# Predictive uncertainty underlies auditory boundary perception

# Predictive uncertainty underlies auditory boundary perception

## Abstract

Anticipating the future is essential for efficient perception and action planning. Yet, the role of anticipation in event segmentation is understudied because empirical research has focused on retrospective cues such as surprise. We address this question in the context of musical phrase-boundary perception. A computational model of cognitive sequence processing was used to control the information-dynamic properties of tone sequences. In an implicit, self-paced listening task ($n$=38), undergraduates dwelled longer on tones generating high entropy (i.e., ~~low~~ high uncertainty) than those generating low entropy (i.e., ~~high~~ low uncertainty). Similarly, sequences that ended on tones generating high entropy were rated as sounding more complete ($n$=31). These entropy effects were independent of both the surprise (i.e., information content) and phrase position of target tones in the original musical stimuli. Our results indicate that events generating high entropy prospectively contribute to ~~prospective~~ segmentation processes in auditory sequence perception, independent of the properties of the subsequent event.

# Statement of relevance

A significant challenge for the human perceptual system is to promote time-sensitive, context-appropriate responses by predictively processing continuous streams of complex sensory information. A large body of research shows that expectations gleaned from a lifetime of experience guide such processes, which are critical in high-risk environments like traffic or manual labor. Because most studies have focused on the degree of surprise evoked by events, there is little evidence for the role of prospective expectations in perceptual organization. Here, we control entropy in musical tone sequences by using an information-theoretic model that has been shown to reflect listeners' ~~prospective~~ predictive uncertainty. Tones that afforded relatively high uncertainty were found to draw implicit attention and influence explicit ratings of sequence completeness. Focusing attention on instances where upcoming events are statistically unconstrained could contribute to an adaptive mechanism facilitating stream segmentation that leads to efficient learning and information processing in a complex, dynamic world.

# Introduction

Humans make sense of a complex, dynamic world by segmenting sequences of events into manageable units (Zacks & Swallow, 2007; Kurby & Zacks, 2008; Richmond & Zacks, 2017). Past work on segmentation has focused on retrospective cues for boundary identification, often conceptualizing group boundaries as coinciding with instances of increased relative change in stimulus features or low transition probabilities (e.g., speech: Saffran & Kirkham, 2018; action sequences: Hard et al., 2011; music: Hartmann et al., 2017; Pearce et al. 2010). However, the sophisticated prediction capabilities of the human mind (Hutchinson & Barrett, 2019) suggest that event boundaries are also anticipated prospectively. For example, in natural conversation, turn-taking happens so rapidly that speakers likely anticipate the end of their conversation partner's sentence (Levinson, 2016). Here we investigate the role of entropy, or degree of prospective uncertainty about an upcoming event, in determining the perception of group boundaries in auditory sequences. We define *prediction* as the psychological processes of generating an expectation about a future event, in terms of how likely the various possible outcomes are. We define *uncertainty* as the imprecision (or extent of equi-probability) of such a prediction.

Though most previous work has focused on retrospective boundary identification of boundaries, anticipatory processing has some preliminary support. Previous work has observed that wWhen self-pacing through sequential images of action sequences, participants tend to "dwell" (or pause) on perceived boundary images (Hard et al., 2011; Hard et al., 2019; Kosie & Baldwin, 2019a, 2019b). Kosie and Baldwin (2019b) proposed that this "dwell time effect" resulted from selective attention to moments of uncertainty afforded by perceiving a goal completion event. No cognitive model was devised to test this theory, however, potentially due to the challenges in modeling expectancy in event processing of action sequences. Indeed, one methodological drawback of this methodology was demonstrated by the finding that participants' dwellinged on boundary slides even

4

55    when those slides were out of order, suggesting that they were responding to conceptual salience

56    rather than to underlying expectancy dynamics (Hard et al., 2011). Cohen et al. (2007) have proposed

57    an entropy-based segmentation model for language, but because it computes statistics from the corpus

58    it is segmenting—including parts it has not yet seen—it does not fully capture segmentation

59    processing in real time (Christiansen & Chater, 2016).

60        Because music is not only hierarchically structured (Lerdahl & Jackendoff, 1983), but also

61    statistically well-defined, it is an ideal domain for testing psychological theories of probabilistic

62    perception (Koelsch, Vuust, & Friston, 2019). As with non-musical sequences (Zacks et al., 2001),

63    there is generally high inter-participant agreement regarding the location of musical phrase

64    boundaries (Deliège, 1987; but see Pearce et al., 2010), and as with action sequences, listeners self-

65    pacing through musical chords "dwell" on boundary chords (Kragness & Trainor, 2016, 2018). Since,

66    however, entropy correlates strongly with phrase boundaries in music (Hansen et al., 2017), previous

67    studies were not optimized to separate prospective effects of expectancy dynamics ~~vs.~~from effects of

68    canonical boundary features on perceptual grouping. ~~The~~ *Information Dynamics of Music* ~~Model~~

69    (IDyOM) (Pearce, 2005) is a computational model of auditory expectation which ~~provides a means~~

70    ~~of~~ enables modelling boundary perception quantitatively using the information-theoretic concepts of

71    entropy and information content, computed in reference to pre-existing long-term knowledge (Hansen

72    & Pearce, 2014; Hansen et al., 2016). Entropy ~~enables~~ facilitates a test of ~~prospective~~ uncertainty as

73    a prospective mechanism for boundary perception which can be pitted directly against information

74    content (a measure of surprise) as a retrospective cue. For example, an individual may form a highly

75    certain ~~prospective~~ prediction ~~for~~ about the next note in a melody but then be surprised when a

76    different note ~~actually~~ follows. Another advantage of ~~using~~ melodic sequences is that~~, unlike images~~

77    ~~of actions,~~ any given note has little intrinsic meaning in isolation from its preceding musical context,

78    ensuring that ~~any~~ observed effects on perception reflect the statistical structure of the sequence and

**5**

not inherent features of the boundary stimulus itself. However, because uncertainty ~~processing~~ is not always available for explicit introspection (Hansen et al., 2016), implicit measures are paramount for investigating the cognitive mechanisms underlying boundary perception.

The present study used ~~the~~ IDyOM ~~model~~ to control the information-dynamic properties of melodic sequences in two experiments ~~that~~ assess~~ing~~ed the role of ~~prospective~~ predictive uncertainty in sequence processing. We measured participants' dwell times (Experiment 1) and explicit ratings of phrase completeness (Experiment 2) for tones that afforded high/low entropy and were phrase-beginning/phrase-ending in the melodies from which they were drawn. We predicted that tones ~~that generated~~ generating high ~~levels of prospective~~ uncertainty would lead to longer dwell times ~~(Experiment 1)~~ and higher ~~explicit~~ ratings of phrase completeness, regardless of original phrase status, ~~(Experiment 2)~~ and that this effect would be independent from ~~that of~~ retrospective surprise.

## Experiment 1: Implicit Self-Pacing Task

**Methods**

*Participants.* Thirty-eight McMaster University undergraduates received psychology course credits for participating in the study ($M_{age}$ = 19.3 years, 1 person declined to report their age, $SD_{age}$ = 3.78, 8 men, 30 women). None of the participants were professional musicians (for more information about musical training levels, see Table S1 in SOM-R2). This sample size exceeds or corresponds to those of previous studies using this methodology to assess comparable effects (e.g., Hard et al., 2011; Kragness & Trainor, 2016, 2018). All participants were fluent in English.

*Stimuli.* Fifty-six monophonic stimulus sequences were selected from the soprano (i.e., highest) part in 370 four-part chorale harmonizations by Johann Sebastian Bach (Dörffel, 1875) (see SOM-R1 for details of the stimulus selection procedure). These chorale melodies are not generally known by present-day listeners in Canada. Unfamiliarity was, moreover, made more likely through

103 complete removal of rhythmic information by granting participants control over tone durations in the

104 self-paced dwell-time paradigm (Experiment 1) or by presenting stimuli with isochronized tone

105 durations (Experiment 2). All chords, interference tones, and self-pacing tones were generated in

106 MaxMSP's grand piano timbre.

107 Each stimulus context contained a full phrase (musical group) of seven to 17 pitches followed

108 by the initial tone of the subsequent phrase in the original chorale melody. Tones associated with

109 phrase beginnings and endings were unambiguously identified from notations in the musical score.

110 This practice seems at least as objective as the reliance on trained "expert coders" to determine event

111 boundaries in research using visual action sequences (e.g., Hard et al., 2019; Kosie & Baldwin, 2019a,

112 2019b). We included both phrase endings and phrase beginnings as target tones to provide a strong

113 test of entropy's role in segmentation, controlling for compositional cues in the melodies that might

114 signal melodic phrase endings in other ways.

115 Fourteen stimulus contexts were selected for each of the four experimental conditions,

116 comprising phrase beginnings with high ("BegHi") or low entropy ("BegLo") and phrase endings

117 with high ("EndHi") or low entropy ("EndLo"). Entropy, in this regard, quantifies the level of

118 uncertainty governing a listener's expectations about what the pitch of the next tone following the

119 relevant phrase beginning or phrase ending would be. Thus, Western-enculturated listeners are

120 expected to be relatively sure about which pitch will follow the target tone in "BegLo" and "EndLo"

121 contexts, but relatively unsure in "BegHi" and "EndHi" contexts. "Target tone", in this respect, refers

122 to the final tone in "BegLo" and "BegHi" contexts and the penultimate tone in "EndLo" and "EndHi"

123 contexts.

124 The entropy level generated by each tone in the corpus was estimated by the *Information*

125 *Dynamics of Music Model* (IDyOM, version 1.3) (Pearce, 2005). This variable-order *n*-gram model

126 uses unsupervised statistical learning to generate probability distributions governing a relevant

127  feature of each tone in a monophonic melody. IDyOM was trained on a large dataset of 5,332 German

128  folk songs (Schaffrath, 1995), 152 Nova Scotian songs and ballads (Creighton, 1966), and 120

129  English hymns (Nicholson et al., 1950)[1]. For each tone in the chorale melody, IDyOM generated a

130  probability distribution (summing to 1) over the 44 pitch values occurring in the training corpus (i.e.,

131  MIDI pitches 45-89 corresponding to A2-F6) by combining *n*-gram models of varying order. Entropy

132  then quantifies the shape of these probability distributions with high entropy for "flat" (relatively

133  uniform) distributions, where there is high uncertainty about the next event, and low entropy for

134  "spiky" (relatively nonuniform) distributions, where one or a small number of continuations are

135  highly probable.

136        The set of 56 stimulus contexts was selected in a way that prioritized extreme high or low

137  entropy values while ensuring that three conditions were met: First, as shown by a non-parametric

138  Kruskal-Wallis test, all four conditions, including EndHi (Median = 2.45, IQR = 1.76), BegHi

139  (Median = 2.69, IQR = 1.78), EndLo (Median = 2.31, IQR = 2.70), and BegLo (Median = 3.37, IQR

140  = 1.83), were matched on information content (i.e., inverse log-probability) for the event of interest, $\chi^2(3)$ = 4.55, *p* = .208; second, as shown by Mann-Whitney *U*-tests, EndHi (Median = 2.97, IQR =

142  0.08) and BegHi (Median = 3.00, IQR = 0.12) stimuli, *U* = 78, *p* = .376, as well as EndLo (Median =

143  1.07, IQR = 0.30) and BegLo (Median = 0.97, IQR = 0.35) stimuli, *U* = 90, *p* = .734, were matched

144  on entropy governing the next event in the sequence. The experimenter selecting these stimuli paid

145  no attention to any other musical features.

146        For the secondary analysis of all tones in the stimulus set, IC and entropy were re-estimated

147  by re-running IDyOM with the same configuration on the final stimulus contexts. This was done

148  because IC and entropy estimates for the initial tones in each stimulus context sometimes relied on

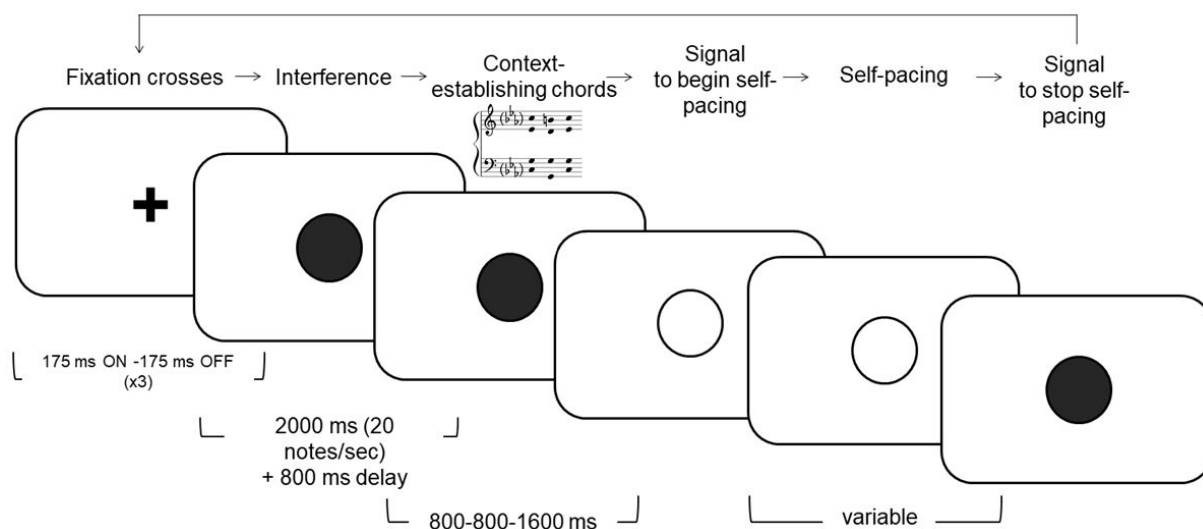149  tones from the preceding phrase in the original chorales, which was excluded from the stimuli used.

---

[1] For more information about the IDyOM implementation and parameters, please see SOM-R.

150    While unproblematic for stimulus selection based on target tones, this presented a problem for tone-

151    level analysis. Note that due to their late position in the tone sequences, target tone entropy and IC

152    values were identical for the two models (one used in stimulus generation and analyses of target tones,

153    the other used in the analysis of all tones).

154          *Procedure.*        The experimental procedures (for Experiment 1 and 2) received prior approval

155    from the McMaster University Research Ethics Board and was carried out in accordance with the

156    provisions of the World Medical Association Declaration of Helsinki. Participants were seated facing

157    a computer screen in a sound-attenuated room. They were instructed to press the spacebar on a

158    computer keyboard with the pointer finger of their dominant hand to elicit the onset of each

159    subsequent tone in the sequence. Tones decayed naturally, but were not terminated until the spacebar

160    was pressed again to initiate the next tone. Participants were instructed to progress as quickly or

161    slowly as they liked while listening carefully, and could not repeat previously heard tones. They were

162    led to falsely believe that their memory for the sequences would be tested afterwards to motivate them

163    to attend to the task (Kragness & Trainor, 2016). No other instructions regarding timing, pacing,

164    rhythmicity, or expressivity were given. If a participant asked for further information, they were told

165    to play through the piece in a way that would maximize their performance in the subsequent memory

166    task.

167          Prior to each trial, participants saw three flashes of a fixation cross, then heard 40 50-ms tones

168    (for a total of 2000 ms) chosen randomly on each trial from range E2 to A5 to minimize carryover

169    from the context of the previous sequence, followed by three context-establishing chords with

170    durations of 800, 800, and 1600 ms (Figure 1). The context-establishing chords were played in the

171    key of the relevant melody. Throughout each trial, a circle on the screen indicated when to begin self-

172    pacing through the melody (light green) and when to stop (dark green).

173

**Figure 1.** Depiction of a trial from Experiment 1. In each trial, participants saw a fixation cross, followed by interference tones, then three context-establishing chords and a signal (white circle) to begin self-pacing. They then self-paced through the tone sequence until the occurrence of a stop signal (black circle). The box depicts examples of tone sequences from each condition containing target tones (boxed) generating relatively uncertain (high entropy) or relatively certain (low entropy) expectations about the pitch of the next tone, matched on IC of the current tone. The double slash indicates whether target tones were phrase beginnings (after double slash) or phrase endings (prior to double slash) in the original notation.

*Data processing and statistical analysis.*     Despite systematic efforts to avoid duplicate stimulus contexts (e.g., multiple occurrences of a repeated phrase from a single melody or identical

185 phrases across melodies), it was discovered after data collection that one melodic context occurred

186 both amongst the "BegHi" and "EndHi" stimulus sets (with different target tones). Given that results

187 did not differ substantially when excluding dwell times for these stimuli, we report statistical analyses

188 including the full dataset here, which included 56 total tone sequences (i.e., 14 per condition).

189 　　To mitigate effects of extreme data points, a minimum dwell-time threshold of 100 ms was

190 adopted for inclusion. Dwell times greater than 3 standard deviations above a participant's own

191 average (across all target and non-target dwell times) were also omitted (Kosie & Baldwin, 2019a,

192 2019b). These exclusion criteria eliminated an average of 1.31% of all tones and 1.70% of target

193 tones per participant (ranging from 0-4 target tones).
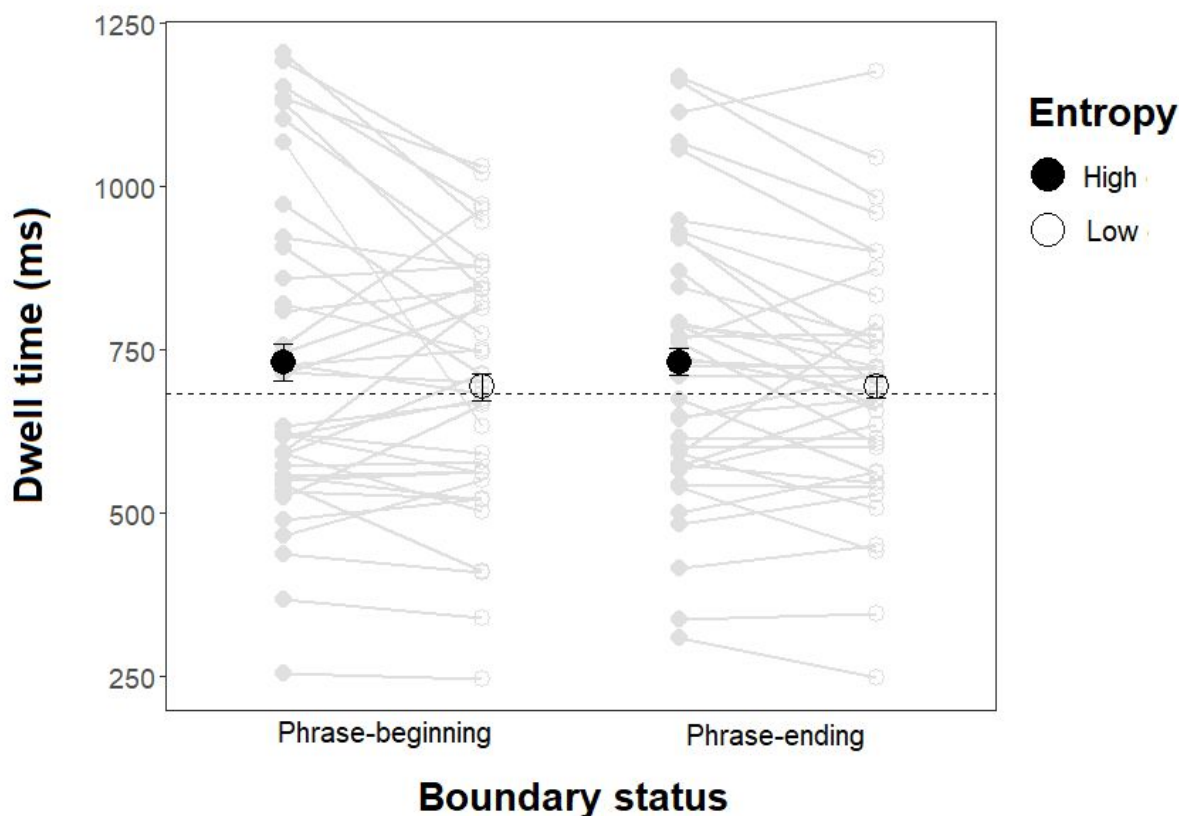
194 　　For the main analysis of target tones, target dwell times were averaged by condition resulting

195 in four condition-wise means per participant. A 2x2 repeated-measures ANOVA (including within-

196 subjects factors boundary status and entropy) was run on target tone dwell times.

197 　　For the secondary analysis of all tones, dwell times were first log-transformed to minimize

198 the positive skew inherent to timing data (cf. Kragness & Trainor, 2018). Subsequently, using the

199 *lmer()* function from the *lme4* package in R (R Core Team, 2019), linear mixed-effects models were

200 fitted with Restricted Maximum Likelihood estimates (REML). Because previous experiments have

201 found that dwell times change systematically throughout trials (Kragness & Trainor, 2016), tone

202 index in the sequence was always included as a predictor. Thus, whereas the null model only included

203 tone index as a fixed effect, two further increasingly complex models added, first, the retrospective

204 cue IC, and, second, the prospective cue entropy. Thereby, we could determine whether prospective

205 predictive processing explained unique variance not already accounted for by retrospective surprise.

206 Random intercepts and slopes of tone number were included for each participant. For all models, this

207 random-effects structure produced the lowest BIC values while avoiding singular fits.

208

**11**

1
2
3
4    209    **Results**
5
6
7    210         *Target tones.*   To examine the effects of boundary status (phrase-ending, phrase-beginning)
8
9    211    and entropy (high, low), a 2x2 repeated-measures ANOVA was run on target tone dwell times.
10
11   212    Whereas no significant interaction ($F(1,37) < 0.01$, $p = .986$, $\eta^2_p < .001$) or main effect of boundary
12
13   213    status ($F(1,37) < 0.01$, $p = .973$, $\eta^2_p < .001$) was found, there was a significant main effect of entropy
14
15   214    ($F(1,37) = 7.24$, $p = .011$, $\eta^2_p = .164$). Thus, as hypothesized, high-entropy target tones were generally
16
17   215    dwelled on longer than low-entropy target tones, regardless of phrase position in the original chorale
18
19   216    melody (Figure 2).
20
21   217         We conducted post-hoc correlational analyses to examine whether participants' musical
22
23   218    sophistication was associated with the magnitude of their dwell time effect. No significant
24
25   219    associations were observed (see SOM-R2 for more details).
26
27
28   220
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

**Figure 2.** Dwell times (ms) for each type of target tone (BegHi, BegLo, EndHi, EndLo) in Experiment 1. The dashed line represents the average dwell time (683 ms) for non-target tones. Error bars represent within-subject 95% confidence intervals (Cousineau, 2005). High-entropy target tones had longer dwell times than low-entropy target tones, and it made no significant difference whether target tones originated from phrase endings or phrase beginnings in the original chorale melody corpus.

*All tones.*        If ~~prospective~~ uncertainty provides a cognitive cue for phrase segmentation, its effect on dwell times should generalize beyond the target tones occupying the extreme ranges of entropy values. Analyzing dwell times for all tones also allowed us to directly compare the effects of prospective entropy vs. retrospective information content (IC). Recall that IC was matched across target tones in the previous analysis.

**13**

233    Model comparisons on models refitted with Maximum Likelihood estimates found that the IC

234    model predicted dwell times significantly better than the null model, $\chi^2(1) = 31.77$, $p < .001$. Adding

235    entropy improved the fit significantly, $\chi^2(1) = 16.64$, $p < .001$. In the full model, log-transformed

236    dwell times increased significantly with IC, $F(1, 19711.3) = 35.26$, $p < .001$, entropy, $F(1, 19711.2)$

237    $= 16.64$, $p < .001$, and marginally non-significantly with tone index in the phrase, $F(1, 37.5) = 3.30$,

238    $p = .077$.

239

# Experiment 2: Explicit completeness ratings

241    In Experiment 1, participants dwelled longer on tones affording high-entropy continuations than on

242    tones affording low-entropy continuations, regardless of whether they were originally phrase

243    beginnings or endings. This suggests that when rhythmic and metrical cues are removed from the

244    musical surface, entropic peaks in prospective pitch expectancy elicit implicit segmentation. Previous

245    dwell-time studies have demonstrated that longer dwell times coincide with perceived boundaries

246    (e.g., Hard et al., 2011), but Experiment 1 did not ~~provide concrete evidence~~guarantee that

247    participants were segmenting the stimuli. Therefore, Experiment 2 was designed to provide

248    converging evidence for effects of ~~prospective~~ prediction on segmentation using an explicit self-

249    report measure of phrase completeness (Palmer & Krumhansl, 1987).

250

**Methods**

252    *Participants.*  Thirty-one McMaster University students (not participants in Experiment 1)

253    took part in Experiment 2. Again, none were professional musicians (see SOM-R2 for more

254    information). This sample size exceeds those from previous studies using this methodology to assess

255    a comparable contrast (e.g., Palmer & Krumhansl, 1987). One participant declined to report their

256    gender and age, but among the remaining participants, the average age was 18.93 years ($SD_{age} = 2.51$

257 years), with 7 men and 23 women. Of the 31 participants, responses from five individuals were

258 omitted due to uninterpretable response sheets (i.e., multiple answers for each sequence, lacking

259 answers for certain sequences).

260 *Stimuli.* Melodic stimulus sequences were identical to those for Experiment 1, except

261 that all notes were played with a constant duration of 400 ms. Unlike in Experiment 1, the target tone

262 was always the final tone in the sequence.

263 *Procedure.* As in Experiment 1, the procedure took place in a sound-attenuating room.

264 Rather than self-pacing through the sequences as in Experiment 1, participants listened to all 56

265 sequences in randomized order. After each sequence, participants rated how complete the sequence

266 sounded (ranging from 1: "totally incomplete" to 7: "totally complete"). If the end of the melody was

267 completely satisfactory, that would constitute a score of 7, but if the melody ended in a way that was

268 implausible and unsatisfactory, that would constitute a score of 1. Participants were encouraged to
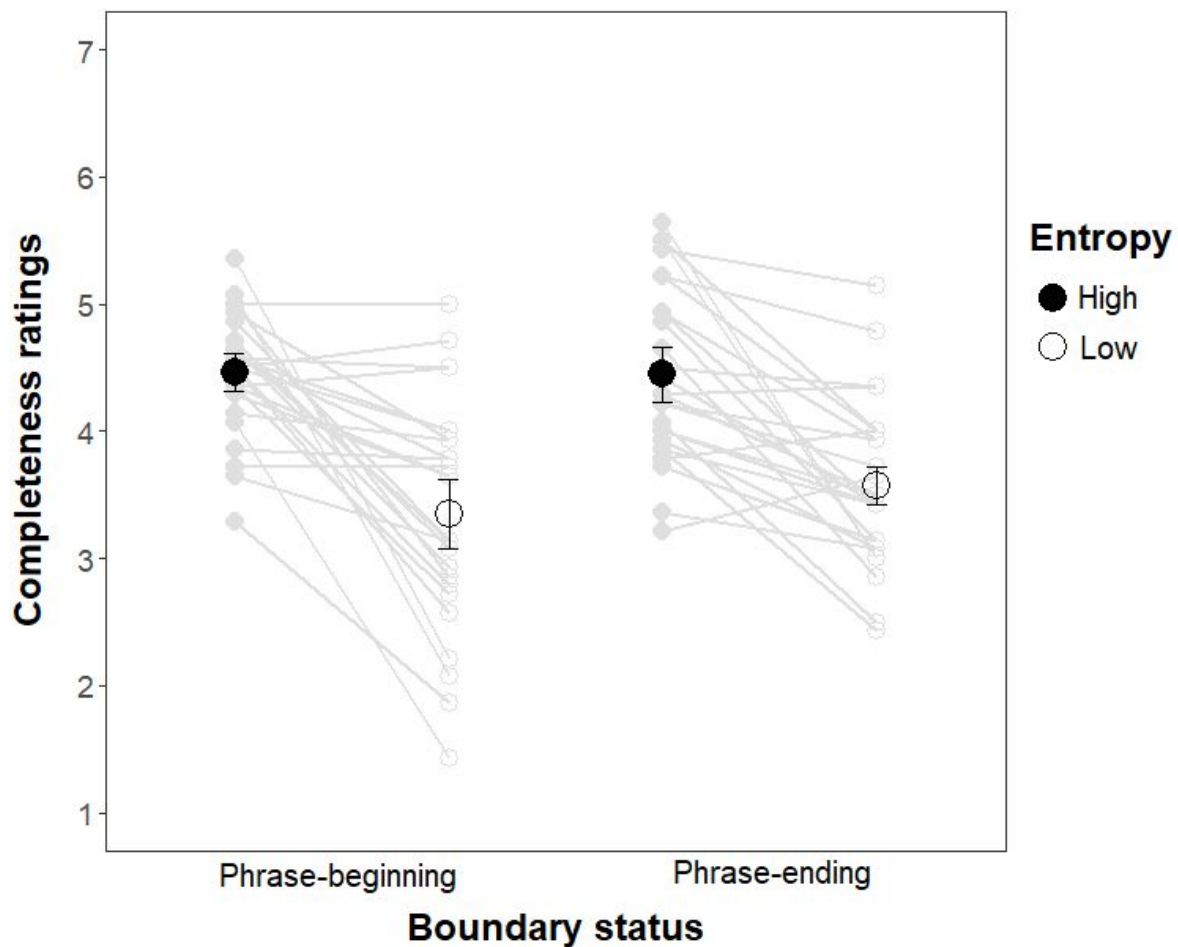
269 use the full range of the scale.

270

**Results**

272 A 2x2 repeated-measures ANOVA with factors boundary status (phrase-ending, phrase-

273 beginning) and entropy (high, low) was run on mean condition-wise ratings. Results were fully

274 consistent with those for Experiment 1. Specifically, no significant interaction ($F(1,25) = 1.80$, $p =$

275 $.192$, $\eta^2_p = .067$) nor main effect of boundary status ($F(1,25) = 0.82$, $p = .373$, $\eta^2_p = .032$) was found,

276 whereas there was a significant main effect of entropy ($F(1, 25) = 44.11$, $p < .001$, $\eta^2_p = .638$). High-

277 entropy target tones were rated as constituting more complete phrase endings than low-entropy target

278 tones, regardless of phrase position in the original chorale melody (Figure 3).

279 Again, no significant associations with musical sophistication were observed (see SOM-R2

280 for more details).

**15**

281



**Figure 3.** Completeness ratings for each type of excerpt (BegHi, BegLo, EndHi, and EndLo) in Experiment 2. Error

bars represent within-subject 95% confidence intervals (Cousineau, 2005). Stimulus sequences with final tones

generating high entropy were generally deemed more complete than those generating low entropy. It made no

significant difference whether tones originated from phrase beginnings or phrase endings in the original chorale

melodies.

## General Discussion

Although prediction is a fundamental component in influential theories of perceptual organization

(Hutchinson & Barrett, 2019), evidence for the role of ~~prospective~~ uncertainty ~~(a prospective measure~~

~~of prediction)~~ remains weak due to the empirical focus on retrospective measures of surprise (Hansen

293    & Pearce, 2014). Here we tested the hypothesis that uncertainty relates to boundary perception in

294    auditory sequences, using stimuli from Western tonal music ~~in which~~with well-defined phrase

295    boundaries ~~are well-defined~~. Sequences ~~that ended~~ ending on tones generating high-entropy

296    expectations were perceived as more complete than those ending on tones generating low-entropy

297    expectations (Experiment 2). This was also indicated by longer dwell times on high-entropy target

298    tones ~~generating high entropy~~ ~~;~~ and, indeed, across all tones in the stimulus sequences, entropy

299    explained unique variance in dwell times not ~~already~~ accounted for by event probability (Experiment

300    1).

301         Our work raises the key question why segmentation follows peaks of ~~statistical~~ uncertainty.

302    Christiansen and Chater's (2016) *Now-or-Never Bottleneck* posits that information ~~currently~~ in

303    working memory needs to be processed ~~here and~~ now or be forever lost. This constraint necessitates

304    "chunk-and-pass" processing whereby fleeting input—such as the content of music, speech, or action

305    sequences—is quickly segmented and encoded as higher-level representational units. Following from

306    ~~Christiansen and Chater's (2016)~~this theory, ~~it is possible that~~ events that afford high-entropy

307    predictions may require more bits to encode and thus may require higher working memory

308    deployment. The likelihood of exceeding memory capacity is higher after high-uncertainty events

309    than after low-uncertainty events, ~~leading to a~~causing higher probability of "chunking" and

310    perceiving a segment boundary.

311         This framework may also explain ~~the~~ previously demonstrated "dwell time" effect~~s~~ ~~observed~~

312    ~~in previous studies~~ (Hard et al., 2011, 2019; Kosie & Baldwin, 2019a, 2019b; Kragness & Trainor,

313    2016, 2018), since there is a time delay associated with segmentation and reintegration into previous

314    knowledge. This reintegration process, however, may have a cost. Specifically, taking in new

315    information is harder while reintegration takes place. Because the human mind aims to be one step

316    ahead, it will attempt to balance this cost optimally. Therefore, pauses in the stimulus stream may

317   induce a chunk to be processed even if it ends on low uncertainty (without fully exceeding working

318   memory capacity). This may constitute one ~~of the~~ potential mechanism~~s~~ explaining why Gestalt-like

319   principles of temporal proximity generally seem to apply to auditory sequence processing (Lerdahl

320   & Jackendoff, 1983).

321         The relatively high working memory capacity required at phrase boundaries may explain

322   previously observed *phrase-final lengthening.* Specifically, across ~~a variety of~~various languages,

323   musical instruments, and performance contexts, speakers and performers tend to slow down at phrase

324   endings (speech: Wightman et al., 1992; music: Palmer, 1989; Repp, 1992). While originally

325   interpreted as a communicative gesture in music (Palmer, 1989), piano performers exhibit phrase-

326   final lengthening even when attempting to play without expression (Penel & Drake, 1998). Combined

327   with the observation that listeners are less prone to detect lengthening on boundary tones than within-

328   phrase tones (Repp, 1992), ~~this led~~ Penel and Drake (1998) ~~to~~ hypothesize~~d~~ that perceptual biases

329   contribute to group-final lengthening, although the source of this bias remained unspecified. ~~We~~

330   ~~propose that O~~one such source could be processing constraints due to ~~predictive~~ uncertainty, which

331   likely apply across ~~multiple~~ domains of sequential perception and production.

332         Here we specifically focused on modelling the uncertainty of a single feature, pitch, as a cue

333   for phrase closure. Of course, the probabilistic characteristics of many other features (for instance,

334   temporal, spectral, syntactic, etc.) might affect ~~the perception of~~ completeness perception. In music,

335   these might include duration, intensity, inter-onset intervals, and performer gestures (Lerdahl &

336   Jackendoff, 1983). Whether ~~predictive~~ uncertainty in temporal features influences musical phrase

337   grouping remains to be tested. However, given that sensory systems prioritize anticipatory ~~processing~~

338   over reactive processing (Christiansen & Chater, 2016; Hutchinson & Barrett, 2019), it seems

339   plausible that our findings should extend to the temporal domain. On the other hand, non-probabilistic

340   and non-pitch-related features may also constrain the statistical learning giving rise to the entropy

341  effects found here, as observed in speech segmentation (Yang, 2004). Incorporating metrical

342  structure, previously heard motives, and limiting the number of accented tones per phrase would, for

343  example, most likely improve the predictive power of our entropy-based model. Future work should

344  more directly contrast the effect of anticipatory vs. adaptive cues and of probabilistic (top-down) vs.

345  Gestalt-related (bottom-up) cues to establish their relative contribution and investigate how this may

346  vary under different experimental conditions.

347      Another concern is whether IDyOM accurately reflects listener expectations. Morgan et al.

348  (2019) found that IDyOM predictions entailed higher entropy than that computed across several

349  participants ~~who~~ provid~~ing~~ed single-tone sung continuations to melodic contexts. Task constraints

350  likely explain this discrepancy as expectations for multiple continuations were not assessed.

351  Furthermore, by manipulating entropy of upcoming events rather than simply analyzing the entropy

352  of instantiated continuations, the present study differs crucially from Morgan et al. (2019). Moreover,

353  whereas they recruited self-identified musicians, who make melodic predictions with demonstrably

354  lower average entropy than non-musicians (Hansen & Pearce, 2014; Hansen, Vuust, & Pearce, 2016),

355  IDyOM was configured to model expectations of the general population. At the same time, Morgan

356  et al. (2019) made an important contribution by demonstrating a greater contribution of statistical

357  learning than of Gestalt-based principles in predicting listener expectations. This supports IDyOM's

358  suitability in predicting auditory boundary perception.

359      The finding that ~~predictive~~ uncertainty influences phrase boundary perception suggests a

360  pertinent role for training effects. Expertise effects may be particularly prominent in the musical

361  domain where skills and experiences differ substantially between individuals. Although some

362  ~~previous~~ studies suggest limited effects of musical expertise on melodic segmentation processes

363  (Palmer & Krumhansl, 1987, but see Hartmann et al., 2017), expertise levels have not always been

364  widely sampled or manipulated systematically. The same limitation applies to the current study where

**19**

365 no significant effects of expertise were seen (see Tables S2 and S3 in SOM-R2 for details). Yet, recent

366 research shows that stylistic specialization results in expectations about melodic continuations that

367 are generally lower in entropy whenever greater confidence is warranted (Hansen & Pearce, 2014;

368 Hansen et al., 2016). The transformation of high-entropy predictions into low-entropy predictions

369 with domain-relevant training or implicit exposure should allow musicians to perceive phrasal

370 coherence across longer timespans. This would be consistent with observations that experts have

371 access to more abstract and deeper levels of hierarchical structure (Chaffin & Imreh, 2002; Chi &

372 Feltovich, 1981) which, in turn, may be associated with larger working memory capacity (Meinz &

373 Hambrick, 2010). While awaiting sampling across more diverse expertise levels in future research,

374 our results relating chunk size to underlying expectancy dynamics enables a novel interpretation of

375 classical findings pertaining to expertise and working memory.

376 By offering an empirical challenge to the view that segmentation primarily relies on

377 retrospective processes, the present work contributes to the emergence of an increasingly coherent

378 model of the human mind as an eager predictive processor of sensory input. Embedded in the constant

379 flux of time, the mind is continually forced to evaluate and recombine retrospective and prospective

380 cues according to their immediate usefulness, and we hypothesize that sequential input in such varied

381 domains as language, music, and visual action sequences are all subject to the constraints arising from

382 this mental machinery.

# References

383     Chaffin, R., & Imreh, G. (2002). Practicing perfection: piano performance as expert memory.

385         *Psychological Science*, *13*(4), 342-349.

386     Chi, M. T., Feltovich, P. J., & Glaser, R. (1981). Categorization and representation of physics

387         problems by experts and novices. *Cognitive Science*, *5*(2), 121-152.

388     Christiansen, M. H., & Chater, N. (2016). The now-or-never bottleneck: a fundamental constraint

389         on language. *Behavioral and Brain Sciences*, *39*, e62. doi:10.1017/S0140525X1500031X.

390     Cohen, P., Adams, N., & Heeringa, B. (2007). Voting experts: an unsupervised algorithm for

391         segmenting sequences. *Intelligent Data Analysis*, *11*(6), 607-625.

392     Cousineau, D. (2005). Confidence intervals in within-subject designs: a simpler solution to Loftus

393         and Masson's method. *Tutorials in Quantitative Methods for Psychology*, *1*(1), 42-45.

394     Creighton, H. (ed.). (1966). *Songs and Ballads from Nova Scotia*. New York, NY: Dover.

395     Deliege, I. (1987). Grouping conditions in listening to music: an approach to Lerdahl &

396         Jackendoff's grouping preference rules. *Music Perception*, *4*(4), 325-359.

397     Dörfell (ed.) (1875). *371 vierstimmige Choralgesänge von Johann Sebastian Bach* (4th ed.).

398         Leipzig, Germany: Breitkopf & Härtel.

399     Hansen, N. C.,& Pearce, M. (2014). Predictive uncertainty in auditory sequence processing.

400         *Frontiers in Psychology* 5, 1052.

401     Hansen, N. C., Vuust, P., & Pearce, M. (2016). "If you've got to ask, you'll never know": Style-

402         congruent musical expertise optimises predictive auditory processing. *PLOS ONE*, *11*(10):

403         e0163584. doi:10.1371/journal.pone.0163584

404     Hansen, N. C. ,Vuust, P., Pearce, M., & Huron, D. (2017, August). *Entropic Ebbs and Flows: The*

405         *Expectancy Dynamics of Musical Phrases*. Paper presented at the Society for Music Perception

406         and Cognition Meeting, San Diego, CA.

Hard, B. M., Meyer, M., & Baldwin, D. (2019). Attention reorganizes as structure is detected in

dynamic action. *Memory & Cognition*, *47*(1), 17-32.

Hard, B. M., Recchia, G., & Tversky, B. (2011). The shape of action. *Journal of Experimental*

*Psychology: General, 140*(4), 586-604. doi:10.1037/a0024310

Hartmann, M., Lartillot, O. & Toiviainen, P. (2017). Interaction features for prediction of

perceptual segmentation: effects of musicianship and experimental task. *Journal of New Music*

*Research, 46*(2), 156-174. doi:10.1080/09298215.2016.1230137

Hutchinson, J. B., & Barrett, L. F. (2019). The power of predictions: an emerging paradigm for

psychological research. *Current Directions in Psychological Science*, *28*(3), 280-291.

Koelsch, S., Vuust, P., & Friston, K. (2019). Predictive processes and the peculiar case of music.

*Trends in Cognitive Sciences*, *23*(1), 63-77.

Kosie, J. E., & Baldwin, D. (2019a). Attention rapidly reorganizes to naturally occurring structure

in a novel activity sequence. *Cognition*, *182*, 31–44. doi:10.1016/j.cognition.2018.09.004

Kosie, J. E., & Baldwin, D. (2019b). Attentional profiles linked to event segmentation are robust to

missing information. *Cognitive Research: Principles and Implications*, *4*(1), 8.

doi:10.1186/s41235-019-0157-4

Kragness, H. E. & Trainor, L. J. (2016). Listeners lengthen phrase boundaries in self-paced music.

*Journal of Experimental Psychology: Human Perception and Performance*, *42*(10), 1676-1686.

doi:10.1037/xhp0000245

Kragness, H. E. & Trainor, L. J. (2018). Young children pause on phrase boundaries in self-paced

music listening: the role of harmonic cues. *Developmental Psychology*, 54(5), 842-856.

doi:10.1037/dev0000405

Kurby, C. A., & Zacks, J. M. (2008). Segmentation in the perception and memory of events. *Trends*

*in Cognitive Sciences*, *12*(2), 72-79.

431  Lerdahl, F., & Jackendoff, R. (1983). *A generative theory of tonal music*. Cambridge, MA: MIT

432      Press.

433  Levinson, S. C. (2016). Turn-taking in human communication: origins and implications for

434      language processing. *Trends in Cognitive Sciences, 20*(1), 6-14. doi:10.1016/j.tics.2015.10.010

435  Meinz, E. J., & Hambrick, D. Z. (2010). Deliberate practice is necessary but not sufficient to

436      explain individual differences in piano sight-reading skill: the role of working memory capacity.

437      *Psychological Science*, *21*(7), 914-919.

438  Morgan, E., Fogel, A., Nair, A., & Patel, A. D. (2019). Statistical learning and Gestalt-like

439      principles predict melodic expectations. *Cognition*, *189*, 23-34.

440  Nicholson, S., Knight, G. H., and Bower, J. D. (Ed.). (1950). *Ancient and Modern Revised*. Suffolk,

441      UK: William Clowes and Sons.

442  Palmer, C. (1989). Mapping musical thought to musical performance. *Journal of Experimental

443      Psychology: Human Perception and Performance*, *15*(2), 331.

444  Palmer, C., & Krumhansl, C. L. (1987). Independent temporal and pitch structures in determination

445      of musical phrases. *Journal of Experimental Psychology: Human Perception and Performance*,

446      *13*(1), 116.

447  Pearce, M. T. (2005). *The construction and evaluation of statistical models of melodic structure in

448      music perception and composition* (Doctoral dissertation). City University, London, UK.

449      Retrieved from https://openaccess.city.ac.uk/id/eprint/8459/1/

450  Pearce, M. T., Müllensiefen, D., & Wiggins, G. (2010). The role of expectation and probabilistic

451      learning in auditory boundary perception: a model comparison. *Perception, 39*(10), 1367-1391.

452      doi:10.1068/p6507

453  Penel, A., & Drake, C. (1998). Sources of timing variations in music performance: a psychological

454      segmentation model. *Psychological Research*, *61*(1), 12-32.

R Core Team (2019). *R: A language and environment for statistical computing*. R Foundation for

Statistical Computing, Vienna, Austria. Retrieved from https://www.R-project.org/.

Repp, B. H. (1992). Probing the cognitive representation of musical time: structural constraints on

the perception of timing perturbations. *Cognition*, *44*(3), 241-281.

Richmond, L. L., & Zacks, J. M. (2017). Constructing experience: event models from perception to

action. *Trends in Cognitive Sciences*, 21(12), 962-980.

Saffran, J. R., & Kirkham, N. Z. (2018). Infant statistical learning. *Annual Review of Psychology*,

*69*, 181-203. doi:10.1146/annurev-psych-122216-011805

Schaffrath, H. (1995). *The Essen Folksong Collection in the Humdrum Kern Format* (D. Huron,

Ed.). Menlo Park, CA: Center for Computer Assisted Research in the Humanities. Retrieved

from https://kern.humdrum.org/cgi-bin/browse?l=essen/europa/deutschl

Wightman, C. W., Shattuck-Hufnagel, S., Ostendorf, M., & Price, P. J. (1992). Segmental durations

in the vicinity of prosodic phrase boundaries. *The Journal of the Acoustical Society of America*,

*91*(3), 1707–1717. doi:10.1121/1.402450

Yang, C. D. (2004). Universal grammar, statistics or both? *Trends in Cognitive Sciences*, *8*(10),

451-456.

Zacks, J. M., & Swallow, K. M. (2007). Event segmentation. *Current Directions in Psychological

Science*, *16*(2), 80-84.

Zacks, J. M., Tversky, B., & Iyer, G. (2001). Perceiving, remembering, and communicating

structure in events. *Journal of Experimental Psychology: General*, *130*(1), 29-58.

**24**

# Predictive uncertainty underlies auditory boundary perception

## Abstract

Anticipating the future is essential for efficient perception and action planning. Yet, the role of anticipation in event segmentation is understudied because empirical research has focused on retrospective cues such as surprise. We address this question in the context of musical phrase-boundary perception. A computational model of cognitive sequence processing was used to control the information-dynamic properties of tone sequences. In an implicit, self-paced listening task ($n$=38), undergraduates dwelled longer on tones generating high entropy (i.e., high uncertainty) than those generating low entropy (i.e., low uncertainty). Similarly, sequences that ended on tones generating high entropy were rated as sounding more complete ($n$=31). These entropy effects were independent of both the surprise (i.e., information content) and phrase position of target tones in the original musical stimuli. Our results indicate that events generating high entropy prospectively contribute to segmentation processes in auditory sequence perception, independent of the properties of the subsequent event.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

16

# Statement of relevance

17 A significant challenge for the human perceptual system is to promote time-sensitive, context-

18 appropriate responses by predictively processing continuous streams of complex sensory information.

19 A large body of research shows that expectations gleaned from a lifetime of experience guide such

20 processes, which are critical in high-risk environments like traffic or manual labor. Because most

21 studies have focused on the degree of surprise evoked by events, there is little evidence for the role

22 of prospective expectations in perceptual organization. Here, we control entropy in musical tone

23 sequences by using an information-theoretic model that has been shown to reflect listeners' predictive

24 uncertainty. Tones that afforded relatively high uncertainty were found to draw implicit attention and

25 influence explicit ratings of sequence completeness. Focusing attention on instances where upcoming

26 events are statistically unconstrained could contribute to an adaptive mechanism facilitating stream

27 segmentation that leads to efficient learning and information processing in a complex, dynamic world.

2

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

28

# Acknowledgments

29    

# Introduction

Humans make sense of a complex, dynamic world by segmenting sequences of events into manageable units (Zacks & Swallow, 2007; Kurby & Zacks, 2008; Richmond & Zacks, 2017). Past work on segmentation has focused on retrospective cues for boundary identification, often conceptualizing group boundaries as coinciding with instances of increased relative change in stimulus features or low transition probabilities (e.g., speech: Saffran & Kirkham, 2018; action sequences: Hard et al., 2011; music: Hartmann et al., 2017; Pearce et al. 2010). However, the sophisticated prediction capabilities of the human mind (Hutchinson & Barrett, 2019) suggest that event boundaries are also anticipated. For example, in natural conversation, turn-taking happens so rapidly that speakers likely anticipate the end of their conversation partner's sentence (Levinson, 2016). Here we investigate the role of entropy, or degree of uncertainty about an upcoming event, in determining the perception of group boundaries in auditory sequences. We define *prediction* as the psychological processes of generating an expectation about a future event in terms of how likely various possible outcomes are. We define *uncertainty* as the imprecision (or extent of equi-probability) of such a prediction.

Though most previous work has focused on retrospective boundary identification, anticipatory processing has some preliminary support. When self-pacing through sequential images of action sequences, participants tend to "dwell" (or pause) on perceived boundary images (Hard et al., 2011; Hard et al., 2019; Kosie & Baldwin, 2019a, 2019b). Kosie and Baldwin (2019b) proposed that this "dwell time effect" resulted from selective attention to moments of uncertainty afforded by perceiving a goal completion event. No cognitive model was devised to test this theory, however, potentially due to challenges in modeling expectancy in event processing of action sequences. Indeed, one methodological drawback was demonstrated by participants' dwelling on boundary slides even when those slides were out of order, suggesting that they were responding to conceptual salience

54    rather than to underlying expectancy dynamics (Hard et al., 2011). Cohen et al. (2007) have proposed

55    an entropy-based segmentation model for language, but because it computes statistics from the corpus

56    it is segmenting—including parts it has not yet seen—it does not fully capture segmentation

57    processing in real time (Christiansen & Chater, 2016).

58      Because music is not only hierarchically structured (Lerdahl & Jackendoff, 1983), but also

59    statistically well-defined, it is an ideal domain for testing psychological theories of probabilistic

60    perception (Koelsch, Vuust, & Friston, 2019). As with non-musical sequences (Zacks et al., 2001),

61    there is generally high inter-participant agreement regarding the location of musical phrase

62    boundaries (Deliège, 1987; but see Pearce et al., 2010), and as with action sequences, listeners self-

63    pacing through musical chords "dwell" on boundary chords (Kragness & Trainor, 2016, 2018). Since,

64    however, entropy correlates strongly with phrase boundaries in music (Hansen et al., 2017), previous

65    studies were not optimized to separate prospective effects of expectancy dynamics from effects of

66    canonical boundary features on perceptual grouping. *Information Dynamics of Music* (IDyOM)

67    (Pearce, 2005) is a computational model of auditory expectation which enables modelling boundary

68    perception quantitatively using the information-theoretic concepts of entropy and information

69    content, computed in reference to pre-existing long-term knowledge (Hansen & Pearce, 2014; Hansen

70    et al., 2016). Entropy facilitates a test of uncertainty as a prospective mechanism for boundary

71    perception which can be pitted directly against information content (a measure of surprise) as a

72    retrospective cue. For example, an individual may form a highly certain prediction about the next

73    note in a melody but then be surprised when a different note follows. Another advantage of melodic

74    sequences is that any given note has little intrinsic meaning in isolation from its preceding musical

75    context, ensuring that observed effects on perception reflect the statistical structure of the sequence

76    and not inherent features of the boundary stimulus itself. However, because uncertainty is not always

77   available for explicit introspection (Hansen et al., 2016), implicit measures are paramount for

78   investigating the cognitive mechanisms underlying boundary perception.

79          The present study used IDyOM to control the information-dynamic properties of melodic

80   sequences in two experiments assessing the role of uncertainty in sequence processing. We measured

81   participants' dwell times (Experiment 1) and explicit ratings of phrase completeness (Experiment 2)

82   for tones that afforded high/low entropy and were phrase- beginning/phrase-ending in the melodies

83   from which they were drawn. We predicted that tones generating high uncertainty would lead to

84   longer dwell times and higher ratings of phrase completeness, regardless of original phrase status,

85   and that this effect would be independent from retrospective surprise.

86

## Experiment 1: Implicit Self-Pacing Task

**Methods**

89          *Participants.*   Thirty-eight McMaster University undergraduates received psychology course

90   credits for participating in the study ($M_{age}$ = 19.3 years, 1 person declined to report their age, $SD_{age}$ =

91   3.78, 8 men, 30 women). None of the participants were professional musicians (for more information

92   about musical training levels, see Table S1 in SOM-R2). This sample size exceeds or corresponds to

93   those of previous studies using this methodology to assess comparable effects (e.g., Hard et al., 2011;

94   Kragness & Trainor, 2016, 2018). All participants were fluent in English.

95          *Stimuli.*          Fifty-six monophonic stimulus sequences were selected from the soprano (i.e.,

96   highest) part in 370 four-part chorale harmonizations by Johann Sebastian Bach (Dörffel, 1875) (see

97   SOM-R1 for details of the stimulus selection procedure). These chorale melodies are not generally

98   known by present-day listeners in Canada. Unfamiliarity was, moreover, made more likely through

99   complete removal of rhythmic information by granting participants control over tone durations in the

100  self-paced dwell-time paradigm (Experiment 1) or by presenting stimuli with isochronized tone

101 durations (Experiment 2). All chords, interference tones, and self-pacing tones were generated in

102 MaxMSP's grand piano timbre.

103       Each stimulus context contained a full phrase (musical group) of seven to 17 pitches followed

104 by the initial tone of the subsequent phrase in the original chorale melody. Tones associated with

105 phrase beginnings and endings were unambiguously identified from notations in the musical score.

106 This practice seems at least as objective as the reliance on trained "expert coders" to determine event

107 boundaries in research using visual action sequences (e.g., Hard et al., 2019; Kosie & Baldwin, 2019a,

108 2019b). We included both phrase endings and phrase beginnings as target tones to provide a strong

109 test of entropy's role in segmentation, controlling for compositional cues in the melodies that might

110 signal melodic phrase endings in other ways.

111       Fourteen stimulus contexts were selected for each of the four experimental conditions,

112 comprising phrase beginnings with high ("BegHi") or low entropy ("BegLo") and phrase endings

113 with high ("EndHi") or low entropy ("EndLo"). Entropy, in this regard, quantifies the level of

114 uncertainty governing a listener's expectations about what the pitch of the next tone following the

115 relevant phrase beginning or phrase ending would be. Thus, Western-enculturated listeners are

116 expected to be relatively sure about which pitch will follow the target tone in "BegLo" and "EndLo"

117 contexts, but relatively unsure in "BegHi" and "EndHi" contexts. "Target tone", in this respect, refers

118 to the final tone in "BegLo" and "BegHi" contexts and the penultimate tone in "EndLo" and "EndHi"

119 contexts.

120       The entropy level generated by each tone in the corpus was estimated by the *Information*

121 *Dynamics of Music Model* (IDyOM, version 1.3) (Pearce, 2005). This variable-order *n*-gram model

122 uses unsupervised statistical learning to generate probability distributions governing a relevant

123 feature of each tone in a monophonic melody. IDyOM was trained on a large dataset of 5,332 German

124 folk songs (Schaffrath, 1995), 152 Nova Scotian songs and ballads (Creighton, 1966), and 120

125    English hymns (Nicholson et al., 1950)[1]. For each tone in the chorale melody, IDyOM generated a

126    probability distribution (summing to 1) over the 44 pitch values occurring in the training corpus (i.e.,

127    MIDI pitches 45-89 corresponding to A2-F6) by combining *n*-gram models of varying order. Entropy

128    then quantifies the shape of these probability distributions with high entropy for "flat" (relatively

129    uniform) distributions, where there is high uncertainty about the next event, and low entropy for

130    "spiky" (relatively nonuniform) distributions, where one or a small number of continuations are

131    highly probable.

132         The set of 56 stimulus contexts was selected in a way that prioritized extreme high or low

133    entropy values while ensuring that three conditions were met: First, as shown by a non-parametric

134    Kruskal-Wallis test, all four conditions, including EndHi (Median = 2.45, IQR = 1.76), BegHi

135    (Median = 2.69, IQR = 1.78), EndLo (Median = 2.31, IQR = 2.70), and BegLo (Median = 3.37, IQR

136    = 1.83), were matched on information content (i.e., inverse log-probability) for the event of interest, $\chi$

137    $^2(3)$ = 4.55, *p* = .208; second, as shown by Mann-Whitney *U*-tests, EndHi (Median = 2.97, IQR =

138    0.08) and BegHi (Median = 3.00, IQR = 0.12) stimuli, *U* = 78, *p* = .376, as well as EndLo (Median =

139    1.07, IQR = 0.30) and BegLo (Median = 0.97, IQR = 0.35) stimuli, *U* = 90, *p* = .734, were matched

140    on entropy governing the next event in the sequence. The experimenter selecting these stimuli paid

141    no attention to any other musical features.

142         For the secondary analysis of all tones in the stimulus set, IC and entropy were re-estimated

143    by re-running IDyOM with the same configuration on the final stimulus contexts. This was done

144    because IC and entropy estimates for the initial tones in each stimulus context sometimes relied on

145    tones from the preceding phrase in the original chorales, which was excluded from the stimuli used.

146    While unproblematic for stimulus selection based on target tones, this presented a problem for tone-

147    level analysis. Note that due to their late position in the tone sequences, target tone entropy and IC
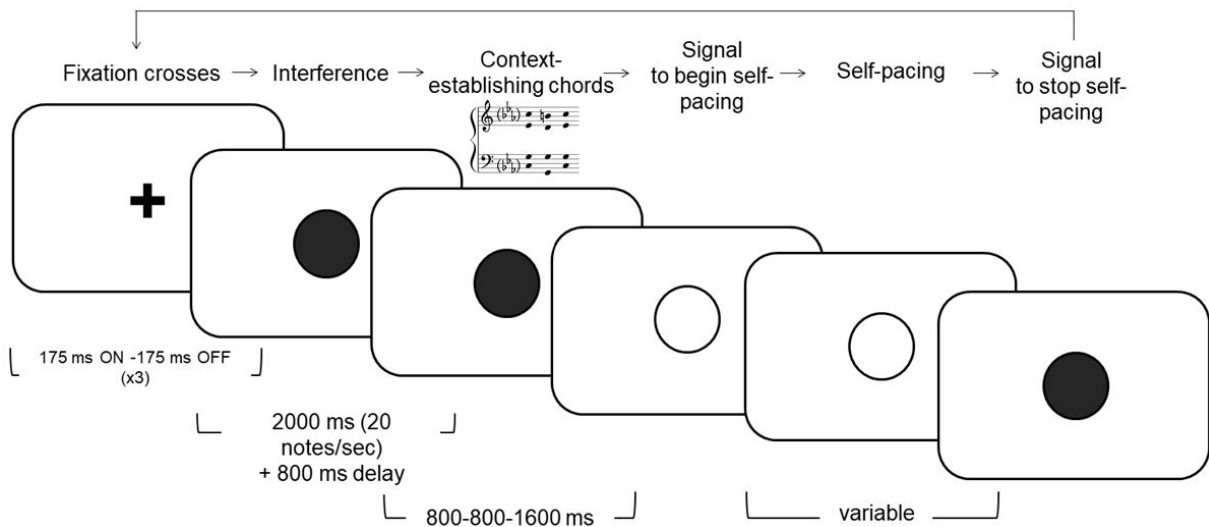
---

[1] For more information about the IDyOM implementation and parameters, please see SOM-R.

**8**

148 values were identical for the two models (one used in stimulus generation and analyses of target tones,

149 the other used in the analysis of all tones).

150   *Procedure.*  The experimental procedures (for Experiment 1 and 2) received prior approval

151 from the McMaster University Research Ethics Board and was carried out in accordance with the

152 provisions of the World Medical Association Declaration of Helsinki. Participants were seated facing

153 a computer screen in a sound-attenuated room. They were instructed to press the spacebar on a

154 computer keyboard with the pointer finger of their dominant hand to elicit the onset of each

155 subsequent tone in the sequence. Tones decayed naturally, but were not terminated until the spacebar

156 was pressed again to initiate the next tone. Participants were instructed to progress as quickly or

157 slowly as they liked while listening carefully, and could not repeat previously heard tones. They were

158 led to falsely believe that their memory for the sequences would be tested afterwards to motivate them

159 to attend to the task (Kragness & Trainor, 2016). No other instructions regarding timing, pacing,

160 rhythmicity, or expressivity were given. If a participant asked for further information, they were told

161 to play through the piece in a way that would maximize their performance in the subsequent memory

162 task.

163   Prior to each trial, participants saw three flashes of a fixation cross, then heard 40 50-ms tones

164 (for a total of 2000 ms) chosen randomly on each trial from range E2 to A5 to minimize carryover

165 from the context of the previous sequence, followed by three context-establishing chords with

166 durations of 800, 800, and 1600 ms (Figure 1). The context-establishing chords were played in the

167 key of the relevant melody. Throughout each trial, a circle on the screen indicated when to begin self-

168 pacing through the melody (light green) and when to stop (dark green).

169

**Figure 1.** Depiction of a trial from Experiment 1. In each trial, participants saw a fixation cross, followed by interference tones, then three context-establishing chords and a signal (white circle) to begin self-pacing. They then self-paced through the tone sequence until the occurrence of a stop signal (black circle). The box depicts examples of tone sequences from each condition containing target tones (boxed) generating relatively uncertain (high entropy) or relatively certain (low entropy) expectations about the pitch of the next tone, matched on IC of the current tone. The double slash indicates whether target tones were phrase beginnings (after double slash) or phrase endings (prior to double slash) in the original notation.

*Data processing and statistical analysis.*    Despite systematic efforts to avoid duplicate stimulus contexts (e.g., multiple occurrences of a repeated phrase from a single melody or identical

**10**

181  phrases across melodies), it was discovered after data collection that one melodic context occurred

182  both amongst the "BegHi" and "EndHi" stimulus sets (with different target tones). Given that results

183  did not differ substantially when excluding dwell times for these stimuli, we report statistical analyses

184  including the full dataset here, which included 56 total tone sequences (i.e., 14 per condition).

185       To mitigate effects of extreme data points, a minimum dwell-time threshold of 100 ms was

186  adopted for inclusion. Dwell times greater than 3 standard deviations above a participant's own

187  average (across all target and non-target dwell times) were also omitted (Kosie & Baldwin, 2019a,

188  2019b). These exclusion criteria eliminated an average of 1.31% of all tones and 1.70% of target

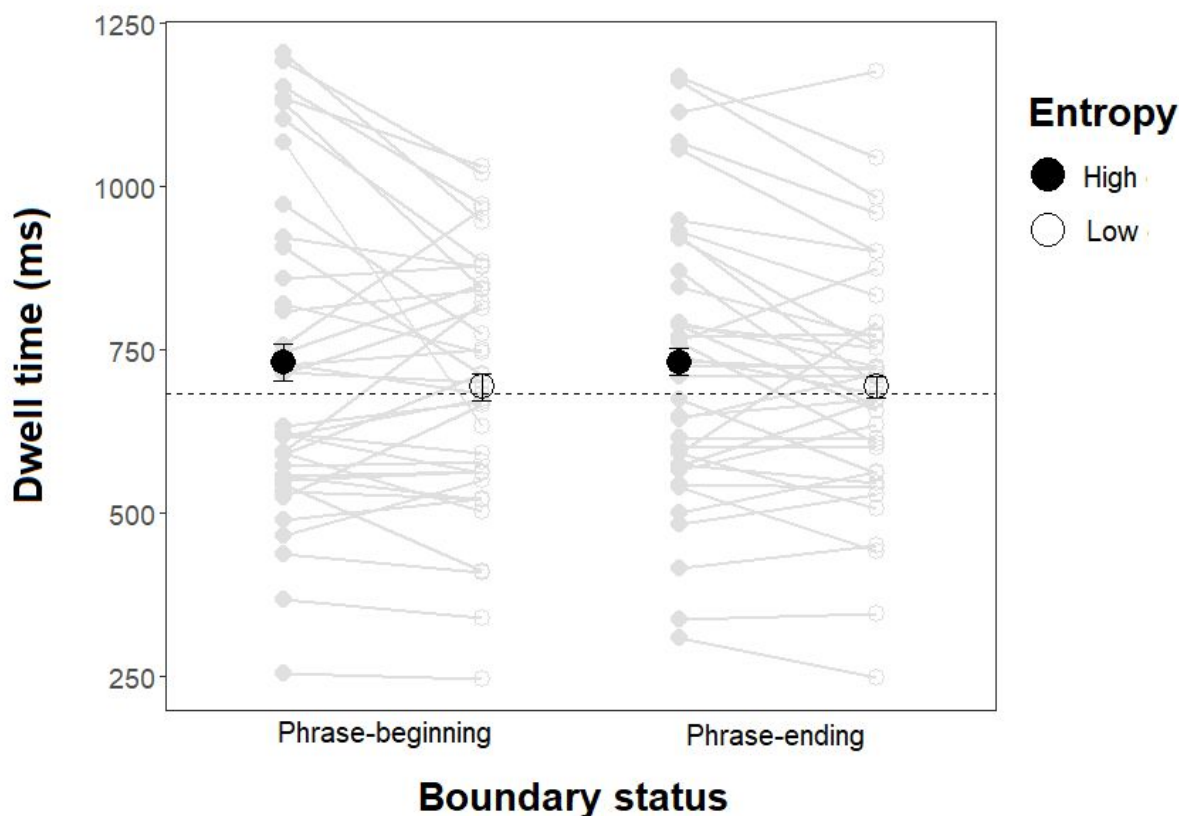189  tones per participant (ranging from 0-4 target tones).

190       For the main analysis of target tones, target dwell times were averaged by condition resulting

191  in four condition-wise means per participant. A 2x2 repeated-measures ANOVA (including within-

192  subjects factors boundary status and entropy) was run on target tone dwell times.

193       For the secondary analysis of all tones, dwell times were first log-transformed to minimize

194  the positive skew inherent to timing data (cf. Kragness & Trainor, 2018). Subsequently, using the

195  *lmer()* function from the *lme4* package in R (R Core Team, 2019), linear mixed-effects models were

196  fitted with Restricted Maximum Likelihood estimates (REML). Because previous experiments have

197  found that dwell times change systematically throughout trials (Kragness & Trainor, 2016), tone

198  index in the sequence was always included as a predictor. Thus, whereas the null model only included

199  tone index as a fixed effect, two further increasingly complex models added, first, the retrospective

200  cue IC, and, second, the prospective cue entropy. Thereby, we could determine whether prospective

201  predictive processing explained unique variance not already accounted for by retrospective surprise.

202  Random intercepts and slopes of tone number were included for each participant. For all models, this

203  random-effects structure produced the lowest BIC values while avoiding singular fits.

204

**11**

1
2
3
4
5

**Results**

6
7

*Target tones.*   To examine the effects of boundary status (phrase-ending, phrase-beginning)

8
9

and entropy (high, low), a 2x2 repeated-measures ANOVA was run on target tone dwell times.

10
11
12

Whereas no significant interaction ($F(1,37) < 0.01$, $p = .986$, $\eta^2_p < .001$) or main effect of boundary

13
14

status ($F(1,37) < 0.01$, $p = .973$, $\eta^2_p < .001$) was found, there was a significant main effect of entropy

15
16

($F(1,37) = 7.24$, $p = .011$, $\eta^2_p = .164$). Thus, as hypothesized, high-entropy target tones were generally

17
18
19

dwelled on longer than low-entropy target tones, regardless of phrase position in the original chorale

20
21

melody (Figure 2).

22
23
24

We conducted post-hoc correlational analyses to examine whether participants' musical

25
26

sophistication was associated with the magnitude of their dwell time effect. No significant

27
28

associations were observed (see SOM-R2 for more details).

29
30

216

31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

**12**

**Figure 2.** Dwell times (ms) for each type of target tone (BegHi, BegLo, EndHi, EndLo) in Experiment 1. The dashed line represents the average dwell time (683 ms) for non-target tones. Error bars represent within-subject 95% confidence intervals (Cousineau, 2005). High-entropy target tones had longer dwell times than low-entropy target tones, and it made no significant difference whether target tones originated from phrase endings or phrase beginnings in the original chorale melody corpus.

*All tones.*       If uncertainty provides a cognitive cue for phrase segmentation, its effect on dwell times should generalize beyond the target tones occupying the extreme ranges of entropy values. Analyzing dwell times for all tones also allowed us to directly compare the effects of prospective entropy vs. retrospective information content (IC). Recall that IC was matched across target tones in the previous analysis.

**13**

229    Model comparisons on models refitted with Maximum Likelihood estimates found that the IC

230   model predicted dwell times significantly better than the null model, $\chi^2(1) = 31.77$, $p < .001$. Adding

231   entropy improved the fit significantly, $\chi^2(1) = 16.64$, $p < .001$. In the full model, log-transformed

232   dwell times increased significantly with IC, $F(1, 19711.3) = 35.26$, $p < .001$, entropy, $F(1, 19711.2)$

233   $= 16.64$, $p < .001$, and marginally non-significantly with tone index in the phrase, $F(1, 37.5) = 3.30$,

234   $p = .077$.

235

# Experiment 2: Explicit completeness ratings

237   In Experiment 1, participants dwelled longer on tones affording high-entropy continuations than on

238   tones affording low-entropy continuations, regardless of whether they were originally phrase

239   beginnings or endings. This suggests that when rhythmic and metrical cues are removed from the

240   musical surface, entropic peaks in prospective pitch expectancy elicit implicit segmentation. Previous

241   dwell-time studies have demonstrated that longer dwell times coincide with perceived boundaries

242   (e.g., Hard et al., 2011), but Experiment 1 did not guarantee that participants were segmenting the

243   stimuli. Therefore, Experiment 2 was designed to provide converging evidence for effects of

244   prediction on segmentation using an explicit self-report measure of phrase completeness (Palmer &

245   Krumhansl, 1987).

246

**Methods**

248        *Participants.*   Thirty-one McMaster University students (not participants in Experiment 1)

249   took part in Experiment 2. Again, none were professional musicians (see SOM-R2 for more

250   information). This sample size exceeds those from previous studies using this methodology to assess

251   a comparable contrast (e.g., Palmer & Krumhansl, 1987). One participant declined to report their

252   gender and age, but among the remaining participants, the average age was 18.93 years ($SD_{age} = 2.51$

**14**

253 years), with 7 men and 23 women. Of the 31 participants, responses from five individuals were

254 omitted due to uninterpretable response sheets (i.e., multiple answers for each sequence, lacking

255 answers for certain sequences).

256 *Stimuli.* Melodic stimulus sequences were identical to those for Experiment 1, except

257 that all notes were played with a constant duration of 400 ms. Unlike in Experiment 1, the target tone
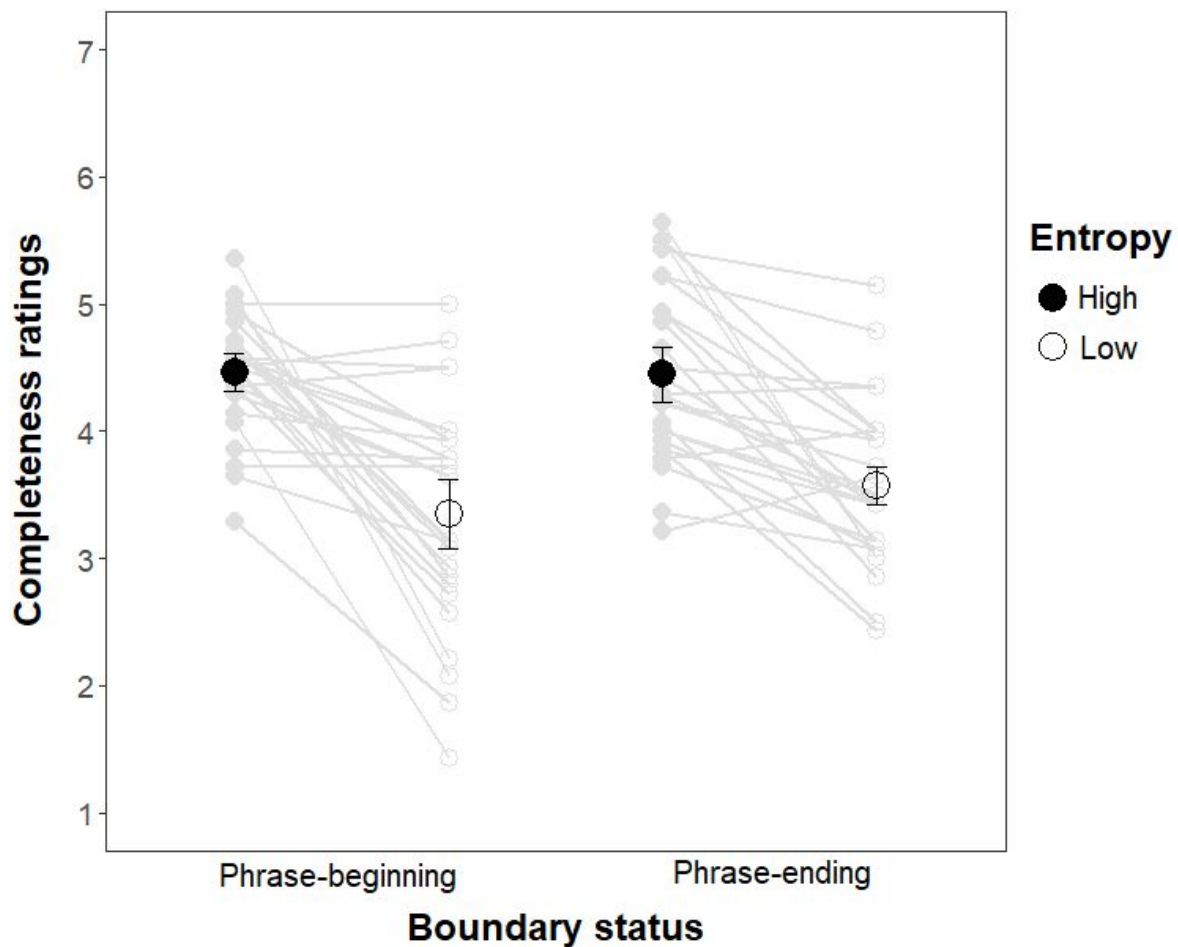
258 was always the final tone in the sequence.

259 *Procedure.* As in Experiment 1, the procedure took place in a sound-attenuating room.

260 Rather than self-pacing through the sequences as in Experiment 1, participants listened to all 56

261 sequences in randomized order. After each sequence, participants rated how complete the sequence

262 sounded (ranging from 1: "totally incomplete" to 7: "totally complete"). If the end of the melody was

263 completely satisfactory, that would constitute a score of 7, but if the melody ended in a way that was

264 implausible and unsatisfactory, that would constitute a score of 1. Participants were encouraged to

265 use the full range of the scale.

266

**Results**

268 A 2x2 repeated-measures ANOVA with factors boundary status (phrase-ending, phrase-

269 beginning) and entropy (high, low) was run on mean condition-wise ratings. Results were fully

270 consistent with those for Experiment 1. Specifically, no significant interaction ($F(1,25) = 1.80$, $p =$

271 $.192$, $\eta^2_p = .067$) nor main effect of boundary status ($F(1,25) = 0.82$, $p = .373$, $\eta^2_p = .032$) was found,

272 whereas there was a significant main effect of entropy ($F(1, 25) = 44.11$, $p < .001$, $\eta^2_p = .638$). High-

273 entropy target tones were rated as constituting more complete phrase endings than low-entropy target

274 tones, regardless of phrase position in the original chorale melody (Figure 3).

275 Again, no significant associations with musical sophistication were observed (see SOM-R2

276 for more details).

**15**

277



**Figure 3.** Completeness ratings for each type of excerpt (BegHi, BegLo, EndHi, and EndLo) in Experiment 2. Error bars represent within-subject 95% confidence intervals (Cousineau, 2005). Stimulus sequences with final tones generating high entropy were generally deemed more complete than those generating low entropy. It made no significant difference whether tones originated from phrase beginnings or phrase endings in the original chorale melodies.

## General Discussion

Although prediction is a fundamental component in influential theories of perceptual organization (Hutchinson & Barrett, 2019), evidence for the role of uncertainty remains weak due to the empirical focus on retrospective measures of surprise (Hansen & Pearce, 2014). Here we tested the hypothesis

**16**

289    that uncertainty relates to boundary perception in auditory sequences, using stimuli from Western

290    tonal music with well-defined phrase boundaries. Sequences ending on tones generating high-entropy

291    expectations were perceived as more complete than those ending on tones generating low-entropy

292    expectations (Experiment 2). This was also indicated by longer dwell times on high-entropy target

293    tones; indeed, across all tones in the stimulus sequences, entropy explained unique variance in dwell

294    times not accounted for by event probability (Experiment 1).

295        Our work raises the key question why segmentation follows peaks of uncertainty. Christiansen

296    and Chater's (2016) *Now-or-Never Bottleneck* posits that information in working memory needs to

297    be processed now or be forever lost. This constraint necessitates "chunk-and-pass" processing

298    whereby fleeting input—such as the content of music, speech, or action sequences—is quickly

299    segmented and encoded as higher-level representational units. Following from this theory, events that

300    afford high-entropy predictions may require more bits to encode and thus may require higher working

301    memory deployment. The likelihood of exceeding memory capacity is higher after high-uncertainty

302    events than after low-uncertainty events, causing higher probability of "chunking" and perceiving a

303    segment boundary.

304        This framework may also explain previously demonstrated "dwell time" effects (Hard et al.,

305    2011, 2019; Kosie & Baldwin, 2019a, 2019b; Kragness & Trainor, 2016, 2018), since there is a time

306    delay associated with segmentation and reintegration into previous knowledge. This reintegration

307    process, however, may have a cost. Specifically, taking in new information is harder while

308    reintegration takes place. Because the human mind aims to be one step ahead, it will attempt to

309    balance this cost optimally. Therefore, pauses in the stimulus stream may induce a chunk to be

310    processed even if it ends on low uncertainty (without fully exceeding working memory capacity).

311    This may constitute one potential mechanism explaining why Gestalt-like principles of temporal

312    proximity generally seem to apply to auditory sequence processing (Lerdahl & Jackendoff, 1983).

313    The relatively high working memory capacity required at phrase boundaries may explain

314 previously observed *phrase-final lengthening.* Specifically, across various languages, musical

315 instruments, and performance contexts, speakers and performers tend to slow down at phrase endings

316 (speech: Wightman et al., 1992; music: Palmer, 1989; Repp, 1992). While originally interpreted as a

317 communicative gesture in music (Palmer, 1989), piano performers exhibit phrase-final lengthening

318 even when attempting to play without expression (Penel & Drake, 1998). Combined with the

319 observation that listeners are less prone to detect lengthening on boundary tones than within-phrase

320 tones (Repp, 1992), Penel and Drake (1998) hypothesized that perceptual biases contribute to group-

321 final lengthening, although the source of this bias remained unspecified. One such source could be

322 processing constraints due to uncertainty, which likely apply across domains of sequential perception

323 and production.

324    Here we specifically focused on modelling the uncertainty of a single feature, pitch, as a cue

325 for phrase closure. Of course, the probabilistic characteristics of many other features (for instance,

326 temporal, spectral, syntactic, etc.) might affect completeness perception. In music, these might

327 include duration, intensity, inter-onset intervals, and performer gestures (Lerdahl & Jackendoff,

328 1983). Whether uncertainty in temporal features influences musical phrase grouping remains to be

329 tested. However, given that sensory systems prioritize anticipatory over reactive processing

330 (Christiansen & Chater, 2016; Hutchinson & Barrett, 2019), it seems plausible that our findings

331 should extend to the temporal domain. On the other hand, non-probabilistic and non-pitch-related

332 features may also constrain the statistical learning giving rise to the entropy effects found here, as

333 observed in speech segmentation (Yang, 2004). Incorporating metrical structure, previously heard

334 motives, and limiting the number of accented tones per phrase would, for example, most likely

335 improve the predictive power of our entropy-based model. Future work should more directly contrast

336 the effect of anticipatory vs. adaptive cues and of probabilistic (top-down) vs. Gestalt-related (bottom-

**18**

337 up) cues to establish their relative contribution and investigate how this may vary under different

338 experimental conditions.

339 Another concern is whether IDyOM accurately reflects listener expectations. Morgan et al.

340 (2019) found that IDyOM predictions entailed higher entropy than that computed across several

341 participants providing single-tone sung continuations to melodic contexts. Task constraints likely

342 explain this discrepancy as expectations for multiple continuations were not assessed. Furthermore,

343 by manipulating entropy of upcoming events rather than simply analyzing the entropy of instantiated

344 continuations, the present study differs crucially from Morgan et al. (2019). Moreover, whereas they

345 recruited self-identified musicians, who make melodic predictions with demonstrably lower average

346 entropy than non-musicians (Hansen & Pearce, 2014; Hansen, Vuust, & Pearce, 2016), IDyOM was

347 configured to model expectations of the general population. At the same time, Morgan et al. (2019)

348 made an important contribution by demonstrating a greater contribution of statistical learning than of

349 Gestalt-based principles in predicting listener expectations. This supports IDyOM's suitability in

350 predicting auditory boundary perception.

351 The finding that uncertainty influences phrase boundary perception suggests a pertinent role

352 for training effects. Expertise effects may be particularly prominent in the musical domain where

353 skills and experiences differ substantially between individuals. Although some studies suggest limited

354 effects of musical expertise on melodic segmentation processes (Palmer & Krumhansl, 1987, but see

355 Hartmann et al., 2017), expertise levels have not always been widely sampled or manipulated

356 systematically. The same limitation applies to the current study where no significant effects of

357 expertise were seen (see Tables S2 and S3 in SOM-R2 for details). Yet, recent research shows that

358 stylistic specialization results in expectations about melodic continuations that are generally lower in

359 entropy whenever greater confidence is warranted (Hansen & Pearce, 2014; Hansen et al., 2016). The

360 transformation of high-entropy predictions into low-entropy predictions with domain-relevant

**19**

361 training or implicit exposure should allow musicians to perceive phrasal coherence across longer

362 timespans. This would be consistent with observations that experts have access to more abstract and

363 deeper levels of hierarchical structure (Chaffin & Imreh, 2002; Chi & Feltovich, 1981) which, in turn,

364 may be associated with larger working memory capacity (Meinz & Hambrick, 2010). While awaiting

365 sampling across more diverse expertise levels in future research, our results relating chunk size to

366 underlying expectancy dynamics enables a novel interpretation of classical findings pertaining to

367 expertise and working memory.

368 By offering an empirical challenge to the view that segmentation primarily relies on

369 retrospective processes, the present work contributes to the emergence of an increasingly coherent

370 model of the human mind as an eager predictive processor of sensory input. Embedded in the constant

371 flux of time, the mind is continually forced to evaluate and recombine retrospective and prospective

372 cues according to their immediate usefulness, and we hypothesize that sequential input in such varied

373 domains as language, music, and visual action sequences are all subject to the constraints arising from

374 this mental machinery.

**References**

375

376 Chaffin, R., & Imreh, G. (2002). Practicing perfection: piano performance as expert memory.

377 *Psychological Science*, *13*(4), 342-349.

378 Chi, M. T., Feltovich, P. J., & Glaser, R. (1981). Categorization and representation of physics

379 problems by experts and novices. *Cognitive Science*, *5*(2), 121-152.

380 Christiansen, M. H., & Chater, N. (2016). The now-or-never bottleneck: a fundamental constraint

381 on language. *Behavioral and Brain Sciences*, *39*, e62. doi:10.1017/S0140525X1500031X.

382 Cohen, P., Adams, N., & Heeringa, B. (2007). Voting experts: an unsupervised algorithm for

383 segmenting sequences. *Intelligent Data Analysis*, *11*(6), 607-625.

384 Cousineau, D. (2005). Confidence intervals in within-subject designs: a simpler solution to Loftus

385 and Masson's method. *Tutorials in Quantitative Methods for Psychology*, *1*(1), 42-45.

386 Creighton, H. (ed.). (1966). *Songs and Ballads from Nova Scotia*. New York, NY: Dover.

387 Deliege, I. (1987). Grouping conditions in listening to music: an approach to Lerdahl &

388 Jackendoff's grouping preference rules. *Music Perception*, *4*(4), 325-359.

389 Dörfell (ed.) (1875). *371 vierstimmige Choralgesänge von Johann Sebastian Bach* (4th ed.).

390 Leipzig, Germany: Breitkopf & Härtel.

391 Hansen, N. C.,& Pearce, M. (2014). Predictive uncertainty in auditory sequence processing.

392 *Frontiers in Psychology* 5, 1052.

393 Hansen, N. C., Vuust, P., & Pearce, M. (2016). "If you've got to ask, you'll never know": Style-

394 congruent musical expertise optimises predictive auditory processing. *PLOS ONE*, *11*(10):

395 e0163584. doi:10.1371/journal.pone.0163584

396 Hansen, N. C. ,Vuust, P., Pearce, M., & Huron, D. (2017, August). *Entropic Ebbs and Flows: The*

397 *Expectancy Dynamics of Musical Phrases*. Paper presented at the Society for Music Perception

398 and Cognition Meeting, San Diego, CA.

399    Hard, B. M., Meyer, M., & Baldwin, D. (2019). Attention reorganizes as structure is detected in

400         dynamic action. *Memory & Cognition*, *47*(1), 17-32.

401    Hard, B. M., Recchia, G., & Tversky, B. (2011). The shape of action. *Journal of Experimental*

402         *Psychology: General, 140*(4), 586-604. doi:10.1037/a0024310

403    Hartmann, M., Lartillot, O. & Toiviainen, P. (2017). Interaction features for prediction of

404         perceptual segmentation: effects of musicianship and experimental task. *Journal of New Music*

405         *Research, 46*(2), 156-174. doi:10.1080/09298215.2016.1230137

406    Hutchinson, J. B., & Barrett, L. F. (2019). The power of predictions: an emerging paradigm for

407         psychological research. *Current Directions in Psychological Science*, *28*(3), 280-291.

408    Koelsch, S., Vuust, P., & Friston, K. (2019). Predictive processes and the peculiar case of music.

409         *Trends in Cognitive Sciences*, *23*(1), 63-77.

410    Kosie, J. E., & Baldwin, D. (2019a). Attention rapidly reorganizes to naturally occurring structure

411         in a novel activity sequence. *Cognition*, *182*, 31–44. doi:10.1016/j.cognition.2018.09.004

412    Kosie, J. E., & Baldwin, D. (2019b). Attentional profiles linked to event segmentation are robust to

413         missing information. *Cognitive Research: Principles and Implications*, *4*(1), 8.

414         doi:10.1186/s41235-019-0157-4

415    Kragness, H. E. & Trainor, L. J. (2016). Listeners lengthen phrase boundaries in self-paced music.

416         *Journal of Experimental Psychology: Human Perception and Performance*, *42*(10), 1676-1686.

417         doi:10.1037/xhp0000245

418    Kragness, H. E. & Trainor, L. J. (2018). Young children pause on phrase boundaries in self-paced

419         music listening: the role of harmonic cues. *Developmental Psychology*, 54(5), 842-856.

420         doi:10.1037/dev0000405

421    Kurby, C. A., & Zacks, J. M. (2008). Segmentation in the perception and memory of events. *Trends*

422         *in Cognitive Sciences*, *12*(2), 72-79.

**22**

1
2
3
4
5    423    Lerdahl, F., & Jackendoff, R. (1983). *A generative theory of tonal music*. Cambridge, MA: MIT
6
7    424        Press.
8
9    425    Levinson, S. C. (2016). Turn-taking in human communication: origins and implications for
10
11   426        language processing. *Trends in Cognitive Sciences, 20*(1), 6-14. doi:10.1016/j.tics.2015.10.010
12
13
14   427    Meinz, E. J., & Hambrick, D. Z. (2010). Deliberate practice is necessary but not sufficient to
15
16   428        explain individual differences in piano sight-reading skill: the role of working memory capacity.
17
18   429        *Psychological Science*, *21*(7), 914-919.
19
20
21   430    Morgan, E., Fogel, A., Nair, A., & Patel, A. D. (2019). Statistical learning and Gestalt-like
22
23   431        principles predict melodic expectations. *Cognition*, *189*, 23-34.
24
25   432    Nicholson, S., Knight, G. H., and Bower, J. D. (Ed.). (1950). *Ancient and Modern Revised*. Suffolk,
26
27   433        UK: William Clowes and Sons.
28
29
30   434    Palmer, C. (1989). Mapping musical thought to musical performance. *Journal of Experimental*
31
32   435        *Psychology: Human Perception and Performance*, *15*(2), 331.
33
34   436    Palmer, C., & Krumhansl, C. L. (1987). Independent temporal and pitch structures in determination
35
36
37   437        of musical phrases. *Journal of Experimental Psychology: Human Perception and Performance*,
38
39   438        *13*(1), 116.
40
41   439    Pearce, M. T. (2005). *The construction and evaluation of statistical models of melodic structure in*
42
43
44   440        *music perception and composition* (Doctoral dissertation). City University, London, UK.
45
46   441        Retrieved from https://openaccess.city.ac.uk/id/eprint/8459/1/
47
48   442    Pearce, M. T., Müllensiefen, D., & Wiggins, G. (2010). The role of expectation and probabilistic
49
50   443        learning in auditory boundary perception: a model comparison. *Perception, 39*(10), 1367-1391.
51
52   444        doi:10.1068/p6507
53
54
55   445    Penel, A., & Drake, C. (1998). Sources of timing variations in music performance: a psychological
56
57   446        segmentation model. *Psychological Research*, *61*(1), 12-32.
58
59
60

447 R Core Team (2019). *R: A language and environment for statistical computing*. R Foundation for

448     Statistical Computing, Vienna, Austria. Retrieved from https://www.R-project.org/.

449 Repp, B. H. (1992). Probing the cognitive representation of musical time: structural constraints on

450     the perception of timing perturbations. *Cognition*, *44*(3), 241-281.

451 Richmond, L. L., & Zacks, J. M. (2017). Constructing experience: event models from perception to

452     action. *Trends in Cognitive Sciences*, 21(12), 962-980.

453 Saffran, J. R., & Kirkham, N. Z. (2018). Infant statistical learning. *Annual Review of Psychology*,

454     *69*, 181-203. doi:10.1146/annurev-psych-122216-011805

455 Schaffrath, H. (1995). *The Essen Folksong Collection in the Humdrum Kern Format* (D. Huron,

456     Ed.). Menlo Park, CA: Center for Computer Assisted Research in the Humanities. Retrieved

457     from https://kern.humdrum.org/cgi-bin/browse?l=essen/europa/deutschl

458 Wightman, C. W., Shattuck-Hufnagel, S., Ostendorf, M., & Price, P. J. (1992). Segmental durations

459     in the vicinity of prosodic phrase boundaries. *The Journal of the Acoustical Society of America*,

460     *91*(3), 1707–1717. doi:10.1121/1.402450

461 Yang, C. D. (2004). Universal grammar, statistics or both? *Trends in Cognitive Sciences*, *8*(10),

462     451-456.

463 Zacks, J. M., & Swallow, K. M. (2007). Event segmentation. *Current Directions in Psychological*

464     *Science*, *16*(2), 80-84.

465 Zacks, J. M., Tversky, B., & Iyer, G. (2001). Perceiving, remembering, and communicating

466     structure in events. *Journal of Experimental Psychology: General*, *130*(1), 29-58.

**24**