

An informal logic of feedback-based temporal control

1 Sam Tilsen^{1*}

2 ¹Cornell Phonetics Lab, Department of Linguistics, Cornell University, Ithaca, NY, USA

3 * Correspondence:

4 tilsen@cornell.edu

5 **Keywords: articulation, articulatory timing, speech rate, motor control, feedback, dynamical**
6 **systems, phonology, prosody**

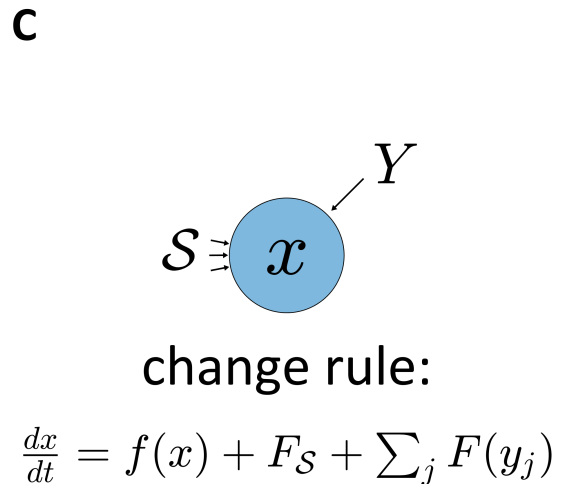
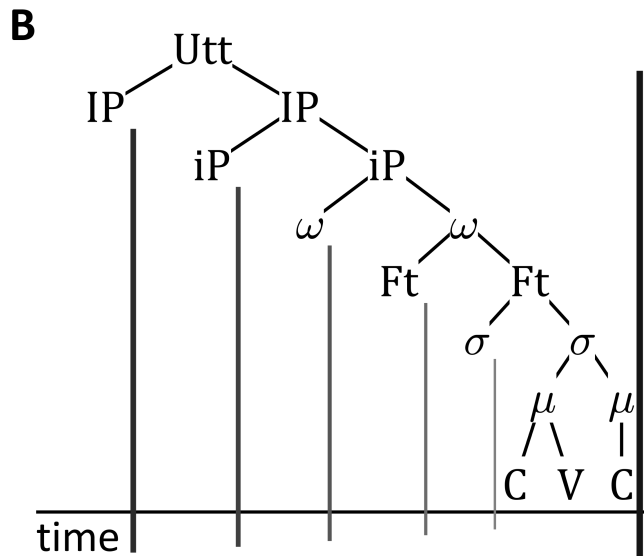
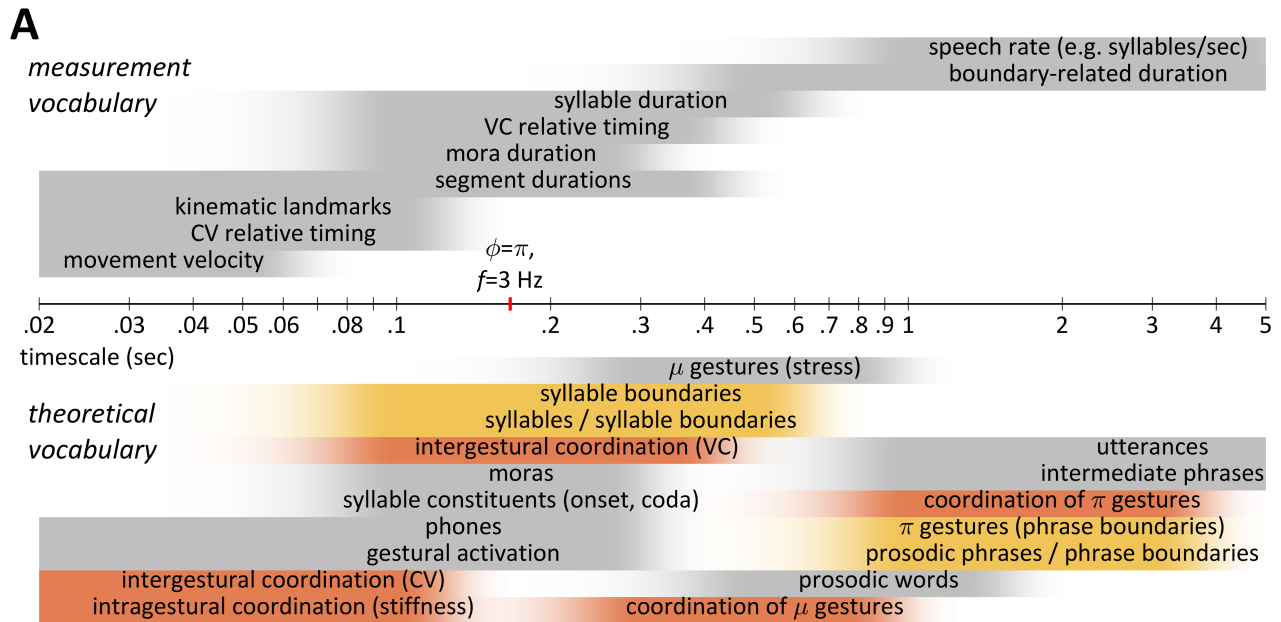
7 **Abstract**

8 A conceptual framework and mathematical model of the control of articulatory timing is presented,
9 in which feedback systems play a fundamental role. The model applies both to relatively small
10 timescales, such as within syllables, and to relatively large timescales, such as multi-phrase
11 utterances. A crucial distinction is drawn between internal/predictive feedback and external/sensory
12 feedback. It is argued that speakers modulate attention to feedback to speed up and slow down
13 speech. A number of theoretical implications of the framework are discussed, including consequences
14 for the understanding of syllable structure and prosodic phrase organization.

15 **1 Introduction**

16 Perhaps you have been in a situation in which it was necessary to *shush* someone. For example,
17 imagine you are reading in a library, when a rude person nearby begins talking on their cell phone.
18 You glare at them and say "shhh", transcribed phonetically as [ʃ:~]. What determines the duration of
19 this sound? Consider now a different situation: in a coffee shop you are ranting to your friend about
20 the library incident, and your friend tells you to slow down because you are talking too fast. You take
21 a deep breath and proceed more slowly. How do you implement this slowing? The focus of this paper
22 is on how variation in the temporal properties of event durations (your "shhh") and variation in event
23 rate (your rapid coffee shop rant) relate to one another. More specifically, what is the mechanistic
24 connection between control of event timing on short timescales and control of speech rate on longer
25 timescales? It is argued that the answer to this question involves a notion of feedback, and that the
26 same feedback mechanisms are involved on both timescales. In other words, control of event timing
27 involves feedback, and control of rate is reducible to control of timing.

28 Temporal patterns in speech are challenging to characterize because they exist across a wide range
29 of analysis scales. Figure 1A shows rough approximations of timescales associated with various
30 measurements and theoretical vocabularies. Even over the modest range of 20 ms to 5,000 ms
31 (shown in a logarithmic axis), there is a diversity of ways to associate time intervals with theoretical
32 constructs. Furthermore, there are certain terms—"coordination", "boundaries"—which reappear
33 across scales, and problematically necessitate different interpretations.



34

35 Figure 1. (A) Comparison of timescales associated with various measurements and theoretical
 36 constructs used to conceptualize temporal patterns. Time axis is logarithmic. Shaded intervals
 37 approximately represent ranges of time in which terminology applied. (B) Hierarchical conception of
 38 prosodic structure and implicit projection of units to boundaries in a temporal coordinate. (C) Generic
 39 system schema, where change in the state variable x is a function of x itself and of forces from the
 40 surroundings S and from other systems Y .

41 It is rarely the case that models of small scale phenomena, such as articulatory timing within syllables,
 42 are integrated with models of larger scale phenomena, such as boundary-related slowing. One
 43 noteworthy exception is the π -gesture model (1), which modulates the rate of a global dynamical
 44 clock in the vicinity of phrase boundaries, thereby slowing the timecourse of gestural activation.
 45 Another example is the multiscale model of (2), where oscillator-based control of gestural timing is
 46 limited to syllable-sized sets of gestures that are competitively selected with a feedback-based

47 mechanism. This early combination of oscillator- and feedback-based control led to the development
 48 of Selection-Coordination theory (3,4), an extension of the Articulatory Phonology framework that
 49 uses feedback control to account for a variety of cross-linguistic and developmental patterns. A recent
 50 proposal in this context is that speech rate is controlled by adjusting the relative contributions of
 51 external (sensory) feedback and internal (predictive) feedback (5). One of the aims of this paper is to
 52 elaborate on this idea, advancing that generalization that temporal control in speech is largely (but
 53 not exclusively) feedback-based.

54 A broader aim is to argue for a worldview in which speech patterns are understood to result from
 55 interactions of dynamical systems. The "informal logic" developed here advocates for new way of
 56 thinking about patterns in speech. It is relevant both for the study of speech motor control,
 57 specifically in relation to feedback and control of timing, and for theories of phonological
 58 representation, sound patterns, and change. The informal logic challenges the prevailing ontologies
 59 of many phonological theories by rejecting the notion that speech is cognitively represented as a
 60 structure of hierarchically connected objects, as in Figure 1B. It also rejects the notion that such units
 61 project "boundaries" onto the temporal dimension of the acoustic signal. Most importantly, the logic
 62 holds that speakers never control event durations directly: rather, durational control is accomplished
 63 via a class of systems which *indirectly* represent time. They do this by integrating the forces they
 64 experience from other systems, or from a surroundings.

65 The systems-oriented approach can provide a more coherent understanding of temporal phenomena
 66 across scales. Its logic is qualified as "informal" because, unlike a formal logic, it does not rely heavily
 67 on symbolic forms; rather, the schemas presented below are iconic and indexical, designed to help
 68 users rapidly interpret complex patterns of system interactions. At the same time, the schemas can
 69 be readily mapped to an explicit mathematical model. All model equations and simulation details are
 70 described in Supplementary Material, and all code used to conduct simulations and generate figures
 71 has been made available in a repository, here: <https://github.com/tilsen/TiR-model.git>. Finally,
 72 although its implications are fairly general, the scope of this paper is narrowly focused on describing
 73 a logic of *temporal* control. Issues related to "spatial" dimensions of feedback or to feedback
 74 modalities are set aside for future extensions of the model.

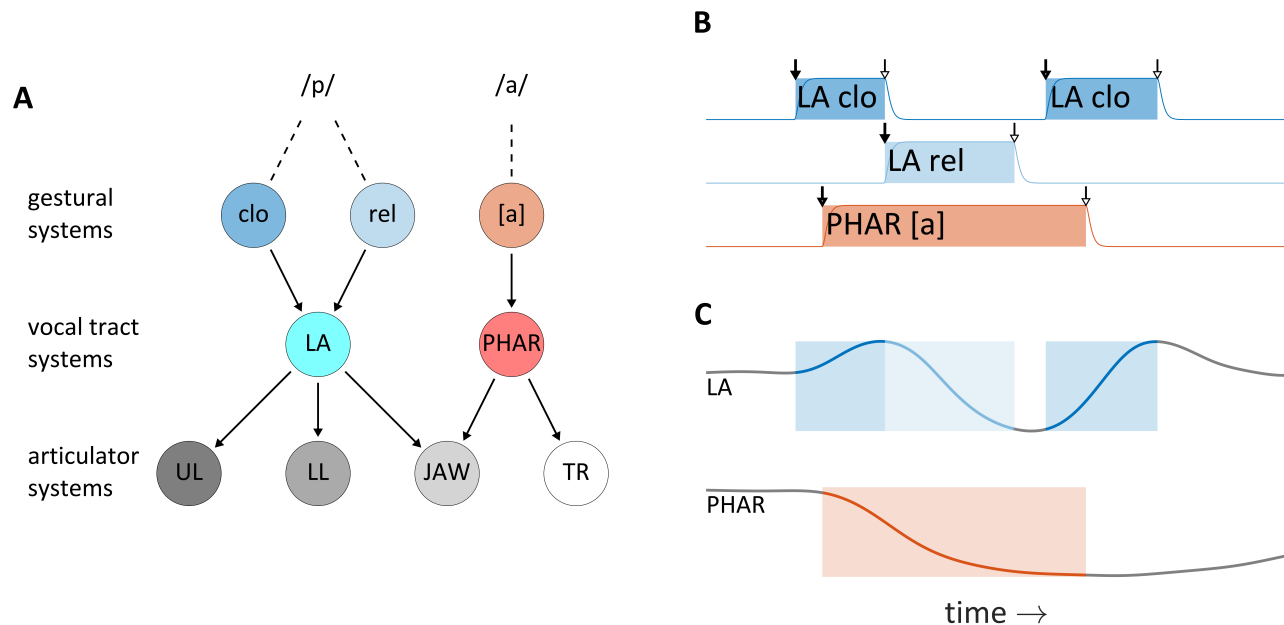
75 **2 Background**

76 In what follows, the objects of our analyses are systems and their relations are interaction forces.
 77 Systems are abstract entities which have time-varying internal states. Our analytical task is to
 78 formulate change rules to describe how the system states evolve over the course of an utterance, as
 79 shown generically in Figure 1C. This setup provides a frame in which to analyze and interpret the
 80 causes of empirical patterns in speech. Moreover, to draw generalizations about systems and their
 81 interactions we must classify them. To accomplish this in the following sections we define terms
 82 below such as *internal*, *external*, *feedback*, and *sensory*. These terms are necessarily relative and
 83 therefore potentially ambiguous out of context, thus the reader should pay careful attention to these
 84 definitions to avoid confusion.

85 **2.1 Gestural systems and control of gestural activation**

86 Before addressing the role of feedback, we describe the understanding of articulatory control
 87 adopted here, which originates from Task Dynamics (6,7). In Task Dynamics (TD), changes in the

88 physical outputs of speech—vocal tract shape and distributions of acoustic energy—are indirectly
 89 caused by systems called *articulatory gestures*. Figure 2A schematizes the organization of system
 90 interactions in the TD model: gestural systems exert driving forces on vocal tract systems, which in
 91 turn exert forces on articulator systems. (As an aside, note that the framework attributes no
 92 ontological status to phones or phonemes—these are merely "practical tools" (8) or inventions of
 93 scientific cultures (9,10)). Gestural system states are defined in normalized activation coordinates
 94 which range from zero to one, and gestures are understood to abruptly become active and
 95 subsequently deactivate, as in Figure 2B. When their activation is non-zero, gestures exert forces on
 96 vocal tract systems, which can lead to movement, as shown in Figure 2C for timeseries of lip aperture
 97 (LA) and pharyngeal constriction (PHAR).



98

99 Figure 2. System organization and interactions in the Task Dynamics model. (A) Organization of
 100 system interactions. (B) Gestural activation intervals for the CVC syllable *pop*. (C) Vocal tract geometry
 101 changes resulting from the actions of gestural systems on vocal tract systems. Lip aperture (LA) and
 102 pharyngeal constriction (PHAR) timeseries are shown.

103 In both a theoretical and technical sense, gestures should be understood as *systems*—entities which
 104 have internal states and which experience and exert forces. Accordingly, gestures are not
 105 movements, nor are they periods of time in which movements occur. To reinforce this point we often
 106 refer to them (redundantly) as *gestural systems*. The distinction is important because it is common
 107 to refer to movements of vocal organs as "gestures"—but this can cause confusion. Similarly, the
 108 periods in which gestural systems obtain states of high activation (shaded intervals in Figure 2B) are
 109 sometimes called "gestures"—these periods are better described as *gestural activation intervals*. The
 110 point here is simply that metonymic extensions of "gesture" to refer to physical movements or
 111 activation intervals should not be conflated with the systems themselves. Furthermore, the vocal
 112 tract and articulator system states of the TD model are nervous system-internal representations of
 113 the physical geometry of the vocal tract/effectors. The actual geometry of the vocal tract is not
 114 modelled explicitly in TD and can in principle diverge from these internal representations.

115 The TD framework is particularly valuable because it clarifies the questions that must be addressed
 116 in order to understand temporal patterns in speech. There are two questions of paramount
 117 importance regarding temporal control: (i) What causes inactive gestural systems to become active?
 118 and (ii) What causes active gestural systems to become inactive? These questions are correspond to
 119 the arrows marking initiations and terminations of the gestural activation in Figure 2B.

120 (i) *What causes the gestures to become active?* In answering this question, we temporarily adopt the
 121 perspective that the entire set of gestures is a "system". In that case, one possible answer is that
 122 there are some *external* systems which exert forces on the gestures. By "external" we mean systems
 123 which are "outside" of the set of gestures, and we refer to such systems as *extra-gestural*. Another
 124 possibility is that the gestural systems experience forces from each other, in which case the activating
 125 forces come from "inside of the system" or are *internal* to the system of gestures, i.e. *inter-gestural*.
 126 Note that the first gesture to become active must necessarily be activated by an extra-gestural
 127 system, because there is presumably no way for a gestural system to spontaneously "activate itself"
 128 or to be activated by inactive gestural systems.

129 (ii) *What causes the gestures to cease to be active?* The extra-gestural and inter-gestural forces
 130 described above are both plausible sources of deactivation. A third possibility, unavailable in the case
 131 of activating forces, is that deactivation is caused by actions of individual gestural systems on
 132 themselves, i.e. *intra-gesturally*. We elaborate below on how this differs from inter-gestural control.

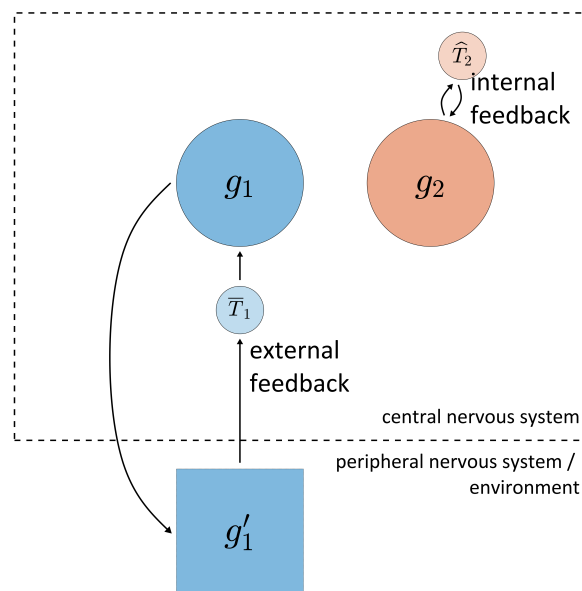
133 The Task Dynamics model of speech production developed by Saltzman and Munhall (7) did not
 134 resolve which of the various sources of initiating and terminating forces are utilized. Saltzman and
 135 Munhall heuristically hand-specified activation intervals to fit empirical data, but they proposed that
 136 the model could be extended with the serial network of (11) to dynamically control gestural
 137 activation. In this serial network, the hidden layers responsible for sequencing might be interpreted
 138 as extra-gestural forces. However, many early descriptions of timing in the TD-based theory of
 139 Articulatory Phonology (12,13)—in particular references to "phasing"—imply that initiating forces
 140 are inter-gestural and that terminating forces are intra-gestural, in line with the explicit
 141 interpretations of phasing in (14). In contrast, later descriptions hypothesize that gestures are
 142 activated by a separate system of gestural planning oscillators (15,16), which are extra-gestural.

143 To summarize, the systems-view of gestural control in the Task Dynamics framework provides two
 144 generic options for what causes gestures to become active or cease to be active—extra-gestural
 145 systems or other gestures (inter-gestural forces)—along with a third option of intra-gestural control
 146 as a form of self-deactivation. There is no theoretical consensus on which of these are actually
 147 involved in control of articulatory timing, or in what contexts they may be utilized.

148 2.2 External feedback vs. internal feedback

149 The term *feedback* has a variety of different uses. Here *feedback* refers to information which—in
 150 either a direct or indirect manner—is produced by some particular system, exists outside of that
 151 system, and subsequently plays a role in influencing the state of that same system. Thus feedback is
 152 always defined relative to a particular reference system. Feedback in this sense is a very general
 153 notion, and does not presuppose that "sensory" organs such as the cochlea or muscle stretch
 154 receptors are involved.

155 For a logic of feedback-based temporal control of speech it is crucial to distinguish between *external*
 156 *feedback* and *internal feedback*, as illustrated in Figure 3. The reference system is the central nervous
 157 system (CNS, consisting of cortex, brainstem, and spinal cord). External feedback involves information
 158 that (i) is originally generated within the CNS, (ii) is transformed to information outside of the CNS,
 159 and (iii) is subsequently transformed back to information within the CNS. For example, activation of
 160 the gestural system g_1 causes the production of various forms of information in the environment
 161 (movement of articulators, generation of acoustic energy), which is in turn transduced in the
 162 peripheral nervous system (depolarization of hair cells in the cochlea and sensory muscle fibers) and
 163 subsequently produces information in cortical systems. For current purposes we draw no distinctions
 164 between various sensory modalities, which are lumped together as system g'_1 in the Figure 3. The
 165 information associated with g'_1 can ultimately influence the state of g_1 , and hence meets our
 166 definition of feedback. Notice that Figure 3 includes a system labeled \bar{T}_1 , which uses the external
 167 feedback from g'_1 to act on g_1 .



168

169 Figure 3. Schematic illustration of distinction between internal and external feedback. The dashed
 170 line represents the boundary of the central nervous system. Systems g_1 and g_2 are gestural systems,
 171 g'_1 is system which represents information associated with g_1 outside of the central nervous system,
 172 and T_1 and T_2 are hypothetical systems which use feedback to act on g_1/g_2 .

173 In contrast to external feedback, internal feedback is information which never exists outside of the
 174 CNS. For example, in Figure 3 the gestural system g_2 generates information that system \hat{T}_2 uses to act
 175 on g_2 . Thus the contrast between external and internal feedback is based on whether the relevant
 176 information at some point in time exists "outside of"/"external to" the central nervous system.
 177 External feedback may be also described as "sensory" feedback, but with a caveat: one could very
 178 well also describe internal feedback as "sensory," in that internal feedback systems experience forces
 179 from other systems, and this property can reasonably be considered a form of *sensation*. The point is
 180 simply that the word "sensory" is ambiguous regarding what is being sensed, and so the qualifiers
 181 *internal* and *external* are preferred, with the CNS being the implied reference system. Internal
 182 feedback can also be described as "predictive", but we should be cautious because this term strongly
 183 evokes an agentic interpretation of systems.

184 The distinction between external and internal feedback is only partly orthogonal to distinction
 185 between extra-gestural, inter-gestural, and intra-gestural control. The full system of gestures is by
 186 definition within the CNS; hence feedback associated with inter-gestural and intra-gestural control is
 187 by definition internal feedback. In contrast, extra-gestural control may involve either external
 188 feedback (e.g. auditory or proprioceptive information) or internal feedback from CNS-internal
 189 systems. This can be confusing because "extra"-gestural control does not entail external feedback—
 190 hence the necessity to keep tabs on the system boundaries to which our vocabulary implicitly refers.
 191 When describing feedback, the reference system is the CNS. When describing control of gestural
 192 activation, the reference system is either the full system of gestures (for extra-gestural control) or
 193 individual gestural systems (for inter- vs. intra-gestural control).

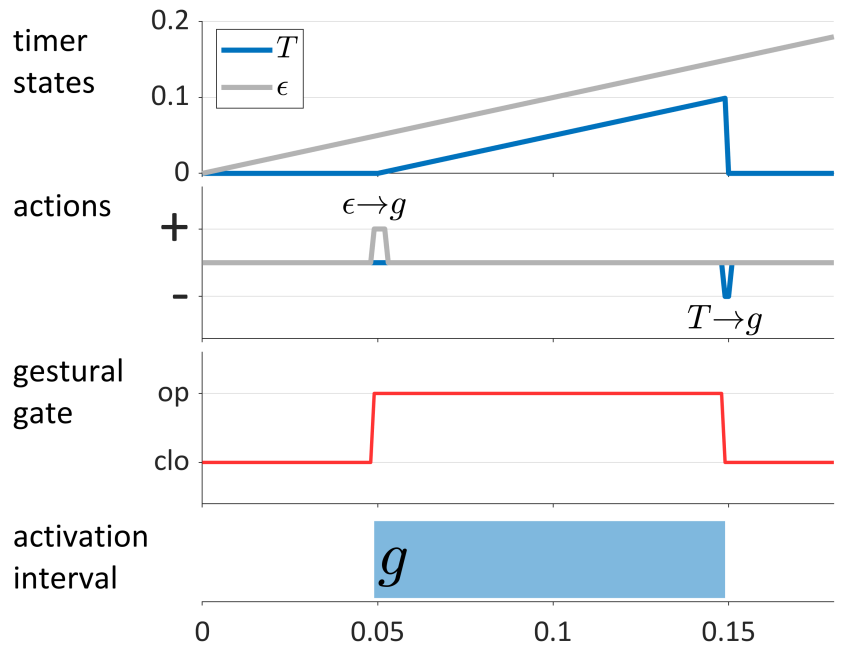
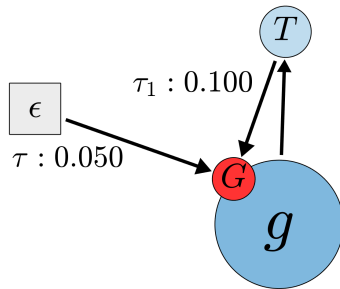
194 The Task Dynamic model incorporates no feedback of any form for gestural systems. Nonetheless,
 195 Saltzman and Munhall cited the necessity of eventually incorporating sensory feedback, stating:
 196 "without feedback connections that directly or indirectly link the articulators to the intergestural
 197 level, a mechanical perturbation to a limb or speech articulatory could not alter the timing structure
 198 of a given movement sequence" (8: p. 360). Note that here Saltzman and Munhall expressed a
 199 concern with the *temporal* effects of perturbation rather than *spatial* effects—in this paper we are
 200 also focused on timing but recognize that a complete picture should incorporate a fully embodied
 201 and sensorially differentiated model of the articulatory and acoustic dimensions of feedback.

202 **2.3 Time-representing systems and timing control**

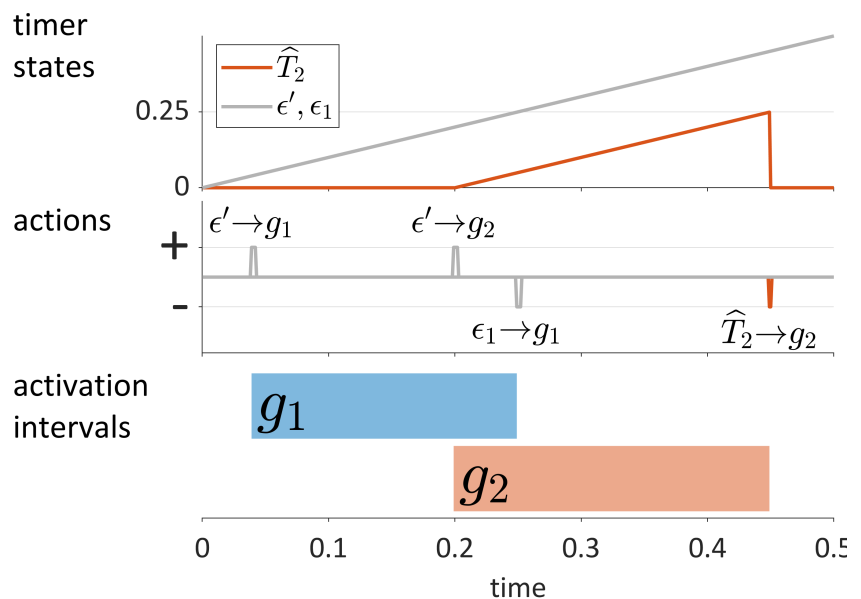
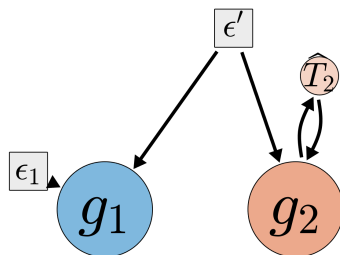
203 To augment our classification of the ways in which gestural systems may be activated or deactivated,
 204 we need to think about how time may be "measured", "estimated", or "represented" by the nervous
 205 system. Researchers have adopted various ways of talking about different types of systems that serve
 206 this function (14,17)—timers, clocks, timekeepers, virtual cycles, etc., with the discussion of (17)
 207 being particularly informative. For current purposes, we describe such systems as "time-
 208 representers" (TiRs) and develop a multidimensional classification. Despite this name, we emphasize
 209 that temporal representations are *always indirect*: the states of the time-representer (TiR) systems
 210 are never defined in units of time.

211 Before classifying TiRs, we make a couple points regarding their interactions with gestures. First, each
 212 gestural system is associated with a gating system, labeled "G" in Figure 4A. The gating system states
 213 are treated as binary: gates are either open or closed. When a gestural gate is open, the activation
 214 state of the associated gestural system transitions rapidly toward its normalized maximum activation
 215 of 1. Conversely, when the gate is closed, the gestural system transitions rapidly toward its minimum
 216 value. For current purposes, transitions in gestural activation states occur in a single time step, as in
 217 (7). Nothing hinges on this simplified implementation and the model can be readily extended to allow
 218 for activation ramping or nonlinearities to better fits of empirical tract variable velocity profiles (18).

A



B



219

220 Figure 4. (A) Model of interactions between gestures and TiRs, with depiction of the gestural gating
 221 system G that TiRs act upon. Panels on the right show timer states, timer actions on gestures, gestural
 222 gating system states, and gestural activation interval. (B) Distinction between autonomous TiRs (ϵ' ,
 223 ϵ_1) and non-autonomous TiRs (\hat{T}_2).

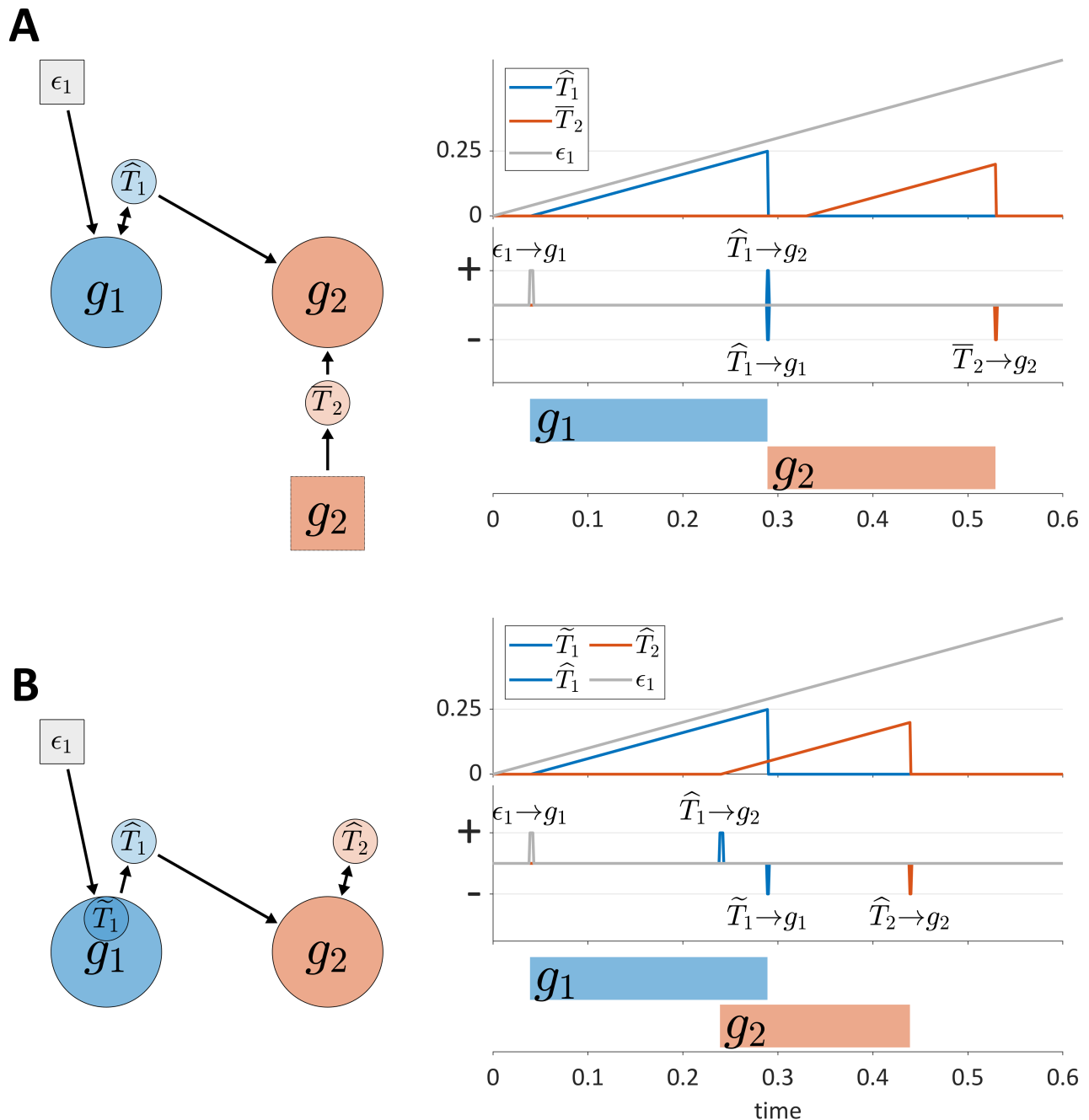
224 Second, TiRs act on gestural gating systems, not directly on gestures, and thus function to
 225 activate/deactivate gestural systems. The actions of TiRs are modeled as brief, pulse-like forces, and
 226 always depend on TiR-internal states: each TiR has threshold parameters (τ) which specify the internal
 227 states (in units of activation) at which the TiR acts on gating systems. The action threshold parameters
 228 are labelled on the arrows of Figure 4A. To reduce visual clutter in model schemas, gating systems
 229 are omitted from subsequent figures.

230 One main dimension of TiR classification involves whether a TiR is autonomous or non-autonomous.
 231 An *autonomous* TiR does not depend on either gestural or sensory system input to maintain an
 232 indirect representation of time. Figure 4B shows two examples of autonomous TiRs. The first is ϵ' ,
 233 which activates gestures g_1 and g_2 . The second is ϵ_1 , which deactivates g_1 . Note that autonomous
 234 TiRs *do* require an external input to begin representing time—they need to be "turned on"/de-
 235 gated—but subsequently their state evolution is determined by a growth rate parameter. This
 236 parameter may vary in response to changes in a hypothesized "surroundings" or contextual factors.

237
 238 In contrast to autonomous TiRs, the states of *non-autonomous* TiRs depend on input from a gestural
 239 or sensory system. Non-autonomous TiRs integrate the forces that they experience from a given
 240 system. An example is \hat{T}_2 in Figure 4B, which receives input from g_2 and deactivates g_2 upon reaching
 241 a threshold state of activation, here $\tau = 0.25$. Non-autonomous TiRs are associated with integration
 242 rate parameters α , which determine how much the forces they experience contribute to changes in
 243 their internal states.

244
 245 The key difference between autonomous TiRs and non-autonomous ones is that the states of the
 246 autonomous TiRs evolve independently from the states of gestures or sensory systems. In the
 247 example of Figure 4B the states of autonomous TiRs ϵ' and ϵ_1 are assumed to be 0 at the beginning
 248 of the simulation and increase linearly in a way that represents the elapsed time. In this example (but
 249 not in general), the growth rates of autonomous TiR states were set to $1/\Delta t$, (where Δt is the
 250 simulation time step); consequently, their activation states exactly correspond to elapsed time. This
 251 is convenient for specifying threshold parameters that determine when TiRs act on other systems.
 252 Similarly, the integration rate parameters of non-autonomous TiRs were parameterized to represent
 253 the time elapsed from the onset of gestural activation. In general, the correspondence between TiR
 254 activation values and elapsed time is neither required nor desirable, and we will see how changes in
 255 TiR growth rates/integration rates are useful for modeling various empirical phenomena.

256
 257 Another dimension of TiR classification involves the sources of input which non-autonomous TiRs
 258 make use of to represent time. Non-autonomous TiRs can be described as *external* or *internal*,
 259 according to whether they integrate external or internal feedback. This distinction is illustrated in
 260 Figure 5A, where the non-autonomous TiR \hat{T}_1 can be described as internal because it integrates
 261 feedback directly from gesture g_1 . In contrast, the non-autonomous TiR \bar{T}_2 is external because it
 262 integrates feedback from sensory systems which encode the actions of g_2 outside of the CNS.



263

264 Figure 5. (A) External vs. internal sources of feedback for non-autonomous TiRs. Panels on the right

265 show timer states, timer actions, and gestural activation intervals. (B) Example of inter-gestural vs.

266 isolated/intra-gestural TiRs.

267

268 Non-autonomous, internal TiRs are further distinguished according to whether they are inter-gestural

269 or intra-gestural (internal to a gesture). Intra-gestural internal TiRs can only act on the particular

270 gestural system that they are associated with, and can integrate forces only from that gesture. Inter-

271 gestural TiRs can act on and experience forces from any gestural system. For example, in Figure 5B,

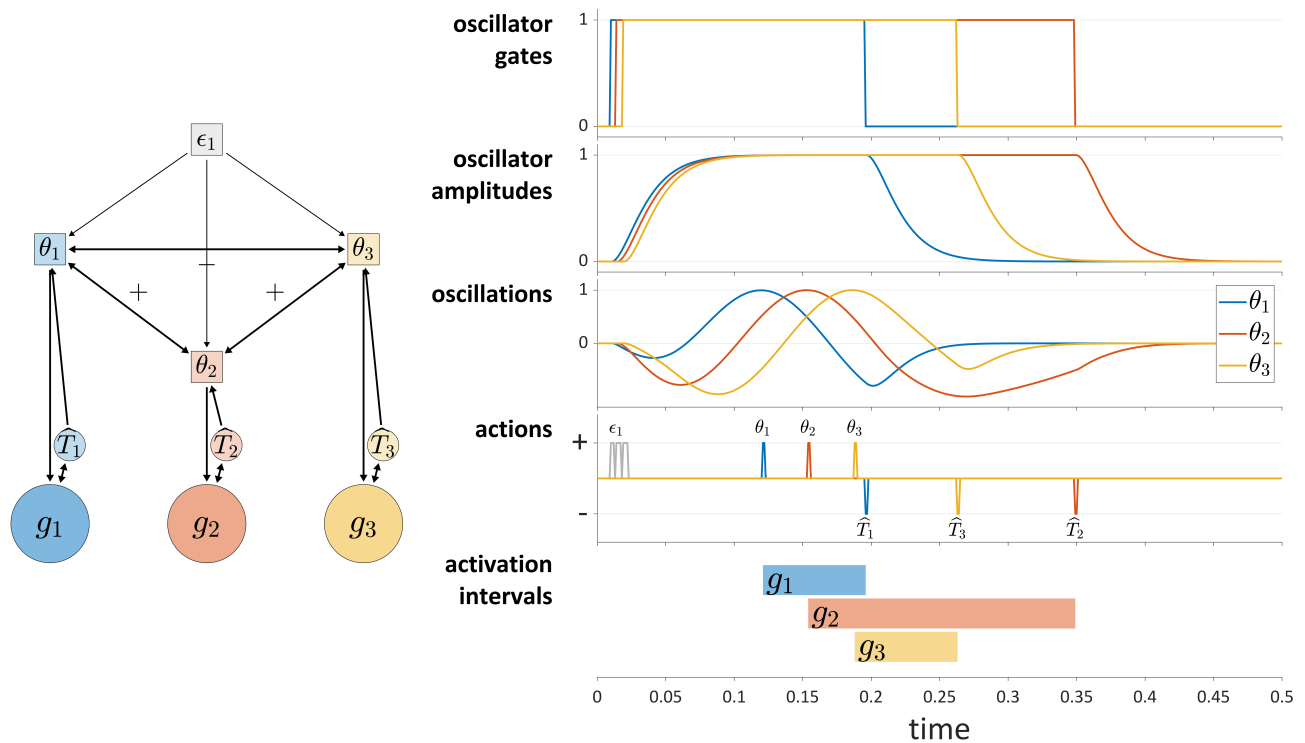
272 the deactivation of g_1 is controlled by an intra-gestural TiR \tilde{T}_1 , but the inter-gestural TiRs \hat{T}_1 and \hat{T}_2 273 activate and deactivate g_2 , respectively. The distinction is useful if we wish to impose the condition

274 that a TiR is isolated from all systems other than a particular gesture.

275

276 The distinction between inter-gestural and intra-gestural TiRs can be viewed in relation to different
277 aspects of the virtual cycles that Tuller and Kelso (14) proposed to govern gestural timing. Tuller and
278 Kelso held that each gesture could be associated with a virtual cycle, which might be described as a
279 "single-shot" oscillation. Different phases of the cycle were hypothesized to correspond to events
280 such as gesture initiation, achievement of maximum velocity, target achievement, and gesture
281 termination. It was suggested in (19) that when a virtual cycle phase of $3\pi/2$ rad (270°) is reached, a
282 gesture is deactivated. In this regard intra-gestural TiRs can implement the functions of virtual cycles:
283 their activation states can be converted to a normalized coordinate that ranges from 0 to 2π , and
284 their growth rates can be adjusted to match the natural frequency of an undamped harmonic
285 oscillator. However, Tuller and Kelso (14) also proposed that intergestural timing might involve
286 specification of the initiation of the virtual cycle of one gesture relative to the virtual cycle of another.
287 Only inter-gestural TiRs can serve this function, because unlike intra-gestural TiRs, they can act on
288 gestural systems that they are not directly associated with. For all of the purposes that follow in this
289 manuscript, intra-gestural TiRs are unnecessary and exclusively use of inter-gestural TiRs.

290 Autonomous TiRs can differ in whether their state evolution is aperiodic or periodic. Periodic (or
291 technically, quasi-periodic) TiRs are used in the coupled oscillators model (15), where each gesture is
292 associated with an oscillatory system called a *gestural planning oscillator*. The planning oscillators are
293 autonomous TiRs because they do not integrate gestural or sensory system states, as can be seen in
294 Figure 6. They are often assumed to have identical frequencies and to be strongly phase-coupled,
295 such that the instantaneous frequencies of the oscillators are accelerated or decelerated as a function
296 of their phase differences. When a given planning oscillator reaches a particular phase, it "triggers"
297 the activation of the corresponding gestural system. The "triggering" in our framework means that
298 the TiR acts upon a gestural system, in the same way that other TiRs act upon gestural systems. The
299 schema in Figure 6 illustrates a system of three periodic TiRs in which θ_1 and θ_3 are repulsively phase
300 coupled to one another while being attractively phase coupled to θ_2 .



301

302 Figure 6. The coupled oscillators model in the TiR framework. Periodic TiRs θ_1 , θ_2 , and θ_3 are phase
 303 coupled as indicated by (+/-) symbols. The oscillator gates, radial amplitudes, and oscillations
 304 (amplitude \times cosine of phase) are shown. Due to the pattern of phase coupling imposed here,
 305 initiation of gestural systems g_1 and g_3 are symmetrically displaced from initiation of g_2 .

306 The phase coupling configuration in Figure 6 generates a pattern of relative phase that—via phase-
 307 dependent actions on gestural systems—leads to a symmetric displacement of initiations of gestures
 308 g_1 and g_3 relative to initiation of g_2 . Statistical tendencies toward symmetric displacement patterns
 309 of this sort are commonly observed in two phonological environments: in simple CV syllables, the
 310 initiations of constriction formation and release are displaced in opposite directions in time from the
 311 initiation of the vocalic gesture (20); in complex onset CCV syllables, the initiations of the first and
 312 second constriction are equally displaced in opposite directions from initiation of the vocalic gesture
 313 (12,21,22).

314 The coupled oscillators model has not been used to govern gestural deactivation. Furthermore, a
 315 gating mechanism is needed to prevent oscillators from re-triggering gestural systems in subsequent
 316 cycles or to prevent them from triggering gestures prematurely. To address this, in the current
 317 implementation each oscillator is described by three state variables: a phase angle, a radial
 318 amplitude, and the derivative of the radial amplitude. Furthermore, each oscillator is associated with
 319 a gating system that controls oscillator amplitude dynamics. These gates are closed by extra-gestural
 320 TiRs, as shown in in Figure 6. Moreover, a condition is imposed such that oscillators can only trigger
 321 gestural activation when their amplitudes are above a threshold value. The "oscillations" panel of
 322 Figure 6 shows a representation of oscillator states that combines phase and amplitude dimensions
 323 (the product of the amplitude and the cosine of phase). Further details are provided in the
 324 Supplementary Material.

325 An important hypothesis is that oscillator frequencies are constrained in a way that aperiodic TiR
 326 growth rates are not. We refer to this as the *frequency constraint hypothesis*. The rationale is that the
 327 oscillator states are believed to represent periodicity in a short-time integration of neuronal
 328 population spike-rates; this periodicity is likely to be band-limited due to intrinsic time-constants of
 329 the relevant neural circuits and neurophysiology. A reasonable candidate band is theta, which ranges
 330 from about 3-8 Hz (23,24), or periods of about 330 to 125 ms. On the basis of these limits, certain
 331 empirical predictions regarding temporal patterns can be derived, which we examine in detail below.

332 Stepping back for a moment, we emphasize that all TiRs can be understood to "represent" time, but
 333 this representation is *not* in units of time. The representation results either (i) from the integration of
 334 gestural/sensory system forces (non-autonomous TiRs), (ii) from a constant growth rate/frequency
 335 (autonomous TiRs) understood to be integration of surroundings forces, or (iii) from a combination
 336 of surroundings forces and forces from other TiRs (as in the case of coupled oscillators). Thus the
 337 systems we hypothesize represent time indirectly and imperfectly, in units of experienced force.

338 The utility of TiRs lies partly their ability to indirectly represent time and partly in their ability to act
 339 on gestures or other systems. Table 1 below summarizes the types of TiRs discussed above. All TiRs
 340 are associated with a parameter vector τ that specifies the activation states at which the TiR acts
 341 upon other systems, along with a parameter vector χ whose sign determines whether actions open
 342 or close gestural gating systems. Autonomous TiRs are associated with a parameter ω which is either
 343 a growth rate (aperiodic TiRs) or angular frequency (periodic TiRs). The latter are also associated with
 344 a phase-coupling matrix. Non-autonomous TiRs are associated with a vector α of integration factors,
 345 which determines how input forces contribute to growth of activation. Additional simulation
 346 parameters and details are described in Supplementary Material.

Table 1. Summary of TiRs

| symbols | autonomous / non-autonomous | feedback source | sub-classes | periodic/ aperiodic | parameters |
|---------------|--------------------------------|--------------------|----------------|------------------------|---------------------------|
| ε | autonomous | | | aperiodic | $\omega, \chi/\tau$ |
| θ | autonomous | | | periodic | $\omega, \chi/\tau, \Phi$ |
| \bar{T} | non-autonomous | CNS-external | extra-gestural | | $\alpha, \chi/\tau$ |
| \hat{T} | non-autonomous | CNS-internal | inter-gestural | | $\alpha, \chi/\tau$ |
| \tilde{T} | non-autonomous | g-internal | inter-gestural | | $\alpha, \chi/\tau$ |

347

348 2.4 Deterministic behavior of TiRs and effects of stochastic forces

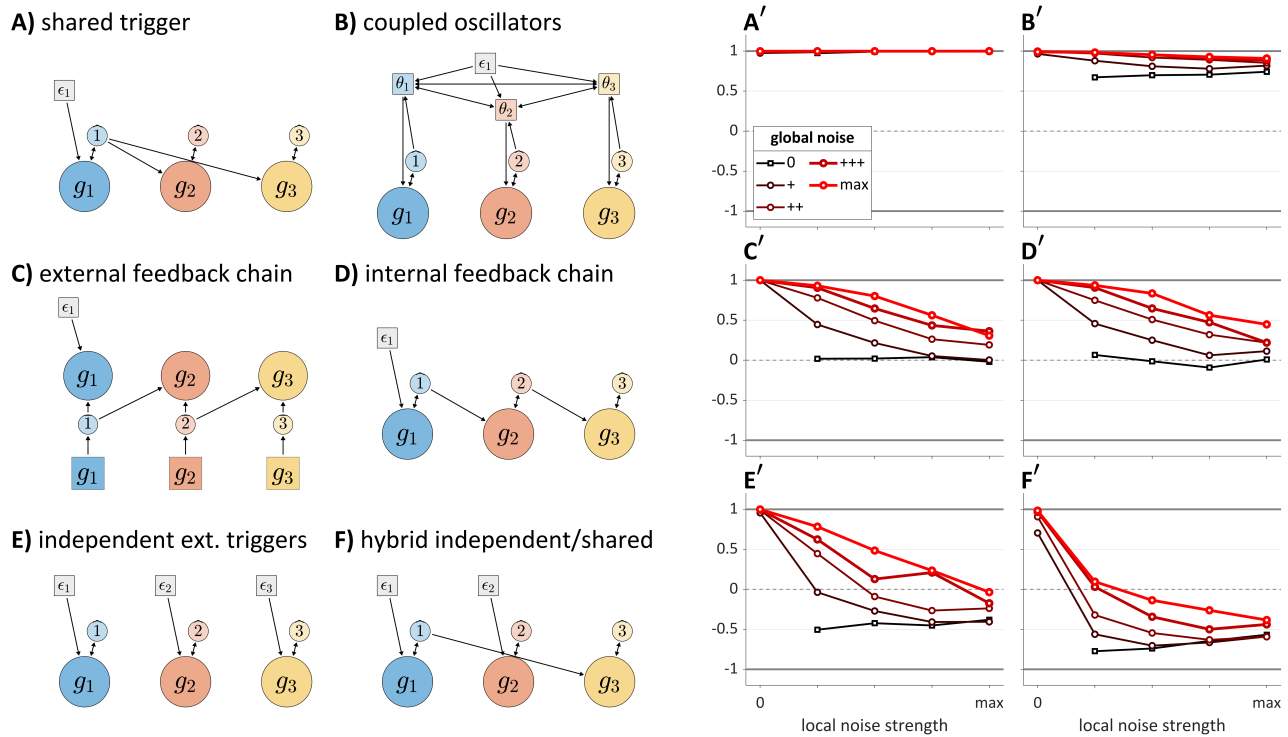
349 Under certain conditions, the time δ when a TiR acts on some other system (δ is relative to when TiR
 350 activation began to grow) is fully determined by its parameters. In the case of autonomous, aperiodic
 351 TiRs, the growth rate ω and action threshold τ determine δ . In two-dimensional ω/τ parameter space,
 352 constant δ are straight lines of positive slope, since increases of ω (which shorten δ) can be offset by
 353 increases of τ (which lengthen δ). Thus either changes in TiR rate ω or in its action threshold τ , or in
 354 some combination of the two, can generate the same change in action timing. This holds for τ and
 355 the integration rate α of non-autonomous TiRs as well, as long as the input force to the TiR is constant.
 356 For coupled oscillator TiRs, δ depends in complicated ways on the initial phases of the systems, the

357 oscillator frequencies, and the strengths of phase coupling forces (putting aside oscillator amplitude
358 dynamics).

359 For even a simple system of three gestures, there is a rich set of possible ways in which temporal
360 control can be organized. How can the organization of control be inferred from empirical
361 observations? What we call "noise" may be quite useful in this regard. An essential characteristic of
362 natural speech is that it is unavoidably stochastic, and as a consequence, no two utterances are
363 identical. We interpret stochastic forces here as variation across utterances in the influence of the
364 surroundings on time-representing systems. Moreover, in modeling noise we distinguish between
365 *global noise*—stochastic variation that affects all TiRs equally—and *local noise*—stochastic variation
366 that differentially affects TiRs. This distinction is important because the relative amplitudes of local
367 and global noise can influence timing patterns.

368 The analysis of stochastic variation below focuses on correlations of successive time intervals
369 between gestural initiations in three-gesture systems. These intervals are referred to as Δ_{12} and Δ_{23} .
370 We examine correlations (henceforth " Δ -correlations") rather than interval durations, because
371 correlations more directly reflect interactions between systems. Five different local and global noise
372 levels were crossed, from 0 to a maximum level (see Supplementary Material: Simulations for further
373 detail). Figure 7 panels A-F show the structures of each model tested, and corresponding panels A'-F'
374 show how Δ -correlation varies as a function of global and local noise levels. Each line corresponds to
375 a fixed level of global noise, and horizontal points represent different local noise levels.

376 The "shared trigger" model (A) shows that if both non-initial gestures are activated by feedback from
377 the initial one, Δ -correlation is trivially equal to 1, regardless of noise. The reason for this is simply
378 that the same TiR (here $\hat{1}$) activates g_2 and g_3 . Note that this trivial correlation occurs for external
379 feedback control as well (not shown). The coupled oscillators model (B) is unique among the systems
380 examined in that it always produces non-trivial positive correlations. The reason for this has to do
381 with phase coupling. Even when oscillator frequencies are heterogenous due to local noise, phase-
382 coupling forces stabilize the oscillators at a common frequency. As long as phase-coupling forces are
383 strong, local noise has relatively small effects on the phase evolution of oscillators. Global frequency
384 noise always leads to positive correlations because it results in simulation-to-simulation variation in
385 frequency that equally influences Δ_{12} and Δ_{23} , causing them to covary positively. However, a more
386 complex analysis of correlation structure in the coupled oscillators model in (20) has shown that when
387 coupling strengths are also subject to noise, the model can generate negative correlations.



388

389 Figure 7. Noise-related correlation patterns for a variety of three-gesture systems. Panels (A-F) show
 390 model schemas and corresponding panels (A'-F') show correlations of intervals between initiation of
 391 gestural systems. Local noise levels increase along the horizontal axes, while global noise levels are
 392 indicated by the lines in each panel. Cases where both global and local noise are zero are excluded.

393 The external and internal feedback "chain models" (C and D) exhibit nearly identical, complex
 394 patterns of correlation that depend on the relative levels of global and local noise. The patterns are
 395 nearly identical because the two models are topologically similar—they are causal chains—differing
 396 only in regard to the temporal delay associated with sensory feedback. When there is no local noise,
 397 these chain models exhibit Δ -correlations of 1, since the global noise has identical effects on Δ_{12} and
 398 Δ_{23} . Conversely, when there is no global noise, Δ -correlation is 0, since local noise has independent
 399 effects on Δ_{12} and Δ_{23} . In between those extremes, the correlation depends on the relative levels of
 400 local and global noise: increasing local relative to global noise leads to decorrelation of the intervals.

401 Unlike the other models, the independent extra-gestural triggers model (E) and hybrid model (F) can
 402 generate substantial negative correlations. In particular, negative correlations arise when g_2 is
 403 influenced by local noise. This occurs because whenever the TiR which activates g_2 does so relatively
 404 early or late, Δ_{12} and Δ_{23} will be influenced in opposite ways. Note that the negative correlations are
 405 stronger when the activation of g_1 and g_3 are caused by the same TiR, as is the case for the hybrid
 406 model (F). At the same time, global noise induces positive Δ -correlation, counteracting the negative
 407 correlating effect of local noise. When we examine speech rate variation below, we will see that the
 408 opposing effects of global and local noise are not specific to "noise" per se: any source of variation
 409 which has similar effects on all TiRs tends to generate positive interval correlations, while the absence
 410 of such variation can lead to zero or negative correlation.

411 **3 A hybrid model of gestural timing and speech rate control**

412 Equipped with a new logic of temporal control, we now develop a hybrid model of gestural timing
413 which is designed to accommodate a wide range of empirical phenomena. The primary requirement
414 of the model is that for each gesture which is hypothesized to drive articulatory movement in an
415 utterance, the model must generate commands to activate and deactivate that gesture.

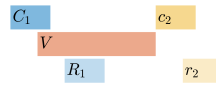
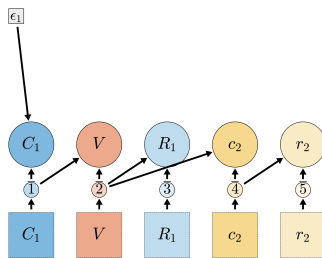
416 **3.1 Model space and hypotheses**

417 For even a single CVC syllable, the set of all logically possible models is very large. Nonetheless, there
418 are a number of empirical and conceptual arguments that we make to greatly restrict this space.
419 Below we consider various ways in which gestural activation might be controlled for a CVC syllable
420 uttered in isolation. Note that we adopt the modern "split-gesture" analysis in which constriction
421 formation and constriction release are driven by separate gestural systems; this analysis has been
422 discussed and empirically motivated in (20,25,26). With that in mind we use the following gestural
423 labeling conventions: C/c and R/r correspond to constriction formation and release gestures,
424 respectively; upper case labels C/R correspond to pre-vocalic gestures (or, gestures associated with
425 syllable onsets); lower case labels c/r correspond to post-vocalic gestures (or, gestures associated
426 with syllable codas); and gestures/gesture pairs are subscripted according to the order in which they
427 are initiated.

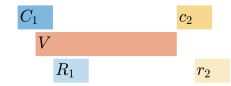
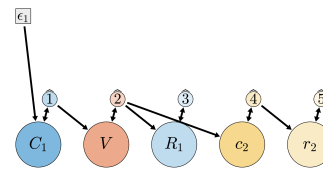
428 The schemas in Figure 8 (A-C) show "extreme" models which—though logically possible—are
429 conceptually and empirically problematic. (A) shows a "maximally sensory" model, where all gestural
430 activation/deactivation is controlled by external feedback systems. This model is problematic because
431 the time delay between efferent motor signals and afferent feedback is too long to be useful for some
432 relative timing patterns, such as the relative timing of consonantal constriction and release in normal
433 speech. (B) shows a "maximally internal" model, where all gestural activation and deactivation is
434 induced by inter-gestural TiRs (keeping in mind that initiation of activation of the first gesture in an
435 utterance is always external). The maximally internal model is problematic because it has no way of
436 allowing for external/sensory feedback to influence timing.

437

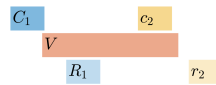
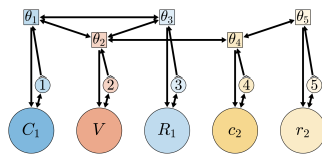
A) maximally sensory



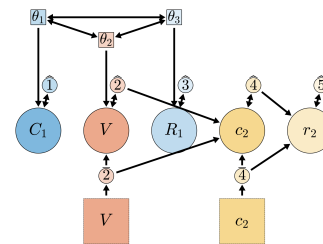
B) maximally internal



C) oscillator triggered



D) hybrid model



438

439 Figure 8. Candidate models of CVC syllables. (A) Maximally sensory model where all activation and
 440 deactivation is controlled by external sensory feedback. (B) Maximally internal model where all
 441 control is governed by internal feedback. (C) Fully oscillator-triggered model where all gestures are
 442 initiated by oscillators. (D) Hybrid model in which pre-vocalic gestural activation is oscillator-governed
 443 while post-vocalic activation is governed by either internal or external feedback.

444 Schema (C) shows an "oscillator triggered" model, where all gestures are activated by coupled
 445 oscillators. Under standard assumptions, this model is problematic because it cannot generate some
 446 empirically observed combinations of pre-vocalic and post-vocalic consonantal timing, as discussed
 447 in (5). The "standard" assumptions are: (i) that all oscillators have (approximately) the same
 448 frequency; (ii) that all oscillators trigger gestural initiation at the same phase of their cycle; and (iii)
 449 that only in-phase and anti-phase coupling are allowed. With these constraints, the model cannot
 450 generate empirically common combinations of pre-vocalic and post-vocalic temporal intervals, where
 451 prevocalic CV intervals are generally in the range of 50-100 ms (20) and post-vocalic VC intervals—
 452 periods of time from V initiation to post-vocalic C initiation—are in the range of 150-400 ms.
 453 Moreover, relaxing any of the three assumptions may be undesirable. Allowing oscillators to have
 454 substantially different frequencies can lead to instability and chaotic dynamics, unless coupling forces
 455 are made very strong. Allowing oscillators to trigger gestures at arbitrary phases is inconsistent with
 456 the neurophysiological interpretation: presumably one particular phase of the cycle represents
 457 maximal population spike rate and should be associated with the strongest triggering force. Allowing
 458 for arbitrary relative phase coupling targets, such as a relative phase equilibrium of $3\pi/2$, may not be
 459 well-motivated from a behavioral or neurophysiological perspective.

460 Although the relatively extreme/monolithic models of Figure 8 (A-C) are individually problematic, the
 461 mechanisms that they employ are practically indispensable for a comprehensive understanding of
 462 timing control. External feedback control is necessary to account for common observation that
 463 segmental durations are lengthened in the presence of feedback perturbations (27–32). Internal

464 feedback is necessary to allow for control under circumstances in which external feedback is not
 465 available, for example during loud cocktail parties, for speakers with complete hearing loss, or during
 466 subvocal rehearsal (internal speech) with no articulatory movement. Finally, oscillator-triggered
 467 control is currently the only known mechanism which adequately explains symmetric displacement
 468 patterns (5,20). Given the utility of these mechanisms it is sensible to adopt a hybrid model which
 469 combines them, as in Figure 8D. The hybrid model of (D) represents the following two hypotheses.

470 *Pre-vocalic coordinative control hypothesis.* Control of the activation of pre-vocalic consonantal
 471 constriction formation (C), release (R), and vocalic initiation (V) is governed by a system of coupled
 472 oscillators.

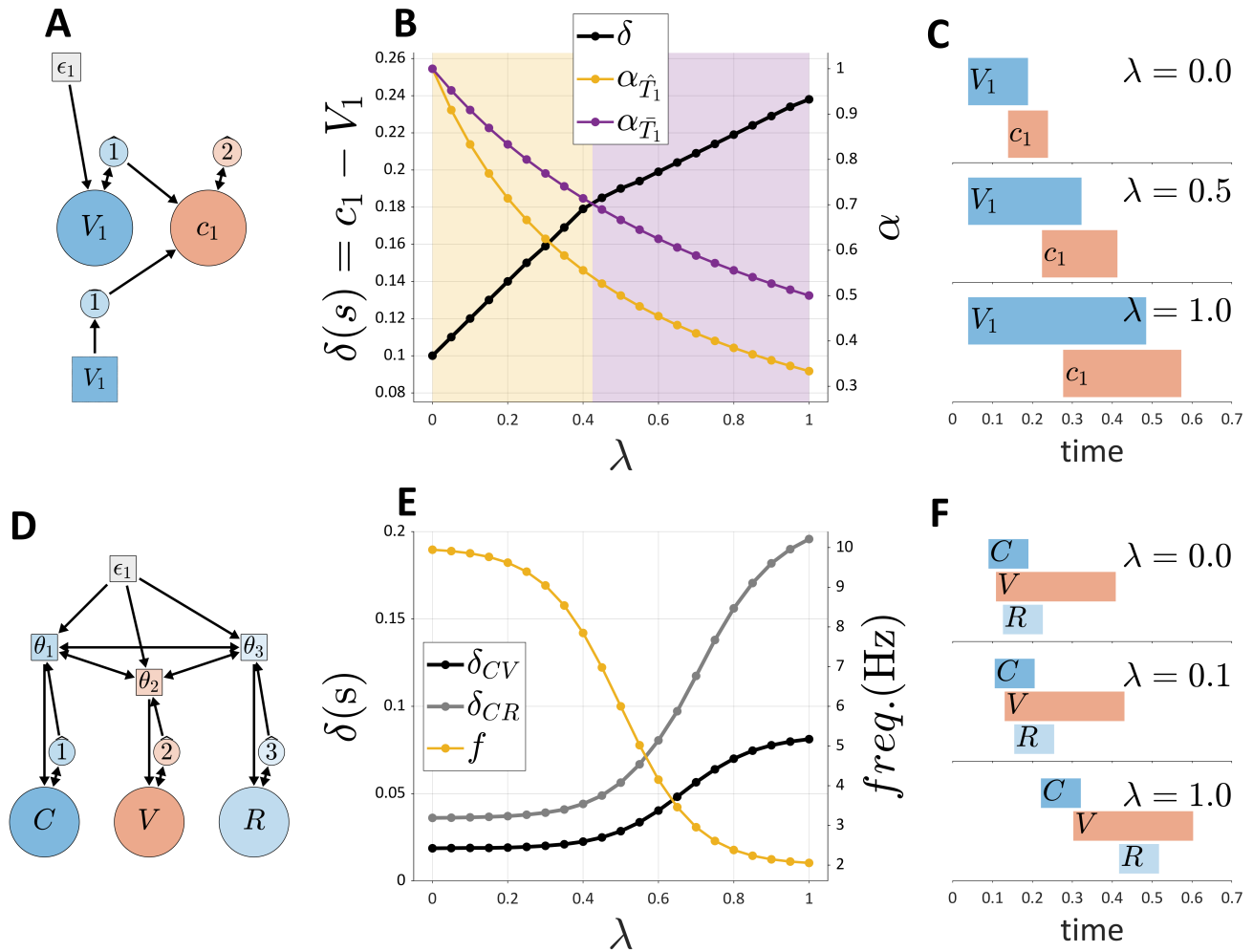
473 *Vocalic/post-vocalic feedback control hypothesis.* The deactivation of vowel gestures and the
 474 activation/deactivation of post-vocalic constriction (c) and release (r) gestures is governed by either
 475 internal or external feedback.

476 Together these hypotheses are referred to as the *hybrid control model*. The specific predictions of the
 477 hypotheses are best considered in light of how interval durations change in response to other sources
 478 of variation, which we examine below.

479 3.2 External influences on parameters

480 The parameters of TiRs are context-dependent: they vary in ways that are conditioned on factors
 481 associated with their surroundings, so-called "external factors". Here we demonstrate two ways in
 482 which external factors may influence timing. An innovation of the model is the idea that these factors
 483 can have differential influences on external vs. internal TiR parameters.

484 Figure 9 (A-C) demonstrates the effects of variation in a hypothetical contextual factor of *self-*
 485 *attention*, or "attention to one's own speech". The figure summarizes simulations of the system
 486 shown in panel (A), where activation of a post-vocalic constriction gesture c_1 is potentially caused by
 487 an internal or external TiR representing feedback from the vocalic gesture V_1 . This is the hypothesized
 488 organization of post-vocalic control in the hybrid model. An external variable λ is posited to represent
 489 self-attention. By hypothesis, the force integration rates of internal and external TiRs are differentially
 490 modulated by λ , such that $\alpha = \alpha' / (1 + \beta\lambda)$, where $\beta_{\text{internal}} < \beta_{\text{external}}$. This reflects the intuition that
 491 when one attends to feedback more closely, feedback-accumulation (i.e. force-integration) rates of
 492 TiR systems are diminished, so that TiRs take longer to act on gestures. This diminishing effect applies
 493 more strongly to internal feedback than external feedback. As a consequence, there is a value of λ
 494 such that as λ is increased, initiation of g_2 switches from being governed by the internal TiR to the
 495 external one. In the example the transition occurs around $\lambda = 0.425$, where a change is visible in the
 496 slope relating the control parameter λ and the interval δ (the time between initiation of V_1 and c_1).
 497 Gestural activation intervals associated with three values of λ are shown in panel (C).



498

499 Figure 9. Simulations of external influences on parameters. (A) Schema for post-vocalic control with
 500 both internal and external TiRs. (B) Dual axis plot showing how δ (left side) and integration rates α
 501 (right side) change with self-attention parameter λ . (C) Gestural activation intervals for several values
 502 of λ . (D) Model schema of pre-vocalic coordinative control. (E) Dual axis plot showing effect of rate
 503 parameter λ on δ -values (left side) and frequencies (right side). (F) Gestural activation intervals for
 504 several values of λ .

505 Panel (B) shows that when TiR parameters are differentially modulated by an external influence,
 506 transitions between internal and external feedback control can occur. In the above example, the
 507 external influence was posited to represent "self-attention" and its state was encoded in the variable
 508 λ ; this variable was then hypothesized to differentially adjust external vs. internal non-autonomous
 509 TiR growth rates. An alternative way in which the same effect might be derived is by allowing the
 510 external variable λ to differentially adjust TiR action-thresholds. Realistically, external variables of this
 511 sort may influence both growth rate and threshold parameters.

512 Another parameter that can respond to external factors is the frequency of the coupled oscillators
 513 which are hypothesized to govern prevocalic gestural initiation. Suppose that the external factor here
 514 is a something novel that we call "pace" and that pace influences oscillator frequencies. However,
 515 because of the frequency constraint hypothesis, we cannot simply allow the oscillator frequencies to
 516 respond linearly to changes in pace. Instead, we impose soft upper and lower frequency bounds by

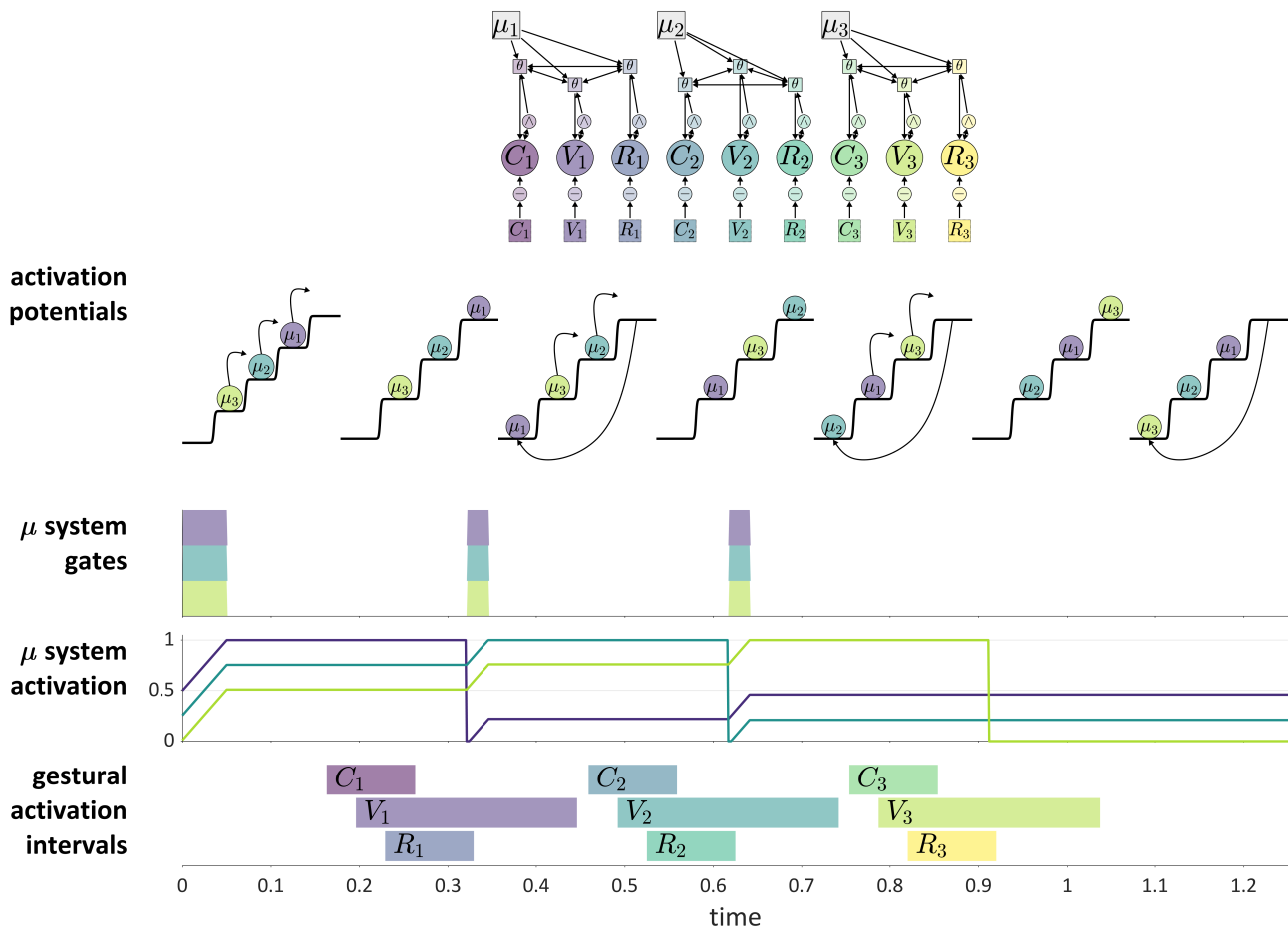
517 attenuating the effect of the pace parameter λ on frequency f . This is accomplished by making the
 518 effective frequency a nonlinear function of λ , as shown in Figure 9E (right side). The consequence of
 519 this limitation on f is that intervals which are governed by coordinative control are predicted to exhibit
 520 nonlinear responses to variation in the external factor: here we can see that the δCV and δCR plateau
 521 at extreme values of λ .

522 In section 3.4 we combine the above effects of self-attention and pace into a general model of the
 523 control of speech rate. But first we introduce another important mechanism, which allows the model
 524 to organize the subsystems of larger utterances.

525 **3.3 Parallel domains of competitive selection**

526 Competitive selection (or competitive queuing) is a dynamical mechanism that, given some number
 527 of actions, iteratively selects one action while preventing the others from being selected. The concept
 528 of competitive selection of actions originates from (33), and many variations of the idea of have been
 529 explored subsequently, both within and outside of speech (2,34–39). One of the key ideas behind the
 530 mechanism is that a serial order of actions is encoded in an initial activation gradient, such that prior
 531 to the performance of an action sequence, the first action in the sequence will have the highest
 532 relative activation gradient, the second action will have the next highest activation, and so on. The
 533 growth of activation is a "competition" of systems to be selected, and selection is achieved by
 534 reaching an activation threshold. Moreover, action selection is mutually exclusive, such that only one
 535 action can be selected at a time.

536 Figure 10 shows how these ideas are understood in the current model. The "actions" which are
 537 competitively selected in this example are three CV syllables, and the selection of these actions is
 538 governed by systems that we refer to as μ -systems. As shown in the model schema, each μ -system
 539 de-gates a system of coupled oscillators, which in turn activate gestures. Each of the μ -systems is
 540 associated with a μ -gating system that—when open—allows the corresponding μ -system activation
 541 to grow. Notice that at time 0 (before the production of the sequence), the pattern of relative
 542 activation of μ -systems corresponds to the order in which they are selected. When μ -system gates
 543 are open, μ -system activations grow until one of the systems reaches the selection threshold. At this
 544 point, all μ -gating systems are closed, which halting growth of μ -system activation. The selected μ -
 545 system is eventually suppressed (its activation is reset to 0) by feedback—specifically by the inter-
 546 gestural TiR associated with the last gesture of the syllable, in this case the vowel gesture. This causes
 547 all μ -systems to be de-gated, allowing their activations to grow until the next most highly active μ -
 548 system reaches the selection threshold. This three-step process—(i) de-gating and competition, (ii)
 549 selection and gating of competitors, and (iii) feedback-induced suppression of the selected system—
 550 iterates until all of the μ -systems have been selected and suppressed. See Supplementary Material:
 551 Model details for further information regarding the implementation.



552

553 Figure 10. Illustration of competitive selection for a sequence of three CV syllables. Top: model
 554 schema. Activation potentials with arrows show transitions between states, and potentials without
 555 arrows shown quasi-steady states. μ -gating system states are shown (shaded intervals are open
 556 states). Bottom: gestural activation intervals.

557 A more abstract depiction of a competitive selection trajectory is included in the activation potentials
 558 of Figure 10. The potentials without arrows are relatively long epochs of time in which μ -systems
 559 exhibit an approximately steady-state pattern of activation. The potentials with arrows correspond
 560 to abrupt intervening transitions in which the relative activation of systems is re-organized by the
 561 competitive selection/suppression mechanism. Along these lines, the dynamics of competitive
 562 selection have been conceptualized in terms of operations on discrete states in (40,41).

563 There are two important questions to consider regarding the application of a competitive selection
 564 mechanism to speech. First, exactly what is responsible for suppressing the currently selected μ -
 565 system? In the example above, which involves only CV-sized sets of gestures, it was the internal TiR
 566 associated with the last gesture of each set. Yet a more general principle is desirable. Second, what
 567 generalizations can we make about the gestural composition of μ -systems? In other words, how is
 568 control of gestural selection organized, such that some gestures are selected together (*co-selected*)
 569 and coordinatively controlled, while others are competitively selected via feedback mechanisms? This
 570 question has been discussed extensively in the context of the Selection-coordination theory of speech
 571 production (3–5), where it is hypothesized that the organization of control follows a typical
 572 developmental progression. In this progression, the use of external sensory feedback for

573 suppression/de-gating is replaced with the use of internal feedback, a process called *internalization*
574 *of control*.

575

576 There are two important points to make about internalization. First, internalization of control is partly
577 optional, resulting in various patterns of cross-linguistic and inter-speaker variation which are
578 detailed in (3) and which we briefly discuss in section 4.1. Second, internalization is flexible within
579 and across utterances, such that various contextual factors (e.g., self-attention) can influence
580 whether external or internal feedback TiRs are responsible for suppressing selected μ -systems.

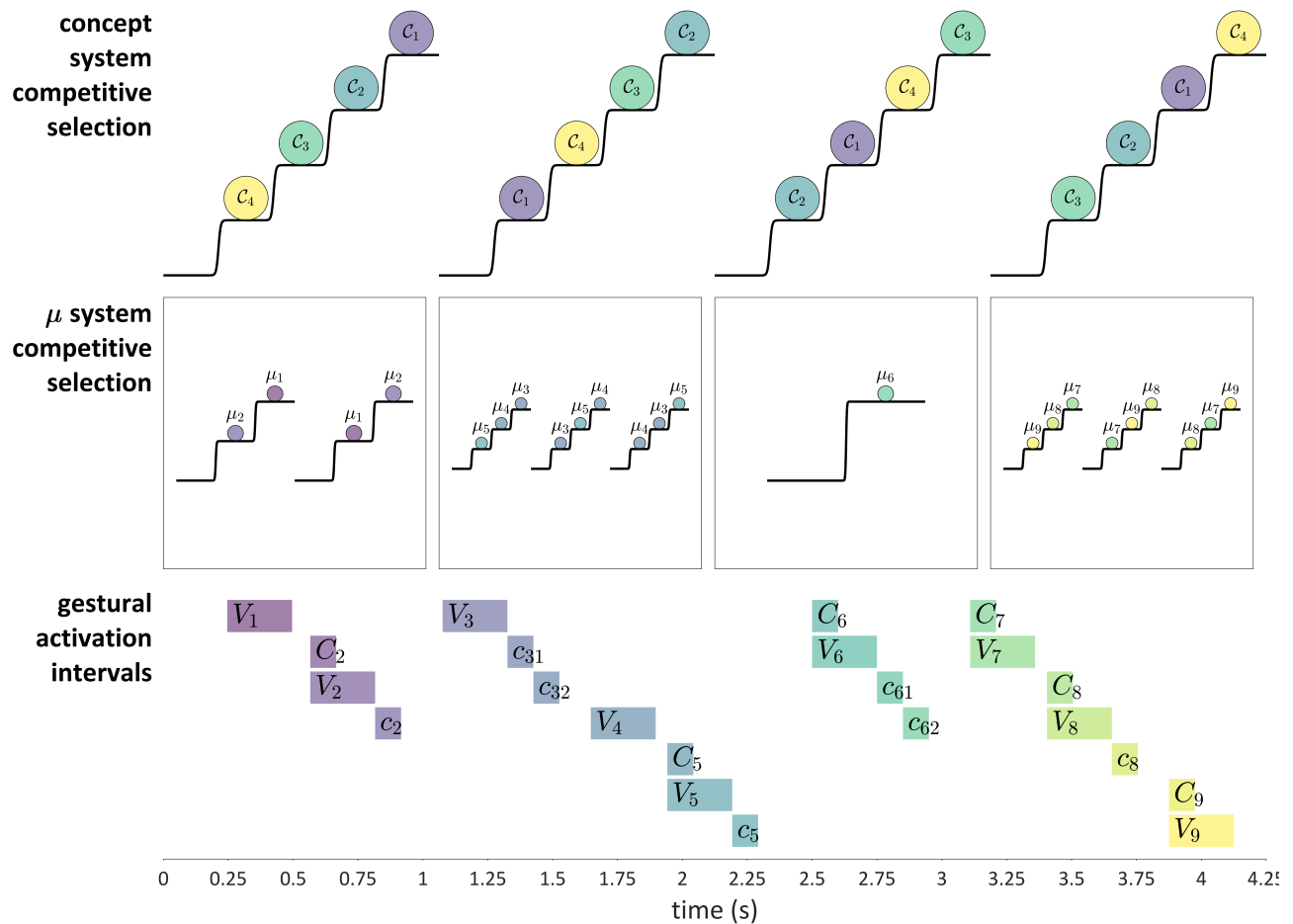
581

582 Furthermore, a recently developed theory of syntactic organization in speech (40) argues that there
583 are two interacting domains of competitive selection. This is known as the *parallel domains*
584 *hypothesis*. One of these domains involves "gestural-motoric" organization of the sort illustrated
585 above, where gestures are organized into competitively selected sets (μ -systems). The other involves
586 "conceptual-syntactic" organization in which concept systems are organized into competitively
587 selected sets. The hypotheses advanced in (40) hold that sets of co-selected conceptual systems
588 correspond loosely to the prosodic unit called the *phonological word* (a.k.a. pwrđ, or ω), which has
589 the property that there is a single accentual gesture associated with set of co-selected conceptual
590 systems. Moreover, under normal circumstances speakers do not interrupt (for example by pausing)
591 the gestural competitive selection processes which are induced by selection a phonological word.

592

593 These parallel domains of conceptual-syntactic and gestural-motoric competitive selection are
594 illustrated Figure 11 for an utterance which would typically be analyzed as four prosodic words, such
595 as [*a dog*] [*and a cat*] [*chased*] [*the monkey*]. Note that to conserve visual space release gestures have
596 been excluded. The top panel shows the sequence of epochs in competitive selection of concept
597 systems \mathcal{C} . Each of these could in general be composed of a number of co-selected subsystems (not
598 shown). For each epoch of concept system selection, there is a corresponding series of one or more
599 epochs of competitive selection of gestural systems. The model accomplishes this by allowing the
600 concept systems to de-gate the corresponding sets of μ -systems. Within each of these sets of μ -
601 systems, the appropriate initial activation gradient is imposed. Further detail on the implementation
602 is provided in the Supplementary Material.

603



604
 605 Figure 11. Illustration of parallel domains of competitive selection for an utterance with the structure.
 606 Top: concept systems \mathcal{C} are competitive selected. Middle: selection a concept system de-gates
 607 corresponding μ -systems which themselves are competitively selected. Bottom: gestural activation
 608 intervals generated by the model.
 609

610 Although there is no *a priori* constraint on the number of domains of competitive selection that might
 611 be modelled, the parallel domains hypothesis that we adopt makes the strong claim that only two
 612 levels are needed—one for conceptual-syntactic organization and one for gestural-motoric
 613 organization. We examine some of the important consequences of these ideas in section 4.2,
 614 regarding phrasal organization. One aspect of prosodic organization which we do not elaborate on
 615 specifically in this paper involves the metrical (stress-related) organization of gestures, but see (42)
 616 for the idea that the property of "stress" relates to which sets of co-selected gestures (μ -systems)
 617 may include accentual gestures, which in turn are responsible for transient increases in self-attention.

618 3.4 A model of speech rate control with selectional effects

619 When given verbal instructions to "talk fast" or "talk slow", speakers are able to produce speech that
 620 listeners can readily judge to be relatively fast or slow. To quantify this sort of variation, speech rate
 621 is often measured as a count of events per unit time, e.g., syllables per second or phones per second.
 622 There are several important points to consider about these sorts of quantities. First, in order to be
 623 practically useful, an event rate must be measured over a period of time in which multiple events
 624 occur. As the size of the counting window decreases, eventually only one full event is included.

625 Second, there is no consensus on which events are the appropriate ones to count—phones, syllables,
 626 words, or something else? In the current framework, many commonly used units do not even have
 627 an ontological status. Third, even if we ignore the above problems, the resulting rate measure cannot
 628 be assumed to be a very good reflection of what speakers are controlling at any particular instant.
 629 There is no evidence to my knowledge that speakers directly control rate quantities such as
 630 syllables/second or phones/second. If we infer that speakers do not in fact control speech rate as an
 631 event rate *per se*, then what are speakers controlling in order to speak fast or slow?

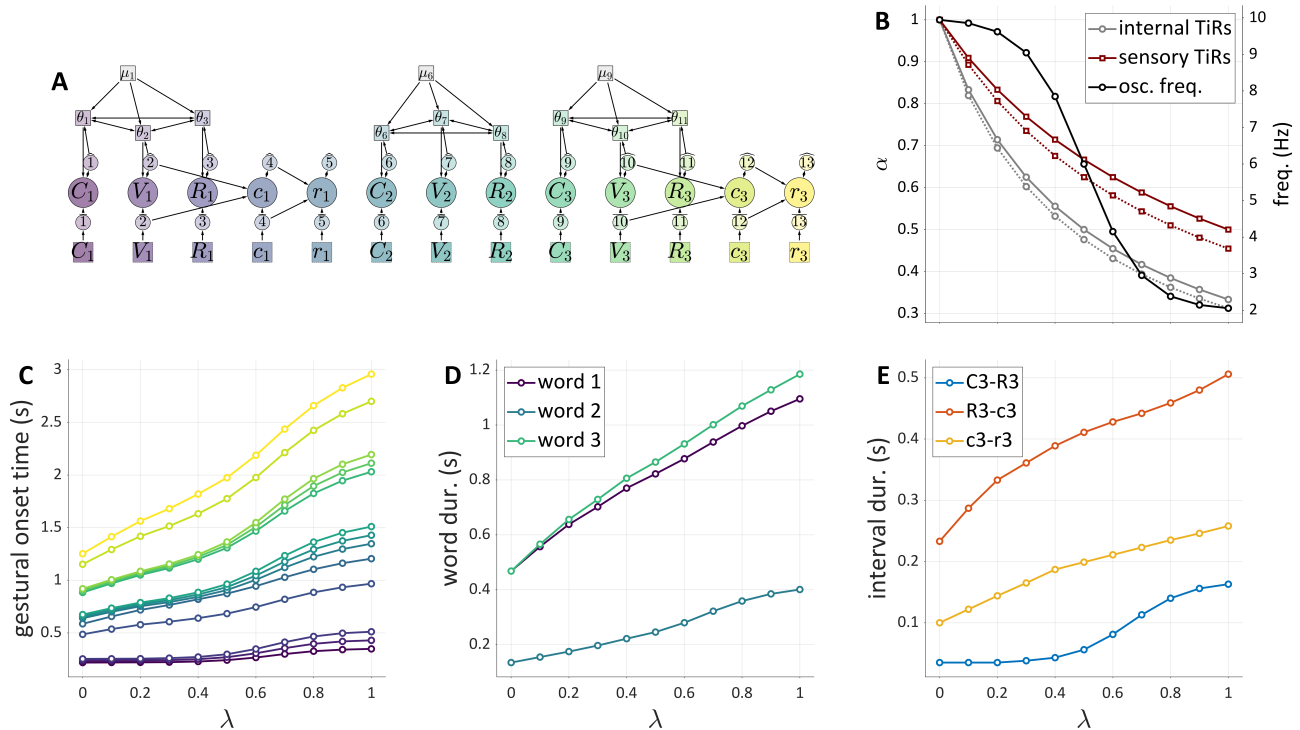
632 The *attentional modulation hypothesis* (5) holds that speakers control rate by modulating their
 633 attention to feedback of their own speech (*self-attention*), and specifically do so in a way that, as self-
 634 attention increases, prioritizes external/sensory feedback over internal feedback. Furthermore, this
 635 hypothesis holds that along with modulating self-attention, speakers may adjust pacing, that is, the
 636 frequencies of gestural planning oscillators. The separate effects of varying these external factors
 637 were already demonstrated in Section 3.2.

638 In addition, a mechanism is need to account for the phenomenon of boundary-related lengthening.
 639 Many empirical studies have shown that speech slows down as speakers approach the ends of
 640 phrases, with greater slowing and increased likelihood of pausing statistically associated with "higher-
 641 level" phrase boundaries (1,43–48). One approach to understanding the mechanism responsible for
 642 such effects is the π -gesture model of (1), in which it was hypothesized that boundary-related
 643 lengthening is caused by a special type of clock modulating system, a " π -gesture". This clock-
 644 modulating system, when active, slows down the rate of a hypothesized nervous system-internal
 645 global clock, relative to real time. Gestural activation dynamics evolve in the internal clock coordinate,
 646 and so gestural activation intervals are extended in time when a π -gesture is active. Furthermore, it
 647 was suggested in (1) that the degree of activation of a π -gesture varies in relation to the strengths of
 648 prosodic boundaries, such that stronger/higher-level boundaries are associated with greater π -
 649 gesture activation and hence more slowing.

650 How can the phenomenon of boundary-related lengthening be conceptualized in the current
 651 framework, where there is no global internal clock for gestural systems? A fairly straightforward
 652 solution is to recognize that in effect, each gestural system has its own "local clocks", in the form of
 653 the internal and external feedback TiRs, whose integration rates are modulated by self-attention. In
 654 that light, it is sensible to adapt the π -gesture mechanism by positing that self-attention effects on
 655 TiR parameters tend to be greater not only in the final set of gestures selected in each prosodic word
 656 (i.e. final μ -system), but also in the final set of co-selected conceptual systems (i.e. the final \mathcal{C} -system).
 657 As for why it is the final set of selected systems that induces these effects, we reason that speakers
 658 may attend to sensory feedback to a greater degree when there are fewer systems that remain to be
 659 selected. At the end of an utterance, there are no more systems that remain to be selected, and thus
 660 self-attention is greatest. We refer to this idea as the *selectional anticipation hypothesis*, because
 661 anticipation of upcoming selection events is proposed to distract a speaker from attention to
 662 feedback of their own speech. Although this hypothesis is admittedly a bit ad hoc, and alternative
 663 accounts should be considered, we show below that the implementation of this idea is sufficient to
 664 generate the lengthening that occurs at the ends of phrases.

665 Putting the above ideas together, Figure 12 shows how interval durations change as a function of
 666 attentional modulation. The utterance here is a competitively selected sequence of three syllables
 667 with forms CVC, CV, CVC, as shown in Figure 12A. Note that the organization of each syllable conforms

668 to the hybrid control model, entailing that pre-vocalic timing is coordinative and vocalic/post-vocalic
 669 timing is feedback-based. As in Section 3.2, the integration rates of external (sensory) and internal
 670 TiRs, along with oscillator frequencies, are made to vary in response to changes in a control parameter
 671 λ ; these relations are shown in Figure 12B. In addition, the integration rate parameters associated
 672 with the final set of gestures are even more strongly modulated by λ (dotted lines of Figure 12B), to
 673 implement the selectional anticipation hypothesis. The initiation times of gestures for each of the 11
 674 values of λ that were simulated are shown vertically in Figure 12C.



675

676 Figure 12. Simulation of variation in speech rate, as controlled by correlated changes in self-attention
 677 and pacing, both indexed by λ . (A) Model schema showing three syllables with the forms CVC, CV,
 678 and CVC. (B) Relations between λ and feedback TiR integration rates (α) and oscillator frequencies.
 679 (C) Times of gestural initiation for each value of λ simulated. (D, E) Word durations and interval
 680 durations of the third word.

681 By simulating variation in speech rate, we are able to generate some of the most essential predictions
 682 of the hybrid control model, introduced in Section 3.1. Recall that this model combined two
 683 hypotheses: prevocalic coordinative control and post-vocalic feedback-control. These hypotheses are
 684 associated with the following three predictions:

685 (i) *Prevocalic attenuation*. The prevocalic coordinative control hypothesis holds that initiation of the
 686 prevocalic constriction and release gestures, along with initiation of the vocalic gesture, is controlled
 687 by a system of coupled oscillators. Moreover, the frequency constraint hypothesis was shown in
 688 Section 3.2 to predict that intervals between these initiations attenuate as rate is increased or
 689 decreased. This effect can be seen in Figure 12E for the C₃-R₃ interval, which is the interval between
 690 constriction formation and release. In other words, the prediction is that prevocalic timing is only so
 691 compressible/expandable, no matter how quickly or slowly a speaker might choose to speak.

692 (ii) *Postvocalic expandability*. Conversely, the post-vocalic feedback-control hypothesis holds that
 693 there is a transition from internally to externally governed control, and that there should be no limits
 694 on the extent to which increasing self-attention can increase the corresponding interval durations.
 695 This prediction is shown in Figure 12E for the R_3 - c_3 interval (which loosely corresponds to acoustic
 696 vowel duration) and the c_3 - r_3 interval (related to constriction duration). These intervals continue to
 697 increase as attention to feedback is increased.

698 (iii) *Sensitivity to feedback perturbation*. Finally, a third prediction of the model is that, when external
 699 feedback governs post-vocalic control (as is predicted for slow rates), perturbations of sensory
 700 feedback will influence post-vocalic control but not prevocalic control.

701 How do these predictions fare in light of current evidence? The ideal tests of predictions (i) and (ii)
 702 require measurements of temporal intervals produced over a wide range of variation in global speech
 703 rate. Unfortunately, most studies of the effects of speech rate do not sufficiently probe extremal
 704 rates, since many studies use categorical adverbial instructions (e.g. *speak fast vs. speak normally vs.*
 705 *speak slowly*). One exception is a recent study using an elicitation paradigm in which the motion rate
 706 of a visual stimulus iconically cued variation in speech rate (49). Utterance targets were words with
 707 either intervocalic singleton or geminate bilabial nasals (/ima/ and /imma/). The study observed that
 708 the timing of constriction formation and release of singleton /m/ exhibited a nonlinear plateau at
 709 slow rates, similar to the prediction for the c_3 - r_3 interval in Figure 12E. This is expected given the
 710 assumption that the formation and release gestures are organized in onset of the second syllable of
 711 the target words. In contrast, the constriction formation-to-release intervals of geminate /mm/ did
 712 not attenuate: they continued to increase in duration as rate slowed. This is expected if the initiation
 713 of the geminate bilabial closure is associated with the first syllable and its release with the second.
 714 Although the dissociation of effects of rate on singletons vs. geminates is not the most direct test of
 715 the hybrid model hypothesis, it shows that more direct tests are warranted.

716 Regarding prediction (iii), a recent study has indeed found evidence that post-vocalic intervals
 717 respond to temporal perturbations of feedback and that pre-vocalic intervals do not (50). This study
 718 found that subtle temporal delays of feedback imposed during a complex onset did not induce
 719 compensatory timing adjustments, while the same perturbations applied during a complex coda did.
 720 This dissociation in feedback sensitivity is a basic prediction of the hybrid model. Another recent study
 721 (51) has found that temporal perturbations induced compensatory adjustments of vowel duration
 722 but not of onset consonant duration (codas were not examined). There may be other reasons why
 723 temporal feedback perturbations have differential effects on prevocalic and vocalic/post-vocalic
 724 intervals, and certainly there is much more to explore with this promising experimental paradigm.
 725 Nonetheless, effects that have been observed so far are remarkably consistent with the predictions
 726 of the hybrid control model.

727 **4 General discussion**

728 The informal logic developed here has many consequences for phonological theories. Below we
 729 discuss three of the most important ones. First, the framework does not allow for direct control over
 730 the timing of articulatory target achievement, and we will argue that this is both conceptually
 731 desirable and empirically justified. Second, structural entities such as syllables and moras can be re-
 732 interpreted in relation to differences in the organization of control. Third, there is no need to posit
 733 the existence of different types of phrases, nor a hierarchical organization of phrases: the appearance

734 of prosodic "structure" above the phonological word can reinterpreted more simply as variation in
735 self-attention conditioned on selection of prosodic words.

736 **4.1 No direct control of target achievement**

737 Some researchers in the TD/AP framework have explicitly hypothesized that control of timing of
738 target achievement is a basic function available in speech (52), or have implicitly assumed such
739 control to be available (53). More generally, outside of the AP/TD framework, it has been argued that
740 speakers prioritize control of the timing of articulatory and acoustic target events over control of the
741 initiation of very same actions that are responsible for achieving those targets (48,54,55). "Target
742 achievement" is defined here as a event in which the state of the vocal tract reaches a putative target
743 state that is associated with a gestural system.

744 Direct control of the timing of gestural target achievement is prohibited by our logic because TiRs
745 control when gestural systems become active and cease to be active, and neither of these events fully
746 determines the time at which targets are achieved. The TiR framework of course allows for *indirect*
747 control of target achievement timing, via the trivial fact that target achievement depends in part on
748 when a gesture is activated. Yet other factors, which are outside the scope of the TiR model, play a
749 role as well. In standard Task Dynamics (7) these factors include the strengths of the forces that
750 gestural systems exert on a tract variable systems—both driving forces and dissipative damping
751 forces—as well as how these forces are blended when multiple gestural systems are active. Or, in an
752 alternative model of how gestures influence tract variable control systems (41), the relevant factors
753 are the strengths, timecourses, and distributions of inhibitory and excitatory forces that gestural
754 systems exert on spatial fields that encode targets. In either case, target achievement cannot be
755 understood to be controlled directly by TiRs.

756 A major conceptual issue with direct control of target achievement is that it requires an unrealistically
757 omniscient system which also has accurate knowledge of the future. In order to control exactly when
758 a target is achieved, a control system must initiate a movement at precisely the right time, which in
759 turn requires that the system is able to anticipate the combined influences on the vocal tract state of
760 all currently active subsystems and all subsystems which might become active in the near future. This
761 all-knowing planner must accomplish these calculations before the critical time at which the
762 movement must be initiated. While such calculations are not in principle impossible, they do require
763 a system which has access to an implausibly high degree of information from many subsystems.

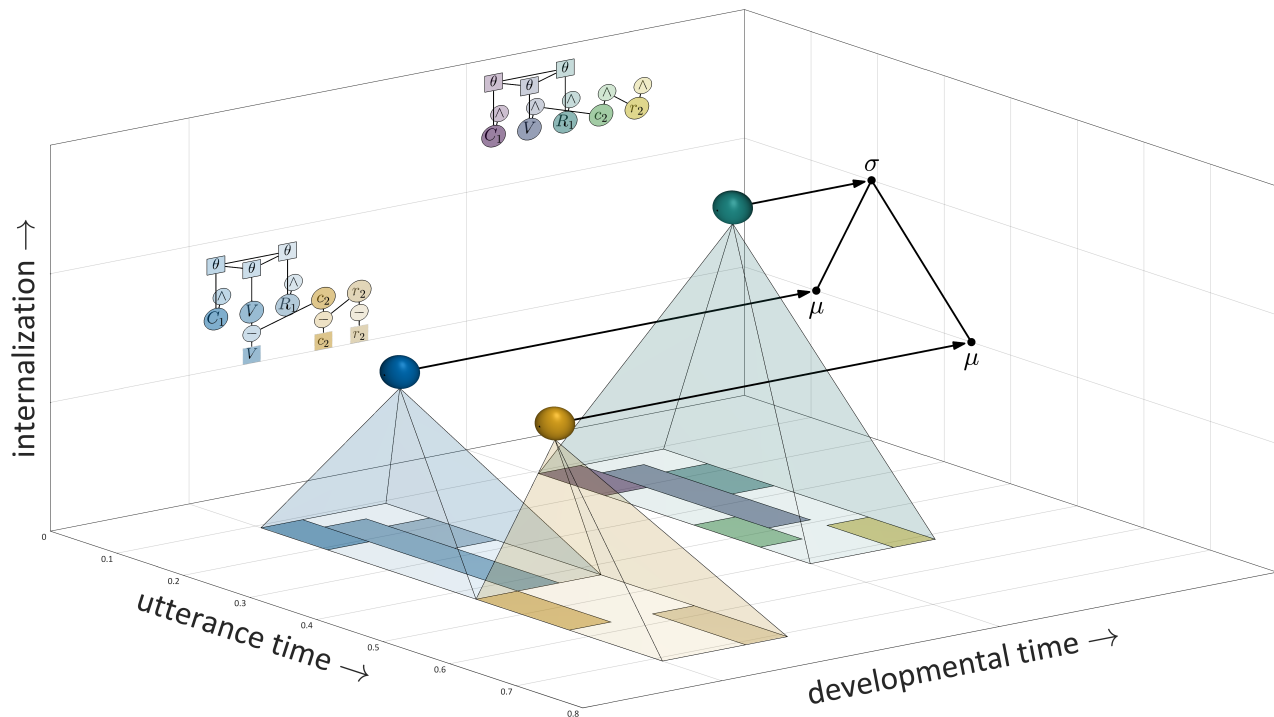
764 A primary empirical argument for direct control of target achievement is premised on the claim that
765 there is less variability associated with timing of target achievement than variability associated with
766 timing of movement onsets. This is argued in (48,54) to suggest that timing of target achievement is
767 not only independently controlled, but also prioritized over timing of movement initiation. The
768 difference in variability upon which the argument is premised has been observed in non-speech
769 studies in which an actor must hit or catch a moving object. Yet these sorts of non-speech examples
770 do not necessarily translate to speech, because in articulation there are no uncontrolled moving
771 objects that the effectors must collide with at the right place in space and time—speech is simply not
772 like catching a ball. Indeed, only one study of speech appears to have concluded that there is less
773 variability in target vs. initiation timing (56), and this interpretation of the data is highly questionable
774 due to differences in how the two events were measured.

775 Empirically observed phonetic and phonological patterns indeed provide the strongest argument
 776 *against* direct control of target achievement timing. Phonetic reduction of targets, which can arise
 777 from insufficient allotment of time for a target to be achieved, is rampant in speech. The "perfect
 778 memory" example of (8) shows how at fast speech rates the word-final [t] can be not only acoustically
 779 absent but also quite reduced kinematically when the preceding and following velar and bilabial
 780 closures overlap. If speakers prioritized the timing of the [t] target relative to either the preceding or
 781 following targets, this sort of reduction presumably would happen far less often. The prevalence of
 782 historical sound changes which appear to involve deletion of constriction targets, argues against the
 783 notion that speakers are all that concerned with achieving targets. Certainly, the consequences of
 784 failing to achieve a target are usually not so severe: in order to recognize the intentions of speakers,
 785 listeners can use contextual information and acoustic cues that not directly related to target
 786 achievement. Rather than being a priority, our informal logic views target achievement as an indirect
 787 and often not-so-necessary consequence of activating gestural systems.

788 4.2 Reinterpretation of syllabic and moraic structure

789 Many phonological theories make use of certain structural entities—syllables (σ) and moras (μ)—as
 790 explanatory structures for phonological patterns. These entities are viewed as groupings of segments,
 791 with moras being subconstituents of syllables, as was shown in Figure 1B. Selection-coordination
 792 theory (3,4) has argued that these entities, rather than being parts of a structure, should be thought
 793 of as different classes of phonological patterns that are learned in different stages of a particular
 794 developmental sequence, over which the organization of control changes. This idea is referred to as
 795 the *holographic hypothesis*, because it holds that what appears to be a multi-level structure of
 796 syllables and moras is in fact a projection over developmental time of two single-level structures
 797 which do not exist simultaneously. This is loosely analogous to a hologram, which encodes a three-
 798 dimensional image in two dimensions.

799 The holographic hypothesis is exemplified in Figure 13 for a CVC syllable. Early in development, the
 800 post-vocalic constriction gesture is controlled entirely by sensory feedback (i.e., extra-gestural TiRs),
 801 and so phonological patterns learned at this time are associated with a moraic structure, reflecting a
 802 stronger differentiation in control of pre-vocalic and post-vocalic articulation. Subsequently, speakers
 803 learn to activate and deactivate the post-vocalic constriction/release with internal TiRs, process called
 804 *internalization*. This leads to initiation of the post-vocalic constriction before termination of the
 805 vocalic gesture, hence an increase in articulatory overlap/coarticulation. Phonological patterns
 806 learned in conjunction with this internalized organization of control are associated with syllables,
 807 rather than moras. Similar reasoning applies to other syllable shapes such as $\{C\}\{CV\} \rightarrow \{CCV\}$ and
 808 $\{CV\}\{V\} \rightarrow \{CVV\}$, where developmental transitions in the internalization of control can account for
 809 cross-linguistic phonetic and phonological variation (3).



810

811 Figure 13. Visualization of the holographic hypothesis, for a CVC form. In an early stage of
 812 development, control over the post-vocalic constriction is based entirely on sensory feedback.
 813 Phonological patterns learned in this stage of development are described with moraic structure. In a
 814 later stage of development, control has been internalized, and phonological patterns learned in this
 815 stage are described with syllabic structure.

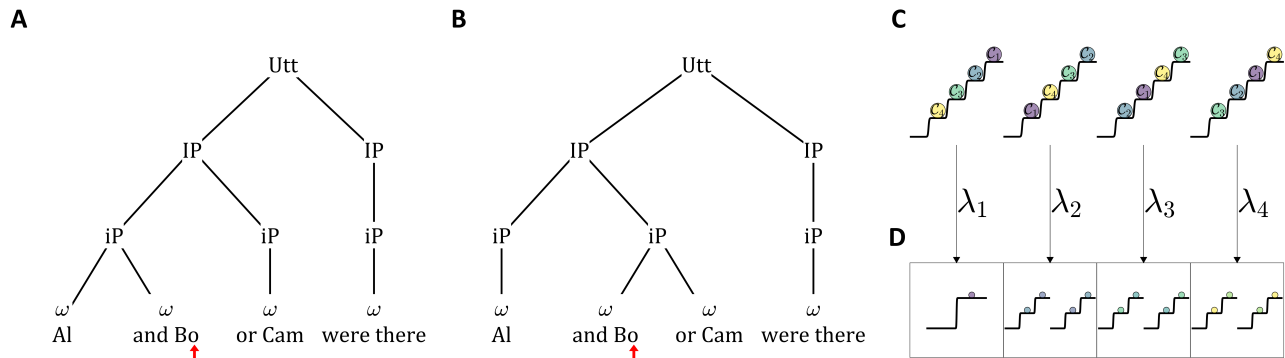
816 Exactly what causes internalization and governs its progression are open questions that presumably
 817 relate to information transmission. More internalization is associated with a greater rate of
 818 information production in speech, or in other words, increased efficiency of communication.
 819 Conversely, too much internalization can result in degrees of articulatory overlap which sacrifice
 820 perceptual recoverability (57–60), reflecting constraints on channel capacity. It is far from clear how
 821 these opposing considerations—information rate vs. channel capacity—might be mechanistically
 822 manifested in a model of utterance-timescale processes. Informational aspects of speech, which by
 823 definition require analysis of the space of possible state trajectories of gestural systems, necessarily
 824 involve attention to patterns on lifespan timescales and speech-community spatial scales. Thus the
 825 challenge lies in understanding how these relatively large timescale informational forces translate to
 826 changes in utterance-scale control.

827 4.3 Reinterpretation of prosodic phrase structure and boundaries

828 There are many prosodic theories in which prosodic words (ω) are understood to be hierarchically
 829 structured into various types of phrases. A "phrase" in this context simply refers to a grouping of
 830 prosodic words. Different types of phrases have been proposed, with two of the most popular being
 831 the "intonational phrase" (IP) and "intermediate phrase" (iP) from (61); these were shown in Figure
 832 1B. Many theories additionally posit that these types of phrases can be recursively hierarchically
 833 structured (62–64), such that a given type of phrase can contain instances of itself. In general, the
 834 motivations for positing phrase structures of this sort are diverse and too complex to address in detail

835 here, but most of them relate either to the likelihood that certain phonological patterns will occur in
 836 some portion of an utterance or to statistical patterns in measures of pitch or duration observed in
 837 longer utterances.

838 To provide an example, consider the question: *Who was in the library?*, answered with the utterance
 839 *Al and Bo or Cam were there*. This utterance has two probable interpretations, and in many theories
 840 these would be disambiguated by the prosodic structures shown in Figure 14 (A vs. B):



841

842 Figure 14. Hierarchical prosodic structure reinterpreted as variation in attentional modulation of
 843 control parameters. (A vs. B): alternative hierarchical prosodic structures purported to encode a
 844 difference in conceptual grouping. Red arrows indicate timepoint discussed in the text. (C, D) In
 845 different epochs of concept system selection, self-attention (λ) may differ, resulting in differences in
 846 temporal control.

847 The motivation for positing the structural distinction between (A) and (B) is that it can account for
 848 certain empirical patterns related to conceptual grouping. Consider specifically the period of time in
 849 the vicinity of the red arrows, near the end of the production of *Bo*, which is often conceptualized as
 850 a phrase "boundary". Here utterance (A), compared to (B), will tend to exhibit a larger fall of pitch,
 851 greater boundary-related lengthening, and a greater likelihood of a pause. The pitch of the following
 852 word may also start at a higher value. Hierarchical structural analyses hold that these differences
 853 occur because there is a "higher-level boundary" here in (A) than in (B), that is, an intermediate phrase
 854 boundary vs. a prosodic word boundary.

855 The logic of multilevel competitive selection makes hierarchical or recursive phrasal structure
 856 unnecessary. If anything, our framework corresponds to a flat, anarchical organization of prosodic
 857 words—though more appropriately it rejects the notion that prosodic words are parts of structures
 858 in the first place, and "boundaries" are seen as wholly metaphoric. How can regularities in
 859 intonational patterns such as in Figure 14 (A vs. B) be understood, without the notions of phrase
 860 hierarchies and boundaries?

861 Recall that each prosodic word is one set of co-selected concept systems, which are associated with
 862 some number of sets of co-selected gestural systems (Figure 11). Furthermore, recall that boundary-
 863 related lengthening was interpreted as a decrease in integration rates of feedback TiRs, and this
 864 parameter modulation is proposed to be greater for the last set of systems in a competitively selected
 865 set (the selectional anticipation hypothesis), as simulated in Figure 12. This reasoning leads to an
 866 alternative understanding of why there exists phonetic and phonological variation that correlates

867 with prosodic organization: rather than being due to "structural" differences, the variation arises from
 868 differences in how TiR parameters are modulated for each prosodic word, as suggested by the arrows
 869 in Figure 14 (C and D). Rather than constructing a structure of prosodic words for each utterance,
 870 speakers simply learn to adjust self-attention in a way that can reflect conceptual relations between
 871 systems of concepts. Presumably many forms of discourse-related and paralinguistic information can
 872 be signaled in this way, including focus phenomena such as emphatic and contrastive focus.

873 5 Conclusion

874 To conclude, we return to the initial questions of this paper: (i) what determines the duration of that
 875 *shush* that you gave to the loud person in the library, and (ii) how do you slow down the rant to your
 876 friend in the coffee shop? According to the feedback-based logic of temporal control, your *shush*
 877 duration is most likely determined by a sensory feedback-based control system (an external, non-
 878 autonomous TiR), and depending upon various factors (how angry you are, how far away the loud
 879 student is), you will diminish the integration rate of the TiR and/or increase its threshold to extend
 880 the duration of the sound. Later on in the coffee shop, you slow down your rant in effect by doing the
 881 same thing: increasing self-attention.

882 There are several important conceptual and theoretical implications of our informal logic. First, all
 883 control of timing must be understood in terms of systems and their interactions, and this
 884 understanding involves the formulation of change rules to describe how system states evolve in time.
 885 Second, the systems which control timing do not "represent" time in any direct sense; the states of
 886 systems are defined in units of activation, and activation is never a direct reflection of elapsed time.
 887 Instead, it is more appropriate to say that timing is controlled via the integration of force, in
 888 combination with thresholds that determine when systems act. Third, the timing of target
 889 achievement is not a controlled event. Finally, much of the theoretical vocabulary that spans the
 890 range of timescales portrayed in Figure 1 is contestable, and new interpretations of empirical patterns
 891 can be derived from our logic. This applies to units such as syllables and moras, and also to hierarchical
 892 and recursive organizations of phrases. Ultimately the logic is useful because it facilitates a unified
 893 understanding of temporal patterns in speech, from the short timescale of articulatory timing to the
 894 large timescale of variation in speech rate.

895 6 Acknowledgments

896 I would like to thank members of the Cornell Phonetics Lab for discussion of the ideas in this
 897 manuscript.

898 7 Data Availability Statement

899 The code for running all simulations and generating all figures in this manuscript can be found on
 900 Github here: <https://github.com/tilsen/TiR-model.git>.

901 8 References

- 902 1. Byrd D, Saltzman E. The elastic phrase: Modeling the dynamics of boundary-adjacent lengthening.
 903 *Journal of Phonetics*. 2003;31(2):149–80.

- 904 2. Tilsen S. A Dynamical Model of Hierarchical Selection and Coordination in Speech Planning. *PLoS*
905 *one*. 2013;8(4):e62800.
- 906 3. Tilsen S. Selection and coordination: The articulatory basis for the emergence of phonological
907 structure. *Journal of Phonetics*. 2016;55:53–77.
- 908 4. Tilsen S. Selection-coordination theory. *Cornell Working Papers in Phonetics and Phonology*,
909 2014. 2014;24–72.
- 910 5. Tilsen S. Three mechanisms for modeling articulation: selection, coordination, and intention.
911 2018. (*Cornell Working Papers in Phonetics and Phonology* 2018).
- 912 6. Kelso JA, Saltzman EL, Tuller B. The dynamical perspective on speech production: Data and theory.
913 *Journal of Phonetics*. 1986;14(1):29–59.
- 914 7. Saltzman E, Munhall K. A dynamical approach to gestural patterning in speech production.
915 *Ecological Psychology*. 1989;1(4):333–82.
- 916 8. Browman C, Goldstein L. Tiers in articulatory phonology, with some implications for casual
917 speech. *Between the grammar and physics of speech*. 1990;341–76.
- 918 9. Ladefoged P. *Vowels and consonants: An introduction to the sounds of the world*. Blackwell
919 Publications; 2001.
- 920 10. Port RF, Leary AP. Against formal phonology. *Language*. 2005;927–64.
- 921 11. Jordan MI. *Serial order: a parallel distributed processing approach*. Technical report, June 1985-
922 March 1986. California Univ., San Diego, La Jolla (USA). Inst. for Cognitive Science; 1986.
- 923 12. Browman C, Goldstein L. Some notes on syllable structure in articulatory phonology. *Phonetica*.
924 1988;45(2–4):140–55.
- 925 13. Browman C, Goldstein L. Articulatory phonology: An overview. *Phonetica*. 1992;49(3–4):155–80.
- 926 14. Kelso JAS, Tuller B. *Intrinsic time in speech production: theory, methodology, and preliminary*
927 *observations*. Sensory and motor processes in language Hillsdale, NJ: Erlbaum. 1987;203:222.
- 928 15. Saltzman E, Nam H, Krivokapic J, Goldstein L. A task-dynamic toolkit for modeling the effects of
929 prosodic structure on articulation. In: *Proceedings of the 4th international conference on speech*
930 *prosody*. Brazil: Campinas; 2008. p. 175–84.
- 931 16. Goldstein L, Byrd D, Saltzman E. The role of vocal tract gestural action units in understanding the
932 evolution of phonology. In: *Action to language via the mirror neuron system*. Cambridge:
933 Cambridge University Press; 2006. p. 215–49.
- 934 17. Schöner G. Timing, clocks, and dynamical systems. *Brain and cognition*. 2002;48(1):31–51.
- 935 18. Sorensen T, Gafos A. The gesture as an autonomous nonlinear dynamical system. *Ecological*
936 *Psychology*. 2016;28(4):188–215.

- 937 19. Browman C, Goldstein L. Gestural syllable position effects in American English. *Producing speech: Contemporary issues*. 1995;19–33.
938
- 939 20. Tilsen S. Exertive modulation of speech and articulatory phasing. *Journal of Phonetics*.
940 2017;64:34–50.
- 941 21. Tilsen S, Zec D, Bjorndahl C, Butler B, L'Esperance M-J, Fisher A, et al. A cross-linguistic
942 investigation of articulatory coordination in word-initial consonant clusters. *Cornell Working
943 Papers in Phonetics and Phonology*. 2012;2012:51–81.
- 944 22. Marin S, Pouplier M. Temporal organization of complex onsets and codas in American English:
945 Testing the predictions of a gestural coupling model. *Motor Control*. 2010;14(3):380–407.
- 946 23. Buzsáki G, Draguhn A. Neuronal oscillations in cortical networks. *science*. 2004;304(5679):1926–
947 9.
- 948 24. Buzsaki G. *Rhythms of the Brain*. Oxford University Press; 2006.
- 949 25. Nam H. Syllable-level intergestural timing model: Split-gesture dynamics focusing on positional
950 asymmetry and moraic structure. In: In J Cole & J I Hualde (Eds), *Laboratory phonology*. Berlin,
951 New York: Walter de Gruyter; 2007. p. 483–506.
- 952 26. Tilsen S. Effects of syllable stress on articulatory planning observed in a stop-signal experiment.
953 *Journal of Phonetics*. 2011;(39):642–59.
- 954 27. Cai S, Ghosh SS, Guenther FH, Perkell JS. Focal manipulations of formant trajectories reveal a role
955 of auditory feedback in the online control of both within-syllable and between-syllable speech
956 timing. *The Journal of Neuroscience*. 2011;31(45):16483–90.
- 957 28. Houde JF, Jordan MI. Sensorimotor Adaptation in Speech Production. *Science*. 1998 Feb
958 20;279(5354):1213–6.
- 959 29. Larson CR, Burnett TA, Bauer JJ, Kiran S, Hain TC. Comparison of voice f responses to pitch-shift
960 onset and offset conditions. *The Journal of the Acoustical Society of America*. 2001;110:2845.
- 961 30. Purcell DW, Munhall KG. Compensation following real-time manipulation of formants in isolated
962 vowels. *The Journal of the Acoustical Society of America*. 2006;119(4):2288–97.
- 963 31. Tourville JA, Reilly KJ, Guenther FH. Neural mechanisms underlying auditory feedback control of
964 speech. *Neuroimage*. 2008;39(3):1429–43.
- 965 32. Villacorta VM, Perkell JS, Guenther FH. Sensorimotor adaptation to feedback perturbations of
966 vowel acoustics and its relation to perception. *The Journal of the Acoustical Society of America*.
967 2007;122(4):2306–19.
- 968 33. Grossberg S. The adaptive self-organization of serial order in behavior: Speech, language, and
969 motor control. *Advances in Psychology*. 1987;43:313–400.

- 970 34. Bullock D. Adaptive neural models of queuing and timing in fluent action. Trends in cognitive
971 sciences. 2004;8(9):426–33.
- 972 35. Bullock D, Rhodes B. Competitive queuing for planning and serial performance. CAS/CNS
973 Technical Report Series. 2002;3(003):1–9.
- 974 36. Bohland JW, Bullock D, Guenther FH. Neural representations and mechanisms for the
975 performance of simple speech sequences. Journal of cognitive neuroscience. 2010;22(7):1504–
976 29.
- 977 37. Glasspool DW. Competitive queuing and the articulatory loop. Psychology Press; 2014.
- 978 38. Bhutani N, Sureshababu R, Farooqui AA, Behari M, Goyal V, Murthy A. Queuing of concurrent
979 movement plans by basal ganglia. Journal of Neuroscience. 2013;33(24):9985–97.
- 980 39. Kristan WB. Behavioral sequencing: Competitive queuing in the Fly CNS. Current Biology.
981 2014;24(16):R743–6.
- 982 40. Tilsen S. Syntax with oscillators and energy levels. Berlin: Language Science Press; 2019. (Studies
983 in Laboratory Phonology).
- 984 41. Tilsen S. Motoric mechanisms for the emergence of non-local phonological patterns. Frontiers in
985 Psychology. 2019;10:2143.
- 986 42. Tilsen S. Space and time in models of speech rhythm. Annals of the New York Academy of
987 Sciences. 2019;1453(1):47–66.
- 988 43. Byrd D, Saltzman E. Intragestural dynamics of multiple prosodic boundaries. Journal of Phonetics.
989 1998;26:173–99.
- 990 44. Byrd D. Articulatory vowel lengthening and coordination at phrasal junctures. Phonetica.
991 2000;57(1):3–16.
- 992 45. Krivokapić J. Gestural coordination at prosodic boundaries and its role for prosodic structure and
993 speech planning processes. Philosophical Transactions of the Royal Society of London B: Biological
994 Sciences. 2014;369(1658):20130397.
- 995 46. Byrd D, Krivokapić J, Lee S. How far, how long: On the temporal scope of prosodic boundary
996 effects. J Acoust Soc Am. 2006 Sep;120(3):1589–99.
- 997 47. Turk AE, Shattuck-Hufnagel S. Multiple targets of phrase-final lengthening in American English
998 words. Journal of Phonetics. 2007;35(4):445–72.
- 999 48. Turk A, Shattuck-Hufnagel S. Timing evidence for symbolic phonological representations and
1000 phonology-extrinsic timing in speech production. Frontiers in psychology. 2020;10:2952.
- 1001 49. Tilsen S, Hermes A. Nonlinear effects of speech rate on articulatory timing in singletons and
1002 geminates. In: Proceedings of the 12th International Seminar on Speech Production. 2020.

- 1003 50. Oschkinat M, Hoole P. Compensation to real-time temporal auditory feedback perturbation
1004 depends on syllable position. *The Journal of the Acoustical Society of America*. 2020;148(3):1478–
1005 95.
- 1006 51. Karlin R, Naber C, Parrell B. Auditory Feedback Is Used for Adaptation and Compensation in
1007 Speech Timing. *Journal of Speech, Language, and Hearing Research*. 2021;64(9):3361–81.
- 1008 52. Gafos AI. A grammar of gestural coordination. *Natural Language & Linguistic Theory*.
1009 2002;20(2):269–337.
- 1010 53. Shaw J, Gafos AI, Hoole P, Zeroual C. Dynamic invariance in the phonetic expression of syllable
1011 structure: a case study of Moroccan Arabic consonant clusters. *Phonology*. 2011;28(03):455–90.
- 1012 54. Turk A, Shattuck-Hufnagel S. Timing in talking: what is it used for, and how is it controlled?
1013 *Philosophical Transactions of the Royal Society B: Biological Sciences*. 2014;369(1658):20130395.
- 1014 55. Turk A, Shattuck-Hufnagel S. *Speech Timing: Implications for Theories of Phonology, Speech
1015 Production, and Speech Motor Control*. Vol. 5. Oxford University Press, USA; 2020.
- 1016 56. Perkell JS, Matthies ML. Temporal measures of anticipatory labial coarticulation for the vowel/u:
1017 Within-and cross-subject variability. *The Journal of the Acoustical Society of America*.
1018 1992;91(5):2911–25.
- 1019 57. Chitoran I, Goldstein L. Testing the phonological status of perceptual recoverability: Articulatory
1020 evidence from Georgian. In: *Proc of the 10th Conference on Laboratory Phonology, Paris, June
1021 29th–July 1st. 2006*. p. 69–70.
- 1022 58. Gick B, Campbell F, Oh S, Tamburri-Watt L. Toward universals in the gestural organization of
1023 syllables: A cross-linguistic study of liquids. *Journal of Phonetics*. 2006;34(1):49–72.
- 1024 59. Liberman AM, Cooper FS, Shankweiler DP, Studdert-Kennedy M. Perception of the speech code.
1025 *Psychological review*. 1967;74(6):431.
- 1026 60. Fowler CA, Rosenblum LD. The perception of phonetic gestures. Modularity and the motor theory
1027 of speech perception. 1991;33–59.
- 1028 61. Beckman ME, Pierrehumbert JB. Intonational structure in Japanese and English. *Phonology*.
1029 1986;3:255–309.
- 1030 62. Féry C. Recursion in prosodic structure. *Phonological Studies*. 2010;13:51–60.
- 1031 63. Ladd DR. Intonational phrasing: the case for recursive prosodic structure. *Phonology*. 1986;3:311–
1032 40.
- 1033 64. Ito J, Mester A. Prosodic subcategories in Japanese. *Lingua*. 2013;124:20–40.
- 1034