

2014-07

Predictive regularity representations in violation detection and auditory stream segregation: from conceptual to computational models.

Schroger, E

<http://hdl.handle.net/10026.1/8210>

10.1007/s10548-013-0334-6

Brain Topogr

All content in PEARL is protected by copyright law. Author manuscripts are made available in accordance with publisher policies. Please cite only the published version using the details provided on the item record or document. In the absence of an open licence (e.g. Creative Commons), permissions for further reuse of content should be sought from the publisher or author.

Brain Topography

Predictive regularity representations in violation detection and auditory stream segregation: from conceptual to computational models

--Manuscript Draft--

| | |
|--|--|
| Manuscript Number: | BTOP-D-13-00085R1 |
| Full Title: | Predictive regularity representations in violation detection and auditory stream segregation: from conceptual to computational models |
| Article Type: | S.I. : Mismatch Negativity |
| Section/Category: | Cognitive neuroscience |
| Keywords: | audition; cognition; auditory object; auditory scene analysis; deviance detection; predictive modelling; mismatch negativity (MMN) |
| Corresponding Author: | Erich Schroger GERMANY |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | |
| Corresponding Author's Secondary Institution: | |
| First Author: | Erich Schroger |
| First Author Secondary Information: | |
| Order of Authors: | Erich Schroger Alexandra Bendixen Susan L Denham Robert W Mill Tamás M Bóhm István Winkler |
| Order of Authors Secondary Information: | |
| Abstract: | <p>Predictive accounts of perception have received increasing attention in the past twenty years. Detecting violations of auditory regularities, as reflected by the Mismatch Negativity (MMN) auditory event-related potential, is amongst the phenomena seamlessly fitting this approach. Largely based on the MMN literature, we propose a psychological conceptual framework called the Auditory Event Representation System (AERS), which is based on the assumption that auditory regularity violation detection and the formation of auditory perceptual objects are based on the same predictive regularity representations. Based on this notion, a computational model of auditory stream segregation, called CHAINS, has been developed. In CHAINS, the auditory sensory event representation of each incoming sound is considered for being the continuation of likely combinations of the preceding sounds in the sequence, thus providing alternative interpretations of the auditory input. Detecting repeating patterns allows predicting upcoming sound events, thus providing a test and potential support for the corresponding interpretation. Alternative interpretations continuously compete for perceptual dominance. In this paper, we briefly describe AERS and deduce some general constraints from this conceptual model. We then go on to illustrate how these constraints are computationally specified in CHAINS.</p> |

Dear Dr. Murray, dear Dr. Sussman (dear Micah and Elyse):

Thank you and the reviewers for the valuable suggestions to improve the manuscript. We took them seriously and made plenty of modifications in the revised manuscript. Please, find below a detailed response to the various points (our text is in blue). In the manuscript, the new text is marked.

Best regards,
Erich (also on behalf of all authors).

CC: christoph.michel@unige.ch, elyse.sussman@einstein.yu.edu

Ref.: Ms. No. BTOP-D-13-00085

Predictive regularity representations in violation detection and auditory stream segregation: from conceptual to computational models Brain Topography

Dear Prof. Schröger,

Many thanks for your submission to this forthcoming special issue of Brain Topography. Prof. Sussman, one of the guest editors of the issue, and I have now received comments from two external experts. Likewise, Prof. Sussman has herself provided comments. You will see that they are all generally favourable, but would propose some minor revisions to further improve the work. I would therefore encourage you to prepare a revision that would undergo re-review by Prof. Sussman and the original reviewers.

For your guidance, Prof. Sussman's and the reviewers' comments are appended below.

If you decide to revise the work, please submit a list of changes or a rebuttal against each point which is being raised when you submit the revised manuscript.

Your revision is due by 23-10-2013.

To submit a revision, go to <http://btop.edmgr.com/> and log in as an Author. You will see a menu item call Submission Needing Revision. You will find your submission record there.

Yours sincerely,
Micah

Micah M. Murray, Ph.D.
Editor-in-Chief
Brain Topography

Guest Editor's Comments:

The reviewers' comments are largely focused on bringing greater clarity to the paper in describing the two models of auditory perception. Reviewer 1 notes that the AERs model is under review. The references should only include manuscripts that are published or have been accepted for publication ('in press'). Manuscripts under review or in preparation can be mentioned in the text but are not to be included in the reference list. Please also check the text for language use and try to avoid colloquialisms (e.g., 'in a nutshell').

As the AERS manuscript is not yet published, we refrain from referring to an unpublished manuscript. This means that in the revised manuscript, the description of the AERS model should stand on its own.

Reviewers' comments:

Reviewer #1: This paper proposes an integration of two separate models on auditory perception: one "verbal-boxological" one (in author's words) on predictive modeling of the acoustic environment (the "AERS" model), and another computational one on auditory stream segregation (the "CHAINS" one); the two of them developed by the authors. In that regard, and although the AERS model is only under revision elsewhere, the author's proposal might be regarded as not really novel. However, the present paper attempts to relate these two models by proposing a series of constraints for computational models. The paper evolved in two clear separate sections, one revising the AERS model, one revising the CHAINS model -in the light of the constraints derived from AERS.

As it stands, the present manuscript will likely provide the first description of the AERS model. Although it cannot include all its facets and all the empirical evidence it is based upon, the revised manuscript should be able to stand on its own, allowing readers to understand the main aspects of AERS.

In my view -and in the lack of the fully developed version, still under review, it is not very clear how the AERS model departs from previous models on the MMN (except for the use of new terminology, for instance borrowed from predictive coding perspectives), and it is not clear how can predictions be derived to implement psychophysiological experiments (i.e., to record the MMN, earlier correlates of deviance detection/regularity encoding, or both). Maybe a comment on this direction might help to improve the outreach of the theoretical (verbal-boxological) proposal. On the other hand, the CHAINS model has been described and validated elsewhere. Also, a comment on potential or future experiments, eventually at electrophysiological level to validate it in the light of the AERS constraints might improve the paper.

A) Separation from previous models of MMN. Except for Winkler's (2007) model of MMN, no previous MMN model posits that the role of the processes reflected by MMN is to support the formation of auditory streams. (I.e., previous models of MMN suggest that the role of the processes reflected by MMN is to call for further processing of the deviant stimulus, ultimately to recruit attentional resources for examining the deviant stimulus in more detail.) AERS has its roots in Winkler's MMN model, but specifies it in a number of ways; for example by providing algorithms for the initial creation of regularity representations, assuming competition between alternative regularity representations, linking the processes of deviance detection to attentive processes, etc.

B) Although AERS, as any model involving prediction and its comparison with the actual input, is generally compatible with the principles of predictive coding, it is not an instantiation of the predictive coding principle for auditory deviance detection. It differs from predictive coding in its roots (i.e., conceptualizing MMN and stream segregation data) as well as in its aims (i.e., focusing on the formation and maintenance of the memory representations underlying the two sets of data). This is now made explicit on p7. A further important difference between AERS and other predictive coding accounts is mentioned on page 15: Unlike other predictive coding models, AERS distinguishes between different types of prediction errors with qualitatively different implications.

C) Throughout the "Auditory Event Representation System (AERS)" section, we added pointers for which assumptions/predictions of AERS would benefit from testing by future experiments (pages 10-15).

D) Future directions of expanding the CHAINS model are mentioned on page 22.

I have some additional minor comments:

1. In abstract, but mainly in page 3 (lines 12-13) the concept of "Auditory violation detection (AVD)" is introduced. The authors should consider revising this concept, as what it is violated is the regularity established on the acoustic environment, but nothing "auditory" itself. A proper term might be "Auditory-regularity violation detection".

Yes, the term proposed by the reviewer is more appropriate. We changed "Auditory violation detection" to "Auditory regularity violation detection".

2. Page 3 (line 27). The adjective "unique" is used to qualify the MMN as psychophysiological indicator of registration of new information. However, as the authors acknowledge later on, there are other (earlier) indicators of the registration of new information, and therefore the MMN is not the "unique" (in the sense of sole) indicator anymore.

Yes, we fully agree! As revealed by Escera, Grimm and other members of this group, the classic MMN has a precursor at a considerably earlier latency. This may also be regarded as a type of MMN. However, to our knowledge, it seems that this precursor of MMN is not that sensitive for complex, abstract rules, for example, for feature conjunctions (Althen, Grimm, Escera, 2013, E J Neurosci). In this sense the classical MMN is still superior. We modified the text accordingly, that is, we omitted the term "unique" on p.3 and we shortly mention that the MMN seems to be more sensitive to more complex, abstract regularities (page 10).

3. Page 9 (lines 52-60). It is claimed that the update of regularity representations is reflected by the MMN. I wonder whether early correlates of deviance detection (as mentioned in the text a few lines above) could also be considered as such correlates of regularity updating. In any case, the authors should comment on this possibility.

We agree. Accordingly, we add this option to the text (pp9-10). However, whether these early detection processes are based on checking predictions against the actual input remains to be seen. This is now mentioned as one of the issues for future research.

4. Page 25 (line 12-13): typo "vying".

"vying" is a synonym for competing.

Reviewer #2: Comments on the manuscript: 'Predictive regularity representations in violation detection and auditory stream segregation: from conceptual to computational models', co-authored by Schröger E., Bendixen A., Denham S.L., Mill R.W., Böhm T.M., & Winkler I.

General comment:

In this review paper, the objective is to bring together within a common conceptual and computational framework two highly active and important research fields in auditory perception: Auditory Violation Detection (AVD, a research field mostly based on the MMN component of the auditory ERP), and Auditory Scene Analysis (ASA, which has been widely studied with behavioural approaches but comparatively rarely from a neuroscience perspective). This most interesting endeavour is motivated by the fact that both processes are related to regularity extraction from the auditory input.

The paper describes two models: firstly, a "boxological" model of ASA and AVD, called Auditory Event Representation System (AERS), and secondly a computational model of ASA (CHAINS), which is based on principles derived from AERS.

Although I fully agree with "the big picture" put forward by the paper, I am a bit more reserved concerning some specific points within AERS (see below). I also think that the paper could be more pedagogical if it more clearly specified which aspects of the models are backed up with experimental evidence (and which points are more speculative), as well as which aspects of the models are new conceptualisations (and which aspects are in fact common with previous attempts at theorizing from the MMN or ASA literature). The links between the AERS model and more "classical" views of the MMN involving sensory memory representations could be made more evident.

Thank you for this comment. We tried to be more explicit in the revised version. However, due to space limitations we cannot elaborate too much on that and this is not the focus of the current ms. Here is the gist of the difference between AERS and previous models of MMN:

Except for Winkler's (2007) model of MMN, no previous MMN model posits that the role of the processes reflected by MMN is to support the formation of auditory streams. (I.e., previous models of MMN suggest that the role of the processes reflected by MMN is to call for further processing of the deviant stimulus, ultimately to recruit attentional resources for examining the deviant stimulus in more detail.) AERS has its roots in Winkler's MMN model, but specifies it in a number of ways; for example by providing algorithms for the initial creation of regularity representations, assuming competition between alternative regularity representations, linking the processes of deviance detection to attentive processes, etc.

We believe that a full comparison amongst the various models of MMN would detract readers from the focus issues of the current ms.

Major issues

1. The AERS model postulates a number of representation levels: Auditory perceptual objects, Auditory predictive regularity representation (proto-object), Auditory stimulus event, Auditory sensory stimulus event representation, Auditory (perceptual) event representation (so five levels of representation according to Table 1 - but only four are displayed Figure 1, and not always with the same terminology, which complicates things rather unnecessarily). There is a definite need to better explain why the model postulates these different levels of representation, whether it is grounded on experimental evidence or speculative, how these different representations differ from one another, etc.

It seems that we were not sufficiently clear in our description. However, the number and definition of the concepts listed in the Table correspond to the ones in the Figure. Please, consider that "Auditory stimulus event" is not a representation (level), it is the stimulus (as defined in the Table and illustrated in Figure 1). In Figure 1 three of these representations defined in the Table are used: "Auditory predictive regularity representation or proto-object" (is used identical in Table and Figure); "Auditory stimulus event representation" (is used identical in Table and Figure; please not that we replaced the previous term "Auditory sensory stimulus event representation" by "Auditory stimulus event representation" as sensory is redundant); "Auditory perceptual event representation" (is used identical in Table and Figure); not used in the Figure is the concept of "Auditory perceptual object (or auditory object representation)" which corresponds to a stream.

With these in mind, there are actually only three levels of representations. 1) An auditory stimulus event representation is an internal representation of the stimulus event – mainly equal to what has been termed as auditory sensory memory in the literature. 2) An auditory perceptual event representation is derived from the auditory stimulus event representation by adding the relation of the auditory event to the context (both the auditory context – i.e., it is a regular sound or a deviant – and the larger context – e.g., this is a target in some task). 3) An auditory predictive regularity representation or proto-object is a representation of a detected regularity which predicts upcoming sounds that would continue the sequence so far in a regular manner (so, this is a predictive model). An auditory perceptual object is the currently dominant proto-object (thus it is not a separate

representational level). We tried to add further clarifying text to the text where these are first described and also modified the figure legend.

2. It is overall difficult to follow exactly the logic behind the AERS model, in part because the text (notably pages 8-10) does not exactly match with Figure 1. For example, the text mentions "auditory sensory memory representations" but they are not present in Figure 1. Furthermore, Figure 1 is difficult to understand. A few examples of the difficulties I encountered trying to understand the model: first, one does not understand from the figure where the information about the "auditory stimulus event" enters in the AERS (a similar ambiguity arises for the output of AERS).

Indeed, "auditory sensory memory representations" are not linked to the terminology used in AERS. This has been resolved in the revision. In footnote 1 we explain how sensory memory representations relates to the corresponding term used in AERS (*"These [auditory stimulus event representations] correspond to auditory sensory memory representations of the classical MMN model (e.g. Näätänen, 1990). We prefer the term auditory stimulus event representations as a stimulus event is represented."*).

Yes, the arrows in Figure 1 did not optimally indicate "where the information about the "auditory stimulus event" enters in the AERS (a similar ambiguity arises for the output of AERS)" (quote from reviewer#2). This has been changed in the revised version.

Second, the "old" or "new" status of a particular event was in my understanding the output of the "comparison" process, but "old" seems to be an input of this "comparison" and there is no arrow between "new" and "comparison", so I got confused. Further, the +/- box between "comparison" and "proto-objects" does not make much sense.

AERS assumes that the comparison is done "streamwise", that is, separately for each established ("old") stream. "Old" and "new" refers to Bregman's "old+new strategy"- This is a heuristic the auditory system utilizes to detect the emergence of new streams. When AERS decides that a "new" stream comes into play, then a new predictive regularity representation may be established (but no comparison is made). That is why the "the +/- box" (quote from reviewer#2) is exactly where it should be. It denotes the output of the comparison process, the auditory sensory stimulus event representation is compared to the prediction and either yields a match ("+") or a mismatch ("-") result. This, in turn, has consequences for the predictive regularity representations as well as for the evaluation (as illustrated in the Figure). We made this more explicit in the revised text and we modified the line from the +/- box to the evaluation in order to avoid misunderstanding.

3. To better understand the link suggested by the authors between ASA and AVD, it would be interesting to know how the AERS model can explain (or not) previously published results suggesting that ASA precedes AVD (e.g. Müller et al., 2005, Psychophysiology; Sussman, 2005, JASA).

AERS itself is not a model of auditory stream segregation. It sets the scene for certain types of models of auditory stream segregation, and, indeed, the focus of the ms. is to show how a model of auditory stream segregation can be based on AERS. Thus much of the literature of auditory stream segregation is trivially compatible with AERS, simply because in this sense AERS is underspecified. However, the Reviewer is right in pointing out that there is an apparent conflict between the main idea of the AERS model (that the processes indicated by MMN serve the formation of auditory streams) and the assumption that auditory stream segregation **precedes** the processes indicated by MMN. We addressed this conflict in a recent experimental paper (Bendixen, Schröger, Ritter & Winkler 2012) where we show on the basis of MMN data that the strict temporal order ("ASA always precedes AVD") does not hold: There are also cases where AVD precedes ASA. Furthermore, we discuss why previous studies have not found similar evidence. In the revised manuscript (page 13),

we briefly mention this issue and refer the reader to Bendixen et al. (2012) for a detailed line of arguments.

Minor issues (by order of appearance in the text):

1. Page 3, first paragraph. For clarity, the two concepts (AVD and ASA) should be mentioned in the same order in the two sentences (it is ASA then AVD in the first sentence and AVD then ASA in the second). With the same logic, perhaps that the order of the two subsequent paragraphs could be inverted.

Done! AVD then ASA order.

2. Page 4, lines 37-42. The hypothesis that links together in a single representation the following elements: "sensory description of incoming sounds", "relationship to the context", "current goals of the listeners", and "consciously perceived" seems a rather strong one. This hypothesis should be justified with respect to previous experimental studies on the topic, or clearly presented as speculative.

Please note that in this Introduction section the topic of the paper is outlined. The description and (sparse) justification is given in the subsequent sections (e.g. in "Auditory Event Representation System (AERS)", and – to a smaller extent - in "Competition and Cooperation between Proto-Objects in a Model of Auditory Scene Analysis (CHAINS)").

3. Page 6, line 54. It is the first time that a direct reference to "memory" is made (it comes back page 8 as "auditory sensory memory representations"), and it should be explained why. As mentioned above, the relationship between the AERS model and sensory memory models of the MMN should be made clear.

See our response to the major comment #2 of reviewer#2. Moreover, we omitted "memory" in this sentence in the revision.

4. Page 10, line 54. Perhaps replace "the separation between A and B" by "the frequency separation between A and B" (the same applies to page 13, line 44).

As a result of other modifications to the manuscript, we have deleted this phrase.

5. Page 11, first paragraph. The various claims made here should be justified with references to the relevant literature. It is also a bit biased to state that the AVD research field might provide new insights for ASA research but not to explore the other direction of cross-fertilisation between research fields (from ASA to AVD).

We have included an example in which ASA influences AVD (pp 11-12).

6. Page 16, line 15. Replace "van Norden" by "van Noorden".

Done!

7. Page 18. If I understood correctly, CHAINS works for sounds which are presented "one at a time". Would it be feasible to extend the model to accomplish segregation of simultaneous sounds, or is it a limit of the model? (My question is definitively not a request for this extension to be done now and integrated in the paper but rather an attempt to understand how general the model could be with some limited modifications).

Yes, in its present state CHAINS works for consecutive sounds. This is a limitation as in real life sounds overlap in time (and ASA may well work in such situations), which will hopefully be taken care of in future research. Furthermore, the current implementation of CHAINS is based on deterministic regularities; this is in conflict with the stochastic nature of many environmental regularities and clearly an issue for refining the model in the future. In the manuscript (page 22) we list some current limitations of CHAINS which can be taken as challenges for future developments. We also refer to the limitations and challenges discussed at length in the paper by Mill et al. (2013) describing the CHAINS model.

8. Page 19, lines 5-11. I did not really understand what "how many times this event has been predicted by already existing chains" means.

Several different chains may predict the same sound. If this is the case, there is no need for "inventing" a new chain (or adding an additional element to an existing chain). Consider the simplest case of repeating the same stimulus A again and again. The chains A, AA, AAA will already predict A. Thus the system would be increasingly reluctant to generate AAAA and AAAAA chains.

9. The bibliography should be checked for accents and umlauts in authors' names.

Thanks. Done.

1
2 Predictive regularity representations in violation detection and auditory stream segregation:
3
4 from conceptual to computational models
5
6
7

8
9 Erich Schröger¹, Alexandra Bendixen^{1,2}, Susan L. Denham³, Robert W. Mill⁴, Tamás M.
10
11 Böhm⁵, & István Winkler^{5,6}
12
13
14
15
16

17 ¹Institute of Psychology, University of Leipzig, Leipzig, Germany

18
19 ²**Auditory Psychophysiology Lab**, Department of Psychology, Cluster of Excellence

20
21 “Hearing4all”, European Medical School, Carl von Ossietzky University, Oldenburg,
22
23
24 Germany

25
26
27 ³Cognition Institute and School of Psychology, University of Plymouth, Plymouth, UK

28
29 ⁴MRC Institute of Hearing Research, Nottingham, United Kingdom

30
31 ⁵Institute of Cognitive Neuroscience and Psychology, Research Centre for Natural Sciences,
32
33 Hungarian Academy of Sciences, Budapest, Hungary

34
35
36 ⁶Institute of Psychology, University of Szeged, Szeged, Hungary
37
38
39
40

41 This research was supported by the Hungarian Academy of Sciences (Lendület
42
43 project, LP2012-36/2012 to IW), by the Reinhart Koselleck grant of the German Research
44
45 Foundation (Deutsche Forschungsgemeinschaft, DFG, SCH 375/20-1 to ES), by the DFG
46
47 Cluster of Excellence 1077 “Hearing4all”, by the German Academic Exchange Service
48
49 (Deutscher Akademischer Austauschdienst, DAAD, Project 56265741), and by the Hungarian
50
51 Scholarship Board (Magyar Ösztöndíj Bizottság, MÖB, Project 39589).
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Abstract

1
2 Predictive accounts of perception have received increasing attention in the past twenty years.
3
4 Detecting violations of auditory regularities, as reflected by the Mismatch Negativity (MMN)
5
6 auditory event-related potential, is amongst the phenomena seamlessly fitting this approach.
7
8
9 Largely based on the MMN literature, **we propose** a psychological conceptual framework
10
11 called the Auditory Event Representation System (AERS), **which is based on the assumption**
12
13 **that auditory regularity violation detection** and the formation of auditory perceptual objects
14
15 are based on the same predictive regularity representations. Based on this notion, a
16
17 computational model of auditory stream segregation, called CHAINS, has been developed. In
18
19 CHAINS, the auditory sensory event representation of each incoming sound is considered for
20
21 being the continuation of likely combinations of the preceding sounds in the sequence, thus
22
23 providing alternative interpretations of the auditory input. Detecting repeating patterns allows
24
25 predicting upcoming sound events, thus providing a test and potential support for the
26
27 corresponding interpretation. Alternative interpretations continuously compete for perceptual
28
29 dominance. In this paper, we briefly describe AERS and deduce some general constraints
30
31 from this conceptual model. We then go on to illustrate how these constraints are
32
33 computationally specified in CHAINS.
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

51
52 **Keywords:** audition, cognition, auditory object, auditory scene analysis, deviance detection,
53
54 predictive modelling, mismatch negativity (MMN)
55
56
57
58
59
60
61
62
63
64
65

1 The processing of auditory information serves to discover the distal sources of sensory input
2 and to detect potentially important events in the environment. To date, these two important
3 functions have been studied relatively independently of each other in the fields of *auditory*
4 *regularity violation (deviance) detection* (Näätänen 1990) and *auditory scene analysis*
5 (Bregman 1990).
6
7
8
9
10

11
12
13
14 *Auditory regularity violation detection* (AVD) is concerned with identifying new information
15 in a given context, which is of potential interest to the listener. The basic idea is that new
16 information requires detailed evaluation because we do not know about it yet (as opposed to
17 the redundant repetition of old information). If we learn the regularities inherent in the
18 dynamic sensory input, we can readily know what is “old” and detect what is “new”. The
19 classic orienting reaction approach (Sokolov 1963) inspired this extremely fruitful field of
20 irregularity detection, which established a psychophysiological indicator of the successful
21 registration of new information, the Mismatch Negativity (MMN; Näätänen et al. 2011).
22
23
24
25
26
27
28
29
30
31
32

33
34
35
36 *Auditory scene analysis* (ASA) is concerned with the problem of identifying the concurrently
37 active sound sources (Bregman, 1990). This is a considerable challenge for the information
38 processing system, because the travelling waves emitted by the different sound sources and
39 their echoes mix together before they reach our ears. The crucial task consists in disentangling
40 this mixture by grouping (*integrating*) information that belongs together and separating
41 (*segregating*) information that belongs to different sources. The classic Gestalt principles
42 (Köhler 1947) such as the laws of common fate, similarity, and continuity have been
43 successfully used as a guideline to study the segregation of the auditory *input into coherent*
44 *sound sequences, termed “auditory streams”*. Auditory streams are important units of auditory
45 perception and, as argued by Winkler and colleagues (Winkler et al, 2009), they also serve as
46 units of information storage possessing the defining characteristics of perceptual object
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 representations. That is, representations of auditory streams 1) are temporally persistent; 2)
2 encode conjoined auditory features; 3) are separable from other streams; 4) can absorb natural
3 variability of the input; and 5) predict upcoming sound events (Winkler et al., 2009). Note
4 that auditory streams do not always correspond to a single sound source (Bregman, 1990).
5
6 However, as has often been noted, sound patterns, such as a melody, can also be regarded as
7 perceptual objects (Kubovy and Van Walkenburg, 2001; Griffith and Warren, 2004).
8
9

10
11
12
13
14
15
16
17 Because new information must either be related to previously detected sound sources, or to
18 decisions on the presence of new sources, in a general theoretical sense one can assume that
19 the detection of new information must be an important factor in auditory scene analysis. A
20
21 more specific link between the two sets of phenomena can be formed by noticing that both
22
23 **AVD and ASA** require knowledge about the history of stimulation. That is, they both utilize
24
25 some form of representation of the auditory context. The assumption that **AVD and ASA**
26
27 utilize at least partly overlapping representations of the auditory context is the first
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
61
62
63
64
65

66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

1 various cognitive operations. Although AERS is a “verbal-boxological” model (cf. Jacobs and
2 Grainger 1994) and lacks detailed computational specificity, it organizes a vast amount of
3 literature and yields various useful constraints for computational models. We will demonstrate
4 this in the second half of the paper by describing a specific computational implementation of
5 AERS, the CHAINS model (Mill et al. 2011; Mill et al. 2013). Please see Table 1 for a
6 definition of key terms and concepts in the two models.
7
8
9
10
11
12
13
14
15
16

17 **The relation between AERS and the predictive-coding approach to perception**

18
19
20
21
22 In general, the predictive-modeling account of perception explains perception as the result of
23 an interaction between the current sensory information and a model of what we already know
24 about the world. This idea is by no means novel. In fact, it originates from Helmholtz’s notion
25 of unconscious inferences (Helmholtz 1867), extended in two important ways:
26
27
28
29
30
31

32
33
34 1) At any given moment of time, there are multiple models potentially applicable to the
35 current sensory input. Therefore, some mechanism for selecting the optimal model should
36 exist.
37
38
39

40
41 2) In order to efficiently represent an ever-changing environment and the arrival of
42 unexpected new information, the perceptual system must monitor how well its current
43 representations suit the environment, improving them when possible, initiating new
44 representations if necessary, and adjusting its selection of dominant representations
45 appropriately.
46
47
48
49
50
51
52
53
54
55

56 With these two principles in mind, the nature of a predictive (or generative) representation can
57 be summarized as follows: The representation generates predictions that take into account
58 prior experience and provide probabilistic information about what is likely to appear next in a
59
60
61
62
63
64
65

1 given context. Consider, for instance, that you have chosen a whistling voice as the ringtone
2 for your mobile phone. Now, while you are in a crowd of people, you suddenly perceive a
3 whistling voice. At least two explanatory models for this sensory input will be formed by your
4 perceptual system: one of them suggesting that your mobile phone is ringing, the other one
5 suggesting that someone is actually whistling in the crowd. Let us assume that the “mobile
6 phone” model wins the competition initially. This will lead to the prediction that you will
7 immediately hear the whistling voice again, which might be so strong that you misperceive a
8 second sound event that is only vaguely similar to a whistling sound. To put this in theoretical
9 terms, the prediction constrains and biases the interpretation of the sensory input. The actual
10 input is compared with the prediction and the difference between the two is computed as
11 prediction error. This prediction error, in turn, can be used for adjusting current
12 representations and/or for selecting a dominant representation (interpretation) from those
13 available (e.g., someone on the street was whistling).

14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34 The most complete variant of predictive coding, Friston’s free energy principle (Friston and
35 Kiebel 2009b), suggests that a) model selection is based on Bayesian inference and b) there is
36 a hierarchy of models with increasing generality, with prediction errors from each level
37 propagating upwards in the hierarchy (bottom-up) and models for each level being selected by
38 the next higher level (top-down). The goal of the system is to minimize the overall prediction
39 error formulated as an entropy type of measure, the *free energy*. In this hierarchical predictive
40 coding model, no level within the hierarchy is of special relevance. This approach aims at
41 describing the overall functioning of the perceptual system and can be validated by comparing
42 the behavior of computational implementations of the model with human perceptual decisions
43 and brain activation measures (in other words, with the *outcome* of perceptual processes). In
44 contrast, psychological descriptions of ASA and AVD focus on the *processes* leading to the
45 conscious perception of sounds.

1
2 Some recent computational models of MMN (Garrido et al. 2009; Lieder et al. 2013;
3
4 Wacongne et al. 2012) have provided a link between predictive coding and the brain's
5
6 response to auditory regularity violations. Winkler and Czigler (2012) suggested that the
7
8 representations of auditory regularities involved in AVD may be mapped to models of an
9
10 intermediate level in a predictive coding hierarchy. However, predictive processing has not
11
12 yet been applied systematically to explain ASA. We will do this by specifying the nature of
13
14 the predictive regularity representations that compete to explain the auditory sensory input,
15
16 and by considering how they are initially formed and maintained. This approach fills
17
18 important holes in both bodies of literature:
19
20
21
22
23
24
25

26
27 1) In general, predictive coding studies have rarely addressed the issue of how models are
28
29 initially formed or maintained (see, however, Kiebel et al. 2009).
30

31
32 2) Models of ASA have largely disregarded bi-/multi-stable perceptual phenomena (for a
33
34 general review, see Schwartz et al. 2012), which has led to the underestimation of the role of
35
36 competition between alternative representations in ASA (see, however, Mill et al. 2013).
37
38
39
40

41 In summary, although AERS is generally compatible with predictive coding models (though
42
43 not necessarily with any particular model), it differs from predictive coding in its roots as well
44
45 as in its aims. AERS is not an instantiation of the predictive coding principle for auditory
46
47 deviance detection; instead, AERS describes the common basis for auditory deviance
48
49 detection and stream segregation in terms of the formation and maintenance of the memory
50
51 representations underlying these processes. By AERS, we lay the groundwork for
52
53 computational models describing ASA in terms of continuous competition between
54
55 regularity-based descriptors of auditory event sequences.
56
57
58
59
60
61
62
63
64
65

Table 1: Definition of terms

| | |
|----|--|
| 1 | <i>Auditory event representation system (AERS)</i> : A cognitive system producing <i>auditory</i> |
| 2 | <i>perceptual event representations</i> from <i>auditory stimulus events</i> . |
| 3 | |
| 4 | |
| 5 | <i>Auditory perceptual object (or auditory object representation)</i> : A member of the set of |
| 6 | currently dominant <i>proto-objects</i> that occupies perceptual awareness. ‘Auditory stream’ is a |
| 7 | similar term but one which is not explicitly identified in terms of predictive representations. |
| 8 | |
| 9 | |
| 10 | <i>Auditory predictive regularity representation</i> or <i>proto-object</i> : Terms used to describe the |
| 11 | representation of a sequence of sounds linked together by some detected rule or repeating |
| 12 | pattern in the form of a generative model. Incoming sounds are checked against the |
| 13 | predictions of current <i>proto-objects</i> that compete to ‘explain’ them. These representations |
| 14 | have the potential to become the perceptual objects in conscious awareness, if and when they |
| 15 | are dominant (i.e., they are in a highly activated state, assumed to indicate their selection as |
| 16 | the most likely description of the current input). |
| 17 | |
| 18 | |
| 19 | |
| 20 | <i>Auditory stimulus event</i> : A discrete sound, localized in time and generated by some source |
| 21 | in the external world; i.e. the physical input to our sensory systems (e.g., each of the sounds |
| 22 | in the particular sequence of sounds generated by our mobile phone when a text message |
| 23 | arrives). |
| 24 | |
| 25 | |
| 26 | <i>Auditory stimulus event representation</i> : The integrated description of the perceived features |
| 27 | of an <i>auditory stimulus event</i> ; is shaped by the predictions from <i>AERS</i> . |
| 28 | |
| 29 | <i>Auditory perceptual event representation</i> : The description of an <i>auditory stimulus event</i> in |
| 30 | the brain; an <i>auditory sensory stimulus representation</i> , which is linked to a perceptual object, |
| 31 | and expanded with the description of its relationship to the auditory and the general context |
| 32 | (e.g. the text-message sound as it appears in our perception). Auditory perceptual event |
| 33 | representations are the output of <i>AERS</i> . |
| 34 | |
| 35 | |
| 36 | |
| 37 | <i>CHAINS</i> : A computational model (and its implementation into a Matlab/C-based computer |
| 38 | program) that incorporates aspects of <i>AERS</i> . It can be used to simulate possible <i>perceptual</i> |
| 39 | <i>organizations</i> for specific parameters, which – in turn – can be tested experimentally. |
| 40 | |
| 41 | <i>Chains</i> : A formal description (within <i>CHAINS</i>) of a sequence of <i>auditory stimulus events</i> |
| 42 | including their timing. |
| 43 | |
| 44 | |
| 45 | <i>Initial sound analysis</i> : The early (bottom-up) activation patterns in the auditory system |
| 46 | caused by each <i>auditory stimulus event</i> ; this analysis is not regarded as part of <i>AERS</i> . |
| 47 | |
| 48 | |
| 49 | <i>Perceptual organization</i> : A complete description of the auditory environment, in terms of |
| 50 | auditory object representations, as it appears in perception. |
| 51 | |
| 52 | |
| 53 | |
| 54 | |
| 55 | |
| 56 | Auditory Event Representation System (AERS) |
| 57 | |
| 58 | |
| 59 | |
| 60 | |
| 61 | |
| 62 | |
| 63 | |
| 64 | |
| 65 | |

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
61
62
63
64
65

AERS is characterized by four major constituents, (1) the formation of **auditory stimulus event representations**, (2) the **formation of regularity representations that predict** subsequent sensory input, (3) comparison between the predictions and the sensory input, and (4) evaluation of the relevance of the relationship between the incoming sound events and the context (Figure 1). AERS is assumed to receive its input from subcortical and cortical levels, for example, in the form of spectrotemporal response patterns encoding features such as spectral energy maxima. These sound features have to be combined into a unitary auditory **stimulus event** representation that is held accessible for some time. The auditory N1 ERP response may (partly) reflect processes engaged in this function (Näätänen and Winkler 1999). However, the formation of **stimulus event** representations does not only rely on input but also on the “bias” exerted by the prior context. This context supports the formation of predictive representations that are used to compute a-priori probabilities **for events embedded in the** sensory input. These representations capture current auditory regularities such as a pitch alternation regularity of two tones (A and B) differing in frequency presented in a regular pace (ABABABAB...).

Regularity representations are generative models in the sense that they produce predictions for future expected parts of the pattern (i.e. **upcoming** sounds, such as that the tone following an A will be a B). These predictions strongly guide the formation of **auditory stimulus event representations**¹ of the incoming sounds. As predictive regularity representations have been **normally** active before the occurrence of the current sound, the auditory **stimulus event** representations are always shaped by the a-priori probabilities. Regularity representation is a concept that is well known and frequently used in AVD research. **Winkler and colleagues (2009) suggested that** the concept of a *stream* essentially corresponds to a regularity

¹ These correspond to auditory sensory memory representations of the classical MMN model (e.g. Näätänen, 1990). We prefer the term auditory stimulus event representations as a stimulus event is represented.

1 representation, although the notion of a *regularity* is seldom mentioned in this context
2 (instead, the term *coherence* is used to refer to the principles holding together the sounds that
3 form a stream) and streams are primarily regarded as perceptual, not as encoding units
4 (however, these two aspects are obviously not contradictory). Within AERS, streams are
5 regarded as generative models based on detected regularities. Any incoming sound receives
6 an interpretation biased towards being a new token of the currently dominant stream (such as
7 the continuation of the voice of a speaker). So far, few MMN studies have attempted to
8 distinguish predictive processing from alternative explanations (e.g., reevaluating the
9 immediate history of stimulation at the arrival of each new sound event; but see Paavilainen et
10 al., 2007; Bendixen et al., 2008, 2009). Further, if regularity representations lie at the heart of
11 auditory streams, then MMN should only be elicited by sounds belonging to the stream whose
12 regularity is violated. Ritter and colleagues' (2000, 2006) results appear to support this
13 assumption. These assumptions of AERS should be further tested by future research.

14 Predictions are compared with the emerging auditory stimulus event representations created
15 for the auditory input. This comparison is not a single unitary process. It is computed at
16 multiple anatomical and temporal levels of sound processing. Recent studies have revealed
17 that even subcortical areas of the brain can be involved and that some form of deviance
18 detection can take place as early as 30 ms from the onset of the violation (for a review, see
19 Grimm and Escera 2012). A regularity representation needs to be updated when the incoming
20 sound mismatches its prediction; this updating process is reflected by the MMN (Winkler
21 2007; Winkler and Czigler 1998; Winkler et al. 2009). According to the proposal of Grimm
22 and Escera (2012) of a hierarchical novelty system (see also Escera et al., this issue),
23 respective updating processes can possibly also be initiated by the precursors of MMN.
24 However, there is evidence that more complex regularities such as feature conjunctions are
25 not encoded at these early levels (Althen et al., 2013).

1
2 If the input matches the prediction, no updating is required. **Instead**, confidence in the model
3
4 might be strengthened. The repetition positivity (RP) component (Haenschel et al. 2005;
5
6 Costa-Faidella et al. 2011) and the induced gamma-band response (Herrmann et al. 2004)
7
8 might be candidates for brain signals reflecting this matching process. **So far no study tested**
9
10 **whether the early deviance-related responses (as reviewed by Grimm and Escera, 2012) can**
11
12 **be regarded as signs of prediction error. Likewise, no study tested whether RP is only elicited**
13
14 **by true repetition or also by non-repeating, but predictable sounds (although Bendixen et al.,**
15
16 **2008, noted an effect similar to RP in a complex MMN paradigm based on predictive**
17
18 **regularities). Future studies may shed light on whether or not these ERP responses reflect**
19
20 **processes assumed by AERS.**
21
22
23
24
25
26
27

28
29 The outcome of the comparison describes the relation between the **auditory stimulus event**
30
31 representation of the incoming sound and the prediction(s) stemming from previously
32
33 detected regularities. In fact, several predictive representations may coexist, providing
34
35 alternative descriptions of the auditory scene. Let us again consider the alternating pitch
36
37 regularity ABABAB... Alternation is only one of several possible descriptions of the tone
38
39 sequence. It has been termed the *local* rule in the literature (Horváth et al. 2001) as it makes
40
41 local predictions regarding the next expected sound (sound $n+1$); that is, after a tone A tone B
42
43 is predicted, and after B, A is predicted. Another possible regularity that can be derived from
44
45 this sequence is that every second tone is A, while every other sound is B. This regularity
46
47 generates the same sound sequence, but makes its predictions with regard to the second
48
49 upcoming sound (sound $n+2$), thus termed the *global* rule. Horváth and colleagues (2001)
50
51 showed **that auditory predictive regularity** representations for both local and global
52
53 regularities are active in parallel.
54
55
56
57
58
59
60
61
62
63
64
65

1 The alternation regularity example can also be used to illustrate the conceptual similarity of
2 streams and predictive regularity representations. When presented with an ‘ABABAB...’
3
4 stimulus, participants often report hearing an integrated percept; that is, the perception that
5 one sound source has produced all the tones by regularly alternating between A and B tones.
6
7 This corresponds to the “n+1” (local) description of the alternation regularity. However,
8
9 participants may also report hearing two separate sound sources: an A-stream, consisting only
10 of the ‘A’ sounds, and a B-stream, consisting only of the ‘B’ sounds. This corresponds to the
11
12 “n+2” (global) description of the alternation regularity. Few MMN studies addressed the issue
13
14 whether MMN is only related to the currently dominant (perceived) stream/regularity-
15
16 representation or also to the currently non-dominant ones (e.g., Szalárdy et al., in press;
17
18 Winkler et al., 2005, 2006) and the results are somewhat equivocal. This issue requires further
19
20 research.
21
22
23
24
25
26
27
28
29
30

31 Having realized this similarity between regularity representations (in AVD research) and
32 streams (in ASA research), how can the two perspectives inform each other in a fruitful
33
34 manner? One aspect of the AVD field that can provide new insights for ASA research is the
35
36 issue of how predictive models are formed. The two research fields have opposite approaches
37
38 here: while AVD assumes that a mixture of sounds comes as an incoherent series of events in
39
40 which the regularities must first be discovered, ASA research typically assumes that a mixture
41
42 is by default interpreted as a series of events belonging together (*integrated*), while the
43
44 formation of more than one representations to describe the input sequence (i.e., *segregated*)
45
46 only happens after the accumulation of corresponding evidence. However, this integration-by-
47
48 default view has recently been challenged by a number of groups (e.g., Deike et al. 2012;
49
50 Denham et al. 2013) and thus the ASA research field must face the question of how any
51
52 stream representation initially develops (Winkler et al. 2009, 2012). It may prove highly
53
54 fruitful to borrow paradigms and findings from AVD research here (e.g., roving standard
55
56
57
58
59
60
61
62
63
64
65

1 paradigms as put forward by Bendixen et al. 2007; Cowan et al. 1993; Haenschel et al. 2005;
2 Sussman et al. 2007; Winkler et al. 1996). At the same time, one aspect in which AVD may
3 benefit from ASA concerns the notion of perceptual organization and the fact that any
4 sequence can have multiple interpretations. This bi- or multistability is rarely considered in
5 AVD research; usually the sequence of sounds is assumed to be processed as one stream
6 throughout. Finally, the notion of a joint representational basis of AVD and ASA has already
7 led to a re-consideration of how the two processes are arranged in time. While previous
8 studies had concluded that ASA precedes AVD (e.g., Müller et al. 2005; Sussman 2005),
9 more recent evidence suggests that this temporal relation is more flexible and depends on the
10 strength of the acoustic and regularity cues that are available to the auditory system (Bendixen
11 et al. 2012).

12 After a predictive representation (a regularity representation or stream) has been set up, its
13 validity is tested by the occurrence of every new sound event that it predicts. The information
14 about whether the prediction was met – and if not, how far the incoming sound deviated from
15 it – is passed on to the evaluation process. Here, representations of incoming sounds are
16 related to what is currently known about the environment; i.e., the relationship between the
17 incoming sound and the context (including the current goals of the listener) is evaluated. We
18 propose that the P3a ERP response reflects the outcome of this evaluation process and acts as
19 a kind of “significance” marker of sensory events (Horváth et al. 2008; Rinne et al. 2006).

20 The resulting information package is the primary output of AERS that is available for further
21 processing. We term this an auditory perceptual event representation, because it describes the
22 sound event together with its relation to both the auditory and the general context. This is a
23 more realistic conceptualization of the minimal “units” of auditory perception than one
24 restricted to the physical stimulus. In contrast to classical accounts of sensory memory such as
25 echoic memory (Neisser 1967) or (short and long) auditory stores (Cowan 1984; Massaro

1 1972), AERS emphasizes that auditory **perceptual** event representations usually are more than
2 a mental “echo” of the auditory stimulus event and incorporate our knowledge **regarding** the
3 context, our intentions, and even affective affordances (“answer the mobile phone”).
4
5
6
7

8
9
10 Central to AERS is the proposal that the evaluation process also participates in the search for
11 new regularities. If a prediction is not met, this may be due to the fact that the existing
12 predictive model is basically valid but an exception with respect to the regularity occurred
13 (e.g., one footstep of a walking person sounded a bit different from the ones experienced in
14 the past because the person stepped **onto some object**). **In such a case, although the predictive**
15 **regularity representation may need some modification (updating) it should not be discarded as**
16 **a valid description (i.e., the model that we are hearing a series of footsteps can be**
17 **maintained)**. However, it may also be the case that a new regularity started, that is, a new
18 stream came into play, which does not require the updating of the existing regularity
19 representation but rather the creation of a new one (e.g., another walking person approaches).
20
21
22 The former case (i.e. a mismatch between the prediction and the actual continuation of the
23 stream) corresponds to a prediction error in predictive coding models. However, the latter (i.e.
24 the residue that cannot be explained by current predictive representations) corresponds to
25 what Bregman (1990) captured in his “old+new strategy”, a heuristic the auditory system
26 utilizes to detect the emergence of new streams. The information that cannot be accounted for
27 by the existing streams (i.e. the residue) can be assessed at this stage and the presence of a
28 new sound source can be considered. **As noted earlier, comparisons are only done within**
29 **streams (see Ritter et al., 2000, 2006). Sounds that do not belong to the given stream are not**
30 **compared and no deviance (error) signal is generated. This is marked by the “old” (i.e.,**
31 **belonging to one of the detected streams) information entering the comparison, whereas**
32 **“new” (belonging to no previously detected stream) information initiates the formation of a**
33 **new regularity representation (see Figure 1).**
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
61
62
63
64
65

1
2 At this point, the description given by AERS diverges from existing predictive coding models,
3
4 which lump together deviation from a prediction and the residue (the ‘old’ and the ‘new’)
5
6 under a single error signal, whereas AERS distinguishes these two prediction errors and
7
8 assigns different follow-up actions to them (i.e., updating an existing vs. forming a new
9
10 auditory predictive regularity representation). The distinction between processing prediction
11
12 errors and the residue may be reflected in ERPs: the former is assumed to elicit the MMN
13
14 (and possibly earlier deviance-detection-related ERP components), whereas the latter may be
15
16 reflected by components notably sensitive to large acoustic changes, such as the P1 and N1.
17
18 Although most known ERP data are compatible with this assumption, it has not been directly
19
20 tested.
21
22
23
24
25
26
27
28

29 **Specifying AERS and extracting some computational principles**

30 *Formation of proto-objects*

31
32
33
34 AERS provides a general scheme for forming predictive representations of repeating sound
35
36 patterns. However, it makes no suggestion about how distinct sounds are linked together into
37
38 a coherent representation. Thus the first issue to be addressed by a computational model based
39
40 on AERS is how associations between temporally separate sounds are formed. Intuitively, it
41
42 should be easier to connect similar than highly dissimilar sounds. This principle has been
43
44 termed the law of similarity by the Gestalt school of psychology. However, the Gestalt school
45
46 focused on vision and space, where display items are present side by side. In contrast, in the
47
48 auditory modality, similarity is mediated primarily by time. Thus the principle of similarity
49
50 translates to a temporal version of smooth continuation. That is, similarity is better expressed
51
52 in terms of temporal rate of feature change (Jones 1976; Winkler et al. 2012). This modified
53
54 definition of similarity receives support from numerous studies of auditory stream segregation
55
56
57
58
59
60
61
62
63
64
65

1 (for reviews see Moore and Gockel 2002; Moore and Gockel 2012). These results show that
2 sounds with even moderate frequency separation may segregate when presented in close
3 succession (high rate of change due to the short inter-stimulus interval), whereas sounds with
4 much higher frequency separation may be perceptually grouped together if there are longer
5 time intervals between them (low rate of change due to the longer inter-stimulus interval). The
6 interplay between featural and temporal separation is limited by the temporal constraints of
7 the underlying memory processes as well as by the organization of feature spaces in the
8 auditory system. Thus a computational implementation must choose parameters in accordance
9 with the known perceptual and neurophysiological properties. Recent studies suggest that the
10 initial formation of predictive representations (termed “proto-objects” in Mill et al. 2011; cf.
11 also Rensink 2000) may primarily rely on the above notion of similarity. Similarity (rate of
12 feature change) may also determine the time needed to establish a proto-object, if we assume
13 that links are formed with a probability related to their similarity; i.e., more similar sounds are
14 more likely to be associated and, therefore, such groups of sounds (and the corresponding
15 regularities) are found more quickly. In response to a sequence of sounds, the proto-object
16 discovered first will emerge first in perception, and will remain there without competition
17 until at least one more alternative proto-object is discovered (cf. Winkler et al. 2012).

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 The formation of a proto-object, however, needs an additional step beyond establishing links
45 between sounds. Sounds linked together by similarity can only affect the processing of
46 upcoming sounds when one can draw predictions from them. That is, the building of a proto-
47 object is only complete once it has shown the potential to predict upcoming sound events
48 (because only then can new events be “absorbed” by this proto-object). The simplest way this
49 can happen is when a repeating pattern is detected. AERS suggests a more general
50 formulation: repetition of an inter-sound relationship, with predictions taking the form of
51 value distributions in the parameter space. In implementations, one needs to carefully choose
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 realistic parameters both for what kind of inter-sound relations are handled by the model and
2 for limiting the possible length of proto-objects. The auditory system appears to show quite
3 surprising constraints in terms of the length of the patterns (number of items within and/or
4 duration spanned by the pattern) that can be extracted (e.g., Sussman and Gumenyuk 2005;
5 **Boh et al. 2011**). These constraints need to be (even) more systematically investigated within
6 the field of AVD to permit their inclusion within computational models of ASA.
7
8
9

10
11
12
13
14
15
16
17 Finally, it is important to note that temporal adjacency is not a necessary prerequisite for
18 linking sounds together, as was shown by Bendixen and colleagues (2012b) for AVD and by
19 numerous streaming experiments for ASA (e.g., Müller et al. 2005). This is important as it
20 allows the auditory system to form parallel representations of alternative proto-objects for the
21 same sequence of sounds. Evidence for the auditory system maintaining alternative regularity
22 representations for describing the same sound sequence has been obtained in MMN studies
23 (Horváth et al. 2001).
24
25
26
27
28
29
30
31
32

33 34 35 36 *Maintenance of proto-objects*

37
38 Predictive processes may have a dual role in maintaining proto-object representations. Firstly,
39 they may help to improve internal cohesion, by **strengthening** links between the elements of a
40 proto-object. This has been suggested by the results of studies showing that proto-objects
41 within which individual sounds or sound features can be predicted by some simple rule (e.g. a
42 regularly repeating pattern) are more likely to emerge in perception for longer periods of time
43 compared with proto-objects within which sounds are only linked by more diffuse rules (e.g. a
44 predictable feature distribution) (Bendixen et al. 2010, 2013). A proto-object with an internal
45 predictive rule can be regarded as a proto-object with an internal structure, which allows it to
46 provide more precise predictions compared to similar proto-objects with no internal structure.
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

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
61
62
63
64
65

Second, successful predictions should help to maintain a proto-object, whereas failures should decrease its chances of survival. AERS suggests that predictive failures reduce the effectiveness of the given proto-object in the competition for perceptual dominance. Further, it suggests that below some threshold, proto-objects become ineffective and stop affecting the processing of incoming sounds. However, they remain in an inactive state, from which they can be reactivated for rather long periods of time (Winkler et al. 2002).

Internal cohesion (the strength of associations between its elements) and predictive success determine the competitiveness of a proto-object, i.e., its effect on other proto-objects within the competition. The moment-to-moment activation levels of the competing proto-objects (resulting from these factors) determine which of them (possibly more than one) are part of conscious awareness at any given time.

Competition and the emergence of perceptual organizations

A perceptual organization is a complete description of the auditory environment as it appears in perception. For example, the repetitive ‘ABA_’ sequence (van Noorden, 1975) is most commonly heard either as a repeating three-tone pattern (i.e., all sounds appearing as a single integrated unit) or as two parallel streams of sound, one consisting of the A, the other of the B sounds, with one of them appearing in the foreground, the other in the background. Whereas in the first case, perceptual organization consists of a single sound object, in the second case, perceptual organization consists of two sound objects and the assignment of the foreground. (Note that in real-life situations, there are almost always multiple sound objects with some of them falling to the background.) As these are alternative perceptual organizations of the same sequence, only one of them can appear in perception in any given time. The questions to be addressed by a computational model implementing the AERS principles are:

1) How are perceptual organizations formed from proto-objects?

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
61
62
63
64
65

2) What is the unit that enters the competition: the proto-objects or the alternative full perceptual organizations?

3) What form does the competition take? How do the competitors affect each other?

Compatibility/incompatibility between proto-objects has been discussed by Winkler et al. (2012; see also Mill et al. 2013), who defined any pair of proto-objects as incompatible if they predicted the same sound event (termed collision). This definition is based on the principle of exclusive allocation (i.e., that any given sound can be part of only one percept at a time). Exclusive allocation mostly holds for auditory perception (see, e.g., Winkler et al. 2006), although there are examples of duplex perception (e.g., Fowler and Rosenblum 1990). An analysis of the possible forms of competition led to the suggestion of collisions as the basic building block of competition between proto-objects: two proto-objects compete with each other, when, and only when, they collide. It can be shown that competition based on this simple principle implicitly leads to the emergence of perceptual organizations as reported by human participants (Mill et al., 2013).

These are then the computational principles extracted from AERS, which underlie the development of the computational model called CHAINS (Mill et al. 2011). We now go on to describe the principles of CHAINS (for detailed accounts, see Mill et al. 2011, 2013). We would like to note that many other computational models of ASA have been formulated (Beauvois and Meddis 1991; McCabe and Denham 1997; Wang and Chang 2008), some of them also incorporating the concept of predictive processing (e.g., Elhilali and Shamma 2008; Grossberg et al. 2004). Nevertheless, we limit ourselves to the description of CHAINS here because it is intimately connected to the MMN-based AERS framework.

Competition and Cooperation between Proto-Objects in a Model of Auditory Scene

Analysis (CHAINS)

Formation of proto-objects

CHAINS is a computational model that allows one to flexibly implement aspects of the AERS conceptual framework already outlined. In keeping with the AERS schematic, CHAINS receives a temporal pattern of pre-analyzed sound events as input, and explores ways to form representations of the embedded regularities. Sounds are encoded throughout CHAINS as discrete tokens, which represent a single point in feature space at a specific instant in time. In our terminology, tokens within CHAIN relate to auditory stimulus events within AERS. The CHAINS **algorithm** does not access the absolute features of a token. Instead, it measures the distance between a token and an incoming stimulus event, and in this way links together “constellations” of stimulus events in an unfolding time-feature space to form *chains* (Figure 2); i.e. its representations are based on relative representations of feature distributions.

The likelihood of a pattern of stimulus events coalescing into a chain in the first place is determined probabilistically according to the interaction of two functions that serve complementary roles. One function specifies the probability that an incoming event is added to a chain; the other specifies the probability that an incoming event is left out, i.e., skipped over. The CHAINS model does not specify what form these functions are to take: this decision is deferred to the modeler. Nevertheless, their general influence will be heavily informed by the empirical data considered earlier, namely, that it is difficult to link events whose features change abruptly; and, conversely, it is difficult *not* to link events whose features change gradually (Figure 2A,B). It is important to emphasize that these two outcomes are not mutually exclusive. On the contrary, it is an essential feature of CHAINS that, when presented with an input event, each chain has the possibility of splitting into *two* parallel

1 variants, one that includes the event, and one that excludes it. This principle leads to an
2 exponentially-growing set of chains, the proliferation of which is to some extent constrained
3 by the low probabilities of perceptually unreasonable linkages.
4
5
6
7
8
9

10 The simple chain-building scheme we have described lays a concrete groundwork for two key
11 principles that appear in the AERS framework, namely, predictive regularity and residual
12 input. A predictive regularity is established when a repeating pattern appears in a chain.
13
14

15 Specifically, if a chain is found to consist of a repeating sequence, it closes to form a loop,
16 and thereupon ceases to grow by incorporating incoming events. This is when it becomes a
17 proto-object and starts to predict events according to the regularity it encodes (Figure 2C).
18
19
20
21

22 The chain will persist in some form as long as its predictions are correct. At the same time,
23 correctly predicting a given input event has immediate implications for the formation of
24 further alternative chains: The probability of adding an event to a chain is reduced in
25 accordance with how many times this event has been predicted by already existing chains.
26
27
28
29
30
31

32 This can be implemented, for instance, by modifying the link probability functions described
33 above. The probability reduction naturally gives rise to a graded interpretation of residual
34 input: events that are predicted by fewer chains are more likely to be built into existing chains,
35 or seed new ones. On the other hand, one can tailor the exclusion probability function to make
36 it easier for a chain, when building, to skip over an event that has been predicted (i.e.,
37 accounted for) by many other chains.
38
39
40
41
42
43
44
45
46
47
48
49
50

51 *Maintenance of proto-objects*

52 Not all features of AERS related to the maintenance of proto-objects have been implemented
53 in CHAINS up to now. For example, in the current state of CHAINS an incoming sound
54 violating a predictive regularity prediction erases the respective chain. In AERS, however, it
55 is assumed that it takes some time before a proto-object becomes ineffective. This is
56
57
58
59
60
61
62
63
64
65

1 suggested from AVD research, where it has been shown that it takes several repetitions of a
2 violation in order to abolish the MMN elicited by the violations (Winkler et al. 1996).
3

4 Moreover, AERS claims that regularity representations that no longer affect the processing of
5 the incoming sounds can remain in an inactive, “dormant” state, and can be reactivated by a
6 single “reminder” (Cowan et al. 1993). This dormant state does not (yet) have an explicit
7 computational analogue in CHAINS; it may, however, map onto the chain’s continuously
8 varying level of excitation described in the next section. Further aspects that are not yet
9 incorporated in the CHAINS model despite existing evidence from MMN/AVD research are
10 summarized in the computational description of CHAINS (Mill et al. 2013). Notwithstanding
11 these future challenges, we now go on to describe the already implemented aspects of the
12 CHAINS model.
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

29 The level of excitation of each chain depends on its predictive success and the presence of
30 collisions with concurrent chains (proto-objects). CHAINS simulates the dynamics of the
31 changing and inter-related excitation levels of all existing chains, which determines its
32 predictions of the perceptual organization of the scene.
33
34
35
36
37
38
39
40

41 *Competition and the emergence of perceptual organizations*

42 From the fact that CHAINS is permitted to skip over input events when building chains, it
43 follows that most individual chains do not describe the input in its entirety. Moreover, the lack
44 of a strict principle of mutual exclusion during the building process implies that the same
45 input event can be incorporated into many chains, and the population of chains at a given
46 moment will contain a degree of redundancy. A single chain, considered in isolation,
47 therefore makes a good candidate for a proto-object, in that it may predict only a fragment of
48 the auditory scene, with incomplete coverage but no internal inconsistencies. At the same
49 time, the population of chains, considered as a whole, provides an overly exhaustive
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

predictive account of a single scene, with complete coverage but many colliding predictions.

The ideal circumstances lie between these two extremes: insofar as a scene is predictable, a subpopulation of the chains should interleave their predictions so as to account for every event exactly once (and the remainder should acquiesce). We refer to these subpopulations as *perceptual organizations*. Of course, for any given scene, there may be more than one organization latent within a population. This provides a natural basis for reasoning about perceptual multi-stability: if the elements within a population compete directly with each other to predict events, then subsets of chains that together predict disjoint aspects of a sensory scene will (implicitly) cooperate to form organizations. This is the process by means of which CHAINS discovers and maintains organizations. We now examine the details of this process.

At the outset, we introduce the notion of a chain's *excitation*, a quantity denoted E_i , which can fall between zero (not excited) and one (fully excited). In a sense, all chains in the population are predictive, but it is the excited chains' predictions that are taken to account for the events in a scene. There are many conceivable contributions to the excitation of chain i , three of which are essential. Firstly, the predictive success of the chain, defined as the rate at which it makes successful predictions (S_i), increases its excitability. Secondly, collisions with any other chain, j , defined as the rate at which predictions of chain i collide with those of chain j (C_{ij}) multiplied by the latter chain's excitation (E_j), reduces the chain's excitability. Thirdly, there is a persistent noise term (U_i), which perturbs the chain's excitability over time.

This behaviour described above can be captured in a simple system of first-order non-linear equations. For each i ,

$$\tau_m \frac{dE_i}{dt} = -E_i + \phi(\alpha_S S_i - \alpha_C \sum_{j \neq i} C_{ij} E_j + \alpha_U U_i + \text{const.})$$

where $\phi(x) = (1 + e^{-x})^{-1}$ is a sigmoid function. Alternatively, the effect of collisions can be

mediated indirectly via an *inhibitory* variable, I_i :

$$\tau_m \frac{dE_i}{dt} = -E_i + \phi(\alpha_S S_i - \alpha_{IE} I_i + \alpha_U U_i + \text{const.})$$

$$\tau_m \frac{dI_i}{dt} = -I_i + \phi(\alpha_C \sum_{j \neq i} C_{ij} E_j + \text{const.})$$

The benefit of the latter scheme is that it limits the effect of collisions with many chains by introducing saturation. In either case, the α parameters control how successes, collisions and noise contribute to the chain's excitation, and τ_m is a time constant that controls how rapidly excitation evolves. (The terms S_i and C_{ij} are dynamic state variables found by leaky integration, though for a repeating stationary sequence such as 'ABA_', we can treat them as though they are constant.)

Consider now how CHAINS might respond to the repeating 'ABA_' sequence. Firstly, we assume that the building process outlined earlier has assembled three chains: one that predicts all three tones (the 'ABA' chain), a second that predicts only the As (the 'A' chain), and a third that predicts only the Bs (the 'B' chain). Principally, there are two ways into which these chains can organise themselves, given the dynamical system mentioned.

In the first scenario, the excitation of chain 'ABA', E_{ABA} , is initially high and the excitation of the other two chains is low. Chain 'ABA' issues three correct predictions per stimulus cycle, whereas chains 'A' and 'B' issue only two and one, respectively. Consequently, chain 'ABA' will be more excited due to more successful predictions than 'A', and 'A' more excited than 'B' in turn. In addition, the predictions of 'ABA' will regularly collide with those of 'A' and 'B', tending to reduce the excitation of the latter even further. This process will stabilise, with E_{ABA} near to one, and E_A and E_B near to zero. In this *integrated* organization, the dominant chain 'ABA' predicts the input events by itself, and chains 'A' and 'B' are non-dominant.

1
2 In the second scenario, the excitation of chains 'A' and 'B' (E_A and E_B , respectively) are
3
4 initially high, and that of chain 'ABA' is low. Here, the contributions due to successful
5
6 predictions are the same as those in the integrated scenario. However, the excitation of chain
7
8 'ABA' is substantially reduced by its collisions with chains 'A' and 'B', whereas chains 'A'
9
10 and 'B' are relatively uninhibited: their own predictions do not collide with each other at all,
11
12 and E_{ABA} is low, so the impact of collisions with chain 'ABA' is small. This process will also
13
14 stabilise, with E_{ABA} nearer zero, and E_A and E_B nearer one. In this *segregated* organization,
15
16 the dominant chains, 'A' and 'B', alternately predict the stimulus events, and chain 'ABA' is
17
18 non-dominant.
19
20
21
22
23
24
25
26

27 The *integrated* and *segregated* organizations of chains are both stable with respect to the
28
29 CHAINS dynamics. If it were not for the noise terms, U_i , the competition would settle into
30
31 one of these two states and remain there. However, the addition of a moderate level of noise
32
33 leads to transitions back and forth between one organization and the other. To adequately
34
35 model perceptual multi-stability, one must choose the α and τ parameters to ensure an
36
37 appropriate balance in the contribution of success, collisions and noise. There is a broad range
38
39 of parameter sets that lead to multi-stability, and we can briefly summarise their respective
40
41 influences as follows. In general increasing α_S (the effect of successes) promotes integration,
42
43 increasing α_C (or α_{IE} , the effect of collisions) promotes segregation, and increasing α_U (the
44
45 effect of noise) increases the rate of switching and promotes segregation to a small extent.
46
47
48
49
50

51 Figure 2D presents example time courses of the excitations of the 'ABA', 'A', and 'B' chains,
52
53 as they compete with each other over a 240 second period to explain an 'ABA_' input
54
55 sequence. The ABA chain is discovered initially, and the 'A' and 'B' chains are discovered
56
57 somewhat later (~25, 52 sec, respectively). The emergence of perceptual organizations is
58
59 evident in these series: either the 'ABA' chain inhibits the 'A' and 'B' chains to produce an
60
61
62
63
64
65

1 integrated phase (e.g., 90–170 sec), or the ‘A’ and ‘B’ chains together inhibit the ‘ABA’ chain
2 to produce a segregated phase (e.g., 60–90 sec). In this example, the noise contribution to
3 each chain (U) suffices to produce phase durations on the order of tens of seconds.
4
5
6
7

8
9 The basic chains dynamics set out above can be augmented in a straightforward manner by
10 adding additional terms inside the sigmoid functions of the excitation and inhibition
11 equations, $\varphi(\cdot)$. For example, a feedback term that causes a chain to excite itself promotes the
12 stability of organizations. (Adding adaptation to this feedback term promotes the stability of
13 organizations for only a limited period after they become dominant.) In addition, one can add
14 a rediscovery term, which excites a chain every time another version of it is rebuilt from the
15 input events. For example, if the parameters of an ‘ABA_’ sequence favor integration (i.e.,
16 smooth changes; see Figure 2A), then the ‘ABA’ chain will be rebuilt or rediscovered quite
17 often. If the chains required to form the segregated organization have already been built (‘A’
18 and ‘B’), the frequent rediscovery of the ‘ABA’ chain will promote integration during the
19 competition. The converse applies if the stimulus parameters promote segregation (Figure
20 2B). Other terms are also conceivable, for instance, those which encode the effort made to
21 attend to a particular sound event or organization.
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

44 The most important feature of the CHAINS dynamics is that they arise naturally from the
45 predictive successes of chains and the collisions between predictions—there is no special
46 effort to predefine *integrated* or *segregated* percepts, as they exist with respect to an ‘ABA_’
47 sequence. Consequently, CHAINS is also bistable when presented with an alternating tone
48 (‘ABAB...’), or a sequence containing three tones. Furthermore, CHAINS exhibits
49 *multistability* when more than two possible stable states exist. For example, organizations can
50 arise in which all three events in an ABA triplet break into three separate chains, or the first
51 two tones in a triplet form one chain, and the third tone breaks off into its own chain.
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Ultimately, which chains participate in the competition will depend on the parameters that
2 govern the probability of their formation in the first place, and which groups of chains
3 subsequently coalesce to form an organization will depend on their compatibility (collisions).
4
5 Because CHAINS makes no assumptions concerning the form the perceptual organizations
6 should ultimately take, it provides a flexible starting point from which to explore multistable
7 perception driven by ambiguous sequences more complex than the classical alternating and
8 galloping tones. As denoted above, the strength of this approach can be further increased by
9 shaping the probability of chain formation in accordance with the vast body of AVD results
10 based on the MMN, resulting in a fruitful integrated AERS-CHAINS framework.
11
12
13
14
15
16
17
18
19
20
21
22
23

24 **Summary**

25
26
27
28
29 We started from the Auditory Event Representation System (AERS), a conceptual framework
30 linking **auditory regularity violation detection** and auditory scene analysis largely based on the
31 MMN (Näätänen et al. 1978) research (Winkler and Schröger submitted). The notion of
32 predictive processes underlying the elicitation of MMN (Winkler et al. 1996) has gained
33 momentum in the last couple of years (e.g., Baldeweg 2007; Bendixen et al. 2012a; Garrido et
34 al. 2009; Näätänen et al. 2011; Schröger 2007; Wacongne et al. 2011; Winkler 2007), partly
35 because of its compatibility with predictive coding theories (e.g., Friston and Kiebel 2009a;
36 Mumford 1992; Rao and Ballard 1999). AERS takes the next step by linking **auditory**
37 **regularity violation detection** and auditory scene analysis through predictive representations
38 of the regularities detected from the sound input, which then serve as proto-objects
39 continuously vying for the possibility of appearing in conscious perception. From AERS, we
40 extracted some theoretical requirements for computational models of auditory stream
41 segregation, many of which have been implemented in the CHAINS model (Mill et al., 2013).
42 CHAINS further specifies the formation, and competition between proto-objects. Applying
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 CHAINS to the auditory streaming paradigm (Van Noorden, 1975) the time-course of the
2 excitation of the three typical proto-objects has been shown, demonstrating that CHAINS can
3 model the dynamics of the competition and the emergence of perceptual organizations in
4 multistable auditory stimulus configurations **in a way that closely resembles** perceptual
5 reports of human listeners. Thus CHAINS **demonstrates** that the principles of AERS provide a
6 viable basis for computational models of auditory stream segregation.
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
61
62
63
64
65

References

- 1
2
3
4
5 **Althen H, Grimm S, Escera C (2013) Simple and complex acoustic regularities are**
6 **encoded at different levels of the auditory hierarchy. *European Journal of***
7 ***Neuroscience*. doi: 10.1111/ejn.12346. [Epub ahead of print]**
- 8 **Baldeweg T (2007) ERP repetition effects and mismatch negativity generation - A**
9 **predictive coding perspective. *Journal of Psychophysiology* 21 (3-4):204-213**
- 10 **Beauvois MW, Meddis R (1991) A computer-model of auditory stream segregation.**
11 ***Quarterly Journal of Experimental Psychology Section a-Human Experimental***
12 ***Psychology* 43 (3):517-541**
- 13 **Bendixen A, Bohm TM, Szalárdy O, Mill R, Denham SL, Winkler I (2013) Different**
14 **roles of similarity and predictability in auditory stream segregation. *Learning***
15 ***and Perception* 5:37-54**
- 16 **Bendixen A, Denham SL, Gyimesi K, Winkler I (2010) Regular patterns stabilize**
17 **auditory streams. *Journal of the Acoustical Society of America* 128 (6):3658-3666**
- 18 **Bendixen A, Prinz WG, Horváth J, Trujillo-Barreto NJ, Schröger E (2008) Rapid**
19 **extraction of auditory feature contingencies. *Neuroimage*, 41(3): 1111-1119.**
- 20 **Bendixen A, Roeber U, Schröger E (2007) Regularity extraction and application in**
21 **dynamic auditory stimulus sequences. *Journal of Cognitive Neuroscience* 19**
22 **(10):1664-1677**
- 23 **Bendixen A, SanMiguel I, Schröger E (2012a) Early electrophysiological indicators for**
24 **predictive processing in audition: A review. *International Journal of***
25 ***Psychophysiology* 83 (2):120-131.**
- 26 **Bendixen A, Schröger E, Ritter W, Winkler I (2012b) Regularity extraction from non-**
27 **adjacent sounds. *Frontiers in Psychology* 3:143.**
- 28 **Bendixen A, Schröger E, Winkler I (2009) I heard that coming: event-related potential**
29 **evidence for stimulus-driven prediction in the auditory system. *The Journal of***
30 ***Neuroscience*, 29 (26): 8447-8451.**
- 31 **Boh B, Herholz SC, Lappe C, Pantev C (2011) Processing of complex auditory patterns**
32 **in musicians and nonmusicians. *PLoS One*, 6 (7): e21458.**
- 33 **Bregman AS (1990) Auditory scene analysis. The perceptual organization of sound. MIT**
34 **Press, Cambridge, MA**
- 35 **Costa-Faidella J, Grimm S, Slabu L, Diaz-Santaella F, Escera C (2011) Multiple time**
36 **scales of adaptation in the auditory system as revealed by human evoked**
37 **potentials. *Psychophysiology* 48 (6):774-783**
- 38 **Cowan N (1984) On short and long auditory stores. *Psychological Bulletin* 96 (2):341-**
39 **370.**
- 40 **Cowan N, Winkler I, Teder W, Näätänen R (1993) Memory prerequisites of mismatch**
41 **negativity in the auditory event-related potential (ERP). *Journal of Experimental***
42 ***Psychology-Learning Memory and Cognition* 19 (4):909-921**
- 43 **Deike S, Heil P, Böckmann-Barthel M, Brechmann A (2012) The build-up of auditory**
44 **stream segregation: a different perspective. *Frontiers in Psychology* 3.**
- 45 **Denham SL, Gyimesi K, Stefanics G, Winkler I (2013) Perceptual bistability in auditory**
46 **streaming: How much do stimulus features matter? *Learning & Perception* 5**
47 **(2):73-100.**
- 48 **Elhilali M, Shamma SA (2008) A cocktail party with a cortical twist: How cortical**
49 **mechanisms contribute to sound segregation. *Journal of the Acoustical Society of***
50 ***America* 124 (6):3751-3771.**
- 51 **Escera C, Leung S, Grimm S (in press) Deviance detection based on regularity encoding**
52
53
54
55
56
57
58
59
60
61
62
63
64
65

along the auditory hierarchy: electrophysiological evidence in humans. *Brain Topography*, this issue.

- Fowler CA, Rosenblum LD (1990) Duplex perception - a comparison of monosyllables and slamming doors. *Journal of Experimental Psychology-Human Perception and Performance* 16 (4):742-754**
- Friston K, Kiebel S (2009a) Cortical circuits for perceptual inference. *Neural Networks* 22 (8):1093-1104**
- Friston K, Kiebel S (2009b) Predictive coding under the free-energy principle. *Philos Trans R Soc Lond B Biol Sci* 364 (1521):1211-1221.**
- Garrido MI, Kilner JM, Stephan KE, Friston KJ (2009) The mismatch negativity: A review of underlying mechanisms. *Clinical Neurophysiology* 120 (3):453-463.**
- Griffiths TD, Warren, JD (2004) Opinion: What is an auditory object? *Nature Review Neuroscience* 5: 887-892**
- Grimm S, Escera C (2012) Auditory deviance detection revisited: Evidence for a hierarchical novelty system. *International Journal of Psychophysiology* 85 (1):88-92.**
- Grossberg S, Govindarajan KK, Wyse LL, Cohen MA (2004) ARTSTREAM: a neural network model of auditory scene analysis and source segregation. *Neural Networks* 17 (4):511-536.**
- Haenschel C, Vernon DJ, Dwivedi P, Gruzelier JH, Baldeweg T (2005) Event-related brain potential correlates of human auditory sensory memory-trace formation. *Journal of Neuroscience* 25 (45):10494-10501**
- Helmholtz Hv (1867) *Handbuch der physiologischen Optik. Allgemeine Encyclopädie der Physik*, vol Bd 9. Voss, Leipzig**
- Herrmann CS, Munk MHJ, Engel AK (2004) Cognitive functions of gamma-band activity: memory match and utilization. *Trends in Cognitive Sciences* 8 (8):347-355.**
- Horváth J, Czigler I, Sussman E, Winkler I (2001) Simultaneously active pre-attentive representations of local and global rules for sound sequences in the human brain. *Cognitive Brain Research* 12 (1):131-144**
- Horváth J, Winkler I, Bendixen A (2008) Do N1/MMN, P3a, and RON form a strongly coupled chain reflecting the three stages of auditory distraction? *Biological Psychology* 79 (2):139-147**
- Jacobs AM, Grainger J (1994) Models of visual word recognition - sampling the state-of-the-art. *Journal of Experimental Psychology-Human Perception and Performance* 20 (6):1311-1334.**
- Jones MR (1976) Time, our lost dimension: toward a new theory of perception, attention, and memory. *Psychological Review* 83 (5):323-355.**
- Kiebel SJ, von Kriegstein K, Daunizeau J, Friston KJ (2009) Recognizing sequences of sequences. *Plos Computational Biology* 5 (8):e1000464**
- Köhler W (1947) *Gestalt Psychology: An introduction to new concepts in modern psychology*. Liveright Publishing, New York**
- Kubovy M, Van Valkenburg D (2001) Auditory and visual objects. *Cognition* 80: 97-126**
- Lieder F, Daunizeau J, Garrido MI, Friston KJ, Stephan KE (2013) Modelling Trial-by-Trial Changes in the Mismatch Negativity. *Plos Computational Biology* 9 (2).**
- Massaro DW (1972) Preperceptual images, processing time, and perceptual units in auditory perception. *Psychological Review* 79:124-145**
- McCabe SL, Denham MJ (1997) A model of auditory streaming. *Journal of the Acoustical Society of America* 101 (3):1611-1621.**
- Mill R, Bohm T, Bendixen A, Winkler I, Denham SL CHAINS - Competition and cooperation between fragmentary event predictors in a model of auditory scene**

- analysis. In: *Information Sciences and Systems (CISS), 2011 45th Annual Conference on*, Baltimore, 2011. IEEE, pp 1-6.
- 1
2
3 **Mill RW, Bohm TM, Bendixen A, Winkler I, Denham SL (2013) Modelling the**
4 **Emergence and Dynamics of Perceptual Organisation in Auditory Streaming.**
5 **Plos Computational Biology 9 (3).**
- 6 **Moore BCJ, Gockel H (2002) Factors influencing sequential stream segregation. Acta**
7 **Acust United Acust 88 (3):320-333**
- 8 **Moore BCJ, Gockel HE (2012) Properties of auditory stream formation. Philosophical**
9 **Transactions of the Royal Society B-Biological Sciences 367 (1591):919-931.**
- 10 **Müller D, Widmann A, Schröger E (2005). Deviance-repetition effects as a function of**
11 **stimulus feature, feature value variation, and timing: a mismatch negativity**
12 **study. Biological Psychology, 68(1): 1-14.**
- 13
14 **Mumford D (1992) On the computational architecture of the neocortex II. The role of**
15 **cortico-cortical loops. Biol Cybern 66 (3):241-251.**
- 16
17 **Näätänen R (1990) The Role of attention in auditory information-processing as revealed**
18 **by event-related potentials and other brain measures of cognitive function.**
19 **Behavioral and Brain Sciences 13 (2):201-232**
- 20
21 **Näätänen R, Gaillard A, Mäntysalo S (1978) Early selective-attention effect on evoked**
22 **potential reinterpreted. Acta Psychologica 42:313-329**
- 23
24 **Näätänen R, Kujala T, Winkler I (2011) Auditory processing that leads to conscious**
25 **perception: A unique window to central auditory processing opened by the**
26 **mismatch negativity and related responses. Psychophysiology 48 (1):4-22.**
- 27
28 **Näätänen R, Winkler I (1999) The concept of auditory stimulus representation in**
29 **cognitive neuroscience. Psychological Bulletin 125 (6):826-859**
- 30 **Neisser U (1967) Cognitive Psychology. Appleton-Century-Crofts, New York**
- 31 **Paavilainen P, Arajärvi P, Takegata R (2007). Preattentive detection of nonsalient**
32 **contingencies between auditory features. Neuroreport 18(2): 159-163.**
- 33
34 **Rao RPN, Ballard DH (1999) Predictive coding in the visual cortex: a functional**
35 **interpretation of some extra-classical receptive-field effects. Nature Neuroscience**
36 **2 (1):79-87.**
- 37
38 **Rensink RA (2000) Seeing, sensing, and scrutinizing. Vision Research 40 (10-12):1469-**
39 **1487.**
- 40
41 **Rinne T, Särkkä A, Degerman A, Schröger E, Alho K (2006) Two separate mechanisms**
42 **underlie auditory change detection and involuntary control of attention. Brain**
43 **Research 1077:135-143**
- 44
45 **Ritter W, De Sanctis P, Molholm S, Javitt DC, Foxe JJ (2006) Preattentively grouped**
46 **tones do not elicit MMN with respect to each other. Psychophysiology 43(5): 423-**
47 **430.**
- 48
49 **Ritter W, Sussman E, Molholm S (2000) Evidence that the mismatch negativity system**
50 **works on the basis of objects. Neuroreport 11(1): 61-63.**
- 51
52 **Schröger E (2007) Mismatch negativity - A microphone into auditory memory. Journal**
53 **of Psychophysiology 21 (3-4):138-146**
- 54
55 **Schwartz JL, Grimault N, Hupé J-M, Moore BC, Pressnitzer D (2012) Multistability in**
56 **perception: binding sensory modalities, an overview. Philos Trans R Soc Lond B**
57 **Biol Sci 367 (1591):896-905**
- 58
59 **Sokolov EN (1963) Higher nervous functions: The orienting reflex. Annual Review of**
60 **Physiology 25:545-580.**
- 61
62 **Sussman ES (2005) Integration and segregation in auditory scene analysis. Journal of**
63 **the Acoustical Society of America 117(3 Pt 1):1285-98.**
- 64
65 **Sussman ES, Gumenyuk V (2005) Organization of sequential sounds in auditory**
memory. Neuroreport 16 (13):1519-1523

- 1 Sussman ES, Horváth J, Winkler I, Orr M (2007) The role of attention in the formation
2 of auditory streams. *Perception & Psychophysics* 69 (1):136-152
- 3 Szalárdy O, Winkler I, Schröger E, Widmann A, Bendixen A (in press) Foreground-
4 background discrimination indicated by event-related brain potentials in a new
5 auditory multistability paradigm. *Psychophysiology*. doi: 10.1111/psyp.12139.
6 [Epub ahead of print]
- 7 van Noorden LPAS (1975) *Temporal Coherence in the Perception of Tone Sequences*.
8 Technical University, Eindhoven
- 9 Wacogne C, Changeux JP, Dehaene S (2012) A Neuronal Model of Predictive Coding
10 Accounting for the Mismatch Negativity. *Journal of Neuroscience* 32 (11):3665-
11 3678.
- 12 Wacogne C, Labyt E, van Wassenhove V, Bekinschtein T, Naccache L, Dehaene S
13 (2011) Evidence for a hierarchy of predictions and prediction errors in human
14 cortex. *Proceedings of the National Academy of Sciences of the United States of*
15 *America* 108 (51):20754-20759.
- 16 Wang DL, Chang P (2008) An oscillatory correlation model of auditory streaming.
17 *Cognitive Neurodynamics* 2 (1):7-19.
- 18 Winkler I (2007) Interpreting the mismatch negativity. *Journal of Psychophysiology* 21
19 (3-4):147-163
- 20 Winkler I, Czigler I (1998) Mismatch negativity: deviance detection or the maintenance
21 of the 'standard'. *Neuroreport* 9 (17):3809-3813
- 22 Winkler I, Czigler I (2012) Evidence from auditory and visual event-related potential
23 (ERP) studies of deviance detection (MMN and vMMN) linking predictive coding
24 theories and perceptual object representations. *International Journal of*
25 *Psychophysiology* 83 (2):132-143.
- 26 Winkler I, Denham S, Mill R, Bohm TM, Bendixen A (2012) Multistability in auditory
27 stream segregation: a predictive coding view. *Philosophical Transactions of the*
28 *Royal Society B-Biological Sciences* 367 (1591):1001-1012.
- 29 Winkler I, Denham SL, Nelken I (2009) Modeling the auditory scene: predictive
30 regularity representations and perceptual objects. *Trends in Cognitive Sciences*
31 13 (12):532-540
- 32 Winkler I, Karmos G, Näätänen R (1996) Adaptive modeling of the unattended acoustic
33 environment reflected in the mismatch negativity event-related potential. *Brain*
34 *Research* 742 (1-2):239-252
- 35 Winkler I, Korzyukov O, Gumenyuk V, Cowan N, Linkenkaer-Hansen K, Ilmoniemi
36 RJ, Alho K, Näätänen R (2002) Temporary and longer term retention of acoustic
37 information. *Psychophysiology* 39 (4):530-534
- 38 Winkler I, Takegata R, Sussman E (2005). Event-related brain potentials reveal multiple
39 stages in the perceptual organization of sound. *Cognitive Brain Research*, 25 (1),
40 291-299.
- 41 Winkler I, van Zuijen TL, Sussman E, Horváth J, Näätänen R (2006) Object
42 representation in the human auditory system. *European Journal of Neuroscience*
43 24 (2):625-634.
- 44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Figures

AERS

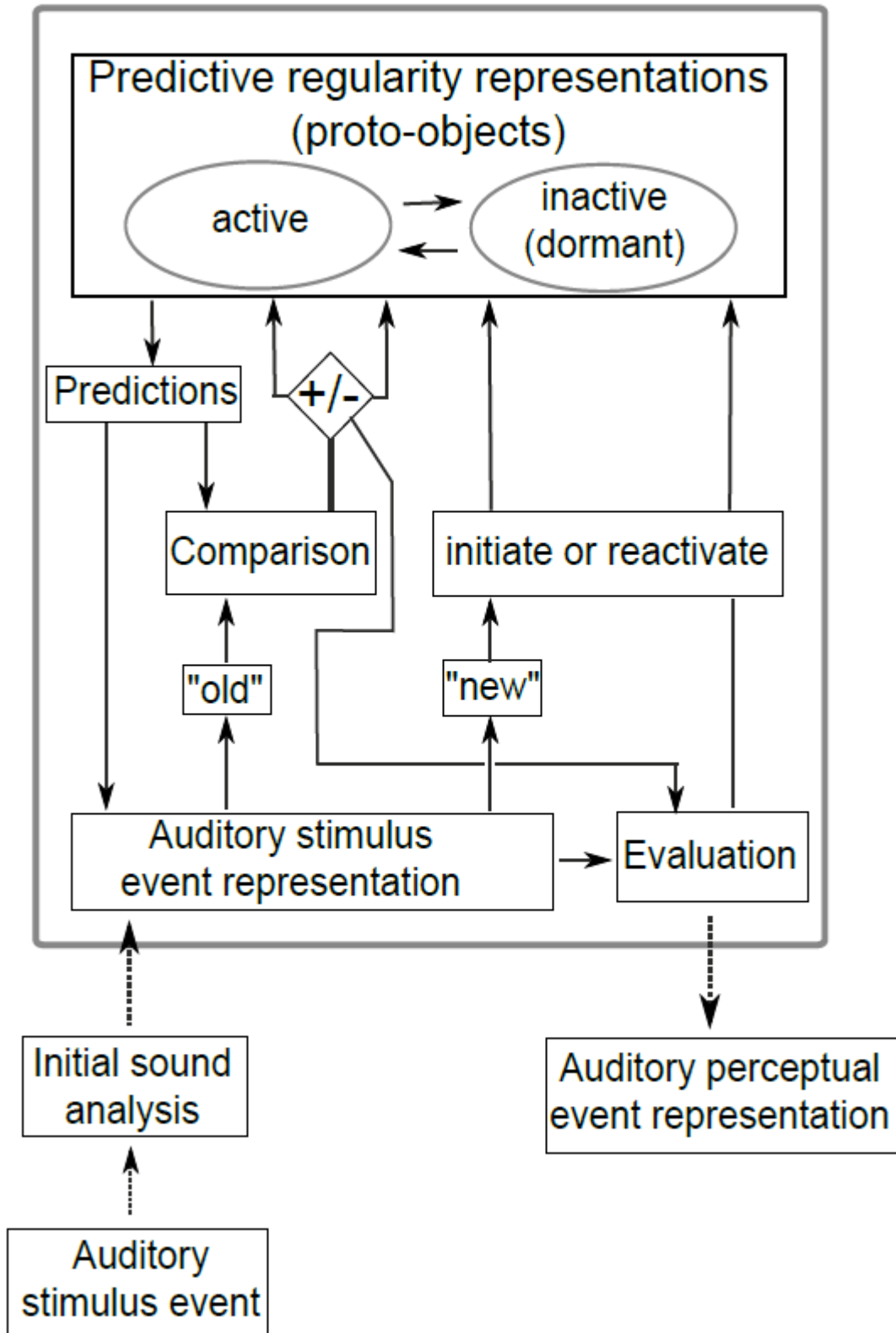


Figure 1. AERS model. The primary input to AERS are the incoming auditory stimulus

1 events with their basic features established by processes concerned with the initial analysis of
2 the sound. Predictive regularity representations encode detected regularities and predict the
3 upcoming sounds. The established auditory stimulus event representations (which, in turn, are
4 biased by the predictions) are compared with the predictions. The outcome of this comparison
5 is used for updating the predictive regularity representations and for the subsequent evaluation
6 process. There, the auditory stimulus event representations are related to auditory context and
7 to the current goals of the organism. The output of AERS is an auditory perceptual event
8 representation (e.g. a particular tone of a flute). They can enter various mental operations and
9 be consciously perceived. Please note, that this event representation is linked to the respective
10 auditory object representation (e.g. the flute). In addition, the evaluation process can initiate
11 the building of new or reactivate old but inactive regularity representations; for more details
12 see text, for the definition of terms see Table 1.
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
61
62
63
64
65

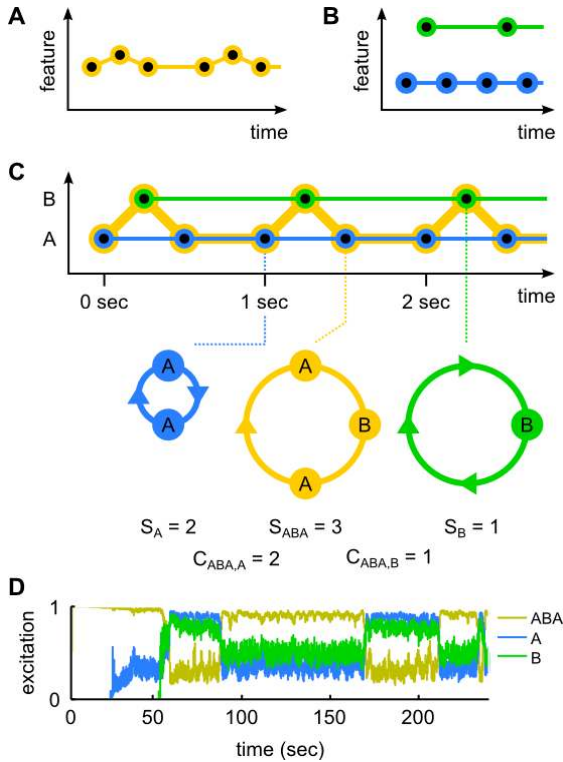


Figure 2. CHAINS model. A) Schematic illustration of circumstances that favor integration.

Events plotted in a time-feature space (solid black markers) vary gradually in feature distance and consequently form into a single chain ('ABA', yellow). Links from A to A are difficult to form, however, because it is difficult to skip over B when it forms a smooth continuation with the As, and the same is true for linking the Bs. B) Schematic illustration of circumstances that favor segregation. Events vary abruptly in feature distance. It is therefore easy to form a chain consisting solely of As (green), because it is probable that 'B' will be successfully skipped.

For the same reason, a complementary chain consisting solely of Bs is probable (blue).

However, the 'ABA' chain is difficult to form, owing to the improbability of building many abrupt links. C) Illustration showing how various aspects of a single 'ABA_' sequence are explained by three chains, with some overlap. (The example assumes that all links that can form, do so with certainty.) At the point where a chain contains a repeating subsequence, it closes to form a predictive loop (shown below). D) Time-varying excitations of the 'ABA', 'A' and 'B' chains (E_{ABA} , E_A and E_B) in competition. Successes and collisions between the chains define a competition that gives rise to organizations that are stable for short periods.

Either 'ABA' dominates alone (integration), or 'A' and 'B' dominate together.

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
61
62
63
64
65

Figure1
[Click here to download high resolution image](#)

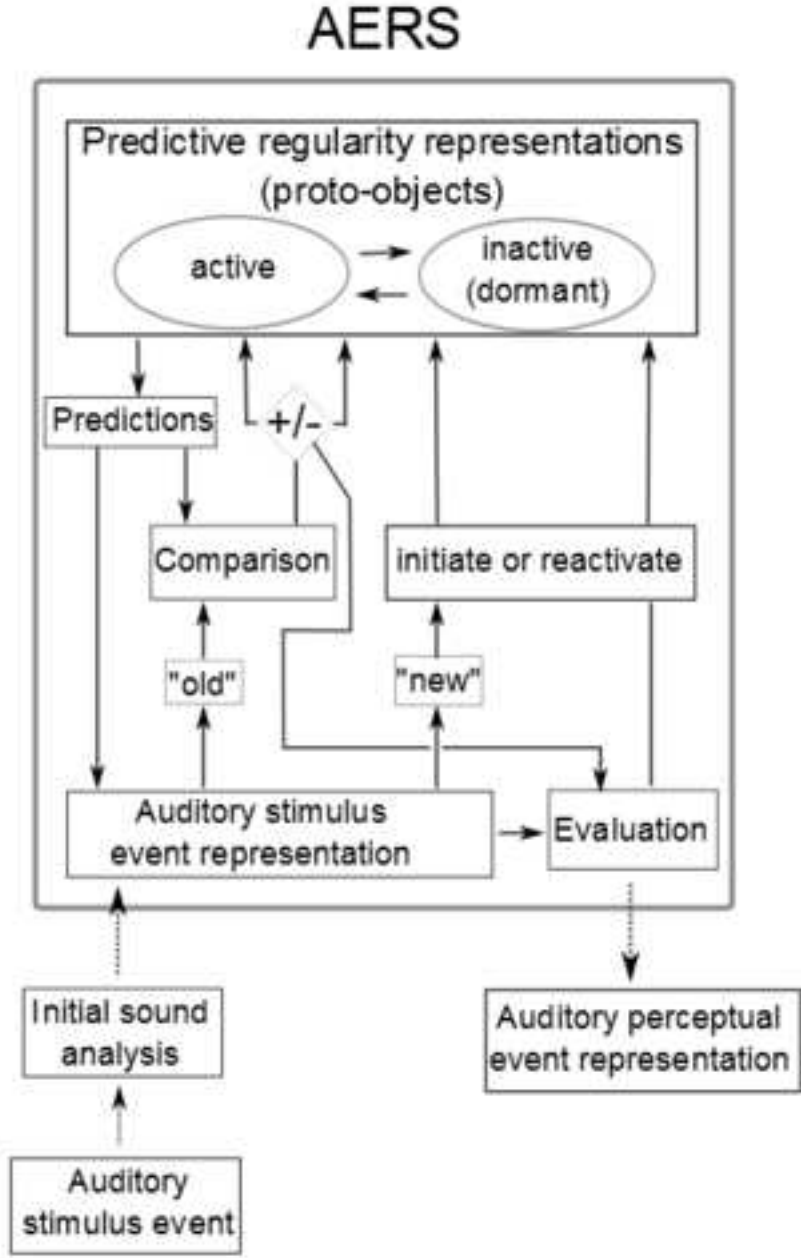


Figure2
[Click here to download high resolution image](#)

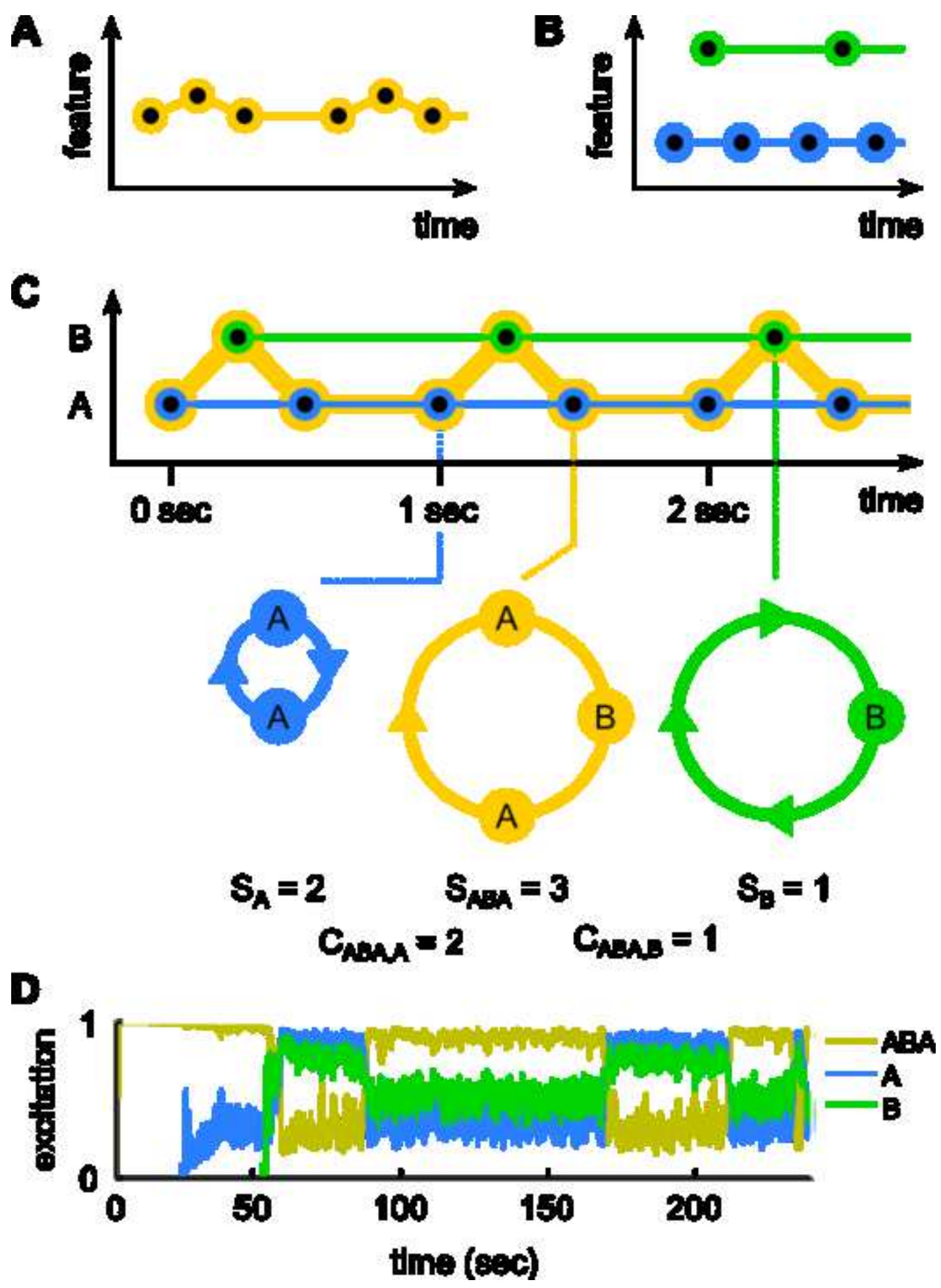


Table 1: Definition of terms

Auditory event representation system (AERS): A cognitive system producing *auditory perceptual event representations* from *auditory stimulus events*.

Auditory perceptual object (or auditory object representation): A member of the set of currently dominant *proto-objects* that occupies perceptual awareness. ‘Auditory stream’ is a similar term but one which is not explicitly identified in terms of predictive representations.

Auditory predictive regularity representation* or *proto-object: Terms used to describe the representation of a sequence of sounds linked together by some detected rule or repeating pattern in the form of a generative model. Incoming sounds are checked against the predictions of current *proto-objects* that compete to ‘explain’ them. These representations have the potential to become the perceptual objects in conscious awareness, if and when they are dominant (i.e., they are in a highly activated state, assumed to indicate their selection as the most likely description of the current input).

Auditory stimulus event: A discrete sound, localized in time and generated by some source in the external world; i.e. the physical input to our sensory systems (e.g., each of the sounds in the particular sequence of sounds generated by our mobile phone when a text message arrives).

Auditory stimulus event representation: The integrated description of the perceived features of an *auditory stimulus event*; is shaped by the predictions from *AERS*.

Auditory perceptual event representation: The description of an *auditory stimulus event* in the brain; an *auditory sensory stimulus representation*, which is linked to a perceptual object, and expanded with the description of its relationship to the auditory and the general context (e.g. the text-message sound as it appears in our perception). Auditory perceptual event representations are the output of *AERS*.

CHAINS: A computational model (and its implementation into a Matlab/C-based computer program) that incorporates aspects of *AERS*. It can be used to simulate possible *perceptual organizations* for specific parameters, which – in turn – can be tested experimentally.

Chains: A formal description (within *CHAINS*) of a sequence of *auditory stimulus events* including their timing.

Initial sound analysis: The early (bottom-up) activation patterns in the auditory system caused by each *auditory stimulus event*; this analysis is not regarded as part of *AERS*.

Perceptual organization: A complete description of the auditory environment, in terms of auditory object representations, as it appears in perception.