrell *et al.*, 2006). Building on TD and SFC, FACTS combines elements of feed-

739  back control and model predictive control. FACTS is an attempt to combine the strengths

740  of each model, while addressing the major shortcomings of each. Specifically, the Task

741  Dynamics model includes a well-developed control law that relates the movements of the

742  speech articulators to high-level tasks, but assumes perfect knowledge of the state of the

743  vocal tract. Conversely, State Feedback Control focuses principally on how the state of the

744  plant can be estimated from sensory information given the noise and time delays inherent

745  in auditory and somatosensory perception, but has to date only been used to control a very

746 simplistic one-dimensional model of pitch. FACTS combines the concept of state prediction

747 and estimation from SFC with the planning model and vocal tract control of TD.
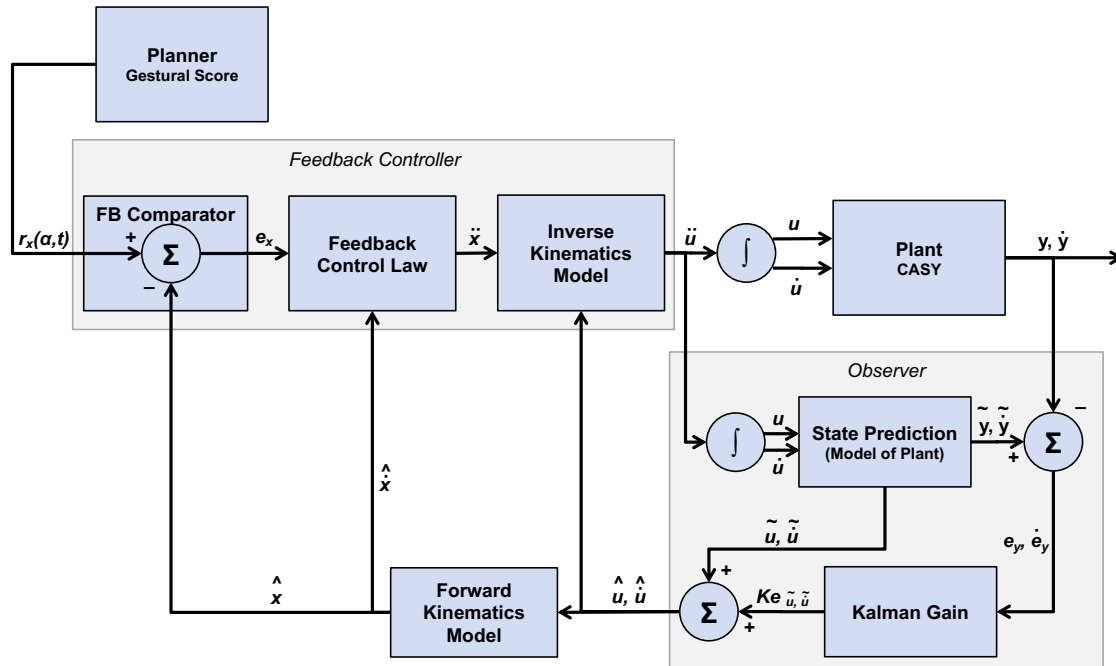


FIG. 9. Control architecture of the Feedback Aware Control of Tasks in Speech (FACTS) model. FACTS builds upon the architecture of the Task Dynamics model by substituting an estimate of the mobility-space state for the true state through an observer module. The observer generates this mobility state estimate through a combination of an internal mobility state prediction and multisensory feedback. As such, FACTS is an implemenation of an integrated model predictive controller, like SFC.

748 The architecture of FACTS is shown in Figure 9. The control component of the model is

749 the same as that for the Task Dynamic model, with a planner generating a gestural score,

750 which is passed to a feedback controller to generate changes at the task ($\ddot{\mathbf{x}}$) and mobility

751 ($\ddot{\mathbf{u}}$) levels. This final motor command, $\ddot{\mathbf{u}}$, is passed to the plant to produce articulator

752  movements as in Task Dynamics. However, where Task Dynamics passes the current plant

753  and tasks states directly back to the feedback controller, FACTS uses an observer to estimate

754  the task and plant states, as in the earlier SFC model. The final motor command $\ddot{\mathbf{u}}$ is passed

755  to an internal model of the plant to generate predicted articulator positions $((\tilde{\mathbf{u}}, \dot{\tilde{\mathbf{u}}}))$, as well

756  as auditory and somatosensory feedback $((\tilde{\mathbf{y}}, \dot{\tilde{\mathbf{y}}}))$. The estimated sensory feedback is then

757  compared with sensory feedback from the plant to generate a sensory error $((\mathbf{e_y}, \dot{\mathbf{e_y}}))$. The

758  estimated mobility state is generated from the predicted mobility state and the sensory

759  error via an unscented Kalman filter, an extension of the linear Kalman filter to nonlinear

760  systems (Wan and Van Der Merwe, 2001). The estimated mobility state is then converted to

761  an estimated task state, needed by the feedback controller to generate the motor command

762  at the next time point, via the same forward kinematics function as in Task Dynamics.

763  The FACTS model is relatively new, and so remains mostly untested. However, the model

764  is able to qualitatively reproduce human responses to external perturbations, including full

765  compensation for mechanical perturbations and partial compensation for auditory pertur-

766  bations (Parrell et al., 2006). This partial compensation is a function of both auditory and

767  somatosensory acuity. One of the features of FACTS is that it builds on the successes of the

768  Task Dynamics model. Since many of the mechanisms of the controller are shared between

769  the two models, FACTS can reproduce the successes of the Task Dynamics model, including

770  coarticulatory effects.

**E. ACT**

The primary focus in the the ACTion-based model of speech production, speech perception, and speech acquisition (ACT) model is the acquisition and development of speech motor control. Kröger *et al.* (2009) introduced ACT as a neurocomputational model that draws elements from both DIVA and Task Dynamics. The architecture of ACT, shown in Figure 10, is essentially a feedforward controller when viewed between the motor plan and the plant. DIVA-style dual auditory/somatosensory feedback pathways are also part of the model. However, those pathways feed indirectly to the planner, by way of high-level comparisons against abstract phoneme templates. Within the present framework, information used to modify the motor plan is considered to be part of the planner, and is therefore outside the scope of low-level control, as defined here. This pathway is indicated by an open, labelled arrow in Figure 10. The plant in ACT is a three-dimensional kinematic model with articulatory control parameters similar to the Maeda and CASY models (Birkholz *et al.*, 2006).

The planner in the ACT model relates to both the speech sound map of DIVA and the gestural score in the Task Dynamics model. Like in DIVA, the basic unit of speech is assumed to be the syllable, and each syllable is represented by a model neuron in the phonemic map (cf. the speech sound map in DIVA). As in DIVA, these abstract syllable representations are linked to specific motor and sensory plans. This is accomplished in ACT through the phonetic map. Unlike in DIVA, where the motor plan is represented as a time-varying desired articulatory position signal, the motor plans in ACT are defined in
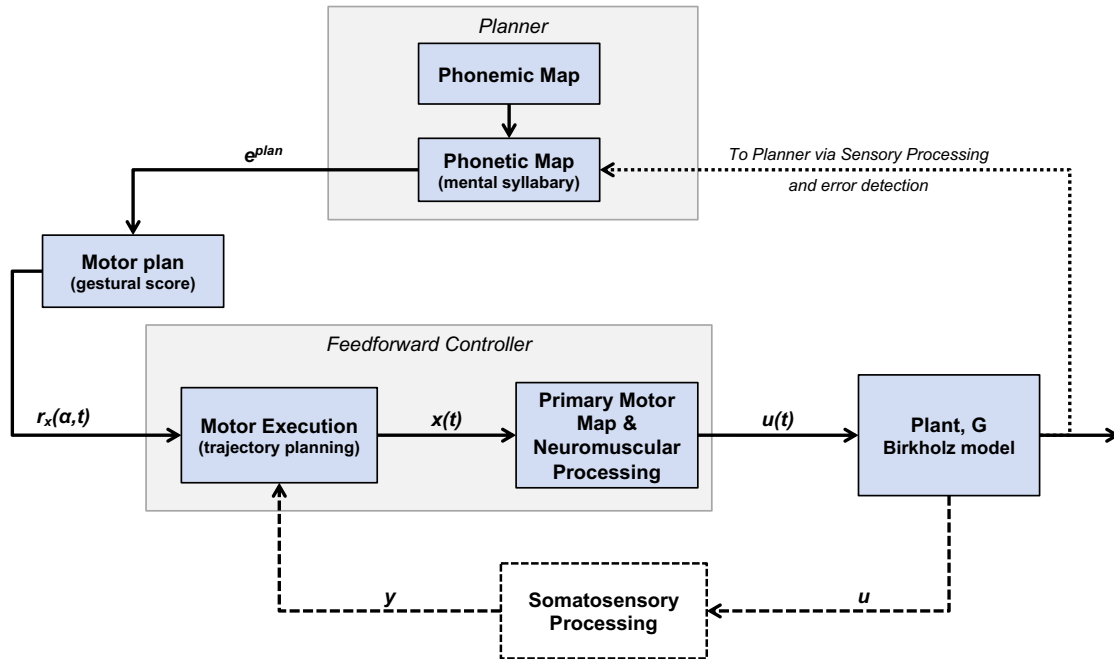
FIG. 10. Control architecture of the ACTion-based model of speech production, speech perception, and speech acquisition (ACT model). ACT draws from both DIVA and Task Dynamics for its architecture, with the model comprising both feedforward and feedback pathways (both somatosensory and auditory), but relying on point-attractor dynamics for its reference signal.

terms of high-level dynamic tasks (or gestures) as in Task Dynamics. Each motor plan is, in effect, a gestural score, which defines the activation levels and temporal extent of each speech gesture, with each speech gesture being defined as a dynamical point-attractor ($\mathbf{r}_x$).

The phonetic map, in addition to linking the syllable to the motor plan, also links the syllable to associated sensory (auditory and somatosensory) expectations. One conceptual difference between ACT and DIVA is that DIVA views the sensory plans as the targets of speech that have associated motor plans, while in ACT the targets are the high-level task gestures with associated sensory expectations. This conceptual difference is reflected principally in terms of how the models are trained (an issue not taken up within the scope

801 of the present review), but the basic architecture of the models is essentially the same: a

802 high-level syllable activates a motor plan used for feedforward or model-predictive control

803 and a sensory plan which can be compared against afferent sensory information.

804 The core control architecture in the ACT model borrows ideas from Task Dynamics, but

805 is quite distinct. As discussed above, TD makes use of task-space comparisons between

806 a reference, derived from the task-based gestural score, and the current (somatosensory)

807 system state to control task-space movements given a control law that is consistent with

808 damped oscillator dynamics. ACT, on the other hand, uses the reference, similarly derived,

809 to directly drive motor action in a feedforward fashion. This is accomplished by the motor

810 execution module, which uses the reference $\mathbf{r}_x(\alpha, t)$ to generate a trajectory in task space

811 ($\mathbf{x}(t)$) that is consistent with damped oscillator dynamics. The task-space trajectory must

812 be transformed into a mobility-space trajectory ($\mathbf{u}(t)$) that can be used as a control signal

813 to drive movements of the plant. This transformation is accomplished by the primary motor

814 map. A subsequent neuromuscular processing step exists in the model, and is presently

815 implemented as a direct, linear mapping. Plans exist for this component to eventually

816 map control signals onto individual and/or combined muscle groups in a neuromuscular

817 model. An additional pathway for somatosensory feedback processing is also planned. This

818 is indicated by dashed lines in Figure 10. This feedback pathway, included in published

819 figures representing ACT, would be used to "control motor execution", presumably in a

820 fashion similar to DIVA. This pathway has not yet been implemented, and the details of its

821 properties have not been fully developed.

822   Like DIVA, the ACT model also has dual somatosensory and auditory feedback pathways.

823   The principal way these feedback pathways are used in the model is to compare the current

824   state of the plant against pre-learned templates representing the desired somatosensory and

825   auditory states. A crucial difference between ACT and other models is that this error signal

826   is used to influence the motor plan, rather than as part of the controller. That is, sensory

827   feedback is used to detect sensory errors for updating the phonetic map to drive trial-to-trial

828   adaption, a model of development and learning.

829   One difference between the ACT model and others is that the mappings that relate the dif-

830   ferent signals (syllables$\mapsto \mathbf{r}_x$, syllables$\mapsto \mathbf{r}_y$, $\mathbf{r}_x \mapsto \mathbf{r}_u$, $r_u \mapsto \dot{\mathbf{u}}$, etc.) are implemented via tunable

831   neural networks rather than as closed-form mathematical expressions. These networks are

832   tuned during a learning phase. Some versions of DIVA presented in the literature, especially

833   earlier in DIVA's development, had neural networks involved in these mechanisms (Guenther,

834   1994). The use of trained neural network models for these transformations allows for flexi-

835   bility in the form of the transformations. It opens the possibility that the transformations

836   might take forms that deviate in unexpected, and potentially even biologically-plausible,

837   ways when compared to mathematically-driven transformations typically adopted. The use

838   of neural networks also makes it likely, however, that key transformations, such as the control

839   law and the inverse kinematic transformations, cannot be easily written down analytically

840   in closed form.

841   The ACT model is able to produce motor equivalence in articulators linked to the same

842   gesture due to the use of high-level tasks rather than articulatory positions as the basic unit

843   of the motor plan (Kröger *et al.*, 2009). The model is also capable of adaptive learning

844    based on high-level auditory errors or somatosensory perturbations, by changing the motor

845    plan. However, the lack of feedback pathways in the controller means online compensation

846    to these perturbations is not accounted for. The ACT model includes hypotheses about the

847    neural structures that underlies the different components but to date has not been used to

848    generate simulated neural activity to compare to neural data.

849    **F.    GEPPETO**

850    The GEPPETO (GEstures shaped by the Physics and by a PErceptually oriented Targets

851    Optimization) model (Patri *et al.*, 2015; Perrier *et al.*, 1996, 2005; ) is a model of speech

852    control based on the equilibrium point hypothesis (Feldman, 1986). The primary focus of

853    GEPPETO has been to investigate the hypotheses that 1) targets for speech production

854    are discrete and phonemic, 2) biomechanics plays a non-trivial role in speech motor control,

855    and 3) speech motor control employs optimal planning principles. In GEPPETO, as in the

856    equilibrium point hypothesis, control occurs at the level of individual muscle lengths. Thus,

857    the mobility space in GEPPETO is composed of lengths, $u_k$, of individual muscles $k$. The

858    command generated by the central controller is a muscle length, or threshold, above which

859    motor neurons will be recruited to contract the muscle. This threshold length is known as

860    the equilibrium point or $\lambda$. Afferent feedback from the muscle about the current muscle

861    length is compared against the current $\lambda$, and contractile force is generated if the muscle

862    length is above the threshold. In GEPPETO, the activation ($A$) of each muscle at time

863    t is based on both the current muscle length $u$ and the current change in muscle length

864    $\dot{u}$: $A_{(k,t)} = [u_k(t) - \lambda_k(t) + \gamma_k \dot{u}_k(t)]^+$, where $\gamma$ is a damping parameter that stabilizes the

865 system. Muscle activation is only generated when the muscle length is greater than $\lambda$:

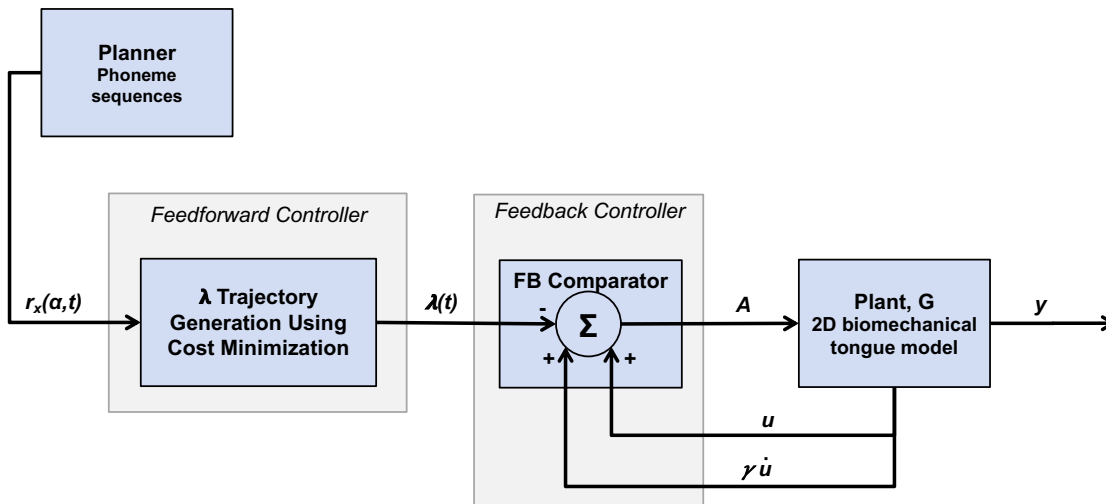866 $[A]^+ = \{A, if A >= 0; 0, otherwise\}$.



FIG. 11. Control architecture of the GEPPETO model. GEPPETO is based on the equilibrium point hypothesis, employing feedback control at the level of individual muscles, with relatively realistic biomechanics to move the speech articulators. Control is mediated by a feedforward process that transforms acoustic speech targets into equilibrium point values.

867 The muscle activation generated by the feedback controller then leads to the generation of

868 force ($f$) in the individual muscles at the level of the plant: $f_k(\lambda_k, t) = \rho_k[exp(c_k A_k(\lambda_k, t) -$

869 $1)]$,where $\rho$ is a magnitude parameter related to the cross-sectional area of the muscle and

870 $c$ is a feedback gain. In this feedback control architecture, force can be generated either

871 by changes in the current $\lambda$ or by changes in the length of the muscles. Importantly in

872 this approach, the ultimate position of the plant results from a combination of descending

873 control ($\lambda$ values), plant biomechanics, and physical constraints.

874      The GEPPETO model, shown in Figure 11, combines the low-level feedback control

875   structure of an equilibrium point model with a high-level feedforward controller that takes

876   acoustic speech targets, defined as convex regions in acoustic (F1-F2-F3) space, as input and

877   output $\lambda$ values that are passed to the feedback controller. Thus, the task space for GEP-

878   PETO is acoustic in nature (though see  for a recent extension of the model to additionally

879   include somatosensory targets). Critically, given the emphasis on the physics of the speech

880   plant, GEPPETO uses a dynamical biomechanical model of the plant with control occur-

881   ring at the level of muscles rather at the level of geometric model parameters/articulators

882   as in the Maeda or CASY plant models. Most published papers on GEPPETO include only

883   the tongue as a controllable articulator. It is modeled as a finite-element model with six

884   muscles whose lengths can be independently controlled. The other vocal tract surfaces and

885   articulators are fixed.

886      The output of the planner in GEPPETO is a series of $n$ acoustic speech targets

887   $(\phi^1, \ldots, \phi^n)$, each of which has an intended duration $(T^1, \ldots, T^n)$. This duration can

888   be affected by variables such as speech rate or stress. An additional constraint sets the

889   amount of effort to be used for each speech target $(w^1, \ldots, w^n)$ , where effort is based on

890   the amount of force that will be generated to produce that target across all the muscles of

891   the plant, categorized into three levels: $w \in \{\text{``}weak\text{''}, \text{``}medium\text{''}, \text{``}strong\text{''}\}$.

892      This time series of targets $\mathbf{r}_x(t) = \{(\phi^1, w^1, T^1), \ldots, (\phi^n, w^n, T^n)\}$ is then passed to the

893   feedforward controller to generate a time series of $\lambda$ values for each of the six muscles in the

894   plant, $(\lambda_1(t), \ldots, \lambda_6(t))$. These $\lambda$ trajectories are generated for each utterance using an op-

895   timization procedure that minimizes displacements in mobility space (i.e. changes in muscle

896 lengths) while producing tongue movements that will achieve the required acoustic targets

897 at the required time with the required amount of effort. In this optimization process, learned

898 internal models are used to estimate the amount of force and acoustic signal generated for

899 any given motor command.

900 GEPPETO shares certain characteristics with other models. First, speech goals are de-

901 fined as regions in acoustic space (F1-F2-F3), as in DIVA. Second, feedback signals are

902 never directly compared against the output of the planner, as in ACT. GEPPETO differs

903 in key ways from other models, however. First, the speech targets in GEPPETO are hy-

904 pothesized to be discrete in time, rather than time-varying regions as in DIVA. Second,

905 the feedforward and feedback controllers in GEPPETO are arranged in a unique, serial or

906 hierarchical arrangement, such that the output of the feedforward controller is used as the

907 input to the lower-level feedback controller. Third, unlike the preplanned trajectories in

908 DIVA, GEPPETO generates new movement plans for each utterance. Finally, it is notable

909 that GEPPETO is unique in the fact that the plant's inputs are not given in mobility-space

910 variables.

911 The largest success of the GEPPETO model has been to replicate many of the character-

912 istic kinematic patterns of speech movements, including velocity profiles (Payan *et al.*, 1997),

913 tongue loops in velar stops (Perrier *et al.*, 2003), and the relationship between velocity and

914 movement curvature (Perrier *et al.*, 2008). This work shows that many of these phenomena

915 need not be directly controlled, since in GEPPETO they are emergent properties of linear

916 changes in $\lambda$ values. One of the drawbacks of the optimization approach in GEPPETO is

917 that it produces identical trajectories each time the same utterance is produced, unlike the

918 variability seen in natural speech. Recently, however, the GEPPETO model was expanded

919 by implementing it in a probabilistic Bayesian framework (B-GEPPETO) that is able to ac-

920 count for token-to-token variability (Patri *et al.*, 2015; ). This newer model also incorporates

921 somatosensory phonemic targets in addition to auditory targets.

### G. Other models

923 All the above models include, at a minimum, the ability to generate motor commands

924 based on some motor plan. These motor commands are then used to move a vocal tract

925 model of some kind. While such complete models are the primary focus of the current review,

926 it is important to also mention more conceptual models which have not been implemented to

927 the same degree. The Hierarchical State Feedback Control model (HSFC) (Hickok, 2012a,b,

928 2014) is an attempt to combine speech motor and psycholinguistic approaches to speech pro-

929 duction. It is a version of an integrated predictive/feedback controller, sharing some aspects

930 with the State Feedback Control model of speech production (Houde and Nagarajan, 2011).

931 Tian & Poeppel (Tian and Poeppel, 2010) propose a hybrid model predictive/feedback con-

932 trol model of speech motor control. The overall architecture is also very similar to the State

933 Feedback Control model (Houde and Nagarajan, 2011).

934 A few other models of speech motor control have been proposed that have focused primar-

935 ily on the biomechanical properties of the plant rather than on the control architecture per

936 se (Dang and Honda, 2002, 2004; Laboissiere *et al.*, 2018; Ostry *et al.*, 1996; Perrier *et al.*,

937 1996; Sanguineti *et al.*, 1990). While these models do not relate control to linguistic speech

938 targets (i.e. describe how or why certain muscle contraction patterns would be used), the

success of these models in recreating measured articulatory trajectories deserves mention in the context of the present review.

One class of these models (reviewed in Sanguineti *et al.* (1998), is based on the equilibrium point control. While this is the same general approach as taken by the GEPPETO model, the focus of this work differs. Rather than implementing control of the speech motor system in terms of higher-level linguistic or task-directed (auditory, articulatory) control, these models focus principally on how muscle forces are generated to move the speech articulators. Typically, the goal is to drive movements to match measured human speech kinematics. These models essentially implement a feedback controller, albeit one that functions entirely at the level of the plant without any distinction between task and mobility space. A separate set of biomechanical models assumes that muscle activations are the output of the controller, rather than equilibrium points (Dang and Honda, 2002, 2004). This is a purely feedforward control architecture.

Both the equilibrium point models (Sanguineti *et al.*, 1998) as well as the direct activation models (Dang and Honda, 2004) have been shown to fit articulatory data well using similar biomechanical models. Interestingly, results from both models suggest that motor commands to certain muscles (or muscle groups) will drive the speech articulators towards a similar location regardless of their initial position. This suggests that speech motor control may be simplified by the use of muscle synergies that will drive the system to a target spatial configuration without the need for complex inverse dynamics models that calculate the precise muscle activations needed for each individual movement.

960  One important thing to note is that, because they focus on the generation of muscle forces

961  given some given motor commands, this class of models is generally complementary to and

962  compatible with control models that output motor commands as articulatory positions, and

963  ignore the generation of muscle activations (such as DIVA, TD, ACT, and FACTS). With

964  some modifications, the output of these models could serve as the input to an equilibrium

965  point model or the Dang & Honda model. In fact, equilibrium point control has been

966  implemented within the DIVA architecture (Zandipour *et al.*, 2004).

967  **IV.  DISCUSSION**

968  The primary goal of the current paper has been to clearly lay out the architectures of a

969  crucial component of existing speech motor control models: the control layer (see Figure 1),

970  that attempts to produce accurate tracking of speech articulation kinematics given a motor

971  plan. Common terminology and basic principles of motor control were used to describe each

972  model, to understand the commonalities between these models, as well as how they differ.

973  It was shown that these models can be cast as special cases of generalized feedforward (Fig

974  4a), feedback (Fig 4b), and model predictive (Fig 4c) controllers. The models discussed

975  here differ in which of these components are used (e.g., some are lacking either feedforward

976  or feedback elements of control), and in the detailed implementation of these mechanisms.

977  These differences are summarized in Table I. Speech production is, however, a complex

978  process with many additional and important considerations, including higher-level motor

979  planning, linguistic, communicative and even social considerations, as well as learning and

980  developmental aspects, all of which contribute to the wide variety of speaking styles observed

| | DIVA | TD | SFC | FACTS | ACT | GEPPETO |
|---|---|---|---|---|---|---|
| Feedback Pathway | Y | Y* | Y | Y | Y | Y |
| Feedforward Pathway | N | N | N | N | Y | Y |
| Internal Prediction/State Estimation | Y | N | Y | Y | N | N |
| Principal Reference | Tourville and Guenther (2011) | Saltzman and Munhall (1989) | Houde and Nagarajan (2011) | Ramanarayanan et al. (2016) | Kröger et al. (2009) | Perrier et al. (2005) |

TABLE I. *Summary of which aspects of motor control modeling are present in each model.*

in real human speech. These aspects are beyond the scope of the present paper, but would make an interesting subject future reviews.

There are clear differences among models in terms of how their final execution of speech motor control is influenced by feedback signals originating from the plant. ACT, for instance, incorporates no explicit feedback into its control mechanisms. SFC implements proportional control, meaning that the motor commands are linearly proportional to the feedback error. DIVA's also implements proportional control which, for its hybrid architecture means that motor commands are linearly proportional to both the error (in the feedback pathway) and the reference (in the feedforward pathway) signals. The simplicity of these designs relative to common engineering approaches is notable. As mentioned above, and by way of example, engineering control systems often take information about the integral or derivative of the error signal into account in order to provide quicker convergence to the target and to deal with persistent errors, respectively. TD – as well as FACTS, by way of adopting key control elements from TD – provides slightly more complexity through a form of PD control, albeit not strictly in the traditional engineering sense of PD control.

A related distinction between the models under consideration is how they function in the absence of feedback. TD, for instance, is solely a feedback architecture, and cannot function

998 in the absence of feedback signals. Similarly, GEPPETO would not be able to function in

999 the absence of proprioceptive feedback about muscle length. Other models could continue

1000 to function without feedback. DIVA is a hybrid feedback/model predictive architecture that

1001 could rely exclusively on its model predictive mechanisms to generate motor commands in

1002 the absence of explicit feedback. With the presence of feedback signals, SFC and FACTS

1003 can utilize that feedback to produce optimal or near-optimal state estimates (under certain

1004 strong assumptions, such as linearity of the plant (Kalman *et al.*, 1960)), but in the absence

1005 of feedback can still rely on the internal state prediction component of their broader state

1006 estimation process to continue functioning through model predictive control. ACT is a

1007 purely feedforward architecture that can function as designed in the absence of sensory

1008 feedback. However, this also means that it is not sensitive to sensory feedback, unlike the

1009 human speech motor control system.

1010 Among models that incorporate feedback, one of the most basic differences is whether

1011 certain feedback signals are treated as idealized signals that are directly and instantaneously

1012 observable, or whether they are treated as true sensory signals that may be potentially

1013 noisy/delayed, subject to conditioning by internal models and that correspond with known

1014 neurological signals. While it seems intuitively correct that any model of biological motor

1015 control should focus on the latter, the former has been sometimes intentionally chosen

1016 in specific aspects of the models, in the interest of focusing on other aspects of control.

1017 TD provides only an idealized view of feedback concerning the positions and velocities of

1018 articulators that does not model the sensory processes in any meaningful way. DIVA, TD

1019 and FACTS also make simplifying assumptions about the somatosensory feedback signal,

which is assumed to be more or less equivalent to the plant's mobility variables. DIVA's auditory and somatosensory feedback are slightly less idealized in that they correspond to known, independent neurological pathways and can incorporate delays associated with sensory transduction and processing. SFC and FACTS begin with the assumption that sensory feedback will be noisy and/or inaccurate, and use that assumption to motivate the well-elaborated integration of sensory feedback with internal model predictions to provide more accurate estimates of the state of the plant. GEPPETO provides perhaps the most realistic implementation of somatosensory feedback given that the feedback in the model (muscle length and change in muscle length) corresponds to well-known afferent signals from muscle spindles. However, no current models seriously attempt to model the sensory system itself – they take it as given that critical information (e.g., formants, articulatory positions) can be extracted from the raw sensory input.

Most models are purely kinematic in how they approach control, in that motor commands are stated in kinematic terms (i.e., as state configurations, and not as forces) and do not account for dynamical considerations related to the effects of inertia, centrifugal and centripetal forces, and the effects of gravity. Control systems that are strictly linear, rigid and slow-moving, highly damped, or that have specialized designs can sometimes operate purely kinematically. It seems likely, however, given existing literature (e.g., Derrick *et al.* (2015); Ostry *et al.* (1996)) that such considerations may be non-negligible for speech production in the biological case. A kinematic approach can be, in the opinion of the authors, partially attributed to models of the plant used in most speech motor control models, which are nearly all kinematic in nature. It is worth noting that other plant models are attempting to provide

enhanced biomechanics (Derrick *et al.*, 2015; Gick *et al.*, 2011; Lloyd *et al.*, 2012) as well, even if a full review of biomechanical vocal tract models is beyond the scope of the present review. GEPPETO represents a notable effort to move beyond kinematic treatment of control, and of the plant, by incorporating a mobility space that represents muscle lengths, as well as motor commands that represent muscle activations that are used to generate muscle forces in a relatively realistic biomechanical model of the tongue.

All architectures rely on a motor plan of some kind – whether an explicitly planned trajectory or a gestural score – that is formed at a higher level of motor processing, and which is issued to the controller in order to be carried out. SFC is a partial exception to this general statement in that, as mentioned above, that model does not explicitly mention the incorporation of a plan, even though the generalized structure of its controller would be able to incorporate a planning module if more detailed specification required it (a specification which FACTS has subsequently elaborated upon). Models of speech motor planning have been discussed and elaborated upon in the literature (Bohland *et al.*, 2010; Byrd *et al.*, 2009; Civier *et al.*, 2013; Saltzman *et al.*, 2008), and display a surprising amount of variety. Although the planning level is beyond the scope of this paper, it is worth noting the variety of planning mechanisms that have been proposed in order to help narrow some of the longest-standing debates concerning speech motor control. In particular, drawing a clear distinction between control architectures and planning mechanisms, as this review has attempted to do, makes it apparent that much of the debate over the quality of competing models of speech production would appear to be concentrated at the planning level, and not at the level of control. For instance, issues surrounding the nature of production goals (e.g., acoustic vs.

articulatory) and the composition of those goals into utterance-size units would primarily be a concern of the planning level. Any role for muscle synergies (Ramanarayanan *et al.*, 2014) and motor primitives would be most naturally incorporated into the planning level, and not the level discussed in this review. The nature of speech production goals has been the subject of particularly strong debate for decades, and is reflected in the nature of the feedback and reference signals in the models, which may be auditory and/or somatosensory, as in DIVA and SFC, or articulatory, as in Task Dynamics. Interestingly, the nature of the feedback signals would appear to have little bearing on the specific architectural choices of the models – the architectures being general enough to handle a range of signals without substantial changes to their configuration.

The parameters that determine the overall characteristics of control are time-invariant in most current models, thereby limiting the models in their ability to capture specific aspects of behavior that require those parameters to change over time. Models may struggle, for instance, to account for interspeaker differences, or long-term changes in speech motor control that occur during development and aging, that could be modeled by adjustments to control parameters. Controllers that adapt their parameters over time are the subject of adaptive control (Åström and Wittenmark, 2013). This well-studied branch of control theory may provide a foundation for models of speech production to incorporate such parameter adjustments as a way to represent the mechanisms of differences or changes mentioned above. A full treatment of adaptive control is outside the defined scope of the present paper, as are issues surrounding speech development. Nonetheless, it should be noted that inroads into adaptive control have been made by some of the models discussed here. ACT allows

1086 for motor planning to be adapted based on sensory feedback errors. DIVA, too, adapts

1087 planned trajectories based on the feedback controller output. This adaptation is of primary

1088 importance during development, but can lead to changes at any time.

1089 Shorter time-scale cognitive and physiological factors – for instance, due to attention,

1090 fatigue and motivation – as well as stochastic variability (Munhall *et al.*, 1994; Saltzman

1091 *et al.*, 1995; Tilsen, 2017) may also most naturally be handled through adjustments to

1092 control-level parameters. Efforts have been made to model learning and adaptation at

1093 the planning level (e.g., GODIVA). However, the value of the proportional gain in DIVA's

1094 controller, as well as the weights assigned to the feedback and model predictive pathways

1095 in their contribution to the motor command, are assumed to be fixed in fully adult speech.

1096 Similarly, the damping and stiffness parameters of the controller in TD are fixed in value. A

1097 notable counterexample to this generalization comes from Kalman filter-based architectures,

1098 such as SFC and FACTS, which change the weight assigned to sensory feedback and internal

1099 model predictions, toward combining them into a single state estimate, based on the degree of

1100 statistical reliability of those two pathways. Such adaptation may be useful in modeling the

1101 impact of sensory feedback impairment on speech motor control. Another notable example

1102 of this type is DIVA's GO signal, which can be adjusted by higher-level processes in order

1103 to control the initiation of movement and overall speaking rate.

1104 A clear understanding of how the various models are structured can aid in clearly defining

1105 theoretical questions of interest. For instance, the many similarities of the models discussed

1106 in this review naturally raise questions about what is gained by allowing the remaining

1107 model dissimilarities to persist, and whether the models can converge to a single, unified

1108 model of the control layer in speech motor control. There is no mathematical reason why the

1109 feedforward/feedback pathways embodied by DIVA couldn't be combined with the forward

1110 dynamic control of TD, as well as the feedback/internal model-based state estimation in

1111 SFC. Indeed, FACTS, as a combination of complementary elements of TD and SFC, has

1112 already taken a step toward beginning these potentially useful combinations. Whether

1113 such a unification is sensible from a theoretical point of view, and precisely what form

1114 that unification might take, can be stated very precisely in mathematical terms using the

1115 model architectures. In general, models can help in defining and circumscribing the space

1116 of possible architectures and solutions to a specified biological control problem (Schaal and

1117 Schweighofer, 2005).

1118 A related, empirical question is whether a model unification is useful for explaining ob-

1119 servations from human speech data. Among the many benefits of developing formal models

1120 of speech motor control is that models can be used to make specific, quantitative predictions

1121 about human speech behavior that are testable in light of data. The predictive capabilities

1122 of formal models can also guide the design of new experiments to test specific aspects of

1123 theory and modeling, inspired by the behavioral predictions of the model, and perhaps pi-

1124 loted *in silico*. Empirical questions regarding the models need not be limited to observable

1125 behaviors, either. Models can also facilitate clearer connections to be drawn between specific

1126 model mechanisms and their observed neurological counterparts, either through structural or

1127 functional neuro-imaging. The connection between engineering and biological mechanisms

1128 has been well developed in several domains of motor control, including speech motor control

(Guenther *et al.*, 1998) and oculomotor control (Lisberger, 1988; Robinson, 1981; Shibata and Schaal, 2001).

The utility of speech motor control models additionally extends beyonds clarifying and formalizing our understanding of speech motor control itself. Models can also be useful for practical applications in speech synthesis. Control models, coupled with faithful mechanical models of the vocal tract, hold promise for applications in flexible and expressive speech synthesis. This kind of synthesis is typically called *articulatory synthesis*. Shadle and Damper (2002) outlined several complementary advantages that articulatory synthesizers should have over now widely adopted data-driven approaches like concatenative synthesis (Black, 2002) and Hidden Markov Model-based synthesis (Schroeter, 2006). Among these advantages are (a) the promise of producing speech associated with extraordinary speakers (e.g., an exceptional opera singer) or hypothetical speakers, from whom data can be difficult or impossible to collect, (b) the promise of changing the quality or type of speaker without having to perform additional statistical training of the synthesizer, (c) the promise of having meaningful parameters that can be helpful in fixing or adjusting the synthesizer output, in addition to providing insights into human speech production.

The models discussed here, in addition to being formal and mechanistic, are also causal, by intention of their development and by virtue of their historical context. Causal models can, as such, serve to encapsulate current theoretical understanding of the mechanisms underlying speech motor control into a compact and rigorous form. Analysis of speech behavior, even in response to challenging or contrived situations, may not always be sufficient for inferring the causal mechanisms of those behaviors. An individual's sensorimotor behav-

1151   ior is, in general, the result of a complex mixture of stable and mature control mechanisms,

1152   learned and adaptive strategies, and possible individual-specific speaking strategies and im-

1153   pairments. Therefore, inferring the underlying mechanisms that contribute to observed

1154   behaviors is exceedingly difficult without an underlying framework. Neurologically relevant,

1155   mechanistic models of sensorimotor control provide a neurocomputational substrate which

1156   can aid in establishing causal relationships among the many component pathways and model

1157   parameters. By modeling and resynthesizing human behaviors, mechanistic models can infer

1158   the mechanisms underlying observed responses, including both impairment mechanisms and

1159   neural adaptation to those impairments. This process is termed *system identification* in an

1160   engineering context, and recent advances in methods for system identification have facili-

1161   tated application to biological multivariate, closed-loop control systems (Engelhart *et al.*,

1162   2016) and human sensorimotor control systems (Boonstra *et al.*, 2013; Engelhart *et al.*,

1163   2015). Inroads have also recently been made in applying similar approaches in the domain

1164   of typical (Mitra *et al.*, 2010) and pathological (Ciccarelli *et al.*, 2016) speech motor control.

1165   **V.   CONCLUSION**

1166       In scanning the published literature on formal models of speech motor control, it is

1167   perhaps understandable to be left with the impression that a dizzying variety of qualitatively

1168   distinct models have been presented. Among all the models, DIVA and TD stand out as

1169   having a relatively long history of representation in the literature, and the efforts to develop

1170   them have remained almost entirely separate. SFC and FACTS make clear and related

1171   modeling contributions that enable the expressive power to explain specific empirical results

1172  in speech production. ACT is inspired by both DIVA and TD, but has a structure all

1173  its own. GEPPETO is the result of yet another distinct effort at model development; it

1174  is concerned with biomechanical considerations in the plant. Clearly, there is a healthy

1175  amount of variety in the various model architectures, especially in their specific use and

1176  method of combining the three essential functional components: feedforward, feedback and

1177  model predictive. However, it is nonetheless possible to view these models as belonging to a

1178  single, coherent framework. The present paper has attempted to cut through the difficulties

1179  associated with varying presentation and terminology, and to directly compare the models

1180  against the backdrop of such a framework. By presenting a clear comparison of the points

1181  of agreement and disagreement among the various models, as well as establishing areas

1182  where all models can be improved, this work can provide a foundation for future model

1183  development to improve our understanding of the speech motor system.

1184  **RESOURCES**

1185  Several of the models discussed in this paper (DIVA, TD, CASY and the Maeda model)

1186  have been implemented as software tools, and are available for download online. Their

1187  addresses on the World Wide Web are included in the references below.

1188  **ACKNOWLEDGMENTS**

1192 understanding the details of their models, as well as the three thoughtful reviewers of the

1193 manuscript.

## Appendix A

1195   To aid the speech motor control practitioner, this Appendix consolidates the key algo-

1196 rithmic steps of three control architectures: Task Dynamics (TD), Directions into Velocities

1197 of Articulators (DIVA), and State Feedback Control (SFC). Bold lower case letters represent

1198 vectors, and bold upper case letters represent matrices. A single overhead dot represents a

1199 time derivative, and a double dot represents a second order time derivative.

### 1.   Directions Into Velocities of Articulators (DIVA)

1201   The Directions Into Velocities of Articulators (DIVA) model is a control architecture

1202 developed by (Guenther *et al.*, 2006) that uses a hybrid of feedback control and model

1203 predictive control. The model has been realized in software, and is available online (Nieto-

1204 Castanon, 2016).

#### a.   *Algorithm*

1206   In the DIVA model predictive controller, the mobility space, $\mathbf{u}$, and state of the plant, $\mathbf{x}$,

1207 are identical, so $\mathbf{u} = \mathbf{x}$. Table II describes the variables in DIVA.

1208   1. Compute a model-predictive control signal (termed *feedforward* in the published liter-

1209     ature on DIVA).

1210     (a) Compute an error using the reference target in mobility space and the current

1211         predicted state of the plant.

$$\mathbf{e_u} = \mathbf{r_u}\left(\mathbf{t}\right) - \tilde{\mathbf{u}} \tag{1}$$

1212     (b) Compute a feedforward control update by scaling the error signal.

$$\dot{\mathbf{u}}_{mp} = g_{mp}G\mathbf{e_u} \tag{2}$$

2. Compute a feedback-driven control signal using the reference target and the sensed plant output to compute an error in task space. Then, use a pseudoinverse Jacobian to convert the error from task space to mobility space. Do this in both the auditory and somatosensory feedback pathways.

$$\mathbf{e}_{aud} = \mathbf{r}_{aud}\left(t\right) - \mathbf{y}_{aud} \tag{3}$$

$$\dot{\mathbf{u}}_{aud} = g_{aud}\mathbf{J}(\mathbf{u})^{-1}\mathbf{e}_{aud} \tag{4}$$

$$\mathbf{e}_{somat} = \mathbf{r}_{somat}\left(t\right) - \mathbf{y}_{somat} \tag{5}$$

$$\dot{\mathbf{u}}_{somat} = g_{somat}\mathbf{J}(\mathbf{u})^{-1}\mathbf{e}_{somat} \tag{6}$$

3. Combine the feedforward and feedback control updates to determine the new plant state.

$$\mathbf{u} = \int \left(\dot{\mathbf{u}}_{mp} + \dot{\mathbf{u}}_{aud} + \dot{\mathbf{u}}_{somat}\right) dt \tag{7}$$

$$\tilde{\mathbf{u}} = \mathbf{u} \tag{8}$$

TABLE II. DIVA variables.

| Variable | Description |
|---|---|
| $\mathbf{e}_u$ | Error between reference target in mobility space and last command issued to the plant |
| $\mathbf{e}_{aud}$, $\mathbf{e}_{somat}$ | Error between the reference target in task space and sensed task space output |
| $\mathbf{r}_u(t)$ | Reference target in mobility space. Defined at each point in time as a region with a center and bounds of acceptable performance. |
| $\mathbf{r}_{aud}(t)$, $\mathbf{r}_{somat}(t)$ | Reference target in task space. Defined at each point in time as a region with a center and bounds of acceptable performance. |
| $\dot{\mathbf{u}}_{aud}$ $\dot{\mathbf{u}}_{somat}$ | Change in mobility space position based on error in task space. A task space velocity update. |
| $\dot{\mathbf{u}}_{ff}$ | Change in mobility space position based on error in mobility space. A task space velocity update. |
| $\mathbf{u}$ | Mobility space position. Computed by integrating the feedforward and feedback mobility space velocities. |
| $\mathbf{y}_{aud}$, $\mathbf{y}_{somat}$ | Task space output |
| $g_{ff}$ | Gain applied to feedforward velocity update |
| $g_{aud}$, $g_{somat}$ | Gain applied to feedback velocity update |
| $G$ | Gain with a value between 0 and 1 that constrains velocities in mobility space from 0 to their maximum. |
| $\mathbf{J}(\mathbf{u})^{-1}$ | Pseudoinverse of the Jacobian. The pseudoinverse converts errors in task space to changes in velocity in mobility space. The pseudoinverse can be computed as the Moore-Penrose pseudoinverse. |

## 2. Task Dynamics

Task Dynamics is a feedback control architecture developed by (Saltzman and Kelso, 1987; Saltzman and Munhall, 1989). The architecture has been realized in software in the Task Dynamics Application (TADA) (Nam *et al.*, 2006) and available online (Nam, 2012).

### a.  Algorithm

The Task Dynamics algorithm is described below, and all variables are defined in Table III.

1. Compute error in task space. In Task Dynamics, the task space, $\mathbf{y}$, and the state, $\mathbf{x}$, are identical, so $\mathbf{y} = \mathbf{x}$, and the error is

$$\mathbf{e_x} = \mathbf{r_x}\left(\alpha, \mathbf{t}\right) - \mathbf{x}. \tag{9}$$

2. Use a dynamical system description of the controller, a second order ordinary differential equation, to compute the new acceleration state of the plant in task space as

$$\ddot{\mathbf{x}} = -\mathbf{M}^{-1}\mathbf{B}\dot{\mathbf{x}} - \mathbf{M}^{-1}\mathbf{K}\mathbf{e_x}. \tag{10}$$

3. Use a pseudoinverse Jacobian to convert the task space acceleration to mobility space acceleration by

$$\ddot{\mathbf{u}} = \mathbf{J^{-1}}(\mathbf{u})\left[\ddot{\mathbf{x}} - \dot{\mathbf{J}}\left(\mathbf{u}, \dot{\mathbf{u}}\right)(\mathbf{u})\right]. \tag{11}$$

4. Integrate mobility space acceleration to get velocity and position in mobility space, so

$$\dot{\mathbf{u}} = \int \ddot{\mathbf{u}}dt \tag{12}$$

$$\mathbf{u} = \iint \ddot{\mathbf{u}}dt \tag{13}$$

### 3.  State Feedback Control

The State Feedback Control is a hybrid feedback/model-predictive control architecture proposed by (Houde and Nagarajan, 2011). Note that the notation used here follows the

TABLE III. Task dynamic variables.

| Variable | Description |
|----------|-------------|
| $\mathbf{x}$, $\dot{\mathbf{x}}$, $\ddot{\mathbf{x}}$ | Task space position, velocity, and acceleration, m by 1 vectors |
| $\mathbf{u}$, $\dot{\mathbf{u}}$, $\ddot{\mathbf{u}}$ | Mobility space position, velocity, and acceleration, n by 1 vectors |
| $\mathbf{r_x}\left(\alpha\right)$ | Reference target in task space, m by 1 vector |
| $\mathbf{M}$ | Inertial coefficients, m by m diagonal matrix |
| $\mathbf{B}$ | Damping coefficients, m by m diagonal matrix |
| $\mathbf{J}$ | Jacobian transformation from mobility space to task space. An m by n matrix with elements $J_{ij} = \frac{\partial x_i}{\partial u_j}$ |
| $\mathbf{J}^{-1}$ | The (pseudo) inverse of the Jacobian. The Moore-Penrose pseudoinverse may be used, or other constraints can be applied to allow inversion of a non-square Jacobian. |
| $\dot{\mathbf{J}}$ | The time derivative of each element of the Jacobian. |

originally-published notation, and differs slightly from the simplified notation used in the main body of the present paper.

### a. Algorithm

1. Create a control update using the current estimate of the plant state by

$$\mathbf{u}_{t-1} = U_t\left(\hat{\mathbf{x}}_{t-1}\right). \tag{14}$$

2. Create the new, true plant state using the true plant dynamics, $G_{dyn}$, by

$$\mathbf{x}_t = G_{dyn}\left(\mathbf{u}_{t-1}, \mathbf{x}_{t-1}\right). \tag{15}$$

3. Create a new, predicted estimate of the plant state using the previous estimate of the plant state, $\hat{\mathbf{x}}_{t-1}$, the previous control signal, $\mathbf{u}_{t-1}$, and an estimate of the plant dynamics, $\hat{G}_{dyn}$, by

$$\tilde{\mathbf{x}}_{t|t-1} = \hat{G}_{dyn}\left(\mathbf{u}_{t-1}, \hat{\mathbf{x}}_{t-1}\right).\tag{16}$$

X

4. Generate the subsequent plant output using the true plant transformation from plant state to plant output by

$$\mathbf{y}_t = G_{out}\left(\mathbf{x}_t\right).\tag{17}$$

5. Create a correction term to the plant state estimate using the sensed feedback from the true plant by

$$\mathbf{y}_{t-N} = G_{out}\left(\mathbf{x}_{t-N}\right)\tag{18}$$

$$\tilde{\mathbf{y}}_{t-\hat{N}} = \hat{G}_{out}\left(\mathbf{u}_{t-1}, \hat{\mathbf{x}}_{(t|t-1)-\hat{N}}\right)\tag{19}$$

$$\mathbf{e}_{\mathbf{y}_{t-\hat{N}}} = \mathbf{y}_{t-N} - \tilde{\mathbf{y}}_{t-\hat{N}}\tag{20}$$

$$\mathbf{e}_{\tilde{\mathbf{x}}_t} = \mathbf{K}_t\left(\mathbf{e}_{\mathbf{y}_{t-\hat{N}}}\right).\tag{21}$$

6. Combine the initial plant state estimate and the correction term to create the current estimate of the plant state by

$$\hat{\mathbf{x}}_t = \tilde{\mathbf{x}}_{t|t-1} + \mathbf{e}_{\tilde{\mathbf{x}}_t}.\tag{22}$$

TABLE IV. State feedback control variables.

| Variable | Description |
| --- | --- |
| $\mathbf{x}_t$ | True plant state at time $t$. |
| $\hat{\mathbf{x}}_t$ | Estimate of the plant state at time $t$ using both the sensed plant output and the predicted plant state. |
| $\mathbf{e}_{\mathbf{y}_{t-N}}$ | Error between the sensed plant output and the predicted plant output. |
| $\mathbf{e}_{\tilde{\mathbf{x}}_t}$ | Error update applied to the predicted estimate of the plant state to create $\hat{\mathbf{x}}_t$ |
| $\tilde{\mathbf{y}}_{t-N}$ | The predicted plant output, derived from estimates of the plant state, estimates of the feedback delay, and estimate of the plant transform from state to output. |
| $\mathbf{K}_t\left(\mathbf{e}_{\mathbf{y}_{t-\hat{N}}}\right)$ | Transformation (e.g. a Kalman gain) applied to the error between the predicted and sensed plant output. The transformation allows the actual plant output to influence the estimate of the plant state. |
| $\tilde{\mathbf{x}}_{(t\|t-1)-\hat{N}}$ | Predicted estimate of the plant state using only the previous estimate of the plant state, the control signal, and the estimated plant dynamics. |
| $G_{dyn}, \hat{G}_{dyn}$ | True and estimated plant dynamics. |
| $G_{out}, \hat{G}_{out}$ | True and estimated transformation from plant state to plant output. |
| $\mathbf{y}_t$ | True plant output. |
| $\mathbf{u}_t$ | Control update to the plant. |
| $\hat{\mathbf{x}}_{(t\|t-1)}$ | Estimate of the plant state based on the control update to the plant, the estimate of the plant dynamics, and the previous estimate of the plant state. |
| $U_t\left(\mathbf{x}_t\right)$ | Controller that issues a control update based on the current estimated state of the plant. While a reference target is not shown in Houde (2011), presumably the reference is internal to $U_t$. |
| $N, \hat{N}$ | Actual delay and estimated delay between the plant output and the sensing of the plant output. |

1243  **Appendix B**

1244     This appendix presents two articulatory speech synthesizers commonly referenced in the

1245  literature: the Configurable Articulatory Synthesizer (CASY), and the Maeda model. Bold

1246  lower case letters represent vectors, and bold upper case letters represent matrices. A single

1247 overhead dot represents a time derivative, and a double dot represents a second order time

1248 derivative.

## 4. Configurable Articulatory Synthesizer

1250 The Configurable Articulatory Synthesizer (CASY) is a geometric model of the vocal tract

1251 based on the work of Mermelstein (1973) and developed by Rubin *et al.* (1996) and Iskarous

1252 *et al.* (2003). The governing equations are presented below, taken from Lammert (2013).

1253 The "q" variables in Lammert *et al.* (2013), that represent the articulators in mobility space,

1254 have been renamed to "u" in this paper for consistency of notation (see Tables V and VI

1255 for details about the variables/constants).

$$x_{PRO} = u_{lx} \tag{23}$$

$$x_{LA} = l_{ut} \sin(a_{ut}) + l_{lt} \cos(u_{ja}) + u_{uy} - u_{ly} \tag{24}$$

$$a = u_{cl} \sin(u_{ja} + u_{ca}) \tag{25}$$

$$b = -u_{cl} \cos(u_{ja} + u_{ca}) \tag{26}$$

$$x_{TBCL} = acos\left(\frac{a - o_x}{\sqrt{(a - o_x)^2 + (b - o_y)^2}}\right) \tag{27}$$

$$x_{TBCD} = r_{ts} - \left(\sqrt{(a - o_x)^2 + (b - o_y)^2} + r_{tb}\right) \tag{28}$$

$$c = u_{ja} + u_{ta} + s_{tb}(u_{cl} - l_{tb}) \tag{29}$$

$$d = a + r_{tb} \sin(u_{ja} + a_{tc}) + u_{tl} \sin(c) \tag{30}$$

$$e = b - r_{tb} \cos(u_{ja} + a_{tc}) - u_{tl} \cos(c) \tag{31}$$

$$x_{TTCL} = acos\left(\frac{d - o_x}{\sqrt{(d - o_x)^2 + (e - o_y)^2}}\right) \tag{32}$$

$$x_{TTCD} = r_{tb} - \sqrt{(d - o_x)^2 + (e - o_y)^2} \tag{33}$$

## 5.  Maeda Articulatory Synthesizer

The Maeda articulatory speech synthesizer is a variable cross-sectional area, tube model of the vocal tract. Resonances of the tube can be computed, and these resonances are the formants. The formants can then be used to shape a vocal source (voiced or unvoiced) to create speech. A MATLAB instantiation of the Maeda synthesizer was created by Ghosh and available for download (Nieto-Castano, 2017).

TABLE V. CASY variables.

| Variable | Description |
|---|---|
| $x$ | Task space variable |
| $u$ | Mobility space variable |
| LX | Lip protrusion |
| UY | Upper lip vertical displacement |
| UT | Upper teeth |
| LY | Lower lip vertical displacement |
| JA | Jaw angle |
| CA | Tongue body angle |
| CL | Tongue body length |
| TL | Tongue tip length |
| TA | Tongue tip angle |
| LA | Lip aperture |
| PRO | Lip protrusion |
| TBCD | Tongue body constriction degree |
| TBCL | Tongue body constriction location |
| TTCD | Tongue tip constriction degree |
| TTCL | Tongue tip constriction location |

Ciccarelli (Ciccarelli, 2017) created a polynomial approximation to the vocal tract compo-

nent to allow fast formant computation and fast, tractable computation of the pseudoinverse

of the Jacobian. The polynomial approximation was determined by running the Ghosh im-

plementation of the Maeda model across a set of parameters, uniformly sampled from the

mobility space of the model, to create a lookup table of parameters and formant values.

TABLE VI. CASY constants.

| Constant | Value |
|----------|-------|
| $l_{ut}$ | 1.1438 |
| $a_{ut}$ | -0.1888 |
| $l_{lt}$ | 1.1286 |
| $o_x$ | 0.7339 |
| $o_y$ | -0.4562 |
| $r_{ts}$ | 0.4 |
| $r_{tb}$ | 0.02 |
| $a_{tc}$ | 1.7279 |
| $l_{tb}$ | 0.8482 |
| $s_{tb}$ | 4.48 |

Formant points outside the standard vowel quadrilateral as determined by visual inspection were excluded. The remaining pairs of articulator points and formants were then fit using a least squares polynomial approximation. The order of the polynomial was a compromise between the fit to the data and the complexity of the polynomial. It was found that a second order polynomial achieved a reasonable balance between these two requirements. While the mapping from articulators to formants is preserved to within a certain error, it has not been evaluated whether the relationship between articulators encoded by the polynomial fundamental alters the trajectories of articulators in previous implementations of the Maeda model.

[1276] [1]In the speech motor control literature, the term 'articulatory space' is often used instead of 'mobility space'.

[1277] The latter term is adopted from the robotics literature (Sciavicco *et al.*, 2012) here to provide a neutral

[1278] terminology for referring specifically to the configuration of the plant, whereas terminology used in the

[1279] literature often leads to confusion over whether the term 'articulatory' refers to low-level descriptions of

[1280] the plant or high-level tasks spaces defined in articulatory terms.

[1281] [2]For this example, the simplifying assumption is made that the feedback signal is in task space, i.e. $\mathbf{y_x}$

[1282] [3]optimal here means closest to the true state of the plant, where "closest" means having the smallest mean

[1283] squared error

[1284] [4]The description of SFC presented here uses a different notation than in Houde and Nagarajan (2011),

[1285] simplified for clarity of presentation. For a more complete mathematical description, see Appendix A.

[1286]

[1287] Åström, K. J., and Hägglund, T. (**1995**). *PID Controllers: Theory, Design, and Tuning*

[1288] (ISA, Research Triangle Park, NC, USA).

[1289] Åström, K. J., and Wittenmark, T. (**2013**). *Adaptive Control* (Courier Corporation).

[1290] Baraduc, P., and Perrier, P. (**2017**). "Motor control of the tongue during speech: predictions

[1291] of an optimization policy under sensorimotor noise," Neuroscience, 403–408.

[1292] Beamish, D., Bhatti, S. A., MacKenzie, I. S., and Wu, J. (**2006**). "Fifty years later: a

[1293] neurodynamic explanation of Fitts' law," Journal of The Royal Society Interface **3**(10),

[1294] 649–654.

[1295] Bellman, R. (**1957**). *Dynamic Programming* (Princeton University Press, Princeton, NJ,

[1296] USA).

[1297] Birkholz, P., Jackèl, D., and Kroger, B. J. (**2006**). "Construction and control of a three-

[1298] dimensional vocal tract model," in *Acoustics, Speech and Signal Processing, 2006. ICASSP*

1299 *2006 Proceedings. 2006 IEEE International Conference on*, IEEE, Vol. 1, pp. I–I.

1300 Black, A. W. (**2002**). "Perfect synthesis for all of the people all of the time," in *Speech*

1301 *Synthesis, 2002. Proceedings of 2002 IEEE Workshop on*, IEEE, pp. 167–170.

1302 Bohlund, J.W., and Guenther, F.H. (**2006**). "An fMRI investigation of syllable sequence

1303 production," Neuroimage **32**(2), 821–841.

1304 Bohland, J. W., Bullock, D., and Guenther, F. H. (**2010**). "Neural representations and mech-

1305 anisms for the performance of simple speech sequences," Journal of Cognitive Neuroscience

1306 **22**(7), 1504–1529.

1307 Boonstra, T. A., Schouten, A. C., and Van der Kooij, H. (**2013**). "Identification of the

1308 contribution of the ankle and hip joints to multi-segmental balance control," Journal of

1309 Neuroengineering and Rehabilitation **10**(1), 23.

1310 Brainard, M. S., and Doupe, A. J. (**2002**). "What songbirds teach us about learning,"

1311 Nature **417**(6886), 351.

1312 Browman, C. P., and Goldstein, L. M. (**1986**). "Towards an articulatory phonology," Phonol-

1313 ogy **3**(01), 219–252.

1314 Browman, C.P., and Goldstein, L.M. (**1992**). "Articulatory phonology: An overview," Pho-

1315 netica **49**(3–4), 155–180.

1316 Browman, C.P., and Goldstein, L. (**1995**). "Dynamics and articulatory phonology," *Mind as*

1317 *Motion: Dynamics, Behavior, and Cognition*, edited by R. Port and T. van Gelder (MIT

1318 Press, Boston, MA, USA).

1319 Browman, C.P., and Goldstein, L. (**1995**). "Gestural syllable position effects in American

1320 English," *Studies in Speech Production: a festschrift for Katherine Safford Harris*, edited

by F. Bell-Berti and L.J. Raphael (American Institute of Physics, Woodbury, NY, USA).

Buchaillard, S., Perrier, P., and Payan, Y. (**2009**). "A biomechanical model of cardinal vowel production: Muscle activations and the impact of gravity on tongue positioning," The Journal of the Acoustical Society of America **126**(4), 2033–2051.

Bullock, D., and Grossberg, S. (**1988**). "Neural dynamics of planned arm movements: emergent invariants and speed-accuracy properties during trajectory formation," Psychological Review **95**(1), 49–90.

Byrd, D., and Saltzman, E. (**2003**). "The elastic phrase: modeling the dynamics of boundary-adjacent lengthening," Journal of Phonetics **31**(2), 149–180.

Byrd, D., Tobin, S., Bresch, E., and Narayanan, S. (**2009**). "Timing effects of syllable structure and stress on nasals: a real-time MRI examination," Journal of Phonetics **37**(1), 97–110.

Ciccarelli, G. A. (**2017**). "Characterization of phone rate as a vocal biomarker of depression," Ph.D. thesis, Massachusetts Institute of Technology.

Ciccarelli, G. A., Quatieri, T. F., and Ghosh, S. S. (**2016**). "Neurophysiological vocal source modeling for biomarkers of disease,".

Civier, O., Bullock, D., Max, L., and Guenther, F. H. (**2013**). "Computational modeling of stuttering caused by impairments in a basal ganglia thalamo-cortical circuit involved in syllable selection and initiation," Brain and Language **126**(3), 263–278.

Dang, J., and Honda, K. (**2002**). "Estimation of vocal tract shapes from speech sounds with a physiological articulatory model," Journal of Phonetics **30**(3), 511–532.

1342 Dang, J., and Honda, K. (**2004**). "Construction and control of a physiological articulatory

1343    model," Journal of the Acoustical Society of America **115**(2), 853–70.

1344 Derrick, D., Stavness, I., and Gick, B. (**2015**). "Three speech sounds, one motor action:

1345    evidence for speech-motor disparity from english flap production," The Journal of the

1346    Acoustical Society of America **137**(3), 1493–1502.

1347 Diedrichsen, J., Shadmehr, R., and Ivry, R. B. (**2010**). "The coordination of movement:

1348    optimal feedback control and beyond," Trends in Cognitive Sciences **14**(1), 31–39.

1349 Engelhart, D., Boonstra, T. A., Aarts, R. G., Schouten, A. C., and van der Kooij, H. (**2016**).

1350    "Comparison of closed-loop system identification techniques to quantify multi-joint human

1351    balance control," Annual Reviews in Control **41**, 58–70.

1352 Engelhart, D., Pasma, J. H., Schouten, A. C., Aarts, R. G., Meskers, C. G., Maier, A. B., and

1353    van der Kooij, H. (**2015**). "Adaptation of multijoint coordination during standing balance

1354    in healthy young and healthy old individuals," Journal of Neurophysiology **115**(3), 1422–

1355    1435.

1356 Feldman, A. G. (**1986**). "Once more on the equilibrium-point hypothesis ($\lambda$ model) for

1357    motor control," Journal of Motor Behavior **18**(1), 17–54.

1358 Feldman, A., Adamovich, S., Ostry, D., and Flanagan, J. (**1990**). "The origin of electromyo-

1359    grams – explanations based on the equilibrium point hypothesis," in *Multiple Muscle Sys-*

1360    *tems* (Springer), pp. 195–213.

1361 Flash, T., and Hogan, N. (**1985**). "The coordination of arm movements: an experimentally

1362    confirmed mathematical model," Journal of Neuroscience **5**(7), 1688–1703.

1363 Fougeron, C., and Keating, P.A. (**1997**). "Articulatory strengthening at edges of prosodic

1364 domains," Journal of the Acoustical Society of America **101**(6), 3728–3740.

1365 Fowler, C., and Saltzman, E. (**1993**). "Coordination and coarticulation in speech produc-

1366 tion," Language and Speech **36**, 171-195.

1367 Garcia, C. E., Prett, D. M., and Morari, M. (**1989**). "Model predictive control: theory and

1368 practice: a survey," Automatica **25**(3), 335–348.

1369 Ghosh, S. S. (**2005**). "Understanding cortical and cerebellar contributions to speech pro-

1370 duction through modeling and functional imaging," Ph.D. thesis, Boston University.

1371 Ghosh, S.S., Tourville, J.A., and Guenther, F.H. (**2008**). "A neuroimaging study of premotor

1372 lateralization and cerebellar involvement in the production of phonemes and syllables,"

1373 Journal of Speech, Language, and Hearing Research **51**(5), 1183–1202.

1374 Gick, B., Stavness, I., Chiu, C., and Fels, S. (**2011**). "Categorical variation in lip posture is

1375 determined by quantal biomechanical-articulatory relations," Canadian Acoustics **39**(3),

1376 178–179.

1377 Golfinopoulos, E., Tourville, J. A., Bohland, J. W., Ghosh, S. S., Nieto-Castanon, A.,

1378 and Guenther, F. H. (**2011**). "fMRI investigation of unexpected somatosensory feedback

1379 perturbation during speech," Neuroimage **55**(3), 1324–1338.

1380 Guenther, F. H. (**1994**). "A neural network model of speech acquisition and motor equivalent

1381 speech production," Biological Cybernetics **72**(1), 43–53.

1382 Guenther, F.H. (**1995**). "Speech Sound Acquisition, Coarticulation, and Rate Effects in a

1383 Neural Network Model of Speech Production," Psychological Review **102**, 594-621.

1384 Guenther, F. H. (**2016**). *Neural Control of Speech* (MIT Press).

1385 Guenther, F. H., Ghosh, S. S., and Tourville, J. A. (**2006**). "Neural modeling and imaging

1386    of the cortical interactions underlying syllable production," Brain and Language **96**(3),

1387    280–301.

1388 Guenther, F. H., Hampson, M., and Johnson, D. (**1998**). "A theoretical investigation of

1389    reference frames for the planning of speech movements.," Psychological Review **105**(4),

1390    611.

1391 Hickok, G. (**2012**a). "Computational neuroanatomy of speech production," Nature Reviews

1392    Neuroscience **13**(2), 135–145.

1393 Hickok, G. (**2012**b). "The cortical organization of speech processing: Feedback control and

1394    predictive coding the context of a dual-stream model," Journal of Communication Disor-

1395    ders **45**(6), 393–402.

1396 Hickok, G. (**2014**). "The architecture of speech production and the role of the phoneme in

1397    speech processing," Language, Cognition and Neuroscience **29**(1), 2–20.

1398 Hogan, N. (**1984**). "An organizing principle for a class of voluntary movements," Journal of

1399    Neuroscience **4**(11), 2745–2754.

1400 Houde, J. F., and Chang, E. F. (**2015**). "The cortical computations underlying feedback

1401    control in vocal production," Current Opinion in Neurobiology **33**, 174–181.

1402 Houde, J. F., and Nagarajan, S. S. (**2011**). "Speech production as state feedback control,"

1403    Frontiers in Human Neuroscience **5**, 82.

1404 Houde, J. F.., Niziolek, C., Kort, N., Agnew, Z.,and Nagarajan, S. S. (**2014**). "Simulating a

1405    state feedback model of speaking," In 10th International Seminar on Speech Production,

1406    202–205.

1407  Iskarous, K., Goldstein, L., Whalen, D. H., Tiede, M., and Rubin, P. (**2003**). "Casy: The

1408  Haskins Configurable Articulatory Synthesizer," in *International Congress of Phonetic*

1409  *Sciences, Barcelona, Spain*, pp. 185–188.

1410  Kalman, R. E. *et al.* (**1960**). "A new approach to linear filtering and prediction problems,"

1411  Journal of Basic Engineering **82**(1), 35–45.

1412  Kröger, B. J., and Birkholz, P. (**2007**). "A gesture-based concept for speech movement con-

1413  trol in articulatory speech synthesis," in *Verbal and Nonverbal Communication Behaviours*

1414  (Springer), pp. 174–189.

1415  Kröger, B. J., Kannampuzha, J., and Neuschaefer-Rube, C. (**2009**). "Towards a neurocom-

1416  putational model of speech production and perception," Speech Communication **51**(9),

1417  793–809.

1418  Laboissiere, R., Ostry, D.J., and Feldman, A.G. (**1996**). "The control of multi-muscle sys-

1419  tems: human jaw and hyoid movements," Biological Cybernetics **74**(4), 373–384.

1420  Lammert, A. C., Ramanarayanan, V., Proctor, M. I., Narayanan, S. *et al.* (**2013**). "Vocal

1421  tract cross-distance estimation from real-time MRI using region-of-interest analysis.," in

1422  *Interspeech*, pp. 959–962.

1423  Lammert, A.C., Shadle, C.H., Narayanan, S.S., and Quatieri, T.F. (**2018**). "Speed-accuracy

1424  tradeoffs in human speech production," PloS one **13**(9), e0202180.

1425  Lisberger, S. (**1988**). "The neural basis for motor learning in the vestibulo-ocular reflex in

1426  monkeys," Trends in Neurosciences **11**(4), 147–152.

1427  Lloyd, J. E., Stavness, I., and Fels, S. (**2012**). "Artisynth: A fast interactive biomechani-

1428  cal modeling toolkit combining multibody and finite element simulation," in *Soft Tissue*

1429   *Biomechanical Modeling for Computer Assisted Surgery* (Springer), pp. 355–394.

1430   Maeda, S. (**1982**). "A digital simulation method of the vocal-tract system," Speech Com-

1431     munication **1**(3-4), 199–229.

1432   Mermelstein, P. (**1973**). "Articulatory model for the study of speech production," Journal

1433     of the Acoustical Society of America **53**(4), 1070–1082.

1434   Miall, R. C., and Wolpert, D. M. (**1996**). "Forward models for physiological motor control,"

1435     Neural Networks **9**(8), 1265–1279.

1436   Mitra, V., Nam, H., Espy-Wilson, C. Y., Saltzman, E., and Goldstein, L. (**2010**). "Retrieving

1437     tract variables from acoustics: a comparison of different machine learning strategies," IEEE

1438     Journal of Selected Topics in Signal Processing **4**(6), 1027–1045.

1439   Munhall, K. G., Löfqvist, A., and Kelso, J. S. (**1994**). "Lip–larynx coordination in speech:

1440     Effects of mechanical perturbations to the lower lip," Journal of the Acoustical Society of

1441     America **95**(6), 3605–3616.

1442   Nam, H. (**2012**). "TADA: TAsk Dynamic Application" $http://www.haskins.yale.edu/

1443     tada_download/index.php$.

1444   Nam, H., Goldstein, L., Browman, C., Rubin, P., Proctor, M., and Saltzman, E. (**2006**).

1445     *TADA (TAsk Dynamics Application) Manual.*

1446   Nazari, M. A., Perrier, P., Chabanas, M., and Payan, Y. (**2011**). "Shaping by stiffening: a

1447     modeling study for lips," Motor Control **15**(1), 141–168.

1448   Nieto-Castanon, A., Guenther, F.H., Perkell, J.S., and Curtin, H.D. (**2005**). "A model-

1449     ing investigation of articulatory variability and acoustic stability during American En-

1450     glish/r/production," Journal of the Acoustical Society of America **117**(5), 3196–3212.

Nieto-Castano, A. (**2017**). "VTCalcs for Matlab" $http://sites.bu.edu/guentherlab/software/vtcalcs-for-matlab/$.

Nieto-Castanon, A. (**2016**). "DIVA Source Code" $http://sites.bu.edu/guentherlab/software/diva-source-code/$.

Niziolek, C.A., and Guenther, F.H. (**2013**). "Vowel category boundaries enhance cortical and behavioral responses to speech feedback alterations," Journal of Neuroscience **33**(29), 12090–12098.

Ostry, D. J., Gribble, P. L., and Gracco, V. L. (**1996**). "Coarticulation of jaw movements in speech production: is context sensitivity in speech kinematics centrally planned?," Journal of Neuroscience **16**(4), 1570–1579.

Parrell, B., Ramanarayanan, V., Nagarajan, S., and Houde, J. (**2018**). "FACTS: A hierarchical task-based control model of speech incorporating sensory feedback," Proc. Interspeech 2018, 1497–1501.

Patri, J.-F., Diard, J., and Perrier, P. (**2015**). "Optimal speech motor control and token-to-token variability: a bayesian modeling approach," Biological Cybernetics **109**(6), 611–626.

Patri, J.F. (**2018**). "Bayesian modeling of speech motor planning: variability, multisensory goals and perceptuo-motor interactions," Ph.D. thesis, Université Grenoble-Alpes.

Payan, Y., and Perrier, P. (**1997**). "Synthesis of VV sequences with a 2D biomechanical tongue model controlled by the Equilibrium Point Hypothesis," Speech Communication **22**(2-3), 185–205.

Perrier, P., Ostry, D., and Laboissière, R. (**1996**). "The equilibrium point hypothesis and its application to speech motor control," Journal of Speech and Hearing Research **39**, 365–78.

1473  Perrier, P., Payan, Y., Zandipour, M., and Perkell, J (**2003**). "Influences of tongue biome-

1474    chanics on speech movements during the production of velar stop consonants: A modeling

1475    study," Journal of the Acoustical Society of America **114**(3), 1582–1599.

1476  Perrier, P., Ma, L., and Payan, Y. (**2005**). "Modeling the production of VCV sequences via

1477    the inversion of a biomechanical model of the tongue,".

1478  Perrier, P., and Fuchs, S. (**2008**). "Speed–curvature relations in speech production challenge

1479    the 1/3 power law," Journal of Neurophysiology **100**(3), 1171–1183.

1480  Ramanarayanan, V., Goldstein, L., and Narayanan, S. S. (**2014**). "Motor control primitives

1481    arising from a learned dynamical systems model of speech articulation," in *Fifteenth Annual*

1482    *Conference of the International Speech Communication Association.*

1483  Ramanarayanan, V., Parrell, B., Goldstein, L., Nagarajan, S., and Houde, J. (**2016**). "A new

1484    model of speech motor control based on task dynamics and state feedback," Interspeech

1485    2016 3564–3568.

1486  Robinson, D. A. (**1981**). "Models of the mechanics of eye movements," Models of Oculomo-

1487    tor Behavior and Control 21–41.

1488  Rubin, P., Baer, T., and Mermelstein, P. (**1981**). "An articulatory synthesizer for perceptual

1489    research," Journal of the Acoustical Society of America **70**(2), 321–328.

1490  Rubin, P., Saltzman, E., Goldstein, L., McGowan, R., Tiede, M., and Browman, C. (**1996**).

1491    "Casy and extensions to the task-dynamic model," in *Speech Production Seminar*, pp.

1492    125–128.

1493  Saltzman, E. (**1986**). "Task dynamic coordination of the speech articulators: a preliminary

1494    model," US Department of Commerce Report.

Saltzman, E., and Kelso, J. (**1987**). "Skilled actions: a task-dynamic approach.," Psychological Review **94**(1), 84.

Saltzman, E., Löfqvist, A., Kinsella-Shaw, J., Kay, B., and Rubin, P. (**1995**). "On the dynamics of temporal patterning in speech," Studies in speech production: A Festschrift for Katherine Safford Harris. Woodbury, New York: American Institute of Physics 469–487.

Saltzman, E., and Munhall, K. (**1989**). "A dynamical approach to gestural patterning in speech production," Ecological psychology **1**(4), 333–382.

Saltzman, E., Nam, H., Krivokapic, J., and Goldstein, L. (**2008**). "A task-dynamic toolkit for modeling the effects of prosodic structure on articulation," in *Proceedings of the 4th international conference on speech prosody. Brazil: Campinas*, pp. 175–184.

Sanguineti, V., Laboissière, R., and Ostry, D. (**1990**). "A dynamic biomechanical model for neural control of speech production," Journal of the Acoustical Society of America **103**(3), 1615–27.

Sanguineti, V., Laboissière, R., and Ostry, D. J. (**1998**). "A dynamic biomechanical model for neural control of speech production," Journal of the Acoustical Society of America **103**(3), 1615–1627.

Schaal, S., and Schweighofer, N. (**2005**). "Computational motor control in humans and robots," Current Opinion in Neurobiology **15**(6), 675–682.

Schroeter, J. (**2006**). "Text to speech synthesis. new paradigms and advances [book review]," IEEE Signal Processing Magazine **23**(6), 72–73.

Sciavicco, L., and Siciliano, B. (**2012**). "Modelling and Control of Robot Manipulators," (Springer Science & Business Media).

1517 Scott, S. H. (**2004**). "Optimal feedback control and the neural basis of volitional motor

1518 control," Nature Reviews Neuroscience **5**(7), 532.

1519 Shadle, C. H., and Damper, R. I. (**2002**). "Prospects for articulatory synthesis: A position

1520 paper," .

1521 Shadmehr, R., and Krakauer, J. W. (**2008**). "A computational neuroanatomy for motor

1522 control," Experimental Brain Research **185**(3), 359–381.

1523 Shibata, T., and Schaal, S. (**2001**). "Biomimetic gaze stabilization based on feedback-error-

1524 learning with nonparametric regression networks," Neural Networks **14**(2), 201–216.

1525 Shiller, D. M., Laboissière, R., and Ostry, D. J. (**2002**). "Relationship between jaw stiffness

1526 and kinematic variability in speech," Journal of Neurophysiology **88**(5), 2329–2340.

1527 Smith, O. J. (**1959**). "A controller to overcome dead time," ISA Journal **6**(2), 28–33.

1528 Takakusaki, K. (**2017**). "Functional neuroanatomy for posture and gait control," Journal of

1529 Movement Disorders **10**(1), 1.

1530 Tian, X., and Poeppel, D. (**2010**). "Mental imagery of speech and movement implicates the

1531 dynamics of internal forward models," Frontiers in Psychology **1**.

1532 Tilsen, S. (**2016**). "Selection and coordination: The articulatory basis for the emergence of

1533 phonological structure," Journal of Phonetics **55**, 53–77.

1534 Tilsen, S. (**2017**). "Exertive modulation of speech and articulatory phasing," Journal of

1535 Phonetics .

1536 Todorov, E. (**2004**). "Optimality principles in sensorimotor control," Nature Neuroscience

1537 **7**(9), 907.

1538  Todorov, E., and Jordan, M. I. (**1998**). "Smoothness maximization along a predefined path

1539  accurately predicts the speed profiles of complex arm movements," Journal of Neurophys-

1540  iology **80**(2), 696–714.

1541  Todorov, E., and Jordan, M. I. (**2002**). "Optimal feedback control as a theory of motor

1542  coordination," Nature Neuroscience **5**(11), 1226.

1543  Tourville, J.A., Reilly, K.J., and Guenther, F.H. (**2008**). "Neural mechanisms underlying

1544  auditory feedback control of speech," Neuroimag **39**(3), 1429–1443.

1545  Tourville, J. A., and Guenther, F. H. (**2011**). "The DIVA model: a neural theory of speech

1546  acquisition and production," Language and Cognitive Processes **26**(7), 952–981.

1547  Villacorta, V.M., Perkell, J.S., and Guenther, F.H. (**2007**). "Sensorimotor adaptation to

1548  feedback perturbations of vowel acoustics and its relation to perception," Journal of the

1549  Acoustical Society of America **122**(4), 2306–2319.

1550  Wan, E. A., and Van Der Merwe, R. (**2001**). "The unscented kalman filter," in *Kalman

1551  Filtering and Neural Networks*, edited by S. Haykin (Wiley, New York).

1552  Wiener, N. (**1948**). *Cybernetics: Control and communication in the animal and the machine*

1553  (Wiley New York).

1554  Wolpert, D. M., Miall, R. C., and Kawato, M. (**1998**). "Internal models in the cerebellum,"

1555  Trends in Cognitive Sciences **2**(9), 338–347.

1556  Zandipour, M., Guenther, F., Perkell, J., Perrier, P., Payan, Y., and Badin, P. (**2004**).

1557  "Vowel-vowel planning in acoustic and muscle space," Proceedings of "From Sound to

1558  Sense: 50+ years of discoveries in speech communication", C103–C108.

1559  Zsiga, E. (**1994**). "An Acoustic and Electropalatographic Study of Lexical and Post-lexical

1560   Palatalization in American English," Haskins Laboratories Status Report on Speech Re-

1561   search **SR-117/118**, 67–79.