

1 PIPPET: A Bayesian framework for generalized
2 entrainment to stochastic rhythms

3 Jonathan Cannon

4 November 30, 2020

5 Department of Brain and Cognitive Science, Massachusetts Institute of
6 Technology, Cambridge, MA, USA

7 Tel.: +314-749-6902

8 jcan@mit.edu

9 **Abstract**

10 When presented with complex rhythmic auditory stimuli, humans are
11 able to track underlying temporal structure (e.g., a “beat”), both covertly
12 and with their movements. This capacity goes far beyond that of a simple
13 entrained oscillator, drawing on contextual and enculturated timing ex-
14 pectations and adjusting rapidly to perturbations in event timing, phase,
15 and tempo. Here we propose that the problem of rhythm tracking is
16 most naturally characterized as a problem of continuously estimating an
17 underlying phase and tempo based on precise event times and their cor-
18 respondence to timing expectations. We formalize this problem as a case
19 of inferring a distribution on a hidden state from point process data in
20 continuous time: either Phase Inference from Point Process Event Tim-
21 ing (PIPPET) or Phase And Tempo Inference (PATIPPET). This ap-
22 proach to rhythm tracking generalizes to non-isochronous and multi-voice

rhythms. We demonstrate that these inference problems can be approximately solved using a variational Bayesian method that generalizes the Kalman-Bucy filter to point-process data. These solutions reproduce multiple characteristics of overt and covert human rhythm tracking, including period-dependent phase corrections, illusory contraction of unexpectedly empty intervals, and failure to track excessively syncopated rhythms, and could be plausibly approximated in the brain. PIPPET can serve as the basis for models of performance on a wide range of timing and entrainment tasks and opens the door to even richer predictive processing and active inference models of rhythmic timing.

Keywords: Bayesian Inference, Active Inference, Timing, Rhythm, Entrainment

1 Introduction

The human brain is remarkably proficient at identifying and exploiting temporal structure in its environment, especially in the auditory domain. This phenomenon is most easily observed in the case of auditory stimuli with underlying periodicity: humans adeptly and often spontaneously synchronize their movements with such auditory rhythms [1], and human brain activity in auditory and motor regions aligns to auditory stimulus periodicity even in the absence of movement [2]. Both of these phenomena are cases of “entrainment” (sensorimotor and neural, respectively), where we define “entrainment” as in [3]: the temporal alignment of a biological or behavioral process with the regularities in an exogenously occurring stimulus.

A simple sinusoidal phase oscillator can entrain to a periodic stimulus; however, it is difficult to discuss the flexible entrainment of human behavior and cognitive processes to variable and sometimes aperiodic patterns such as speech without invoking the cognitive concept of “temporal expectation.” Expecta-

50 tions for event timing can be used to achieve a range of behavioral goals. They
 51 can help us hone our sensory detection, our sensory discrimination, and our
 52 response time for behaviorally important stimuli at the anticipated time [4, 5,
 53 6]. In some situations, temporal expectations attenuate neural responses [7],
 54 which may help to conserve neural resources. And timing expectations bias
 55 our perception of time, allowing us to use prior experience to supplement noisy
 56 sensory data as we make temporal judgments [8].

57 Entrainment in humans involves an interplay of stimulus and temporal ex-
 58 pectation [9]. Nowhere is this clearer than in interaction with music, hu-
 59 mankind’s playground for auditory temporal expectation and entrainment [10].
 60 But the precise nature of this interplay is an open question. The framework
 61 of Dynamic Attending Theory characterizes temporal expectancy as pulses of
 62 “attentional energy” issued by entrained neural oscillators, and mathematical
 63 models based on these ideas describe bidirectional interactions between tempo-
 64 ral expectation and entrainment that reproduce aspects of human behavior and
 65 perception [11, 12]. But although the behavior of these models may be satis-
 66 fying, the groundwork underlying them is less so: key high-level concepts like
 67 the “attentional pulse” are difficult to define mechanistically, so the implemen-
 68 tations of these concepts in models remain impressionistic. Moreover, recent
 69 results have emphasized the relevance and neural correlates of aperiodic modes
 70 of temporal expectation [13, 6, 14], but dynamic attending models are designed
 71 to describe entrainment to periodicity and cannot account for aperiodic forms
 72 of structured temporal expectation such as entrainment to memorized temporal
 73 patterns, irregular musical meters, and the loose temporal regularities of speech
 74 [15].

75 Here, we propose a normative framework for understanding the interaction
 76 of entrainment and expectation. The goal is to first suggest a formal problem

77 that is being solved by general entrainment – namely, the problem of inferring
 78 the state of the exogenous process giving rise to a series of events in time – and
 79 then use mathematics to describe an optimal solution to that problem. This
 80 teleological approach to entrainment complements previous approaches based on
 81 cognitive constructs like dynamic attending. It brings to the table a concrete and
 82 mathematically precise link between the phenomenon of expectation-informed
 83 entrainment and the statistical structure of the stimuli that entrainment is used
 84 to exploit. If such a solution bears sufficient similarities to observations in
 85 humans, then we can begin to discuss human entrainment as a precise reflection
 86 of the temporal structure of the sensory world. Moreover, this approach is
 87 sufficiently general to describe entrainment to “stochastic” rhythms (rhythms in
 88 which some expected events may omitted) based on either periodic or aperiodic
 89 temporal expectations.

90 In the next section, we discuss previous models of expectation in cognition
 91 and where they fall short for our purposes. We then formulate three versions
 92 of the problem of entrainment that are amenable to precise solutions. In the
 93 first, “Phase Inference from Point Process Event Timing” (PIPPET), a hidden
 94 phase variable advances steadily with added noise, and the observer is tasked
 95 with continuously inferring the phase based on the observation of events emit-
 96 ted probabilistically at certain phases with certain degrees of precision. The
 97 variational Bayesian solution to this inference problem provides a continuous
 98 estimate of phase that entrains to the actual phase, as well as an estimated level
 99 of certainty about that phase. In the second, “Phase And Tempo Inference from
 100 Point Process Event Timing” (PATIPPET), the rate of phase advance (tempo)
 101 is also a dynamic variable with drift, and the solution simultaneously estimates
 102 phase, tempo, and certainty about both. The third (multi-PIPPET) general-
 103 izes the first two to incorporate the observation of multiple types of events, each

104 with distinct characteristic phases and precisions, into the inference process.

105 In the following section, we simulate these solutions, drawing on music as
106 a rich source of intuitive examples of entrainment informed by expectation. In
107 doing so, we provide intuition into the range of behaviors of these solutions,
108 and show how they reproduce key aspects of human sensorimotor entrainment
109 behavior that are not explained by other entrainment models. These include:

- 110 1. Failure to track phase through excessive syncopation (events occurring at
111 weakly expected times but omitted at strongly expected times).
- 112 2. Illusory contraction of intervals when expected events are omitted.
- 113 3. Near-linear corrections to phase after event timing perturbations, with
114 larger (and even over-) corrections for stimulus trains with longer inter-
115 onset intervals.

116 In the final section, we discuss the potential contributions of PIPPET and
117 PATIPPET to our understanding of human entrainment.

118 2 Mathematical framework

119 The framework of “predictive processing” has emerged as the preferred lens for
120 modeling the role of expectations in the brain [16, 17]. According to this con-
121 stellation of ideas, expectations (or, interchangeably, “predictions”) from higher
122 levels of the sensory processing hierarchy are sent to lower levels, where they
123 are compared to incoming sensory information and used to compute “predic-
124 tion errors.” These prediction errors are used to inform dynamic adjustments
125 to the expectations at all levels of processing, as well as slower adjustments to
126 the learned models upon which predictions are based. This is formalized as
127 a process of variational Bayesian inference based on a hierarchical generative
128 model.

Predictive processing would be a natural modeling framework for understanding rhythmic expectation and entrainment as inference [18, 19, 20] except for one key limitation: existing predictive coding models that operate in continuous time are structured to perform inference based on continuous observation, characterizing prediction errors in terms of deviation between a true level of input and a mean expected level of input [21, 22]. They describe predictions about “what” rather than “when,” and are therefore ill-suited to characterizing moment-by-moment errors in *timing* prediction, which arrive sporadically, separated by intervals largely devoid of informative prediction error. This may be a fundamental shortcoming in modeling inference in the brain: behavior and neurophysiology suggests that information about “when” is carried by its own distinctive pathways and represented separately from “what,” both in perceptual and motor tasks [23, 6, 10]. Bayesian methods have been applied to describe inferences about timing in the brain [24, 25, 26], but in these cases the problem the brain solves has been formulated as discrete inferences about consecutive intervals rather than a continuous inference process.

Here, we use event timing to inform a continuous variational inference process using the mathematical tool of point processes. The result approximates an ideal observer with respect to a generative process in continuous time that describes the probabilistic generation of a time series of events.

2.1 Phase Inference from Point Process Event Timing (PIPPET)

PIPPET is a simple generative model of a homogeneous, temporally structured series of instantaneous sensory events. This model consists of a phase $\phi \in \mathbb{R}$

153 that advances as a drift-diffusion process:

$$d\phi = dt + \sigma dW_t \quad (1)$$

154 and an inhomogeneous point process that generates events with probability
 155 $\tau(\phi)$, a function of phase. We will refer to $\tau(\phi)$ as a “temporal expectation
 156 template,” though it can also be understood as a hazard function for events. To
 157 achieve both analytical tractability and flexible descriptive power, we assume
 158 that $\tau(\phi)$ is a sum of a constant τ_0 and a countable set of scaled Gaussian
 159 functions indexed by $i = 1, 2, \dots$ etc. Each Gaussian i is centered at a mean
 160 phase ϕ_i with variance v_i and scale τ_i :

$$\tau(\phi) = \tau_0 + \sum_i \tau_i N(\phi|\phi_i, v_i) \quad (2)$$

161 where $N(x|m, v)$ denotes a normalized Gaussian distribution with mean m and
 162 variance v . Each Gaussian mean ϕ_i represents a phase at which an event is
 163 expected; τ_i represents the strength of that expectation; and v_i^{-1} is the tem-
 164 poral precision of that expectation. $\tau_0 > 0$ represents the rate of events being
 165 generated as part of a uniform noise background unrelated to phase. Together,
 166 $\tau(\phi)$ constitutes a likelihood function for an event occurring at phase ϕ . See
 167 Figure 1 for illustration.

168 Note that ϕ is assumed to be on the real line, not the circle. This design
 169 decision allows PIPPET to entrain to temporally patterned expectations with
 170 or without periodic structure by choosing a periodic or aperiodic temporal ex-
 171 pectation template τ . We discuss this decision further in the Discussion section.

172 Given a series of event times $\{t_n\}$, a temporal expectation template $\tau(\phi)$, and
 173 a prior distribution $p_0(\phi)$ describing the distribution of phase at time $t = 0$, the
 174 observer’s goal is to infer a posterior distribution $p_t(\phi)$ describing an estimate

$$\lambda(\phi) = \lambda_0 + \sum_i \lambda_i N(\phi | \phi_i, v_i)$$

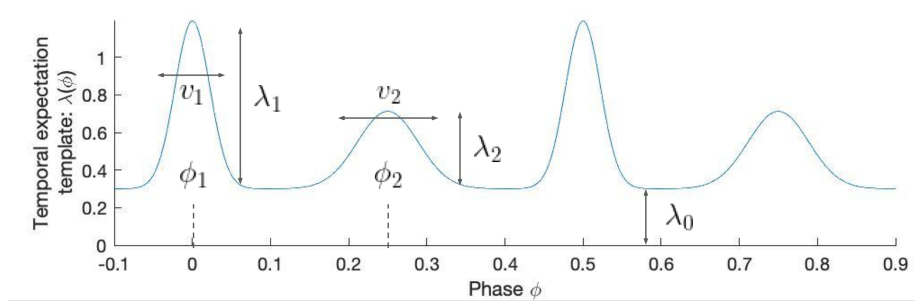


Figure 1: **The temporal expectation template.** In the PIP-PET/PATIPPET generative model, $\tau(\phi)$ represents the instantaneous rate of events occurring when the underlying temporal process is at phase ϕ . This is assumed to be a sum of Gaussian-shaped functions with means ϕ_i representing the phases at which specific events are expected, variances v_i representing the (inverse of) the temporal precision expected of those events, and scales τ_i representing the strength of the expectations. A constant τ_0 is also added, representing the instantaneous rate of events unrelated to the underlying phase.

175 of phase ϕ at every time $t > 0$.

176 In [27], Snyder derives exact equations for the evolution of this posterior
177 distribution over time. Following the predictive processing ansatz of maintaining
178 Gaussian posterior distributions (the Laplace assumption), which provides both
179 computational tractability and neurophysiological plausibility by reducing the
180 representation of the posterior to a mean and a variance, we project the posterior
181 onto a Gaussian at each dt time-step. We do this by moment-matching: we use
182 Snyder’s solution to determine the evolution of the mean and variance of the
183 posterior, and then replace the true posterior with a Gaussian of the same mean
184 and variance. This choice of Gaussian is the choice with minimum KL divergence
185 from the true posterior [28], and therefore also minimizes the free energy of the
186 solution within the family of possible Gaussian posteriors, in accordance with
187 the Free Energy Principle of predictive processing [29].

188 The result of this derivation is a generalization of a Kalman-Bucy filter with
 189 Poisson observation noise. Eden and Brown [30] have derived an explicit form
 190 for this filter, but it relies on a local approximation of the firing rate function
 191 τ that hides some of the interesting effects of events expected at nearby time
 192 points. For τ a mixture of Gaussians, we derive a filter that more accurately
 193 approximates the Bayesian directly from Synder's solution in [27]. Derivation
 194 is presented in Appendix 6.2.

195 **Solution: the PIPPET filter** At any time t , let $\bar{\phi}_t$ denote the mean and V_t
 196 denote the variance of the Gaussian posterior. At each event time t , we let $\bar{\phi}_{t-}$
 197 and V_{t-} denote the left-hand limits of $\bar{\phi}$ and V before the event, and we write
 198 $\bar{\phi}_{t+}$ and V_{t+} to denote their right-hand limit values after the event. $\bar{\phi}_t$ and V_t
 199 evolve according to the ODE

$$\begin{cases} \dot{\bar{\phi}} = 1 - \hat{\Lambda}(\hat{\phi} - \bar{\phi}) \\ \dot{V} = \sigma^2 - \hat{\Lambda}(\hat{V} - V) \end{cases} \quad (3)$$

and at each event $\bar{\phi}_{t+} = \hat{\phi}$ and $V_{t+} = \hat{V}$, where we define

$$\begin{aligned} \hat{\phi} &:= \sum_{i=0,1,\dots} \frac{T_i}{\hat{\Lambda}} \hat{\phi}_i \\ \hat{V} &:= \sum_{i=0,1,\dots} \frac{T_i}{\hat{\Lambda}} \left(K_i + (\hat{\phi}_i - \bar{\phi}_{t+})^2 \right) \end{aligned}$$

(Note that in this formulation, $\bar{\phi}_{t+}$ must be calculated before V_{t+} .)

$$\hat{\phi}_0 := \bar{\phi}_{t-} \text{ and } \hat{\phi}_i := K_i(V_{t-}^{-1}\bar{\phi}_{t-} + v_i^{-1}\phi_i) \text{ for } i > 0.$$

$$T_0 := \tau_0 \text{ and } T_i := \tau_i N(\phi_i | \bar{\phi}_{t-}, v_i + V_{t-}) \text{ for } i > 0.$$

$$K_0 := V_{t-} \text{ and } K_i := \frac{1}{V_{t-}^{-1} + v_i^{-1}} \text{ for } i > 0.$$

$$\hat{\Lambda} := \sum_i T_i$$

Intuitively,

- $\bar{\phi}_t$ is the estimated phase at time t , and V_t is the level of uncertainty about the phase estimate.
- At each event time t , $\tau(\phi)$ serves as a likelihood function for phase, and the role of prior is played by a Gaussian with mean $\bar{\phi}_{t-}$ and variance V_{t-} .
- At any time t , $\hat{\phi}_i$ would be the mean of the posterior if an event occurred and was known to come from Gaussian i . It is a weighted sum of the current mean estimated phase $\bar{\phi}_t$ and the mean ϕ_i of Gaussian i , weighted by the precision $\frac{1}{V_t}$ on estimated phase and the temporal precision $\frac{1}{v_i}$ of the Gaussian generating the event, respectively.
- At any time t , $\hat{\phi}$ and \hat{V} would be the mean and variance of the posterior if an event occurred and its source was not known. These are weighted sums of the influences of each Gaussian, weighted by T_i , the relative likelihood that the event is drawn from Gaussian i .
- Between events, each dt time step is taken as a Bayesian inference with likelihood $1 - \tau(\phi)dt$ and with a Gaussian prior consisting of the posterior of the previous time step carried forward by dt according to the Fokker-Planck evolution associated with the ODE (3).

218 • In the absence of an event, this continuous inference process pushes $\bar{\phi}$ and
 219 V away from $\hat{\phi}$ and \hat{V} with a strength proportionate to $\hat{\Lambda}$, the current
 220 strength of the expectation of an event – thus, the absence of an event
 221 continuously pushes the posterior in the opposite direction as would the
 222 occurrence of an event.

223 **2.2 Phase And Tempo Inference from Point Process Event** 224 **Timing (PATIPPET)**

225 PATIPPET is generative model of homogeneous point process events in time
 226 that extends PIPPET by making the rate of phase advancement itself a noisy
 227 dynamic variable subject to ongoing inference. The dynamic state of the system
 228 is now a two-dimensional vector $\mathbf{x} = \begin{pmatrix} \phi \\ \theta \end{pmatrix}$, where ϕ is the phase as above, θ
 229 is the rate of phase advancement (or tempo), and σ_ϕ and σ_θ are the levels of
 230 phase and tempo noise, respectively:

$$d\mathbf{x} = \begin{pmatrix} \theta \\ 0 \end{pmatrix} dt + \begin{pmatrix} \sigma_\phi dW_t^\phi \\ \sigma_\theta dW_t^\theta \end{pmatrix} \quad (4)$$

231 As above, an inhomogeneous point process generates events with probability
 232 based on a temporal expectation template $\tau(\phi)$, where τ is a sum of Gaussians
 233 and a constant:

$$\tau(\phi) = \tau_0 + \sum_i \tau_i N(\phi | \phi_i, v_i) \quad (5)$$

234 However, in this formulation, we want events to occur with a certain probability
 235 in each $d\phi$ phase bin regardless of tempo, which is not the case if events are
 236 generated with probability $\tau(\phi)dt$; instead, we let events occur with probability

$$\tau(\phi)\mathbb{E}[d\phi] = \tau(\phi)\theta dt.$$

237 Note that this is the same as the PIPPET expression for event rate if we set
 238 $\theta = 1$.

239 Given a series of event times $\{t_n\}$, a temporal expectation template $\tau(\phi)$, and
 240 a prior distribution $p_0(\mathbf{x})$ describing the distribution of phase and tempo at time
 241 $t = 0$, the observer's goal is to infer a posterior distribution $p_t(\mathbf{x})$ describing an
 242 estimate of phase and tempo at every time $t > 0$. A similar derivation provides
 243 a point-process Kalman-Bucy filter that optimally serves this function within
 244 the constraint of Gaussian posteriors, providing a running estimate of a mean
 245 phase and tempo $\hat{\mathbf{x}}_t$ and a phase/tempo covariance matrix \mathbf{V}_t . The solution is
 246 presented in 6.1 and its derivation is presented in 6.2.

247 The resulting PATIPPET filter generalizes the PIPPET filter, and is iden-
 248 tical if the initial tempo distribution is set to a delta distribution at $\theta = 1$ and
 249 σ_θ is set to zero. At each event, the distribution of phase and tempo is dis-
 250 continuously updated to a 2D Gaussian posterior, which evolves continuously
 251 between events. This scheme is similar to [31], which estimates phase and tempo
 252 by updating a 2D Gaussian posterior, but is updated in continuous time and
 253 is significantly more flexible in its capacity to track phase based on arbitrary
 254 temporal expectation templates.

255 **2.3 PIPPET with multiple event streams (multi-PIPPET)**

256 Finally, we generalize PIPPET to include multiple types of events (indexed by
 257 j), each generated as point processes with rates determined by functions $\tau^j(\phi)$
 258 of a single underlying phase:

$$d\phi = dt + \sigma dW_t \quad (6)$$

259

$$\tau^j(\phi) = \tau_0^j + \sum_i \tau_i^j N(\phi | \phi_i^j, v_i^j) \quad (7)$$

260 The Kalman-Bucy estimate of phase for this model is described by mean $\bar{\phi}$
 261 and variance V evolving according to the ODE

$$\begin{cases} \dot{\bar{\phi}} = 1 - \sum_j \hat{\Lambda}^j (\hat{\phi}^j - \bar{\phi}) \\ \dot{V} = \sigma^2 - \sum_j \hat{\Lambda}^j (\hat{V}^j - V) \end{cases} \quad (8)$$

262 and resetting to $\bar{\phi}_{t+} = \hat{\phi}^j$ and $V_{t+} = \hat{V}^j$ when an event occurs in stream j , where
 263 we define $\hat{\Lambda}^j$, $\hat{\phi}^j$, and \hat{V}^j as we defined $\hat{\Lambda}$, $\hat{\phi}$, and \hat{V} above but in reference only
 264 to event stream j .

265 The same adjustment can be made to the PATIPPET generative model, and
 266 the PATIPPET filter can be similarly generalized to account for multiple event
 267 streams.

268 3 Results

269 In this section we conduct a series of simulations to explore parallels between the
 270 behavior of the the PIPPET and PATIPPET filters and human entrainment.
 271 Parameters for these simulations are listed in Appendix 6.3.

272 3.1 Response to events: phase and variance correction

273 We simulated PIPPET filter with simple metronomic expectations to illustrate
 274 its basic behavior. Events occurring near an expected event phase ϕ_i cause the
 275 mean phase estimate $\bar{\phi}$ to shift linearly toward ϕ_i , as indicated by the plateaus
 276 in the phase transition function (Figure 2A). Events occurring far from any
 277 expected event phase ϕ_i caused negligible adjustment in the phase estimate
 278 because they were attributed to the background rate τ_0 of events occurring
 279 unrelated to any specific expectation. This leads to a phase response curve
 280 that crosses zero with negative slope near each expected event phase and sits

281 uniformly near zero away from expected event phases (Figure 2A).

282 If the estimated phase $\bar{\phi}_{t-}$ just before an event time t was very close to an
283 expected event phase ϕ_i , the phase uncertainty V decreased at the event, which
284 effectively “corroborated” the phase estimate (Figure 2B). Events occurring
285 when $\bar{\phi}_{t-}$ was far from any expected event phase had no impact on V , as they
286 were effectively attributed to the background noise rate τ_0 and thus contained
287 no new information about phase. Events occurring in the liminal zone near but
288 not very near an expected event phase ϕ_i caused uncertainty V to increase.

289 **3.2 Stochastic rhythms with uneven subdivision**

290 The PIPPET framework describes entrainment to “stochastic” rhythms in which
291 each expected event phase may or may not be populated by an event. Fur-
292 ther, PIPPET is formulated in sufficient generality to describe entrainment to
293 rhythms based on timing expectations with complex, non-isochronous stress
294 patterns [32] and with non-integer duration ratios using suitably designed (or,
295 presumably, learned) temporal expectation templates $\tau(\phi)$. Such rhythmic pat-
296 terns have been shown to support highly precise synchronization in musicians
297 with appropriate training and enculturated expectations [33], and should there-
298 fore be accounted for by any plausible model of human entrainment. Thus,
299 PIPPET is equipped to model entrainment to a very wide range of rhythmic
300 structures with any degree of predictability.

301 As an example of entrainment to a stochastic rhythm based on a temporal
302 structure with non-integer duration ratios, we simulated entrainment to a swing
303 rhythm. The rhythm is based on an underlying grid of “swung” eighth notes,
304 where the first eighth note of every pair is given a slightly longer duration than
305 the second. Though the “swing” feel is often caricatured using eighth note
306 pairs with a 2:1 duration ratio, this value has been shown to vary by player

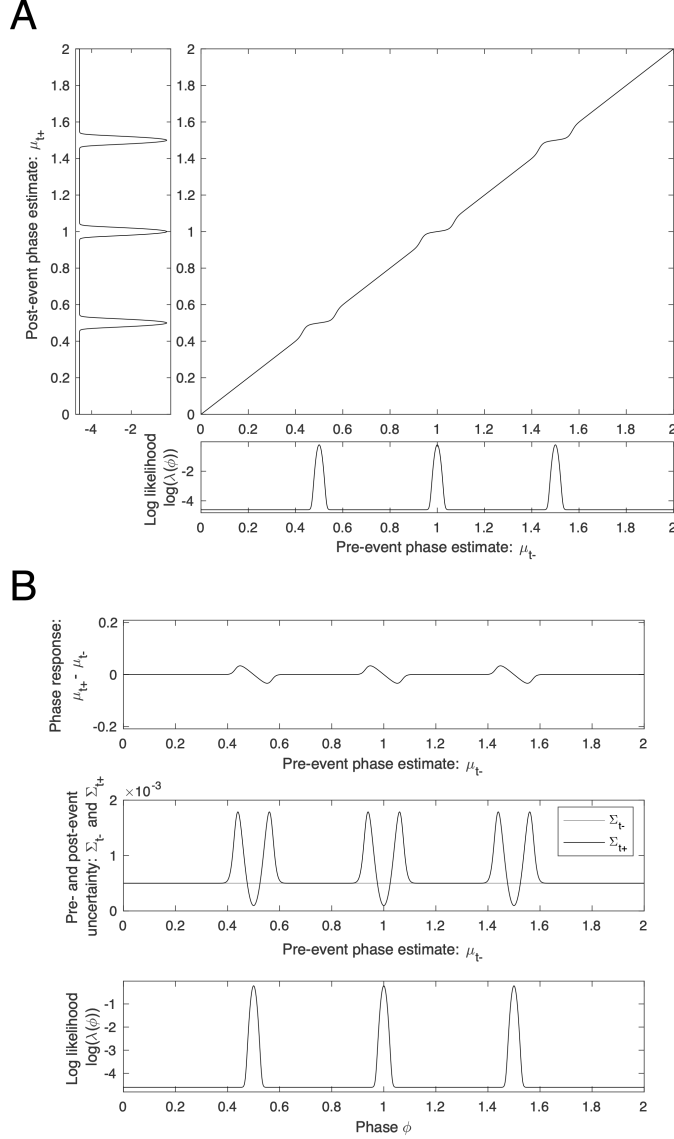


Figure 2: **Characterizing PIPPET's behavior at events** A) Phase transition curve for PIPPET with expectation of three isochronous events. Note that events occurring when the phase estimate $\hat{\phi}_{t-}$ is between expected event phases ϕ_i have little corrective effect on the posterior mean phase $\hat{\phi}_{t+}$, as indicated by a diagonal phase transition curve, whereas events occurring when the estimated phase is near an expected event phase tend draw the phase estimate toward the expected phase, as indicated by plateaus in the phase transition curve. B) Phase and variance response curves. Note that events occurring when estimated phase is very close to an expected event phase cause the variance of the posterior on phase to decrease, whereas events occurring slightly offset from an expected event phase cause the variance to increase. Events occurring far from any expected event phase have little effect on posterior variance.

and tempo and is certainly not limited to small integer ratios [34]. We used a temporal expectation template with a swing ratio close to 3:2 and associated the first eighth note in each pair with a stronger expectation than the second. The simulation entrained to a complex, syncopated rhythm based on this template, and corrected the phase estimate when a phase shift was introduced into the rhythm (Figure 3).

3.3 Failure mode: too much syncopation

Another attractive aspect of the PIPPET framework is that it can account for realistic failures in tracking perfectly timed rhythms. In addition to failures due to time warping described above, failures may occur due to interference between expectations packed closely together in time. Every expected event phase ϕ_i exerts an influence on the evolution of the posterior at all times. This influence is very weak if the current phase estimate is far from ϕ_i . However, if the uncertainty V of the phase estimate is large enough to encompass several expected event phases, or if several events are expected at neighboring phases with insufficient precision, the event may not be fully “attributed” to a single expected event phase. As a result, the adjustment to the phase estimate at an event may reflect an amalgam of these multiple influences, with stronger expectations exerting more influence than weaker ones.

A prime example of this failure mode in human rhythm tracking is tracking overly syncopated rhythms (rhythms with a predominance of events at time points with weaker expectations). Listeners tend to “re-hear” such rhythms by attributing events to metrical positions where events are more strongly expected [35]. Using the expectation template with a swing grid as in the previous section, we simulated a strongly syncopated rhythm (Figure 4). The rhythm’s phase was not tracked successfully due to a convergence of factors. Phase un-

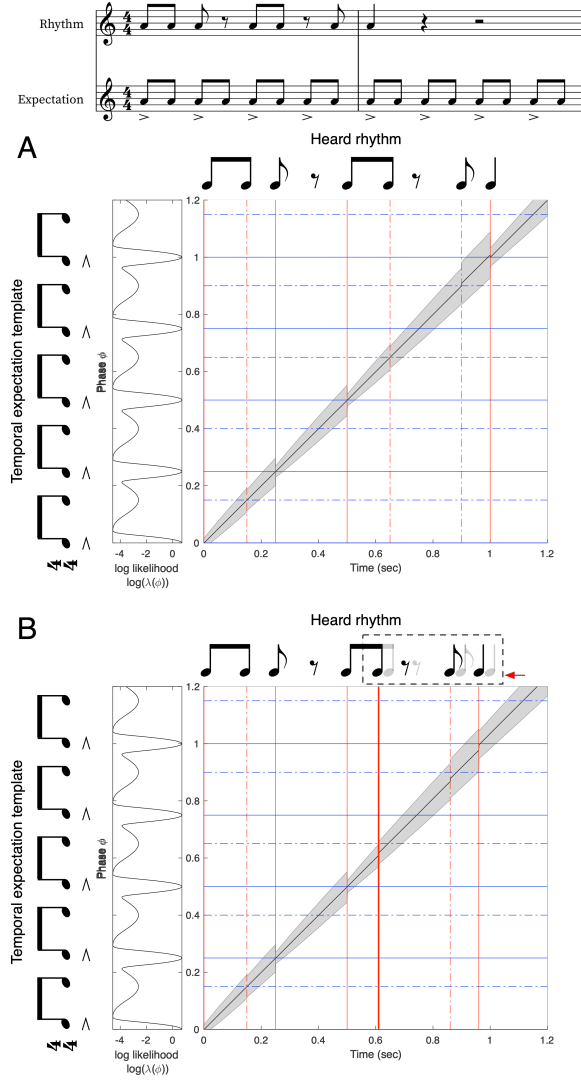


Figure 3: **Tracking phase through swung rhythms.** (Same color key as 5.) A: Phase is estimated over the course of a rhythm. Temporal expectations are not isochronous, but instead represent a swing pattern in which the first eighth note of every pair is slightly longer and more strongly expected than the second. Dotted lines correspond to weak expectations and solid lines correspond to strong expectations. B: A phase shift is introduced into the rhythm, moving all subsequent events earlier in time. When the first early event arrives, uncertainty V increases. Mean estimated phase $\bar{\phi}$ is corrected over the first few events after the shift, and V decreases most substantially when the estimate $\bar{\phi}$ is corroborated by a strongly expected event happening at the appropriate estimated phase.

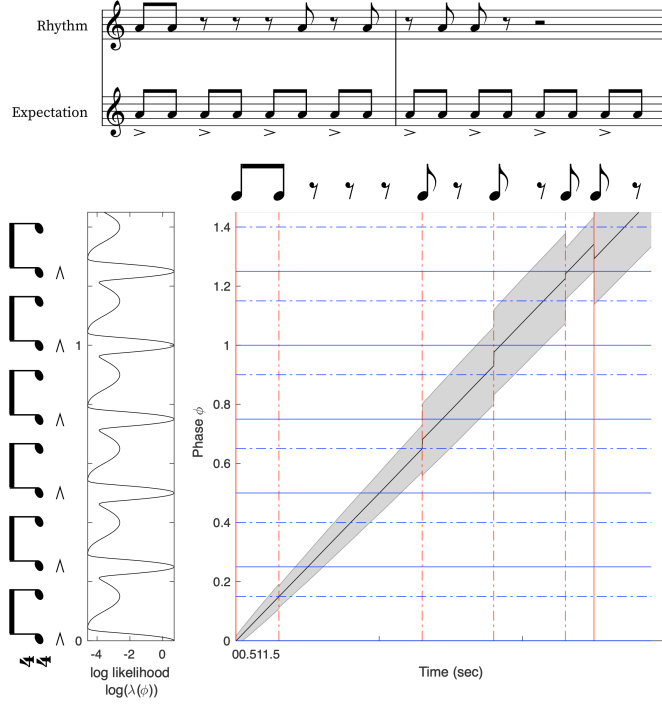


Figure 4: **Too much syncopation causes rhythm tracking failure.** Syncopation combined with imprecise and weak timing expectations on at weak time points can lead to a failure to track phase accurately. In this example, phase uncertainty V increases over a long silence. At the next event, this high uncertainty leads the model to partially attribute a weakly expected event to the nearby phase at which an event is strongly expected. As a result, the model ends up aligning the fifth event with a strong phase rather than a weak one.

333 certainty V was only slightly reduced when events occurred at weakly expected
 334 phases, so it accumulated over the course of the rhythm, and especially during
 335 the long silence. Once V was large, strongly expected event phases ϕ_i began
 336 to exert more influence at each event, until eventually events that should have
 337 been attributed to weak phase points were instead attributed primarily to adja-
 338 cent strong phase points. This type of attribution error in syncopated rhythm
 339 perception is described in [36].

340 3.4 In the absence of events: time warping

341 When an event is strongly expected but no event occurs, an optimal Bayesian
342 observer should initially be biased to believe that in spite of their current esti-
343 mate, the stimulus may not have reached the expected event phase yet. When
344 we stimulated PIPPET with sufficiently strong metronomic expectations by
345 scaling up τ , PIPPET’s behavior at each event was unchanged; however, when
346 strongly expected events were omitted, the mean phase estimate slowed down
347 at each expected event phase, leading to an overall slowing in estimated phase
348 advance (Figure 5).

349 There is evidence of such an effect in human perception. The “filled dura-
350 tion” illusion is the impression that an isochronous sequence has changed tempo
351 when it is initially subdivided by additional predictable events and then sub-
352 divisions are eliminated. According to multiple reports, the magnitude of this
353 effect is reduced or eliminated if the empty intervals precede the filled intervals
354 [37, 38, 39, 40] (though there is some disagreement about this [41]), suggesting
355 that the established expectation of continuing subdivision interferes with per-
356 ceived timing when subdivisions cease. In PIPPET, this effect is created when
357 the slowing of phase advance causes a properly timed event at the end of the
358 empty interval to arrive at an earlier apparent phase than expected, causing the
359 interval to “seem” shorter.

360 A second result that could similarly be accounted for by this aspect of PIP-
361 PET is the surprising finding in [42] that a participant tapping along with a
362 subdivided beat delays their tap following the omission of an expected subdivi-
363 sion. If taps are planned to coincide with the arrival of a specific mean estimated
364 phase, then the slowing of phase induced by an omission of a strongly expected
365 event in PIPPET would delay the subsequent tap.

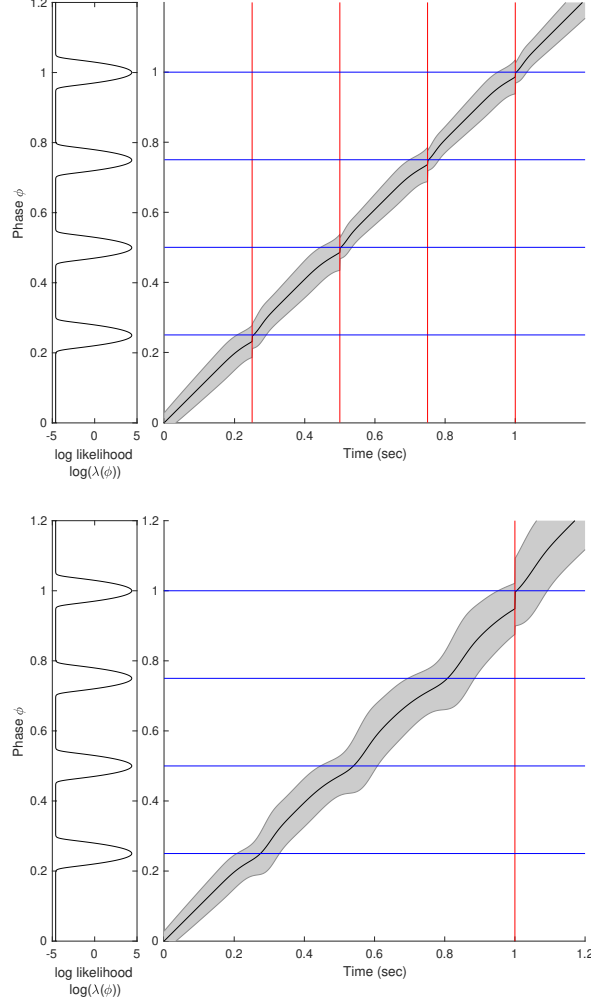


Figure 5: Time warping by the omission of strongly expected events. Black curve tracks the estimated mean phase $\bar{\phi}$ over time. Red lines mark event times; blue lines mark expected event phases. Grey shading represents uncertainty about phase, quantified in the model as variance *Sigma* and displayed by shading two standard deviations up and down. PIPPET is given strong expectations for four isochronous events. Above: when the strongly expected events occur as expected, mean phase stays on track, advancing (on average) at a rate of 1. Below: the first three expected events are omitted. When the strongly expected events do not occur, the advance of $\bar{\phi}$ slows around the expected event phase and then speeds back up. On average over the interval, $\bar{\phi}$ advances at a rate slower than 1. As a results, when the fourth event does occur at time $t = 1$, it occurs when $\bar{\phi}_t$ is still substantially short of $\bar{\phi} = 1$. The event is thus perceived as occurring at an earlier phase than expected.

366 3.5 Tempo inference

367 We simulated the PATIPPET filter with basic metronomic expectations to ob-
368 serve its capacity to infer phase and tempo at once. We gave the model a wide
369 initial range of possible tempi and a simple metronomic stimulus with actual
370 tempo near the upper end of that range. In these conditions and with the pa-
371 rameter set we chose, the model established the appropriate tempo and phase
372 to within a tight range over the course of the first two events (Figure 6).

373 In addition to its value as a model of human rhythmic cognition, the PATIP-
374 PET filter shows promise as a general-purpose tempo tracking algorithm for
375 musical applications. This would require a principled method of choosing val-
376 ues for the various free parameters of the generative model, which might be
377 done a priori based on a labeled corpus, adaptively over the course of listening,
378 or through some combination of the two. We leave a more thorough exploration
379 of the relative performance of this model to future work.

380 3.6 Period-dependent corrections

381 In entrainment literature, finger taps entrained to a metronome generally shift
382 to correct a certain fraction of an event timing perturbation on the next tap.
383 This fraction is called α . In human subjects, α has repeatedly been observed
384 to increase linearly with metronome period (“inter-onset interval,” or IOI), ex-
385 ceeding 1 (i.e., over-correction) for sufficiently long IOIs [43, 44].

386 The PIPPET framework offers a principled explanation for α increasing
387 with IOI. During an event-free interval, phase uncertainty increases over time.
388 When an event does occur, the precision of the prior distribution on phase and
389 tempo is weighed against the precision of the likelihood function associated with
390 the expectation of that event. If the prior is less precise due to accumulated
391 uncertainty, the precision of the likelihood weighs more heavily against it and

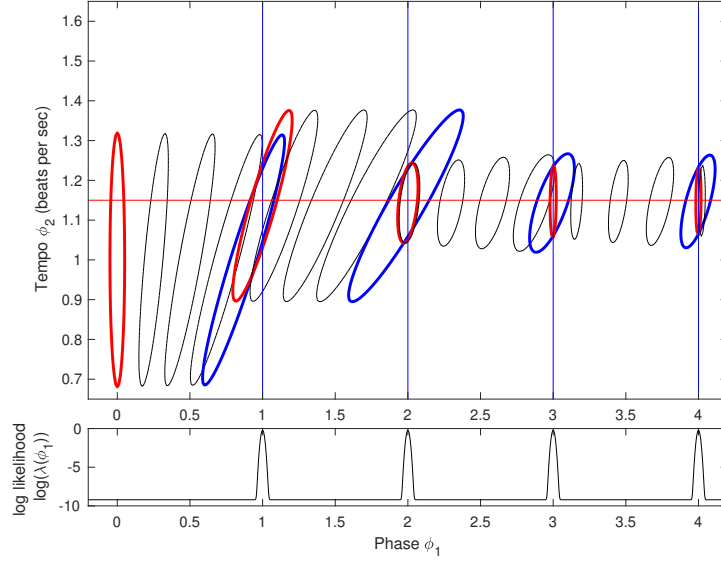


Figure 6: **A point process Kalman-Bucy Filter estimates phase and tempo.** Ellipses trace the contours of the Gaussian posterior distributions on phase and tempo. Black ellipses show a strobed visualization of the evolution of the posterior between events. Blue ellipses are the posterior distributions just before each event, and red ellipses are the posterior distributions just after each event. Here, PATIPPET is initialized with a high variance in its estimate of tempo. The first event occurs relatively early, causing the posterior mean tempo θ to increase. Each subsequent event occurs close to the time expected based on the mean estimated phase $\tilde{\phi}$ and tempo $\tilde{\theta}$, causing, the posterior to contract in both the phase and variance direction as its prediction of event time is fulfilled and its phase and tempo estimates are corroborated. Ultimately, PATIPPET settles on a narrow distribution around the appropriate tempo as it continues to accurately estimate phase.

the adjustment in phase is more thorough. Thus, all else being equal, events spaced more widely apart in time induce more extensive phase corrections.

Since the strongest phase correction PIPPET can make at an event is to fully update the phase estimate to the expected event time, it cannot account for α values above 1. However, it has been previously suggested that α may exceed 1 for long metronome periods due to some period correction occurring in addition to phase correction [43]. We were therefore curious to see whether PATIPPET could reproduce the linear increase of α with increasing IOI up to and beyond $\alpha = 1$.

In Figure 7, we show that with appropriate parameters, PATIPPET can indeed reproduce the experimental observation of a linear increase in α from below to above 1 as IOI increases. In PATIPPET, this phenomenon is a natural consequence of optimal inference in the context of phase and tempo uncertainty that accumulates between observations.

3.7 Multiple event streams

Multi-PIPPET generalizes the PIPPET/PATIPPET framework to cases of multiple distinguishable event types, each with its own set of expectations as a function of phase. One example could be listening, tapping, or dancing to a kit drum track with bass drum, snare, and hi-hat cymbal. Timing perturbations of different instruments in drum rhythms have been shown to differently affect human entrainment [45]. By letting j take values from $\{bass, snare, hihat\}$ and choosing appropriate values for ϕ_i^j , v_i^j , and τ_i^j for each event i on the metrical grid, we can create a set of timing expectations with strength and precision dependent on the specific drum and metrical position that could then be used to optimally track underlying phase and tempo through a complex kit drum rhythm. We illustrate such a template in Figure 8. A similar setup could be

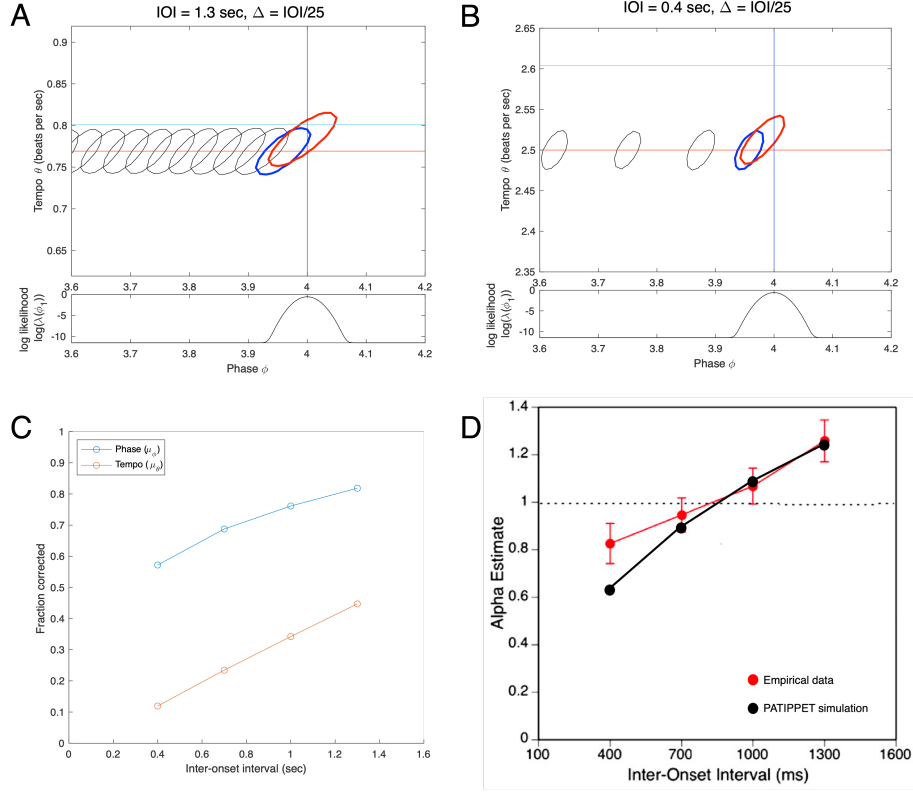


Figure 7: **PATIPPET reproduces human tapping data showing over-correction after timing perturbations to slow metronomes.** A and B) The distribution on phase and tempo leading up to and following a phase shift at the fourth event in an isochronous sequence for two different metronome tempi (i.e., two different inter-onset intervals). See Figure 6 for color key. Note that when the IOI is short, PATIPPET arrives at the phase-shifted event with a high degree of phase and tempo certainty. C) PATIPPET makes a proportionally larger correction to phase and tempo for long IOIs than for short IOIs due to the greater degree of uncertainty preceding each event. D) Alpha (α) is the proportion of a phase shift that is corrected at the next tap time. With this set of parameters, PATIPPET reproduces the empirical observation from [44] that the phase shift is undercorrected when IOIs are short and overcorrected $\alpha > 1$ when IOIs are long.

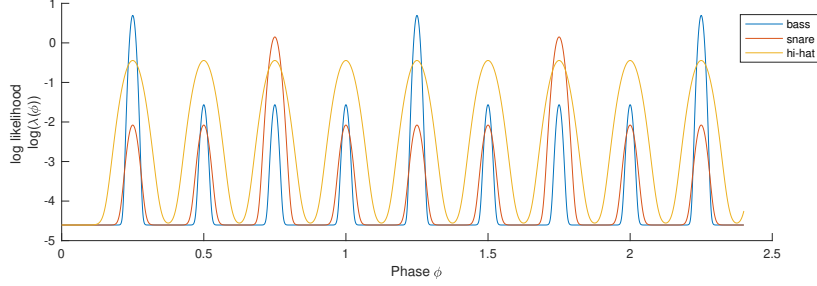


Figure 8: **Example expectation template for a basic rock beat.** In this illustration, bass drum hits are expected more strongly on the first of each cycle of four eighth notes, and are expected with high timing precision such that misplaced bass drum hits will exert a strong influence on phase. Snare drum hits are expected more strongly on the third eighth note of each cycle, and are expected with higher variance such that a misplaced snare hit exerts less influence on estimated phase. Hi-hat hits are evenly expected across all eighth note positions, but they are expected with low precision, so misplaced hi-hat hits will not exert a strong influence on estimated phase.

used to implement the assumption that pitches in a melody match the harmonic context more often in strong metrical positions, allowing event attribution and timing correction during melody listening to be influenced by scale degree.

Multi-PIPPET with $j \rightarrow \infty$ can be used to account for a continuum of event types. Thus, we could create a forward model in which it is more likely for notes played with stronger accents to fall on strong beats, or in which lower pitches are expected with higher timing precision and therefore exert greater influence on synchronization (as observed in [46]).

Multi-PIPPET could also be useful in flexibly modeling tapping data. Experiments have shown that the presence of entrained tapping prior to temporal perturbations in a metronomic stimulus reduces the phase correction response [47], indicating that the estimate of moment-by-moment phase is influenced by the proprioceptive and auditory feedback from tapping. Given working assumptions about how taps are planned and executed based on an underlying phase estimate, the taps themselves could provide a second stream of input to the

ongoing phase estimation that would bias it toward making smaller corrections to timing perturbations.

Importantly, using tap times to inform an estimate of underlying phase challenges our interpretation of this phase representing a purely external source of temporally patterned events. Instead, the inferred phase would be a hybrid of an external phase and the phase of one’s own motor cycle. Functionally, this is similar to the perceptual oscillator forced by both an external stimulus and one’s own periodic action proposed by [48]. This may be an especially useful way to think about synchronization with another agent, where one can adopt strategies ranging from following (assigning high precision to input from the other) to leading (assigning low precision to input from the other, and possibly higher precision to self-generated events). See [49] for a discussion of such a coding strategy as a means of minimizing representational neural resources.

The PIPPET framework could be further generalized to take into consideration additional stream of continuous input. This could be visual input from watching a pendulum, auditory input from a continuously modulated sound, or proprioceptive feedback from continuous entrained motion (as opposed to discrete, timed proprioceptive feedback like tapping). This goes beyond the scope of the mathematics presented here, but is a straightforward application of results proven in [27].

4 Discussion

Here we have presented PIPPET, a framework representing entrainment to a time series of discrete events based on a template of temporal expectations. PIPPET treats the event stream as the output of a point process modulated by the state of a hidden phase variable. The PIPPET filter uses variational Bayes to continuously estimate phase and track phase uncertainty based on

459 this generative model. PATIPPET extends PIPPET to include a generative
 460 model of tempo change, and the PATIPPET filter simultaneously estimates
 461 phase, tempo, and the covariance matrix representing their uncertainty and
 462 their codependence. This framework is intended to serve as a hypothesis for
 463 how the human brain integrates auditory event timing to inform and update an
 464 estimate of the state and rate of an underlying temporal process.

465 Our chosen examples have been auditory rhythms based on cyclical (met-
 466 ric) patterns of temporal expectations. But PIPPET is sufficiently general to
 467 describe entrainment based on non-isochronous and even aperiodic temporal
 468 expectations, an area that has been largely neglected in entrainment model-
 469 ing. Further, it can describe the integration of multiple event streams into an
 470 entrainment process, each with its own associated timing expectations.

471 PIPPET and PATIPPET reproduce several qualitative features of human
 472 entrainment, including realistic failures to track overly perfectly-timed but over-
 473 syncopated rhythms, perceived acceleration of a metronomic pulse when strongly
 474 expected events are omitted, and error correction after metronome timing per-
 475 turbations that increases with increasing inter-onset interval. We show that
 476 these phenomena all follow naturally from our framing of entrainment as a pro-
 477 cess of Bayesian inference based on specific phase-based temporal expectations.

478 4.1 Relationship to other models of timing

479 The dynamics of PIPPET and PATIPPET in response to sensory events are
 480 similar to dynamics of other entrainment models that correct phase and period
 481 based on event timing, e.g., [50, 51]. Models based on dynamic attending the-
 482 ory, e.g., [11, 12], are also similar in explicitly modeling timing expectations
 483 and their effect on phase and period adjustment. Our frameworks differ from
 484 these in three key ways. First, they are derived as optimal solutions to specific

inference problems, and therefore all modeling decisions can be justified within a normative framework. Second, they explicitly track uncertainty in phase and tempo – without this feature, they would not account for observed dependence of phase shift response on inter-onset interval or mimic human failures to track overly-syncopated rhythms. Finally, they allow expectations to influence the inferred phase even in the absence of sensory events, creating the time-warping effect of disappointed expectations evidenced in humans by the “filled duration” illusion.

Bayesian methods have been used elsewhere to analyze rhythmic structure as time series of point events. Some of these are application-focused methods that require offline analyses [52, 53] and therefore do not serve as satisfying models of real-time behavior. Cemgil et al (2000) [31] use a Kalman filter that tracks a distribution on phase and tempo similarly to PATIPPET. However, this model is structured to infer phase and tempo event-by-event rather than in continuous time, and is not equipped to handle stochastic rhythms or temporal structures more complex than approximate isochrony.

Bayesian inference has also been used to model timing estimation in the brain (e.g., [24, 25]), but it is generally used to describe inferences about discrete variables like interval durations and event times, whereas PIPPET describes a continuous inference process underlying predictions about event times. One such model leading to particularly PIPPET-like results was presented in Elliot et al 2014 [26]. The authors created a Bayesian model to explain the results of an experiment that had participants tap along to a stimulus consisting of two jittered metronomes. The model behaves similarly to PIPPET in that it estimates the next event time using a weighted average of previous event times and prior beliefs, with weights informed by expected timing precision. However, like [31], their model infers the anticipated timing of discrete, metronomic events,

512 whereas PIPPET predicts and updates an underlying phase in continuous time
 513 and can therefore generalize to non-isochronous and stochastic rhythms and ac-
 514 count for the effects of event omissions. Additionally, in order to account for
 515 participants ignoring events far from predicted time points, they introduce the
 516 assumption that participants repeatedly test the hypotheses that events come
 517 from one or two separate streams, whereas PIPPET naturally accounts for this
 518 phenomenon by attributing stray events to a background event rate τ_0 .

519 **4.2 Motor, perceptual, and neural entrainment**

520 Throughout this work, we have made mention of perceptual and motor expres-
 521 sions of entrainment, but have remained agnostic as to how we would expect
 522 to observe an expression of phase and tempo inference in humans. These two
 523 readouts sometimes give conflicting results: for example, exposure to musical
 524 performance with expressively irregular timing affects perceptual reports of tim-
 525 ing in subsequent stimuli [54], but does not affect phase correction in tapping
 526 to subsequent stimuli [55].

527 We expect that both physical entrainment and perceptual report are in-
 528 formed by a neural process of estimating underlying phase. Further, principles
 529 of economy suggest that they should share in such an estimate rather than draw-
 530 ing on separately instantiated processes of neural inference. However, neither
 531 motor nor perceptual experiments will necessarily give a straightforward readout
 532 of this inference process. Both readouts may be affected by independent sources
 533 of additional noise, and also potential biases: certain perceptual responses may
 534 be implicitly considered less likely than others, and certain motor errors may be
 535 implicitly considered more costly than others. Thus, an attempt at a normative
 536 Bayesian model at a specific task should be prepared to take into account this
 537 additional layer of complexity.

538 4.3 PIPPET in the brain

539 If the brain is indeed performing an optimal estimation of phase and tempo,
540 then this estimate should be legible in neural activity somewhere in the brain.
541 At the scalp level and in intracortical electrodes, slow electrical oscillations do
542 seem to anticipatorily track the structure of periodic auditory stimuli [56, 57],
543 and this tracking is associated with the subjective passage of time [58]; these os-
544 cillations could be explored as possible estimates of mean underlying phase. In
545 monkeys, the supplementary motor area appears to track the phase underlying
546 periodic visual events [59]; recordings from this region could be another candi-
547 date for reading out mean phase. Nigrostriatal dopaminergic signaling has been
548 identified as a possible marker of timing certainty [60, 61], so those dopaminer-
549 gic populations might be a good place to look for a readout of phase variance.
550 The temporal expectation template is a hazard function, and may therefore be
551 observable by using techniques recently applied to decode the temporal hazard
552 function from EEG data [62], or through its correlation with beta oscillations
553 [63].

554 Though PIPPET and PATIPPET are not committed to a particular brain-
555 based implementation, advances in the brain basis of timing and beat-keeping
556 combined with the hypothesized neural bases of predictive processing suggest
557 the beginnings of a plausible implementation of PIPPET in the brain. A de-
558 tailed discussion of a possible neural basis of beat maintenance is presented in
559 [64]. Briefly, supplementary motor area may maintain an ongoing estimate of
560 mean phase through some combination of intrinsic dynamics and interaction
561 with the basal ganglia, while dopaminergic signaling in striatum may maintain
562 an estimate of phase uncertainty. The phase estimate may be used to inform
563 auditory timing expectancy via learned models in premotor cortex [65]. These
564 expectations may be delivered to the early stages of audition via the top-down

connections along the dorsal auditory pathway, where they can be used to evaluate timing prediction error [66]. These errors, weighted by their precisions, may be transmitted back to the supplementary motor area via the bottom-up connectivity of the dorsal auditory pathway and used to update the estimate of phase.

4.4 Learning and inference outside of PIPPET

If the brain does treat entrainment as a process of inference based on a generative model, this raises the question of how the properties of the generative model are established in the first place. The PIPPET framework does not address this question directly, but by examining the parameters necessary to formulate PIPPET, we can clearly see what components need to be in place before a process of continuous phase and tempo updating can begin.

First, the brain must learn the temporal structures of the expectation template for rhythmic expectation. Learning these underlying structures from an experiential corpus of noisy, stochastic rhythms is not trivial. It seems likely to involve some type of bootstrapping in which a recognition of some degree of temporal structure allows for attribution of events to positions in that structure, allowing for deeper structure learning. Earlier exposure to simpler, less stochastic rhythms would likely help with such a bootstrapping process. For a discussion of the challenges of this type of simultaneous learning and filtering and a proposed solution for non-point-process data, see [67].

The brain must also learn noise and precision parameters for the model. Note that neither the temporal expectation variance parameters v_i nor the noise parameters σ and σ_θ necessarily correspond to the actual precision of the neural or external timing mechanisms in play. The brain may underestimate the noisiness (σ) of the timing process it uses to track underlying phase, leading to under-

adjustment to auditory event timing and minimal time-warping between events,
or do the opposite. Presumably, these parameters must be learned through ex-
perience and prediction error.

The precision parameters v_i may be informed by several factors. First, an
upper bound on the precision of expected event timing is the precision of sensory
timing perception, which is, for example, high for human audition and signifi-
cantly lower for human vision¹. Second, expected event timing precision may
also be informed by the observed relative timing distributions of event streams.
These observations may inform expectations on time scales ranging from a single
sitting to a lifetime of listening. Expected timing may be learned separately for
different sensory modalities, different musical genres (e.g., techno vs. funk), or
even different instruments (e.g., kick drum, snare, hi-hat, as discussed above).
The precision of a beat-based temporal expectation is closely related to the
width of a “beat bin,” the window of time (rather than a single time point) that
is proposed to constitute the “beat” in [68], and to the width of the temporal
“expectancy region” described in dynamic attending theory [11]; in both cases,
this width is increased by imprecision in the immediately preceding stimulus.

When the brain is exposed to a rhythmic stimulus, it must first recognize
that a predictable pattern exists and select an appropriate temporal expectation
template from its learned repertoire. This is its own process of inference, and
may be amenable to a Bayesian description. Since the PIPPET filter maintains
a unimodal posterior, it is not well-suited to model this initial inference process,
which may require maintaining a distribution over multiple distinct possible
starting phases and temporal expectation templates. This problem might be

¹An event can only be experienced after it occurs, so (as pointed out in [25]) the likelihood function on underlying phase associated with this type of uncertainty should be asymmetrical. The analytically tractable incarnation of our framework presented here uses Gaussian likelihood peaks, so cannot account for the effect of asymmetrical likelihoods; however, we could posit a τ function with asymmetrical peaks and use numerical methods rather than the explicit solution derived here to estimate underlying phase at each time step.

615 partially addressed at a modeling level by incorporating a model of meter in-
616 ference based on prior probabilities of hearing specific meters at specific tempi,
617 e.g. [69], as an additional level of inference in parallel with phase and tempo
618 inference.

619 Finally, aspects of the temporal expectation template are likely changing
620 even as a rhythm plays out in time. This is evidenced by the grammar-like
621 structure of music rhythm [70]: certain patterns of events are more expected
622 than others regardless of their metrical positions. PIPPET and PATIPPET take
623 a template of expected event time points as an input, and thus do not take into
624 account immediate stimulus history in creating expectations. However, such
625 effects could be incorporated into a model based on this framework by adding
626 a history dependence to the expectation template τ . The precise details of this
627 history dependence could be based on any suitable formal model for rhythmic
628 grammar (e.g., [71, 72, 70]).

629 4.5 Future directions

630 In evaluating future directions, it is important to be clear that PIPPET and
631 PATIPPET are not “models” but “frameworks.” Directly testing their validity
632 as models of human behavior would require setting values for many free pa-
633 rameters, and it is not yet clear to what extent the parameters of individual
634 expected events should be based on empirical data collected over a lifetime or
635 empirical data collected trial by trial.

636 However, there is a certain extent to which these frameworks can be vali-
637 dated as descriptions of human cognition. First, these models predict certain
638 qualitative effects such as the slowing of perceived phase advance as strong ex-
639 pectations are disappointed. Second, although the parameters in the forward
640 models are not directly empirically measurable values, changes in stimulus his-

641 tory should influence them in predictable ways. For example, if a certain type
 642 of event occurs consistently at a particular metrical position within an extended
 643 stimulus presentation or within the music the listener has experienced in a life-
 644 time of listening, then it should induce stronger phase corrections than an event
 645 that occurs inconsistently as if it has been given a higher value of τ_i . Param-
 646 eters may also be influenced by long term listening experience, but they should
 647 at least respond to recent empirical experience by changing in the direction
 648 predicted by PIPPET.

649 If we find situations in which human behavior differs from solutions to the
 650 inference problems posed by PIPPET and PATIPPET, this suggests that the
 651 tasks being performed in those situations are being performed with a different
 652 objective than optimal inference of phase and tempo based on these generative
 653 models. In this case, we would be challenged to articulate the true nature of
 654 the problem being solved. This might require modifications of the generative
 655 model, e.g., introducing the belief that tempo changes occur in jumps or ramps
 656 rather than as random drift, or modification of the objective of the task, e.g., by
 657 including additional cost functions or priors associated with perceptual report
 658 or motor output as discussed above.

659 Once we are satisfied with the PIPPET framework’s utility in describing
 660 to human behavior, we can use it to model and analyze experimental data.
 661 Given a perceptual or behavioral task, we can suppose that motor or perceptual
 662 human entrainment behavior is optimally solving an inference problem, and
 663 determine the parameters of that problem by fitting them with appropriate
 664 methods. We can study the changes in these parameters over the course of an
 665 experiment, over different variations on the same experiment, over the human
 666 lifespan, across cultures, etc. This approach could add an additional level of
 667 insight to the analysis of a wide range of timing tasks.

668 One specific question that the PIPPET framework might help resolve is how
669 periodic and nonperiodic entrainment differ. PIPPET has no specific machinery
670 to account for ways in which the two situations differ (for neural and behavioral
671 evidence of differences between memory-based and periodicity based entrainment,
672 see, e.g., [14, 6]. However, since it is sufficiently general to model both, it could
673 guide an exploration of parameter differences between the performance of similar
674 tasks in periodic and aperiodic contexts.

675 We can also let the PIPPET framework guide a search for the brain bases
676 of entrainment. Even if perceptual and motor outputs are subject to different
677 biases and costs, they would both be well-served by an optimal estimate of a
678 ground truth, so there is reason to expect to find such an estimate represented in
679 the brain. Such a search could proceed by looking for covariates for PIPPET’s
680 phase and uncertainty estimates in neural data during the performance of tasks
681 that require non-trivial updating of these estimates.

682 Finally, the PIPPET framework can serve as a cog in larger predictive pro-
683 cessing models. The generative models we describe here allow for the evaluation
684 of joint and marginal distributions on specific timing patterns and hidden states
685 underlying them. By introducing additional levels of hidden states and addi-
686 tional sources of sensory input, we can create Bayesian inference models that
687 use event timing to infer higher-order contextual states, e.g. meter, and predict
688 other aspects of sensory input, e.g. pitch, creating a unified picture of human
689 musical expectation.

690 5 Acknowledgments

691 Thanks to Tom Kaplan for extensive discussions and insights motivating this
692 manuscript, and to Darren Rhodes and Nori Jacoby for helpful feedback.

6 Appendix

6.1 The PATIPPET filter

We let $\mathbf{x} = \begin{pmatrix} \phi \\ \theta \end{pmatrix}$ denote the posterior mean and $\mathbf{\Sigma} = \begin{pmatrix} V & \Sigma^{12} \\ \Sigma^{21} & \Sigma^{22} \end{pmatrix}$ denote the posterior covariance. The expressions for the evolution of the PATIPPET filter, which we derive in the following section, are:

$$\begin{cases} d\bar{\mathbf{x}}_t = \begin{pmatrix} \bar{\theta} \\ 0 \end{pmatrix} dt + (\hat{\mathbf{x}} - \bar{\mathbf{x}}_{t-}) \cdot (dN_t - \hat{\Lambda} dt) \\ d\mathbf{\Sigma} = \begin{pmatrix} 2\Sigma^{12} + \sigma_\phi^2 & \Sigma^{22} \\ \Sigma^{22} & \sigma_\theta^2 \end{pmatrix} dt + (\hat{\mathbf{\Sigma}} - \mathbf{\Sigma}_{t-}) \cdot (dN_t - \hat{\Lambda} dt) \end{cases} \quad (9)$$

where we define

$$\begin{cases} \hat{\Lambda} := \sum_{i=0,1,\dots} T_i \hat{\theta}_i \\ \hat{\mathbf{x}} = \frac{1}{\hat{\Lambda}} \sum_{i=0,1,\dots} T_i \begin{pmatrix} K_i^{12} + \hat{\phi}_i \hat{\theta}_i \\ K_i^{22} + \hat{\theta}_i^2 \end{pmatrix} \\ \hat{\mathbf{\Sigma}} := \frac{1}{\hat{\Lambda}} \sum_{i=0,1,\dots} T_i \left(\hat{\theta}_i \mathbf{K}_i + \hat{\theta}_i (\hat{\mathbf{x}}_i - \bar{\mathbf{x}}_{t+}) (\hat{\mathbf{x}}_i - \bar{\mathbf{x}}_{t+})^T \right. \\ \left. + (\hat{\mathbf{x}}_i - \bar{\mathbf{x}}_{t+}) \begin{pmatrix} K_i^{21} & K_i^{22} \end{pmatrix} + \begin{pmatrix} K_i^{12} \\ K_i^{22} \end{pmatrix} (\hat{\mathbf{x}}_i - \bar{\mathbf{x}}_{t+})^T \right) \end{cases}$$

$$\mathbf{K}_0 := \mathbf{\Sigma}, \mathbf{K}_i := (\mathbf{P}_i + \mathbf{\Sigma}^{-1})^{-1} \text{ for } i > 0.$$

K_i^{kl} denotes the entries in \mathbf{K}_i .

$$T_0 := \tau_0, T_i := \tau_i N(\phi_i | \bar{\phi}, v_i^{-1} + V^{-1}) \text{ for } i > 0.$$

$$\hat{\mathbf{x}}_i = \begin{pmatrix} \hat{\phi}_i \\ \hat{\theta}_i \end{pmatrix} := \mathbf{K}_i (\mathbf{P}_i \mathbf{x}_i + \mathbf{\Sigma}^{-1} \bar{\mathbf{x}}) \text{ for } i > 0, \text{ and } \hat{\mathbf{x}}_0 := \bar{\mathbf{x}}.$$

$$P_i := \begin{pmatrix} v_i^{-1} & 0 \\ 0 & 0 \end{pmatrix}$$

6.2 Derivation of differential equations and update equations.

We derive the PATIPPET filter first, and then derive the PIPPET filter as a special case.

Snyder [27] provides a partial differential equation describing the evolution of a probability distribution on a continuously stochastically evolving state that drives the emission of point process events. If the evolution of the underlying state is described by a Gauss-Markov diffusion process:

$$d\mathbf{x} = \mathbf{A}\mathbf{x}dt + \mathbf{B}d\mathbf{W}_t \quad (10)$$

and events are generated at rate $\lambda(\mathbf{x})$, then the evolution of the probability distribution $p_t(\mathbf{x})$ is described by

$$dp_t(\mathbf{x}) = \mathcal{L}[p_{t-}(\mathbf{x})]dt + p_{t-}(\mathbf{x}) \left(\frac{\lambda(\mathbf{x})}{\hat{\Lambda}} - 1 \right) \cdot (dN_t - \hat{\Lambda}dt) \quad (11)$$

where $\hat{\Lambda} := \mathbb{E}[\lambda(\mathbf{x})]$ (with \mathbb{E} denoting expectation under distribution $p_{t-}(\mathbf{x})$), dN_t is the increment in the event count over each dt time step (assumed to be either 1 or 0 with probability 1), and \mathcal{L} is the Kolmogorov forward operator associated with (10):

$$\mathcal{L}[p(\mathbf{x})] = - \sum_i \frac{\partial}{\partial x_i} [\mathbf{A}\mathbf{x}]_i p(\mathbf{x}) + \frac{1}{2} \sum_{i,j} \frac{\partial^2}{\partial x_i \partial x_j} [\mathbf{B}\mathbf{B}^T]_{ij} p(\mathbf{x}) \quad (12)$$

Here we project p onto a Gaussian distribution at each time step by matching mean $\bar{\mathbf{x}}$ and covariance Σ , which is also the projection with minimal KL

720 divergence. We do this by finding the differentials of these moments of p_t and
 721 using them to drive the evolution of these two variables:

$$\begin{aligned}
 d\bar{\mathbf{x}}_t &= \bar{\mathbf{x}}_{t+} - \bar{\mathbf{x}}_{t-} = \int_{\mathbf{x}} \mathbf{x} p_{t+}(\mathbf{x}) d\mathbf{x} - \int_{\mathbf{x}} \mathbf{x} p_{t-}(\mathbf{x}) d\mathbf{x} \\
 &= \int_{\mathbf{x}} \mathbf{x} (p_{t+}(\mathbf{x}) - p_{t-}(\mathbf{x})) d\mathbf{x} = \int_{\mathbf{x}} \mathbf{x} dp_t(\mathbf{x}) d\mathbf{x} \\
 &= \int_{\mathbf{x}} \mathbf{x} \mathcal{L}[p_t(\mathbf{x})] dt d\mathbf{x} + (\hat{\mathbf{x}} - \bar{\mathbf{x}}_{t-}) \cdot (dN_t - \hat{\Lambda} dt) \quad (13)
 \end{aligned}$$

where we define $\hat{\mathbf{x}} := \mathbb{E}[\mathbf{x} \lambda(\mathbf{x})]$, and

$$d\mathbf{\Sigma}_t = \mathbf{\Sigma}_{t+} - \mathbf{\Sigma}_{t-} = \int_{\mathbf{x}} [[\mathbf{x} - \bar{\mathbf{x}}_{t+}]^2] p_{t+}(\mathbf{x}) d\mathbf{x} - \int_{\mathbf{x}} [[\mathbf{x} - \bar{\mathbf{x}}_{t-}]^2] p_{t-}(\mathbf{x}) d\mathbf{x}$$

where $[[\mathbf{x}]]^2$ denotes $\mathbf{x} \mathbf{x}^T$.

$$\begin{aligned}
 d\mathbf{\Sigma}_t &= \int_{\mathbf{x}} [[\mathbf{x} - \bar{\mathbf{x}}_{t+}]^2] (p_{t+}(\mathbf{x}) - p_{t-}(\mathbf{x})) d\mathbf{x} \\
 &\quad + \int_{\mathbf{x}} ([[\mathbf{x} - \bar{\mathbf{x}}_{t+}]^2] - [[\mathbf{x} - \bar{\mathbf{x}}_{t-}]^2]) p_{t-}(\mathbf{x}) d\mathbf{x} \\
 &= \int_{\mathbf{x}} [[\mathbf{x} - \bar{\mathbf{x}}_{t+}]^2] dp_t(\mathbf{x}) - [[\bar{\mathbf{x}}_{t+} - \bar{\mathbf{x}}_{t-}]^2] \\
 &= \int_{\mathbf{x}} [[\mathbf{x} - \bar{\mathbf{x}}_{t+}]^2] \mathcal{L}[p_t(\mathbf{x}|N_t)] dt d\mathbf{x} + (\hat{\mathbf{\Sigma}} - \mathbf{\Sigma}_{t-}) \cdot (dN_t - \hat{\Lambda} dt) \quad (14)
 \end{aligned}$$

722 where we define $\hat{\mathbf{\Sigma}} := \mathbb{E}[[[\mathbf{x} - \bar{\mathbf{x}}_{t+}]^2] \lambda(\mathbf{x})]$.

723 Integrating by parts (or following [30]), we can calculate the appropriate inte-
 724 grals of $\mathcal{L}[p_t(\mathbf{x}|N_t)]$, arriving at a general expression for the variational Bayesian
 725 filter for point process data:

$$\begin{cases} d\bar{\mathbf{x}}_t = \mathbf{A}\bar{\mathbf{x}}_{t-}dt + (\hat{\mathbf{x}} - \bar{\mathbf{x}}_{t-}) \cdot (dN_t - \hat{\Lambda}dt) \\ d\bar{\mathbf{\Sigma}}_t = (\mathbf{A}\bar{\mathbf{\Sigma}}_{t-} + \bar{\mathbf{\Sigma}}_{t-}\mathbf{A}^T + \mathbf{B}\mathbf{B}^T)dt + (\hat{\mathbf{\Sigma}} - \bar{\mathbf{\Sigma}}_{t-}) \cdot (dN_t - \hat{\Lambda}dt) \end{cases} \quad (15)$$

From (4), the PATIPPET generative model is described by the Gauss-Markov diffusion process (10) with

$$\mathbf{x} = \begin{pmatrix} \phi \\ \theta \end{pmatrix} \text{ and } \bar{\mathbf{x}} = \begin{pmatrix} \bar{\phi} \\ \bar{\theta} \end{pmatrix}$$

$$\bar{\mathbf{\Sigma}} = \begin{pmatrix} V & \Sigma^{12} \\ \Sigma^{21} & \Sigma^{22} \end{pmatrix}$$

$$\mathbf{A} := \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \text{ and } \mathbf{B} := \begin{pmatrix} \sigma_\phi & 0 \\ 0 & \sigma_\theta \end{pmatrix}.$$

Plugging into (15), we have

$$\begin{cases} d\bar{\mathbf{x}}_t = \begin{pmatrix} \bar{\theta} \\ 0 \end{pmatrix} dt + (\hat{\mathbf{x}} - \bar{\mathbf{x}}_{t-}) \cdot (dN_t - \hat{\Lambda}dt) \\ d\bar{\mathbf{\Sigma}} = \begin{pmatrix} 2\Sigma^{12} + \sigma_\phi^2 & \Sigma^{22} \\ \Sigma^{22} & \sigma_\theta^2 \end{pmatrix} dt + (\hat{\mathbf{\Sigma}} - \bar{\mathbf{\Sigma}}_{t-}) \cdot (dN_t - \hat{\Lambda}dt) \end{cases} \quad (16)$$

We complete the derivation by calculating $\hat{\Lambda}$, $\hat{\mathbf{x}}$, and $\hat{\mathbf{\Sigma}}$. This proceeds by first deriving a simple expression for $p(\mathbf{x})\tau(\mathbf{x})$ as a sum of scaled normal distributions.

Let $\|x\|_A^2$ denote $x^T A x$. We will make use of the following result, a generalized form of a well-known result about quadratic forms that allows us to write

736 products of multivariate normal distributions as normal distributions (see [73]
 737 for proof and similar application):

$$\|x - a\|_A^2 + \|x - b\|_B^2 = \|a - b\|_{A(A+B)^{-1}B}^2 + \|x - (A+B)^{-1}(Aa + Bb)\|_{A+B}^2 \quad (17)$$

738 In the PATIPPET generative model, events are generated at rate $\lambda(\mathbf{x}) =$
 739 $\theta\tau(\phi) = \theta\tau(\phi)$, where

$$\tau(\phi) = \tau_0 + \sum_{i=1,2,\dots} \frac{\tau_i}{\sqrt{2\pi v_i}} e^{-\frac{1}{2}\|\mathbf{x} - \mathbf{x}_i\|_{P_i}^2}$$

740

$$P_i = \begin{pmatrix} v_i^{-1} & 0 \\ 0 & 0 \end{pmatrix}, \quad \mathbf{x}_i = \begin{pmatrix} \phi_i \\ 0 \end{pmatrix}.$$

741 $p(\mathbf{x})$ is assumed (forced) to be Gaussian, so we can write:

$$p(\mathbf{x}) = \frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}\|\mathbf{x} - \bar{\mathbf{x}}\|_{\Sigma^{-1}}^2}.$$

We calculate:

$$\begin{aligned} p(\mathbf{x})\tau(\mathbf{x}) &= \frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}\|\mathbf{x} - \bar{\mathbf{x}}\|_{\Sigma^{-1}}^2} \left(\tau_0 + \sum_{i=1,2,\dots} \frac{\tau_i}{\sqrt{2\pi v_i}} e^{-\frac{1}{2}\|\mathbf{x} - \mathbf{x}_i\|_{P_i}^2} \right) \\ &= \frac{\tau_0}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}\|\mathbf{x} - \bar{\mathbf{x}}\|_{\Sigma^{-1}}^2} + \sum_{i=1,2,\dots} \frac{\tau_i}{2\pi\sqrt{v_i}|\Sigma|} e^{-\frac{1}{2}\|\mathbf{x} - \mathbf{x}_i\|_{P_i}^2 - \frac{1}{2}\|\mathbf{x} - \bar{\mathbf{x}}\|_{\Sigma^{-1}}^2} \end{aligned}$$

Applying (17),

$$\begin{aligned} p(\mathbf{x})\tau(\mathbf{x}) &= \frac{\tau_0}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}\|\mathbf{x} - \bar{\mathbf{x}}\|_{\Sigma^{-1}}^2} \\ &\quad + \sum_{i=1,2,\dots} \tau_i \left(\frac{1}{\sqrt{2\pi(v_i^{-1} + V^{-1})}} e^{-\frac{1}{2}\|\mathbf{x}_i - \bar{\mathbf{x}}\|_{P_i K_i \Sigma^{-1}}^2} \right) \left(\frac{1}{\sqrt{2\pi \frac{v_i|\Sigma|}{v_i^{-1} + V^{-1}}}} e^{-\frac{1}{2}\|\mathbf{x} - K_i(P_i \mathbf{x}_i + \Sigma^{-1} \bar{\mathbf{x}})\|_{K_i^{-1}}^2} \right) \end{aligned} \quad (18)$$

742 where we define $\mathbf{K}_i := (\mathbf{P}_i + \mathbf{\Sigma}^{-1})^{-1}$.

743 We next calculate:

$$\|\mathbf{x}_i - \bar{\mathbf{x}}\|_{\mathbf{P}_i \mathbf{K}_i \mathbf{\Sigma}^{-1}}^2 = \frac{(\phi_i - \phi)^2}{v_i^{-1} + V^{-1}}$$

744 and

$$|\mathbf{K}_i| = \frac{v_i |\mathbf{\Sigma}|}{v_i^{-1} + V^{-1}}$$

We use these expressions to write (18) in terms of normal distributions:

$$p(\mathbf{x})\tau(\mathbf{x}) = \tau_0 N(\mathbf{x}|\bar{\mathbf{x}}, \mathbf{\Sigma}) + \sum_{i=1,2,\dots} \tau_i N(\phi_i|\bar{\phi}, v_i^{-1} + V^{-1}) N(\mathbf{x}|\mathbf{K}_i(\mathbf{P}_i \mathbf{x}_i + \mathbf{\Sigma}^{-1} \bar{\mathbf{x}}), \mathbf{K}_i) \quad (19)$$

We simplify this expression by defining $T_i := \tau_i N(\phi_i|\bar{\phi}, v_i^{-1} + V^{-1})$ for $i > 0$, and setting $T_0 := \tau_0$ and $\mathbf{K}_0 = \mathbf{\Sigma}$. We define $\hat{\mathbf{x}}_i := \begin{pmatrix} \hat{\phi}_i \\ \hat{\theta}_i \end{pmatrix} := \mathbf{K}_i(\mathbf{P}_i \mathbf{x}_i + \mathbf{\Sigma}^{-1} \bar{\mathbf{x}})$ for $i > 0$ and set $\hat{\mathbf{x}}_0 := \bar{\mathbf{x}}$. This lets us write

$$p(\mathbf{x})\tau(\mathbf{x}) = \sum_{i=0,1,\dots} T_i N(\mathbf{x}|\hat{\mathbf{x}}_i, \mathbf{K}_i) \quad (20)$$

We use this expression and the moments of normal distributions to calculate $\hat{\Lambda}$, $\hat{\mathbf{x}}$, and $\hat{\mathbf{\Sigma}}$:

$$\hat{\Lambda} := \mathbb{E}[\lambda(\mathbf{x})] = \mathbb{E}[\theta\tau(\mathbf{x})] = \sum_{i=0,1,\dots} T_i \hat{\theta}_i \quad (21)$$

$$\hat{\mathbf{x}} := \frac{1}{\hat{\Lambda}} \mathbb{E}[\mathbf{x}\lambda(\mathbf{x})] = \frac{1}{\hat{\Lambda}} \mathbb{E}_p[\mathbf{x}\theta\tau(\mathbf{x})] = \frac{1}{\hat{\Lambda}} \int_{\mathbf{x}} \begin{pmatrix} \phi\theta \\ \theta^2 \end{pmatrix} p(\mathbf{x})\tau(\mathbf{x}) d\mathbf{x}$$

This expression picks out non-central second moment terms of each normal distributions in 20, each of which can be written in terms of the covariance matrix and mean of the distribution. Using K_i^{kl} to denote the entries in \mathbf{K}_i , we can write

$$\hat{\mathbf{x}} = \frac{1}{\hat{\Lambda}} \sum_{i=0,1,\dots} T_i \begin{pmatrix} K_i^{12} + \hat{\phi}_i \hat{\theta}_i \\ K_i^{22} + \hat{\theta}_i^2 \end{pmatrix} \quad (22)$$

The third-order expression for $\hat{\Sigma}$ can also be written in terms of covariance matrices and means since the central third moments of normal distributions are zero.

$$\begin{aligned} \hat{\Sigma} &:= \frac{1}{\hat{\Lambda}} \mathbb{E}_p [[\mathbf{x} - \bar{\mathbf{x}}_{t+}]^2 \lambda(\mathbf{x})] = \frac{1}{\hat{\Lambda}} \mathbb{E}_p [[\mathbf{x} - \bar{\mathbf{x}}_{t+}]^2 \theta \tau(\mathbf{x})] \\ &= \frac{1}{\hat{\Lambda}} \sum_{i=0,1,\dots} T_i \mathbb{E}_p [\theta [[\mathbf{x} - \bar{\mathbf{x}}_{t+}]^2 N(\mathbf{x}|\hat{\mathbf{x}}_i, \mathbf{K}_i)] \\ &= \frac{1}{\hat{\Lambda}} \sum_{i=0,1,\dots} T_i \hat{\theta}_i \mathbb{E}_p [[[\mathbf{x} - \hat{\mathbf{x}}_i]]^2 N(\mathbf{x}|\hat{\mathbf{x}}_i, \mathbf{K}_i)] \\ &\quad + T_i \hat{\theta}_i [[\hat{\mathbf{x}}_i - \bar{\mathbf{x}}_{t+}]^2 \\ &\quad + T_i (\hat{\mathbf{x}}_i - \bar{\mathbf{x}}_{t+}) \mathbb{E}_p [(\mathbf{x} - \hat{\mathbf{x}}_i)^T (\theta - \hat{\theta}_i) N(\mathbf{x}|\hat{\mathbf{x}}_i, \mathbf{K}_i)] \\ &\quad + T_i \mathbb{E}_p [(\mathbf{x} - \hat{\mathbf{x}}_i)(\theta - \hat{\theta}_i) N(\mathbf{x}|\hat{\mathbf{x}}_i, \mathbf{K}_i)] (\hat{\mathbf{x}}_i - \bar{\mathbf{x}}_{t+})^T] \end{aligned} \quad (23)$$

$$\begin{aligned} &= \frac{1}{\hat{\Lambda}} \sum_{i=0,1,\dots} T_i [\hat{\theta}_i \mathbf{K}_i + \hat{\theta}_i [[\hat{\mathbf{x}}_i - \bar{\mathbf{x}}_{t+}]^2 \\ &\quad + (\hat{\mathbf{x}}_i - \bar{\mathbf{x}}_{t+}) \begin{pmatrix} K_i^{21} & K_i^{22} \end{pmatrix} + \begin{pmatrix} K_i^{12} \\ K_i^{22} \end{pmatrix} (\hat{\mathbf{x}}_i - \bar{\mathbf{x}}_{t+})^T] \end{aligned} \quad (24)$$

$$(25)$$

745 These expressions coupled with (28) constitute the PATIPPET filter.

746 The PIPPET filter can be derived as a special case of the PATIPPET filter

by setting $\sigma_\theta = 0$, $\theta_0 = 1$, and all terms in Σ to zero except V . However, this requires finessing various degeneracies, e.g. wherever Σ is inverted. More straightforward is to follow the same process as above, starting from the PIPPET generative model (3). Either way ultimately yields the PIPPET filter (3).

For multiple event streams j ,

$$dp_t(\mathbf{x}) = \mathcal{L}[p_t(\mathbf{x})]dt + p_t(\mathbf{x}) \sum_j (\lambda_j(\phi) - \mathbb{E}_p[\lambda_j(\phi)]) \cdot (\mathbb{E}_p[\lambda_j(\phi)]^{-1} dN_j - dt) \quad (26)$$

This follows directly from application of the derivation above to equation (5) in [74] with a discrete spatial dimension. By the methods above, it yields the multi-PIPPET filter (8) and the multi-PATIPPET filter:

$$\begin{cases} d\mu = dt - \sum_j (\mu^{*j} - \mu)(dN_t^j - \hat{\Lambda}^j dt) \\ d\Sigma = \sigma^2 dt - \sum_j (\Sigma^{*j} - \Sigma)(dN_t^j - \hat{\Lambda}^j dt) \end{cases} \quad (27)$$

and the multi-PATIPPET filter:

$$\begin{cases} d\bar{\mathbf{x}}_t = \begin{pmatrix} \bar{\theta} \\ 0 \end{pmatrix} dt + \sum_j (\hat{\mathbf{x}}^j - \bar{\mathbf{x}}_{t-}) \cdot (dN_t^j - \hat{\Lambda}^j dt) \\ d\Sigma = \begin{pmatrix} 2\Sigma^{12} + \sigma_\phi^2 & \Sigma^{22} \\ \Sigma^{22} & \sigma_\theta^2 \end{pmatrix} dt + \sum_j (\hat{\Sigma}^j - \Sigma_{t-}) \cdot (dN_t^j - \hat{\Lambda}^j dt) \end{cases} \quad (28)$$

6.3 Simulation parameters.

All code used to create figures in this manuscript is available at <https://github.com/joncannon/PIPPET>.

760 PIPPET simulations were conducted by numerical simulation of (1) with
761 $dt = 0.001$ and initialized with $\mu_0 = 0$ and $\Sigma_0 = 0.0002$. Parameters for
762 the simulations shown in each figure are listed below, with t_i used to denote
763 simulated event times. (ϕ_i and t_i are given in units of seconds, and v_i is given
764 in units of s^2 .)

765 *Figure 1:* $\phi_i = t_i = \{0.5, 1, 1.5\}$, $v_i = 0.0001$, $\tau_i = 0.02$, $\tau_0 = 0.01$, $\sigma = 0.05$

766 *Figure 2A:* $\phi_i = t_i = \{0.25, 0.5, 0.75, 1\}$, $v_i = 0.0001$, $\tau_i = 2$, $\tau_0 = 0.01$,
767 $\sigma = 0.05$.

768 *Figure 2B:* Same as Figure 2A, but with $t_i = \{1\}$.

Figure 3A:

$$t_i = \{0, 0.150, 0.25, 0.5, 0.65, 0.9, 1\}$$

$$\phi_i = \{0, 0.15, 0.25, 0.4, 0.5, 0.65, 0.75, 0.9, 1, 1.15\}$$

$$v_i = \{.0001, .0005, .0001, .0005, .0001, .0005, .0001, .0005\}$$

$$\tau_i = \{.05, .01, .05, .01, .05, .01, .05, .01\}$$

$$\tau_0 = 0.01$$

$$\sigma = 0.05$$

769 *Figure 3B:* Same as Figure 3A, but with $t_i = \{0, 0.150, 0.25, 0.5, 0.61, 0.86, 0.96\}$.

770 *Figure 4:* Same as Figure 3A, but with $t_i = \{0, 0.15, .65, .9, 1.15, 1.25\}$.

Figure 5: (No numerical simulation was performed for this figure.)

$$\begin{aligned}
\phi_i^j &= 0.25i \text{ for } j = \text{bass}, \text{snare}, \text{hihat} \\
v_i^{\text{bass}} &= .0001, v_i^{\text{snare}} = .0003, v_i^{\text{hihat}} = .001 \\
\tau_i^{\text{bass}} &= \{.05, .005, .005, .005, \dots\} \\
\tau_i^{\text{snare}} &= \{.005, .005, .05, .005, \dots\} \\
\tau_i^{\text{hihat}} &= \{.05, .05, .05, .05, \dots\} \\
\tau_0 &= 0.01
\end{aligned}$$

771 PATIPPET simulations were conducted by numerical simulation of (4) with
772 $dt = 0.001$. Parameters for the simulations shown in each figure are listed below.

Figure 6:

$$\begin{aligned}
t_i &= \frac{i}{1.15} \\
\phi_i &= i \\
v_i &= \{.0001, .0003, .0001, .0003, .0001, .0003, .0001, .0003\} \\
\tau_i &= \{.02, .01, .02, .01, .02, .01, .02, .01\} \\
\tau_0 &= 10^{-4} \\
\sigma &= 0.05 \\
\sigma_\theta &= 0.05 \\
\mu_0 &= \begin{pmatrix} 0 \\ 1 \end{pmatrix} \\
\Sigma_0 &= \begin{pmatrix} .001 & 0 \\ 0 & .04 \end{pmatrix}
\end{aligned}$$

Figure 7: In four simulations, we set the inter-onset interval Δ to $0.4s$, 0 , $7s$,

1.0s, and 1.3s. In each simulation, we set the perturbation δ to $\frac{\Delta}{25}$.

$$t_i = \{\Delta, 2\Delta, 3\Delta, 4\Delta + \delta\}$$

$$\phi_i = i$$

$$v_i = 0.0002$$

$$\tau_i = \{.02, .01, .02, .01, .02, .01, .02, .01\}$$

$$\tau_0 = 10^{-5}$$

$$\sigma = 0.01$$

$$\sigma_\theta = 0.01$$

$$\boldsymbol{\mu}_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$\boldsymbol{\Sigma}_0 = \begin{pmatrix} 10^{-4} & 0 \\ 0 & 10^{-4} \end{pmatrix}$$

References

1. Repp BH and Su YH. Sensorimotor synchronization: A review of recent research (2006-2012). *Psychonomic Bulletin and Review* 2013; 20:403–52. DOI: 10.3758/s13423-012-0371-2. arXiv: NIHMS150003
2. Merchant H, Grahn J, Trainor L, Rohrmeier M, and Fitch WT. Finding the beat: a neural perspective across humans and non-human primates. *Philosophical transactions of the Royal Society of London. Series B, Biological sciences* 2015; 370. DOI: 10.1098/rstb.2014.0093. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/25646516>
3. Obleser J and Kayser C. Neural Entrainment and Attentional Selection in the Listening Brain. *Trends in Cognitive Sciences* 2019; 23:1–14. DOI:

- 784 10.1016/j.tics.2019.08.004. Available from: [https://doi.org/10.](https://doi.org/10.1016/j.tics.2019.08.004)
785 1016/j.tics.2019.08.004
- 786 4. Lawrance ELA, Harper NS, Cooke JE, and Schnupp JWH. Temporal pre-
787 dictability enhances auditory detection. *The Journal of the Acoustical So-*
788 *ciety of America* 2014; 135:EL357–EL363. DOI: 10.1121/1.4879667.
789 Available from: <http://dx.doi.org/10.1121/1.4879667>
- 790 5. Nobre AC and Van Ede F. Anticipated moments: Temporal structure in
791 attention. *Nature Reviews Neuroscience* 2018; 19:34–48. DOI: 10.1038/
792 **nrn.2017.141**. Available from: [http://dx.doi.org/10.1038/nrn.2017.](http://dx.doi.org/10.1038/nrn.2017.141)
793 141
- 794 6. Morillon B, Schroeder CE, Wyart V, and Arnal LH. Temporal prediction
795 in lieu of periodic stimulation. *Journal of Neuroscience* 2016; 36:2342–7.
796 DOI: 10.1523/JNEUROSCI.0836-15.2016
- 797 7. Lange K. Brain correlates of early auditory processing are attenuated by
798 expectations for time and pitch. *Brain and Cognition* 2009; 69:127–37.
799 DOI: 10.1016/j.bandc.2008.06.004. Available from: [http://dx.doi.](http://dx.doi.org/10.1016/j.bandc.2008.06.004)
800 [org/10.1016/j.bandc.2008.06.004](http://dx.doi.org/10.1016/j.bandc.2008.06.004)
- 801 8. Jazayeri M and Shadlen MN. Temporal context calibrates interval timing.
802 *Nature Neuroscience* 2010; 13:1020–6. DOI: 10.1038/**nn.2590**
- 803 9. Herrmann B, Henry MJ, Haegens S, and Obleser J. Temporal expectations
804 and neural amplitude fluctuations in auditory cortex interactively influence
805 perception. *NeuroImage* 2016; 124:487–97. DOI: 10.1016/j.neuroimage.
806 2015.09.019
- 807 10. Rajendran VG, Teki S, and Schnupp JW. Temporal Processing in Audi-
808 tion: Insights from Music. *Neuroscience* 2018; 389:4–18. DOI: 10.1016/

- 809 j.neuroscience.2017.10.041. Available from: <https://doi.org/10.1016/j.neuroscience.2017.10.041>
- 810
- 811 11. Large EW and Jones MR. The dynamics of attending: How people track
812 time-varying events. *Psychological Review* 1999; 106:119–59. DOI: 10.1037//0033-295x.106.1.119
- 813
- 814 12. Large EW and Palmer C. Perceiving temporal regularity in music. *Cognitive Science* 2002; 26:1–37. DOI: 10.1016/S0364-0213(01)00057-X
- 815
- 816 13. Breska A and Deouell LY. Neural mechanisms of rhythm-based tempo-
817 ral prediction: Delta phase-locking reflects temporal predictability but not
818 rhythmic entrainment. *PLoS Biology* 2017; 15:1–30. DOI: 10.1371/journal.pbio.2001665
- 819
- 820 14. Bouwer FL, Honing H, and Slagter HA. Beat-based and memory-based
821 temporal expectations in rhythm: similar perceptual effects, different un-
822 derlying mechanisms. 2019; 8:55
- 823 15. Rimmele JM, Morillon B, Poeppel D, and Arnal LH. Proactive Sensing of
824 Periodic and Aperiodic Auditory Patterns. *Trends in Cognitive Sciences*
825 2018; 22:870–82. DOI: 10.1016/j.tics.2018.08.003. Available from:
826 <https://doi.org/10.1016/j.tics.2018.08.003>
- 827 16. Friston K. A theory of cortical responses. *Philosophical Transactions of*
828 *the Royal Society B: Biological Sciences* 2005; 360:815–36. DOI: 10.1098/rstb.2005.1622
- 829
- 830 17. Friston K. Does predictive coding have a future? *Nature Neuroscience*
831 2018; 21:1019–21. DOI: 10.1038/s41593-018-0200-7
- 832 18. Vuust P and Witek MA. Rhythmic complexity and predictive coding: A
833 novel approach to modeling rhythm and meter perception in music. *Frontiers in Psychology* 2014; 5:1–14. DOI: 10.3389/fpsyg.2014.01111
- 834

- 835 19. Vuust P, Dietz MJ, Witek M, and Kringelbach ML. Now you hear it: A
836 predictive coding model for understanding rhythmic incongruity. *Annals*
837 *of the New York Academy of Sciences* 2018; 1423:19–29. DOI: 10.1111/
838 *nyas.13622*
- 839 20. Proksch S, Comstock DC, Médé B, Pabst A, and Balasubramaniam R.
840 Motor and Predictive Processes in Auditory Beat and Rhythm Perception.
841 2020; 14. DOI: 10.3389/fnhum.2020.578546
- 842 21. Friston K, Stephan K, Li B, and Daunizeau J. Generalised filtering. *Math-*
843 *ematical Problems in Engineering* 2010; 2010. DOI: 10.1155/2010/621670
- 844 22. Buckley CL, Kim CS, McGregor S, and Seth AK. The free energy principle
845 for action and perception: A mathematical review. *Journal of Mathemati-*
846 *cal Psychology* 2017; 81:55–79. DOI: 10.1016/j.jmp.2017.09.004. arXiv:
847 1705.09156. Available from: [http://dx.doi.org/10.1016/j.jmp.2017.](http://dx.doi.org/10.1016/j.jmp.2017.09.004)
848 09.004
- 849 23. Schwartz M and Kotz SA. A dual-pathway neural architecture for spe-
850 cific temporal prediction. *Neuroscience and Biobehavioral Reviews* 2013;
851 37:2587–96. DOI: 10.1016/j.neubiorev.2013.08.005. Available from:
852 <http://dx.doi.org/10.1016/j.neubiorev.2013.08.005>
- 853 24. Egger SW and Jazayeri M. A nonlinear updating algorithm captures subop-
854 timal inference in the presence of signal-dependent noise. *Scientific Reports*
855 2018 ;18–20. DOI: 10.1038/s41598-018-30722-0
- 856 25. DI Luca M and Rhodes D. Optimal Perceived Timing: Integrating Sensory
857 Information with Dynamically Updated Expectations. *Scientific Reports*
858 2016; 6:1–15. DOI: 10.1038/srep28563
- 859 26. Elliott MT, Wing AM, and Welchman AE. Moving in time: Bayesian causal
860 inference explains movement coordination to auditory beats. *Proceedings*

- 861 of the Royal Society B: Biological Sciences 2014; 281. DOI: 10.1098/rspb.
862 2014.0751
- 863 27. Snyder DL. Filtering and Detection for Doubly Stochastic Poisson Pro-
864 cesses. IEEE Transactions on Information Theory 1972; 18:91–102. DOI:
865 10.1109/TIT.1972.1054756
- 866 28. Oppen M. A Bayesian Approach to On-line Learning. On-Line Learning in
867 Neural Networks 2010 :363–78. DOI: 10.1017/cbo9780511569920.017
- 868 29. Friston K. The free-energy principle: A unified brain theory? Nature Re-
869 views Neuroscience 2010; 11:127–38. DOI: 10.1038/nrn2787
- 870 30. Eden UT and Brown EN. CONTINUOUS-TIME FILTERS FOR STATE
871 ESTIMATION FROM POINT PROCESS MODELS OF NEURAL DATA.
872 Statistica Sinica 2008; 18:1293–310
- 873 31. Cemgil AT, Kappen B, Desain P, and Honing H. On tempo tracking:
874 Tempogram representation and Kalman filtering. Journal of New Music
875 Research 2000; 29:259–73. DOI: 10.1080/09298210008565462
- 876 32. London J, Polak R, and Jacoby N. Rhythm histograms and musical meter:
877 A corpus study of Malian percussion music. Psychonomic Bulletin and
878 Review 2017; 24:474–80. DOI: 10.3758/s13423-016-1093-7
- 879 33. Polak R, London J, and Jacoby N. Both isochronous and non-isochronous
880 metrical subdivision afford precise and stable ensemble entrainment: A
881 corpus study of malian jembe drumming. Frontiers in Neuroscience 2016;
882 10:1–11. DOI: 10.3389/fnins.2016.00285
- 883 34. Friberg A and Sundström A. Swing Ratios and Ensemble Timing in Jazz
884 Performance: Evidence for a Common Rhythmic Pattern. Music Percep-
885 tion 2002; 19:333–49. DOI: 10.1525/mp.2002.19.3.333

- 886 35. Fitch WT and Rosenfeld AJ. Perception and Production of Syncopated
887 Rhythms. *Music Perception* 2007; 25:43–58
- 888 36. Warren RM and Gregory RL. An Auditory Analogue of the Visual Re-
889 versible Figure. *The American Journal of Psychology* 1958; 71:612–3
- 890 37. HALL GS and JASTROW J. STUDIES OF RHYTHM. *Mind* 1886 Jan;
891 os-XI:55–62. DOI: 10.1093/mind/os-XI.41.55. eprint: [https://](https://academic.oup.com/mind/article-pdf/os-XI/41/55/9358438/os-XI\41\55.pdf)
892 [academic.oup.com/mind/article-pdf/os-XI/41/55/9358438/os-](https://academic.oup.com/mind/article-pdf/os-XI/41/55/9358438/os-XI\41\55.pdf)
893 [XI\41\55.pdf](https://academic.oup.com/mind/article-pdf/os-XI/41/55/9358438/os-XI\41\55.pdf). Available from: [https://doi.org/10.1093/mind/os-](https://doi.org/10.1093/mind/os-XI.41.55)
894 [XI.41.55](https://doi.org/10.1093/mind/os-XI.41.55)
- 895 38. Nakajima Y. A psychophysical investigation of divided time intervals shown
896 by sound bursts. *Journal of the Acoustical Society of Japan* 1979; 35:145–
897 51
- 898 39. Meumann E. Beiträge zur Psychologie des Zeitbewußtseins [contributions
899 to the psychology of time consciousness]. *Philosophische Studien* 1896;
900 12:128–254
- 901 40. Grimm K. der einfluß der Zeitform auf die Wahrnehmung der Zeitdauer
902 [the influence of time-form on the perception of duration]. *Zeitschrift für*
903 *Psychologie* 1934; 132:104–32
- 904 41. Repp BH and Bruttomesso M. A filled duration illusion in music: Effects
905 of metrical subdivision on the perception and production of beat tempo.
906 *Advances in Cognitive Psychology* 2009; 5:114–34. DOI: 10.2478/V10053-
907 008-0071-7
- 908 42. Repp B and Jendoubi H. Flexibility of temporal expectations for triple
909 subdivision of a beat. *Advances in Cognitive Psychology* 2009; 5:27–41.
910 DOI: 10.2478/v10053-008-0063-7

- 911 43. Repp BH. Tapping in synchrony with a perturbed metronome: The phase
912 correction response to small and large phase shifts as a function of tempo.
913 *Journal of Motor Behavior* 2011; 43:213–27. DOI: 10.1080/00222895.
914 2011.561377
- 915 44. Repp BH, Keller PE, and Jacoby N. Quantifying phase correction in sen-
916 sorimotor synchronization: Empirical comparison of three paradigms. *Acta*
917 *Psychologica* 2012; 139:281–90. DOI: 10.1016/j.actpsy.2011.11.002.
918 Available from: <http://dx.doi.org/10.1016/j.actpsy.2011.11.002>
- 919 45. Witek MA, Clarke EF, Kringelbach ML, and Vuust P. Effects of Poly-
920 phonic Context, Instrumentation, and Metrical Location on Syncopation
921 in Music. *Music Perception* 2014; 32:201–17
- 922 46. Hove MJ, Marie C, Bruce IC, and Trainor LJ. Superior time perception for
923 lower musical pitch explains why bass-ranged instruments lay down musical
924 rhythms. *Proceedings of the National Academy of Sciences of the United*
925 *States of America* 2014; 111:10383–8. DOI: 10.1073/pnas.1402039111
- 926 47. Repp BH. Phase Correction , Phase Resetting , and Phase Shifts After Sub-
927 liminal Timing Perturbations in Sensorimotor Synchronization. *Journal*
928 *of Experimental Psychology: Human Perception and Performance* 2001;
929 27:600–21. DOI: 10.1037//0096-1523.27.3.600
- 930 48. Heggli OA, Cabral J, Konvalinka I, Vuust P, and Kringelbach ML. A Ku-
931 ramoto model of self-other integration across interpersonal synchronization
932 strategies. *PLoS Computational Biology* 2019; 15:1–17. DOI: 10.1371/
933 journal.pcbi.1007422
- 934 49. Koban L, Ramamoorthy A, and Konvalinka I. Why do we fall into sync
935 with others? Interpersonal synchronization and the brain’s optimization
936 principle. *Social Neuroscience* 2019; 14:1–9

- 937 50. Wing AM and Kristofferson AB. Response delays and the timing of discrete
938 motor responses. *Perception & Psychophysics* 1973; 14:5–12. DOI: 10 .
939 3758/BF03198607
- 940 51. Mates J. A model of synchronization of motor acts to a stimulus sequence
941 - II. Stability analysis, error estimation and simulations. *Biological Cyber-*
942 *netics* 1994; 70:475–84. DOI: 10.1007/BF00203240
- 943 52. Fox C, Rezek I, and Roberts S. Drum ' N ' Bayes : on-Line Variational
944 Inference for Beat Tracking and Rhythm Recognition. *International Com-*
945 *puter Music Conference* 2007. DOI: 10.1016/j.chieco.2016.10.003
- 946 53. Pesek M, Leonardis A, and Marolt M. An Analysis of Rhythmic Pat-
947 terns with Unsupervised Learning. *Applied Sciences* 2019. DOI: 10.3390/
948 app10010178
- 949 54. Repp BH. Obligatory "expectations" of expressive timing induced by per-
950 ception of musical structure. *Psychological Research* 1998; 61:33–43. DOI:
951 10.1007/s004260050011
- 952 55. Repp BH. Compensation for subliminal timing perturbations in perceptual-
953 motor synchronization. *Psychological Research* 2000; 63:106–28. DOI: 10 .
954 1007/PL00008170
- 955 56. Schroeder CE and Lakatos P. Low-frequency neuronal oscillations as in-
956 struments of sensory selection. *Trends in neurosciences* 2009; 32. DOI:
957 10.1016/j.tins.2008.09.012.Low-frequency
- 958 57. Arnal LH and Giraud AL. Cortical oscillations and sensory predictions.
959 *Trends in Cognitive Sciences* 2012; 16:390–8. DOI: 10.1016/j.tics .
960 2012.05.003. Available from: [http://dx.doi.org/10.1016/j.tics.](http://dx.doi.org/10.1016/j.tics.2012.05.003)
961 2012.05.003

- 962 58. Arnal LH and Kleinschmidt AK. Entrained delta oscillations reflect the
963 subjective tracking of time. *Cerebral Cortex* 2017 :e1349583. DOI: 10 .
964 1093/cercor/bhu103
- 965 59. Gámez J, Mendoza G, Prado L, Betancourt A, and Merchant H. The am-
966 plitude in periodic neural state trajectories underlies the tempo of rhythmic
967 tapping. *PLoS biology* 2019; 17:e3000054
- 968 60. Tomassini A, Ruge D, Galea JM, Penny W, and Bestmann S. The Role
969 of Dopamine in Temporal Uncertainty. *Journal of Cognitive Neuroscience*
970 2016. DOI: 10 . 1162/jocn. arXiv: 1511 . 04103. Available from: [http :
971 //dx . doi . org/10 . 1162/jocn%7B%5C_%7Da%7B%5C_%7D00409%7B%5C_
972 %7D5Cnhttp://www.mitpressjournals . org/doi/abs/10 . 1162/jocn%
973 7B%5C_%7Da%7B%5C_%7D00409](http://dx.doi.org/10.1162/jocn%7B%5C_%7Da%7B%5C_%7D00409%7B%5C_%7D5Cnhttp://www.mitpressjournals.org/doi/abs/10.1162/jocn%7B%5C_%7Da%7B%5C_%7D00409)
- 974 61. Sarno S, De Lafuente V, Romo R, and Parga N. Dopamine reward predic-
975 tion error signal codes the temporal evaluation of a perceptual decision re-
976 port. *Proceedings of the National Academy of Sciences of the United States*
977 *of America* 2017; 114:E10494–E10503. DOI: 10.1073/pnas.1712479114
- 978 62. Herbst SK, Fiedler L, and Obleser J. Tracking temporal hazard in the hu-
979 man electroencephalogram using a forward encoding model. *eNeuro* 2018;
980 5:1–17. DOI: 10.1523/ENEURO.0017–18.2018
- 981 63. Tavano A, Schröger E, and Kotz SA. Beta power encodes contextual esti-
982 mates of temporal event probability in the human brain. *PLoS ONE* 2019;
983 14. DOI: 10.1371/journal.pone.0222420
- 984 64. Cannon J and Patel AD. How beat perception coopts motor neurophysiol-
985 ogy: a proposal. *bioRxiv* 2020. DOI: <https://doi.org/10.1101/805838>

- 986 65. Schubotz RI. Prediction of external events with our motor system: towards
987 a new framework. *Trends in Cognitive Sciences* 2007; 11:211–8. DOI: 10.
988 1016/j.tics.2007.02.006
- 989 66. Rauschecker JP. An expanded role for the dorsal auditory pathway in
990 sensorimotor control and integration. *Hearing Research* 2011; 271:16–25.
991 DOI: 10.1016/j.heares.2010.09.001. Available from: [http://dx.doi.](http://dx.doi.org/10.1016/j.heares.2010.09.001)
992 [org/10.1016/j.heares.2010.09.001](http://dx.doi.org/10.1016/j.heares.2010.09.001)
- 993 67. Kneissler J, Drugowitsch J, Friston K, and Butz MV. Simultaneous learn-
994 ing and filtering without delusions: A bayes-optimal combination of pre-
995 dictive inference and adaptive filtering. *Frontiers in Computational Neu-*
996 *roscience* 2015; 9:1–12. DOI: 10.3389/fncom.2015.00047
- 997 68. Danielsen A. Here, There, and Everywhere: three accounts of pulse in
998 D’Angelo’s ‘Left and Right’. 2010 Jan :19–36. DOI: 10.4324/9781315596983-
999 2
- 1000 69. Weij B van der, Pearce MT, and Honing H. A probabilistic model of meter
1001 perception: Simulating enculturation. *Frontiers in Psychology* 2017; 8:1–
1002 18. DOI: 10.3389/fpsyg.2017.00824
- 1003 70. Rohrmeier M. Towards a formalization of musical rhythm. *Proc. of the*
1004 *21st Int. Society for Music Information Retrieval Conf.* 2020
- 1005 71. Pearce MT. The construction and evaluation of statistical models of melodic
1006 structure in music perception and composition. PhD thesis. City Univer-
1007 sity, London, 2005
- 1008 72. Sioros G, Davies ME, and Guedes C. A generative model for the char-
1009 acterization of musical rhythms. *Journal of New Music Research* 2018;
1010 47:114–28. DOI: 10.1080/09298215.2017.1409769. Available from: [http:](http://doi.org/10.1080/09298215.2017.1409769)
1011 [//doi.org/10.1080/09298215.2017.1409769](http://doi.org/10.1080/09298215.2017.1409769)

- 1012 73. Harel Y, Meir R, and Oppor M. A tractable approximation to optimal
1013 point process filtering: Application to neural encoding. Advances in Neural
1014 Information Processing Systems 2015; 2015-Janua:1603–11
- 1015 74. Snyder DL and Fishman P. How to track a swarm of fireflies by observing
1016 their flashes. IEEE Transactions on Information Theory 1975; 21