

ANNALS *of* THE NEW YORK
ACADEMY OF SCIENCES

Neural bases of learning and recognition of statistical regularities

| | |
|-------------------------------|---|
| Journal: | <i>Ann NY Acad Sci</i> |
| Manuscript ID | annals-2000-196.R4 |
| Manuscript Type: | Original Article |
| Date Submitted by the Author: | 18-Dec-2019 |
| Complete List of Authors: | Ordin, Mikhail; Basque Center on Cognition Brain and Language, Consciousness; Ikerbasque, Consciousness research group Polyanskaya, Leona; Basque Center on Cognition Brain and Language, Spoken Language Research Group Soto, David; BCBL, Consciousness research group; Ikerbasque, Basque Foundation for Science |
| Keywords: | statistical learning, segmentation, statistical generalization |
| | |

SCHOLARONE™
Manuscripts

Neural bases of learning and recognition of statistical regularities

Mikhail Ordin^{1,2} and Leona Polyanskaya¹ and David Soto^{1,2}

¹BCBL – Basque Centre on Cognition, Brain and Language and ²Ikerbasque – Basque Foundation for Science, San Sebastián, Spain

Address for correspondence: Basque Centre on Cognition, Brain and Language, Paseo Mikeletegi 69, San Sebastián, 20009 Spain. m.ordin@bcbl.eu

Short title: *Neural bases of statistical learning*

Graphical abstract

Previous studies have only assessed brain responses to trained words and novel non-words, and hence do not provide sufficient information on how the brain mediates the recognition of word-like units versus mere statistical regularities within sequences. The present study addresses this issue, as well as determines whether learning of statistical regularities embedded into a continuous sensory input and discrete constituents comprising this input, and subsequent recognition of extracted constituents, relies on the same mechanisms.

Abstract

Statistical learning is a set of cognitive mechanisms allowing for extracting regularities from the environment and segmenting continuous sensory input into discrete units. The current study used functional MRI (N = 25) in conjunction with an artificial language learning paradigm to provide new insight into the neural mechanisms of statistical learning, considering both the online process of extracting statistical regularities and the subsequent offline recognition of learned patterns. Notably, prior fMRI studies on statistical learning have not contrasted neural activation during the learning and recognition experimental phases. Here we found that learning is supported by the superior temporal gyrus and the anterior cingulate gyrus, while subsequent recognition relied on the left inferior frontal gyrus. Besides, prior studies only assessed the brain response during the recognition of trained words relative to novel non-words. Hence a further key goal of this study was to understand how the brain supports recognition of discrete constituents from the continuous input vs. recognition of mere statistical structure that is used to build new constituents that are statistically congruent with the ones from the input. Behaviorally, recognition performance indicated that statistically congruent novel tokens were less likely to be endorsed as parts of the familiar environment than discrete constituents. fMRI data showed that the left intraparietal sulcus and angular gyrus support the recognition of old discrete constituents relative to novel statistically congruent items, likely reflecting an additional contribution from memory representations for trained items.

Keywords: statistical learning; segmentation; statistical generalization; fMRI; sensory input; information

Introduction

1
2
3 Statistical learning allows agents to detect regularities in the world around them. These
4 statistical cues can be used to split the continuous flow of sensory information (visual, auditory,
5 tactile) into discrete constituents, a process called segmentation. The neural mechanisms
6 underlying segmentation are evolutionary ancient¹⁻⁴ and shared by a diverse range of species⁵⁻
7
8 ⁸. Segmentation based on statistical cues contained in the sensory input operates across
9
10 different domains. In humans, this includes splitting speech into words and phrases⁹, separating
11
12 distinct rhythms and other musical properties [in musical compositions]¹⁰⁻¹², parsing sequences
13
14 of events as well as discerning discrete sequences of actions in a continuous series of human
15
16 activities^{13,14}. For example, while viewing a series of still images representing a continuous
17
18 dynamic activity, viewers' gaze tends to dwell longer on those slides that illustrate the
19
20 boundaries between unfolding events: dwell times are longer for slides that show the grasp of a
21
22 glass has been completed than those that show the grasping action still unfolding. This
23
24 segmentation takes place at various levels: slides that depict boundaries between distinct
25
26 higher-level actions, for example, the boundary between emptying a dishwasher (which includes
27
28 the lower-level action of grasping a glass to take it out of the machine) and starting a new
29
30 sequence of sweeping the floor, attracts even longer gazes¹⁵. It is more difficult to predict
31
32 actions that follow boundaries so they attract more attention. By contrast, when the next action
33
34 can be easily predicted based on previously observed events, less attention is required – and
35
36 dwell times diminish – because the further unfolding of events is highly predictable. These
37
38 results were interpreted within the framework of statistical learning¹⁶. Besides, it has been
39
40 suggested that the segmentation of actions, of continuous sensory input across modalities, and
41
42 segmentation of speech into linguistic constituents like words and phrases – all rely on the same
43
44 cognitive processes related to statistical learning^{14,16}.

41 Statistical learning operates on a variety of cues, including (but not limited to) conditional
42 regularities known as transitional probabilities (TPs). TPs refer to the probability of an event B
43 happening given that event A has occurred. Higher TPs characterize the events that commonly
44 happen sequentially one after another, while lower TPs are aligned with the boundaries
45
46 between the sequences of commonly co-occurring events. Thus, the differences between high
47
48 and low TPs between events allow breaking continuous flow of events into discrete sequences,
49
50 with events within sequences being more predictable than those spanning the edges of these
51
52 sequences. Although this tokenization, or segmentation mechanism has been shown to operate
53
54 across domains (see references above), it has been most extensively studied in the context of
55
56 speech processing. For instance, during speech processing, continuous stream of syllables (i.e.,
57
58 events) can be segmented into separate words by calculating transitional probabilities (TPs)
59
60

1
2
3 between adjacent syllables¹. For example, in the phrase “*pretty baby*”, the probability that the
4 syllable “*ty*” will follow syllable “*pret*” is higher than the probability “*ba*” will follow “*ty*”. Minima in
5 TPs between adjacent syllables, compared to surrounding TPs (i.e., local minima), are aligned
6 with word boundaries. They can be used by infants learning their first language or adults
7 exposed to a foreign language to segment a sequence of syllables into discrete words¹⁷.
8
9

10
11 What is still debated, however, is how new constituents are identified. Do listeners detect
12 word boundaries between consecutive constituents based on lower TPs, or do they merge
13 smaller frequently co-occurring units into a single constituents⁹? Some researchers advocate for
14 clustering mechanisms^{18,19}, while others argue in favor of boundary-finding mechanisms^{20–22}.
15
16 Some studies demonstrate that both human and non-human animals employ various strategies,
17 which might rely on different neural mechanisms, and show that the choice of the strategy is
18 determined by individual preferences, peculiarities of sensory input, and environmental
19 circumstances^{9,23}.
20
21
22
23

24 The use of statistical cues in speech segmentation is usually studied within the artificial
25 language learning paradigm²⁴. A set of artificial words is concatenated into a continuous
26 acoustic stream, with each word in the stream recurring multiple times. The syllable pairs with
27 lower inter-syllabic TPs are more likely to straddle word boundaries than syllable pairs with
28 higher TPs, which, in turn, are more likely to be confined within word boundaries. This enables
29 the segmentation of the continuous syllable stream into words. Performance in speech
30 segmentation tasks is tested by habituating the listener to the constructed acoustic stream, then
31 administering a recognition test, in which participants need to endorse or reject test tokens as
32 word candidates. Test tokens can either be statistical words from the learning stream, or
33 sequences that violate the statistical regularities embedded in the habituation stream. Empirical
34 studies convincingly demonstrate that words are endorsed as legal word candidates, while
35 syllabic sequences that violate statistical regularities are rejected²⁵. We aim to explore the
36 neural bases of statistical learning in the context of this speech segmentation paradigm using
37 the interesting case of so-called phantom words. Phantoms are test tokens that conform to the
38 acquired statistical regularities but never occurred during habituation. For instance, listeners
39 may be exposed to recurrent syllabic triplets, including XYA and BYZ triplets, with 0.5
40 transitional probabilities between syllables within triplets. Although syllabic pairs XY and YZ
41 frequently occur in the familiarization stream, the sequence XYZ is never presented during
42 habituation. It is still not clear if or how the brain differentiates between familiar structural units
43 (i.e. artificial words) and novel units that are structurally congruent with the old ones. Addressing
44 this question promises to provide important insights into the neurocognitive bases of statistical
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 learning, namely, whether the brain relies on structural regularities, memory representations or
4 a mixture of both while endorsing different types of tokens as legitimate candidates as discrete
5 constituents in a continuous environmental input. Observing differences in brain responses to
6 words and phantoms would suggest that different cognitive mechanisms are employed in the
7 recognition of these tokens, notably, even when tokens of both types are not differentiated
8 behaviorally. Phantoms can only be endorsed based on the recognition of their congruence with
9 the statistical regularities embedded in the sensory input, while words may additionally rely on
10 memory representations of whole discrete elements.
11
12
13
14
15

16 To date, there is conflicting evidence as to whether, after exposure to artificial language,
17 phantoms emerge as perceptual units during recognition. Not all researchers have observed
18 that phantoms are confused with holistic triplets (i.e. word-like structures)^{21,26}. Furthermore,
19 individual differences and the native language of the listener can influence whether or not
20 phantoms are confused with words^{22,27}. A key goal of the present study is to provide novel
21 insights into the mechanisms supporting statistical learning using a novel behavioral and
22 neuroimaging protocol to partial out the processing of words, phantoms and pseudorandom
23 sequences during learning and subsequent recognition.
24
25
26
27
28

29 Earlier studies have shown that learning based on transitional probabilities (TPs) in the
30 auditory modality is supported by the superior temporal gyrus (STG)²⁸⁻³¹ and inferior frontal
31 gyrus (IFG)^{28,32,33}, mainly – but not exclusively – in a left-lateralized network. These studies have
32 suggested that the STG is involved in learning TPs in the acoustic input, while the IFG is
33 hypothesized to support the learning of word-like units²⁸. As learning progresses, brain
34 responses to occasional violations of statistical regularities can also be observed in
35 frontoparietal cortex^{34,35}, particularly in the control network, which includes the right angular and
36 bilateral anterior cingulate gyri^{28,29,35,36}. Detection of violations of statistical structure during
37 training is also supported by the temporoparietal junction^{34,37}, in line with its established role in
38 attentional re-orienting to unexpected stimuli³⁸.
39
40
41
42
43
44

45 However, previous studies have only assessed brain responses to trained words and
46 novel non-words. Hence, they do not provide sufficient information on how the brain mediates
47 the recognition of word-like units versus the recognition of mere statistical regularities (i.e., TPs)
48 within sequences. The present study addresses this issue for the first time. A second key
49 objective of this study is to determine whether learning of statistical regularities embedded into a
50 continuous sensory input and discrete constituents comprising this input on the one hand, and
51 subsequent recognition of extracted constituents on the other hand relies on the same
52 mechanisms^{9,28,35}. This is an issue that is hard to tackle using behavioral measures alone.
53
54
55
56
57
58
59
60

1
2
3 Cognitive processes related to memory encoding and retrieval may rely on shared neural
4 mechanisms^{39,40}. However, most prior fMRI studies that used a statistical learning paradigm
5 only recorded brain activity online (during learning) or offline (later, during recognition). Critically,
6 brain responses were not recorded during both learning and recognition within the same
7 experiment. Additionally, the only fMRI study that examined both the learning and recognition
8 stages²⁹ did not test for neural differences between these two phases. Here, we recorded BOLD
9 responses during both the learning and recognition phases in order to determine whether online
10 statistical learning and subsequent offline recognition were supported by similar or different
11 brain substrates.
12
13
14
15
16
17
18
19

20 **Methods**

21 *Participants*

22
23 We analyzed the data from 25 native Spanish participants (11 males between 20 and 33,
24 average age 25.5 years, SD = 3.29). MRI data from one participant was discarded because he
25 did not follow instructions (i.e. during the recognition test he pressed the same button for all
26 responses). All participants had acquired Basque in childhood after the age of two as their
27 second language and were using it daily. We note that prior studies have not revealed
28 differences between monolingual Spanish and bilingual Spanish-Basque participants in the
29 segmentation of statistical units in off-line recognition tests^{41,42}. Also, to mitigate the possibility
30 that individual differences related to bilingualism might influence statistical learning processes,
31 we homogenized our sample by matching participants by proficiency and age of acquisition (2–3
32 years) of the second language, as well as the self-reported extent of their daily exposure to and
33 use of Basque. None of the participants had any prior history of neurological disorders. The
34 experiment was approved by the BCBL ethical committee. All participants provided informed
35 consent.
36
37
38
39
40
41
42
43
44
45
46
47
48
49

50 *Experimental materials and procedures*

51
52 **Learning phase.** The habituation stream was composed of alternating structured and random
53 blocks. Structured blocks consisted of concatenated syllable triplets, so that higher TPs
54 between syllables within these triplets and lower TPs between syllables straddling the triplet
55
56
57
58
59
60

1
2
3 boundaries allow for predicting the following syllable after both triplet-initial and triplet-medial
4 syllables, but not after triplet-final syllables. Random blocks consisted of the same syllables
5 unsystematically concatenated so that the TPs between them were uniform throughout and
6 therefore did not allow for segmentation.
7
8

9
10 The habituation stream was synthesized using MBROLA, with ES1 voice and
11 fundamental frequency invariably set to 110Hz. We used a set of 18 consonant-vowel syllables.
12 The duration of each syllable was 240 ms (100 ms for consonants and 140 ms for vowels).
13 These syllables were used to construct 12 trisyllabic statistical words, with TPs between
14 syllables within triplets set to 0.5. These triplets were randomly concatenated 21 times with the
15 restriction that the same word was never repeated consecutively. TPs between syllables
16 straddling triplet boundaries were approximately 0.16. The difference in TPs between syllables
17 within triplets and TPs straddling triplet boundaries provided statistical cues for the
18 segmentation of the continuous stream of syllables into recurrent trisyllabic sequences (words).
19 This procedure was repeated three times to create three structured blocks of 181.44 sec. We
20 also pseudorandomly concatenated all 18 syllables from the artificial language syllabic inventory
21 six times, such that no syllable was ever repeated consecutively. We prepared three
22 pseudorandom blocks of 25.92 sec each. The same syllable inventories were used during
23 structured and pseudorandom blocks because our goal was to test for differences in brain
24 responses to the presence or absence of statistical structure without introducing any confounds
25 due to differences between the acoustic properties of the stimuli. Also, we elected not to use
26 rest periods as the baseline but rather contrasted activity in structured and pseudorandom
27 blocks during learning. It was likely that overall brain state in pseudorandom blocks would be
28 more similar to that in structured blocks than during rest. Therefore, we could be confident that
29 any differences in neural responses to random and structured blocks were elicited because
30 structured blocks included recognizable structure with extractable and learnable constituents.
31 Contrasting structured blocks to rest would also have introduced the problem of individual
32 variability because mental activity during rest can vary, engaging different mechanisms and
33 different networks⁴³.
34
35
36
37
38
39
40
41
42
43
44
45
46

47
48 The habituation syllabic stream was prepared by alternating structured and
49 pseudorandom blocks three times. At the end of the stream, we added an additional 36
50 randomized syllables (each of the 18 syllables from the inventory repeated twice for a total of
51 8.64 sec) and applied a fade-out filter. A fade-in filter was also applied at the beginning of each
52 stream, thus preventing any potential anchoring effects of stream-initial and stream-final
53 syllables on segmentation. As statistical learning mechanisms are constantly operating on
54
55
56
57
58
59
60

1
2
3 incoming sensory input, it is possible – though unlikely, given the difference in exposure time to
4 structured and random blocks – that participants collapsed conditional statistics across
5 conditions. However, this possibility does not undermine the validity of the cues for statistical
6 segmentation since the syllable pairs with higher TPs were more likely to fall within the recurrent
7 triplets.
8
9
10

11 We prepared two similar habituation streams which had unique orders for recurrent
12 triplets within structured blocks and for syllables within pseudorandom blocks. Streams 1 and 2
13 were used for the first and second runs of the learning session. During the learning phase,
14 participants were asked to listen to an “extraterrestrial language” and to try to detect and
15 memorize the words from that language.
16
17
18
19
20
21

22 **Recognition phase.** The recognition test was comprised of four test runs each comprising 63
23 trials. In each trial, we randomly concatenated either 4 different triplets from the habituation
24 stream, 4 different phantoms – triplets that fit the statistical regularities of the habituation input
25 but had never occurred in the learning stream as whole constituents, or 4 non-words (i.e.,
26 triplets composed of syllables that never occurred consecutively in the habituation stream).
27 Each run included 21 trials of each type. The duration of the stimuli in the recognition test was
28 2880 ms. Each triplet was used an equal number of times across all trials. The stimuli were
29 preceded and followed by 200 ms silence. After each stimulus presentation, participants were
30 asked to decide whether that acoustic sequence had been presented during the learning
31 session, and then to rate their confidence in the given response, on a 4-point scale. The period
32 for each response was fixed to 2000 ms. The trials were separated by a jittered time interval
33 according to a pseudo-exponential distribution from 3000 ms to 5000 ms in steps of 500 ms.
34
35
36
37
38
39
40
41

42 Both the learning and recognition phases were performed inside the scanner. The sound
43 was played via in-ear Sensimetrics S14 headphones. A pair of headphones was placed above
44 the in-ear headphones in order to dampen the noise of the scanner and to enable
45 communication with the experimenter. The stimuli were back projected onto a screen by a
46 mirror on the head coil. The area between the participant’s head and the coil was padded with
47 foam to make the participant more comfortable and to minimize head movements. We asked
48 the participants not to move during scanning.
49
50
51
52

53 To familiarize the participants with the procedure and the experimental protocol and
54 interface, a brief training session was organized outside the scanner, with a 40-second
55
56
57
58
59
60

1
2
3 familiarization stream and 4 recognition trials. The syllables for this training session were
4 different from those used in the actual experiment. The list of statistical words, phantoms and
5 non-words are given in Table 1. The structure of the learning runs and recognition trials is
6 illustrated in Figure 1.
7
8
9

10 11 12 *Functional and structural MRI data acquisition*

13
14 Whole-brain MRI data acquisition was conducted in a 3T MAGNETOM PRISMAfit MR scanner
15 using a 64-channel coil. T1-weighted images were acquired using MP-RAGE sequences with
16 the following parameters: TR = 2530 ms, TE = 2,36 ms, FoV = 256 mm, flip angle = 7 degrees,
17 acquiring 176 contiguous 1 cubic mm slices per run. Functional images were acquired using a
18 multi-band acceleration factor of 6 (multi-slice interleaved mode), with 66 contiguous 2.4 cubic
19 mm slices, TR = 850 ms, TE = 35 ms, flip angle = 56 degrees. We achieved whole-brain
20 coverage.
21
22
23
24
25
26
27

28 29 *fMRI data pre-processing*

30
31 Image pre-processing was performed in FSL 5.0.9 using the FEAT module. First, we used the
32 brain extraction tool (BET⁴⁴) to separate the brain matter from non-brain tissues. The first 11
33 volumes of each run, both in learning and in recognition tests, were discarded to control for
34 magnetic saturation effects and allow for MR signal stabilization. We used a high-pass filter
35 cutoff of 100 sec for the learning runs and of 60sec for the testing runs following FSL manual
36 instructions for blocked and event-related designs. Scans were realigned by using MCFLIRT
37 motion correction (spatial smoothing with a 6-mm FWHM Gaussian kernel applied). Translation
38 parameters did not exceed half a voxel in any direction for any participant in any run. Functional
39 images were registered to T1 structural images (7 degrees of freedom for testing runs and using
40 the boundary based registration BBR algorithm⁴⁵ for the learning runs). Then, the images were
41 registered to the standard MNI152 template using affine registration with 12 degrees of
42 freedom, using full search setting.
43
44
45
46
47
48
49
50
51
52

53 54 *fMRI data analysis*

1
2
3 **Learning phase.** Statistical analysis for the learning stream was performed within the
4 framework of the general linear model in the individual native space first, with statistical maps
5 normalized to the standard space prior to higher-level analyses. Each pseudorandom block was
6 entered as a separate explanatory variable (EV). Critically, chunks with duration of 25.92
7 seconds of structured blocks immediately preceding the pseudorandom blocks (also 25.92
8 seconds) were entered as separate EVs in order to match the two conditions in the number of
9 scans so that the amount of data for the comparison of the BOLD responses during structured
10 and pseudorandom blocks is equated. The numbers of scans were equated in order to have a
11 similar signal for the contrast of blood oxygen level dependent (BOLD) differences between
12 structured and pseudorandom blocks. See Rosenthal *et al.*⁴⁶ for detailed methodological
13 justifications for the necessity to equate the amount of data while comparing the BOLD
14 response between two conditions.
15
16
17
18
19
20
21
22

23 Each EV specified the onset of the pseudorandom block or the onset of the structured
24 chunk. The EVs were introduced in the model along with their temporal derivatives. We applied
25 FILM pre-whitening⁴⁷. Standard and extended motion parameters were introduced in the GLM
26 model as additional regressors of no interest⁴⁸. Furthermore, we used the FSL motion outliers
27 function (<https://fsl.fmrib.ox.ac.uk/fsl/fslwiki/FSLMotionOutliers>) and included the regressors
28 corresponding to the motion outliers in the design matrix in order to deal with the effect of
29 intermediate-to-large motions that could potentially corrupt images beyond the level that the
30 extended motion parameter regression methods could deal with. To detect the volumes
31 containing motion outliers, we calculated the root mean square (RMS) head position difference
32 to the reference volume and compared the 75 percentile +1.5 inter-quartile range of the
33 distribution of RSM values for each run. A confound matrix was generated and used in the GLM
34 to completely remove the effects of these timepoints on the analysis.
35
36
37
38
39
40
41

42 Parameter estimates were calculated for the following contrasts in which brain activity
43 was higher for pseudorandom relative to structured blocks (R1>S1, R2>S2, R3>S3, R4>S4,
44 R5>S5, R6>S6, where “R” stands for pseudorandom, “S” stands for structured, and the number
45 indicates the sequential number of the pseudorandom block or immediately preceding
46 structured chunk, i.e., R1, R2, and R3 are parts of run 1 and R4, R5, and R6 are parts of run 2.
47 The resulting 6 contrasts of parameter estimates reflect acquisition of relevant sequence
48 knowledge.
49
50
51
52

53 We then performed a second level, fixed-effects analysis within each participant using
54 the 6 parameter estimates noted above. Here we tested for different temporal profiles of
55
56
57
58
59
60

1
2
3 learning related activity (S<R and S>R) across the training phase. We assessed a logarithmic
4 trend (specified using the following contrast: -3.125, -1.0, 0.5, 1.0, 1.25, 1.375) and an
5 exponential increase (specified using the following contrast: -1.375, -1.25, -1.0, 0.5, 1.0, 3.125).
6 The logarithmic increase suggests that the difference in BOLD change on pseudorandom and
7 structured blocks is larger at the beginning of the learning phase and attenuates as learning
8 progresses. This trend can be caused by faster learning at the beginning of exposure, with the
9 learning rate then decelerating. Exponential increase is a reciprocal function, which, on the
10 contrary, suggests that learning is slower at the beginning, and accelerates with time, causing
11 the differences in the BOLD response on pseudorandom and structured blocks to increase more
12 rapidly as habituation progresses.
13
14
15
16
17
18

19 The output of the contrasts was fed into a mixed-effect model using the FLAME 1
20 algorithm in FSL⁴⁹ in order to test for the consistent effect across participants (group $Z > 2.3$,
21 cluster significance threshold $P = 0.05$, corrected using Gaussian field theory).
22
23
24
25

26 **Recognition test.** Statistical analysis for the recognition test was first conducted within the
27 framework of the general linear model in native space, with statistical maps normalized to the
28 standard space prior to higher-level, group analyses. We created the EVs for words endorsed
29 as words (*words_acc*), rejected non-words (*nonw_rej*), rejected and endorsed phantoms
30 (correspondingly *phan_rej* and *phan_acc*). We modelled the onset of each EV with durations
31 that corresponded to the length of the stimulus (2.88 secs). Regressors of no interest were
32 introduced into the design matrix as separate EVs to control for variation in decision response
33 time both for the first recognition response (i.e., whether the stimulus was presented during the
34 learning session) and for the confidence rating. The EVs were introduced along with their
35 temporal derivatives. We applied FILM pre-whitening with standard and extended motion
36 parameters and introduced an additional EV for the motion outliers.
37
38
39
40
41
42
43

44 We then estimated contrasts of parameter estimates that were relevant to our study
45 goals. In particular, we assessed the brain substrates that support (1) the recognition of words
46 vs. non-words and critically (2) the recognition of words vs phantoms. In these analyses we
47 used trials with correct responses (i.e. *words_acc* vs. *nonw_rej* for (1) and *words_acc* vs.
48 *phantoms_rej* as well as *words_acc* vs. *phantoms_acc* for (2). The contrasts were estimated in
49 both directions, i.e., for the contrast *words_acc* vs. *nonw_rej* we estimated contrasts of
50 parameter estimates both for *words_acc > acc-nonw_rej* and *words_acc < nonw_rej*. We
51 reasoned that accepting words can rely both on the recognition of word-like structural
52
53
54
55
56
57
58
59
60

1
2
3 information present in the memory traces of the triplets and the statistical structure contained in
4 the lower-level TPs. The comparison of the BOLD signal change between words and phantoms
5 ought to reveal the brain substrates that support the recognition of word-like structural
6 information. The trials with phantoms involve analysis of statistical structure and making
7 decisions based only on statistical congruency. Alternatively, rejecting phantoms may rely on
8 the fact that they are not supported by memory representations because the phantoms were not
9 encountered and extracted as whole units during the learning phase. Hence, we expected the
10 memory network to be more activated for accepted words than both rejected and accepted
11 phantoms.
12
13
14
15
16

17 We performed within-subject, cross-run (fixed effects) analysis to estimate the individual
18 mean of each contrast across all runs of the recognition test. In order to find effects that were
19 consistent across subjects, all contrasts were fed into a mixed effect model for the whole-brain
20 group analysis using FLAME 1, thresholding the statistic images ($Z > 2.3$, $P = 0.05$ corrected
21 using Gaussian Random Field theory^{49,50}.
22
23
24
25
26
27
28

29 **Results**

30 *Behavioral results*

31 Behavioral data was acquired only during the recognition test. We estimated the percentage of
32 endorsed words, phantoms and non-words. Also, we calculated the mean confidence rating
33 assigned to endorsed and rejected words, non-words and phantoms. Although the percentage
34 of correct responses might seem somewhat lower than what is usually reported in artificial
35 language learning experiments, it is not extraordinarily low. The environment, in which
36 participants had to do the task was more challenging (i.e. performed inside the scanner against
37 strong background noise), the number of discrete constituents was larger (12 triplets) than what
38 is usually used in similar experiments (4 triplets), and the differences in TPs between syllables
39 within words and between syllables spanning word boundaries was less pronounced (50% vs.
40 16%) than what is usually used (100% vs. 33%). All these factors increased the difficulty of the
41 task. In addition, the interleaved random blocks could also have had a detrimental effect on
42 learning.
43
44
45
46
47
48
49
50
51
52

53 We performed one-sample t-tests comparing the percentage of endorsed tokens in each
54 condition (words, phantoms and non-words) with chance level performance (50%). The results
55
56
57
58
59
60

(Fig. 2) show that the percentage of endorsed words is significantly above what would be expected by chance, $t(24) = 3.224$, $P = 0.004$, M (mean difference) = 12.08%, 95%CI [4.35: 19.81], $d = 0.645$. The percentage of endorsed non-words is significantly below what would be expected by chance, $t(24) = -6.379$, $P < 0.0005$, $M = -23.63\%$, 95%CI [-34.27: -15.98], $d = 1.276$. The percentage of endorsed phantoms is not significantly different from chance, $t(24) = 1.285$, $P = 0.211$, $M = 10.71\%$, 95%CI [-3.25: 13.95], $d = 0.257$.

Repeated-measures ANOVA revealed significant differences in the percentage of endorsed words, phantoms and non-words, $F(2, 48) = 60.009$, $P < 0.0005$, $\eta_p^2 = 0.714$ (p value corrected with the Greenhouse-Geisser method, df are reported uncorrected). Pairwise comparison (with the Bonferroni correction applied) showed that the proportion of endorsed words was significantly higher than that of endorsed phantoms, $t(24) = 3.371$, $P = 0.008$, $M = 6.72\%$, 95%CI[2.6: 10.84], $d = 0.674$. The percentage of endorsed phantoms was also significantly higher than the percentage of endorsed non-words, $t(24) = 7.245$, $P < 0.0005$, $M = 28.98\%$, 95%CI[20.73: 37.24], $d = 1.449$. Unsurprisingly, the proportion of endorsed words was also significantly higher than the proportion of endorsed non-words, $t(24) = 8.922$, $P < 0.0005$, $M = 35.71\%$, 95%CI[27.45: 43.97], $d = 1.785$.

We then analyzed the confidence ratings. The results showed that endorsed words were assigned significantly higher confidence ratings compared to incorrect, rejected words, $t(23) = -4.253$, $P < 0.0005$, $M = 0.334$, 95%CI [-0.497: -0.172], Cohen's $d = 0.96$. The same pattern was found for non-words, $t(24) = -2.428$, $P = 0.023$, $M = -0.23$, 95%CI [-0.425: -0.035], Cohen's $d = 0.65$. Interestingly, the level of confidence assigned to endorsed phantoms was significantly higher than rejected phantoms, $t(24) = 4.451$, $P < 0.0005$, $M = 0.353$, 95%CI [0.189: 0.517], Cohen's $d = 0.812$. As Figure 3 shows, the level of confidence for accepted phantoms matches that of accepted words, $t(24) = 1.118$, $P = 0.275$. Similar results were found for confidence ratings on trials with rejected words and rejected phantoms, $t(24) = 0.237$, $P = 0.814$.

Taken together, these results show that participants have learnt discrete constituents from the auditory input and reliably endorse them later during the test, while rejecting those sequences which violated the statistical regularities embedded in the familiarization stream. However, when participants encountered phantoms, which were consistent with the statistical probabilities defining the structural constituents but had never been presented as holistic units during the learning phase, participants were at a loss, and could not unambiguously accept or reject them. They responded randomly, at a chance level, which resulted in a significantly higher proportion of accepted words than accepted phantoms. However, once a decision was made, participants assigned higher ratings to accepted than to rejected phantoms, showing that the

1
2
3 metacognitive system treated acceptance of phantoms as a correct response. This lack of
4 difference in confidence ratings assigned to correctly endorsed words and accepted phantoms
5 suggests that metacognitively phantoms were treated as words, even though the cognitive
6 system treated words and phantoms differently.
7
8
9

10 11 *fMRI results*

12
13
14 **Learning phase.** Our goal here was to delineate the neural correlates of learning related
15 changes during the study phase. Accordingly, we tested the effect of training on the differences
16 in BOLD response between the successive structured (S) and pseudorandom (R) chunks (i.e.
17 our index of learning). We compared two types of models in which (1) learning related activity
18 mainly occurred during the first training run and then remained constant during the second
19 training run, and (2) learning related changes emerged in the second training run (see
20 Methods). Based on prior research on perceptual sequence learning, we elected to focus our
21 analyses on the S<R contrast^{46,51–53}. For the sake of completeness, we also conducted a similar
22 analysis based on S>R parameter estimates (as was done in Ref. 28), however, here we did not
23 find any significant results at the group level.
24
25
26
27
28
29
30

31 The significant fit of learning related activity with a logarithmic trend indicates that
32 learning-related brain activity builds faster at the beginning of the exposure and is then
33 attenuated as training progresses. This was found in three clusters: (1) superior-frontal gyrus
34 (SFG) extending to the paracingulate gyrus; (2) right superior temporal gyrus (STG); and (3) left
35 STG. Figure 4 illustrates these results. Table 2 provides information regarding the peak voxels
36 in MNI coordinates of the different contrasts. We did not find any evidence for linear or
37 exponential increases in learning related activity across the two training runs suggesting that
38 learning increases did not continue as training progressed further in the second run of the
39 learning phase.
40
41
42
43
44
45
46
47

48 **Recognition phase.** Following the learning phase, participants were presented with a
49 recognition test. On each trial, previously studied words, phantoms or non-words appeared for
50 an old/new recognition decision followed by confidence ratings. Our key goal was to isolate the
51 neural substrates implementing recognition of word units relative to phantoms (i.e. statistically
52 congruent tokens that were not embedded in the learning input). We therefore ran three
53 contrasts. First, we compared BOLD activity changes when words were accepted relative to
54
55
56
57
58
59
60

1
2
3 when non-words were rejected. This contrast is similar to comparing pseudorandom and
4 structured blocks during the learning phase. Most crucially, we then compared BOLD activity
5 changes when words were endorsed relative to when phantoms were rejected as well as when
6 words were endorsed relative to when phantoms were endorsed. These contrasts aimed to
7 isolate the brain substrates activated by the recognition of word-like units as whole constituents
8 relative to the recognition of merely statistical structure. We reasoned that phantoms may be
9 accepted as legitimate elements of an artificial language due to their statistical congruency with
10 word-like constituents^{22,27}. Rejecting phantoms may therefore rely on the fact that they are not
11 supported by memory representations; phantoms were not encountered and extracted as whole
12 units during the learning phase. Hence, we reasoned that by splitting the tokens into accepted
13 and rejected the chances of dissociating the brain basis of words vs phantom processing would
14 increase.
15
16
17
18
19
20
21
22

23 Endorsed words, compared to rejected non-words (*words_acc > non-words_rej*), elicit
24 BOLD response changes in the left inferior frontal gyrus (LIFG) around BA44, pars opercularis
25 extending to par triangularis in BA45 (see Table 2). Critically, endorsed words relative to
26 rejected phantoms (*words_acc > phan_rej*) elicited BOLD increases in the anterior part of the
27 cingulate cortex, posterior division of the STG (strongly right lateralized), and in the left
28 hemisphere in a cluster that involved the posterior division of the angular gyrus and anterior part
29 of the intra-parietal sulcus (see Table 3). Figure 5 illustrates these results.
30
31
32
33

34 Finally, we report that no differences were found when comparing accepted words vs.
35 accepted phantoms, which suggests that the recognition of equally familiar tokens may have a
36 similar neural underpinning, whether the recognition is based merely on statistical congruency,
37 or strengthened by memory representations of word-like units. This result is in line with the
38 behavioral data on confidence, showing that accepted words are treated as accepted phantoms,
39 and that participants processed the trials, in which they endorsed phantoms, as correct
40 responses.
41
42
43
44
45
46
47

48 **Comparison of learning and recognition phases.** In order to compare neural substrates
49 supporting online statistical learning with offline recognition of holistic constituents, we
50 compared the brain activity maps associated with learning-related activity for structured vs.
51 pseudorandom blocks, and the brain activity maps associated with the recognition of words
52 relative to non-words. Note that structured/pseudorandom blocks during learning equate to the
53 presentation of words/non-words during recognition, both acoustically and statistically. Hence,
54
55
56
57
58
59
60

1
2
3 contrasting the recognition effect with the learning effect in terms of the associated BOLD
4 activity maps will reveal whether these processes are supported by similar neural substrates.
5
6

7 Learning and recognition contrasts of parameter estimates were fed into a whole-brain
8 mixed effects model paired t-test with subjects as random factors. We found that learning-
9 related activity was stronger in the left STG and the right superior frontal gyrus, extending to the
10 anterior cingulate cortex ($Z > 2.3$, $P = 0.05$ with a corrected significance level using Gaussian
11 Random Field theory^{49,50}). Table 4 and Figure 6 illustrates these results.
12
13
14
15

16 17 **Discussion**

18
19
20 The key objectives of the current study were (1) addressing the common and distinct neural
21 substrates that support online statistical learning and the subsequent offline recognition of the
22 learned constituents, and (2) defining the neural substrates that support the recognition of
23 holistic constituents (learned words) as opposed to recognition of merely statistical structure
24 (phantoms). Behaviorally, we found that participants could successfully recognize words vs non-
25 words in terms of discrimination accuracy and response confidence. Higher confidence ratings
26 were assigned to correctly endorsed words and correctly rejected non-words compared to the
27 corresponding incorrect responses. The proportion of accepted words was also significantly
28 higher than that of phantoms. Importantly, rejected phantoms were assigned lower confidence
29 rating compared to endorsed phantoms, showing that on rejected phantoms participants
30 estimate the likelihood of making an error to be higher than on accepted phantoms. At the same
31 time, endorsed phantoms and words were assigned similar confidence. Overall, the results
32 show that accepted phantoms were treated as accepted words, *metacognitively*, although in
33 terms of accuracy (i.e., cognitive decisions), they were not confused with words. One of the
34 driving forces that underlie endorsing sensory input as part of the environment is the
35 concordance with the statistical regularities the individuals are familiar with. This motivates their
36 response to both words and phantoms. Memory representations further strengthen the
37 recognition of word-like units. Sometimes, phantoms are endorsed despite the lack of memory
38 support for these tokens as whole constituents, suggesting that memory does not play a crucial
39 role in the recognition of legitimate constituents^{21,54}. The fact that statistically congruent tokens
40 are sometimes endorsed in the absence of memory representations would allow for the transfer
41 of processing skills from recently encountered to novel situations, as long as these novel
42 situations exhibit recently encountered statistical features.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 We turn now to the neuroimaging results. We found that the bilateral STG supports the
4 online extraction of conditional statistical cues during learning. This is in keeping with a number
5 of earlier studies both in the auditory and visual modalities^{29,31,34} and with prior work
6 demonstrating a role for STG in associative learning and relational memory^{55,56}, here related to
7 statistical structure and the acquisition of the relational positions of syllables in a stream.
8 Learning was also mediated by the cognitive control network. The level of processing load
9 during auditory perception is known to regulate activity in the control network, especially in the
10 paracingulate and anterior cingulate gyri⁵⁷. Accordingly, we found an increased activity in these
11 areas when pseudorandom sequences were presented during learning following structured
12 sequences.
13
14
15
16
17
18

19 Notably, we did not observe that activity in the IFG changed differently for
20 pseudorandom and structured blocks during the learning phase. This result seems at odds with
21 the study of Ablam and Okanoya³², who showed that online segmentation of recurrent tone
22 sequences in a continuous tone stream elicits higher levels of activity in the LIFG. Karuza et
23 al.²⁸ however, found involvement of the LIFG for learning forward speech but not backward
24 speech, and suggested that the results by Ablam and Okanoya³² and their own results²⁸ reveal
25 the role of the LIFG in TP calculation and the formation of structural representations.
26 Importantly, in the Ablam and Okanoya's³² study, participants were first trained on three-tone
27 sequences presented in isolation, and later these tone triplets were concatenated in a
28 continuous stream and presented in alternation with the same tones randomly concatenated
29 (i.e., not built into triplets). As participants were already familiarized with the recurrent triplets
30 prior to the exposure, the activation in the LIFG could actually indicate a neural response to the
31 recognition of the already learnt constituents rather than the process of formation of new
32 representations. The role of the LIFG in the recognition of learned constituents rather than
33 online segmentation of TPs is shown by Turk-Browne *et al.*⁵⁶, who correlated familiarity ratings,
34 assigned to discrete constituents during the recognition test with activity in the LIFG during
35 learning exposure. Higher activity in the LIFG was observed for those constituents which were
36 later rated as more familiar, indicating that neural responses in the LIFG differed for recognized
37 vs. unrecognized tokens embedded into a continuous sensory input. Our findings are also in line
38 with the hypothesis that activation in the LIFG is related to the recognition of discrete
39 constituents rather than to the online segmentation of continuous input. By assessing learning
40 and recognition processes independently within the same paradigm, our study allowed us to
41 isolate the contributions of the left STG and the anterior cingulate cortex to learning, while the
42 left IFG was critically involved in subsequent, offline recognition processes. We believe that the
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 LIFG is implicated in monitoring that statistically congruent sequences are indeed discrete
4 constituents learned from the environment, which is in line with its proposed role as a general
5 domain and modality independent sequence processor²⁸. The right STG seems to be equally
6 involved in learning and recognition. The conclusion that statistical learning and recognition are
7 supported by different neural substrates also agrees with previous studies in other domains and
8 modalities, which have shown differences in the neural networks supporting successful
9 encoding (i.e., learning) and successful retrieval (i.e., recognition) of semantic and perceptual
10 associations⁵⁵. The lack of differential involvement of the IFG in the present study during
11 learning and recognition is in keeping with the proposed account.

12
13
14 For the first time, in a functional MRI of statistical learning, phantom sequences were
15 included in the recognition test phase alongside words. This allowed us to determine whether
16 the processing of phantoms and words is mediated by a similar neurocognitive mechanism. Our
17 behavioral evidence suggests that participants are confused by phantom cases and are as likely
18 to accept as to reject them. Hence, we tested the extent to which processing of words and
19 phantoms could be dissociated in brain responses. Overall, BOLD responses to endorsed
20 words and endorsed phantoms did not differ, confirming the conclusion that accepted phantoms
21 are metacognitively considered to be correct responses, and that once a phantom is accepted, it
22 is processed as a legitimate structural constituent. However, when we compared BOLD
23 responses for accepted words and rejected phantoms, we found stronger responses in the left
24 angular gyrus and intraparietal sulcus associated with endorsed words. This reveals active
25 activation of the memory network³⁹ elicited by retrieving memory representations of words as
26 whole constituents presented during learning. Also, we found that the level of neural activity
27 associated with accepted words and rejected phantoms differed in the anterior division of the
28 cingulate cortex (ACC), which is frequently related to error detection and conflict resolution⁵⁹. A
29 significant difference in neural activation in the ACC for accepted words vs rejected phantoms
30 suggests that on the trials in which phantoms were rejected a competing response was present,
31 and thus the participants estimated the likelihood of making an error on such trials as high. This
32 is manifested in the low confidence ratings assigned to rejected phantoms. Absence of
33 difference in neural activation between accepted words and phantoms confirms our earlier
34 conclusion, based on the behavioral results: accepted phantoms are treated as correctly
35 endorsed words. This pattern of results also invites a tentative explanation, which nevertheless
36 needs to be further empirically tested. We propose that the tokens are accepted mainly based
37 on recognition of statistical structure, while rejection is based on the lack of memory
38 representation. Thus, it is possible that memory representations do not yield additional support
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 for recognizing constituents from the sensory input. It is rather the absence of memory
4 representations that leads to the rejection of some statistically congruent items. We believe
5 these findings are important considering that prior behavioral studies have not consistently
6 observed that phantoms are confused with holistic triplets (i.e. word-like structures). Hence this
7 work lays the foundation for future studies to further explore the conditions in which the brain
8 can distinguish words vs phantoms, for instance, by manipulating the amount of exposure
9 during the training phase. New insights may be provided by exploring the different cognitive
10 mechanisms which underlie rejection and acceptance of novel statistically congruent tokens.
11 This understanding might be important, for example, in the field of language learning, to explain
12 the phenomenon of generalization, when a rule learnt on a small set of examples is generalized
13 over previously unencountered cases, and the phenomenon of fossilization, when the transfer
14 from known to novel situations does not happen and progress in learning is halted.

15
16 As we argued in the introduction, statistical learning mechanisms are evolutionary
17 ancient, operate across a wide variety of taxonomically divergent species, and predate the
18 emergence of language. Hence, although these mechanisms are engaged in speech processing
19 and language acquisition, it is very unlikely that they evolved specifically for these purposes. We
20 suggest that statistical learning mechanisms evolved to detect abrupt changes in the
21 environment. The structure of ecologically relevant natural states is usually relatively stable, with
22 rapid transitions between longer lasting stable states^{60,61}. For survival and reproduction (i.e.,
23 fitness), organisms need to monitor the environment and detect and react to rapid ecologically-
24 relevant changes as they suddenly occur. These fitness needs likely gave rise to the early
25 emergence of statistical learning mechanisms during evolution and explain their spread across
26 different taxa, domains and modalities. Detecting structural regularities is probably more
27 important than detecting recurrent constituents, because any breach in ecological stability
28 signals rapid and fitness-relevant environmental changes. The ability to detect statistical
29 structure presumably predates and underlies the segmentation of the dynamical flow of
30 experience and supports building abstract representations of segmented constituents that
31 reflect the structure of the environment. As phantoms in our study were statistically congruent
32 with words, the cognitive system might have confused some of them as being equally familiar as
33 words at the behavioral level, because both correspond to stable states in the statistical
34 structure of the acoustic environment. Nevertheless, additional support from memory
35 representations (evidence for the memory support of words relative to phantoms is discussed
36 below) affects confidence judgements. However, the presence of this support at the neural level
37 does not override the importance of detecting breaches in statistical structure which signal

1
2
3 environmental changes. In the environment of evolutionary adaptiveness that shaped the
4 functions of the statistical learning mechanisms breaches in the statistical structure rather than
5 recurrent constituents had to be monitored and required a behavioral response because they
6 cued sudden ecological changes.
7
8

9
10 It may be argued that this activity contrast is related to endorsement vs. rejection,
11 irrespective of underlying statistical structure. Several considerations argue against this. First,
12 the BOLD response to the auditory sequence was modelled separately from that associated
13 with the behavioral response in the recognition test that took place 4 seconds later. Critically,
14 the BOLD activity patterns that we report are time locked to the onset of the auditory sequence.
15
16 Second, the patterns of BOLD activity were found in putative substrates of statistical learning.
17
18 Finally, the fact that there is a brain signal that distinguishes words (accepted) vs phantoms
19 (rejected) likely reflects the contribution of memory representations derived throughout the
20 training. Statistically congruent novel items (i.e., phantoms) that did not receive additional
21 support from activation in the memory-related brain areas were rejected. If the level of the
22 memory activation was the same for words and phantoms, then phantoms were accepted as
23 constituents of the previously experienced sensory input. The role of the left intraparietal sulcus
24 in memory is well established^{39,58}. Previous studies have also found post-learning sensitivity in
25 the angular gyrus to the presentation of statistically congruent sequences^{28,35}. The angular
26 gyrus also underlies discrimination of pseudo-words and real words from natural languages⁶².
27 The relations between linguistic experience and the functionality of the left angular gyrus is
28 supported by work by Mechelli *et al.*⁶³, which showed that bilinguals and highly proficient L2
29 learners have stronger grey matter density in the anterior division of the angular gyrus than
30 monolinguals. The important role of the angular gyrus in the recognition of artificial words may
31 be facilitated by its strong connectivity with temporal cortices via the arcuate fasciculus⁶⁴ and
32 also with the inferior frontal gyrus, both BA 44⁶⁵ and BA 45⁶⁶, via the longitudinal fasciculus.
33 These connections are ipsilateral, which explains the simultaneous left-lateralized activation in
34 multiple cortical areas, communicating through a major connection hub of the left angular
35 gyrus⁶⁷. Functional connectivity studies have also revealed a broad cortical network involved in
36 statistical learning, with strong functional connections between both left and right STG, which
37 we also observed here, alongside the LIFG⁶⁸. Overall, our recognition results are in line with the
38 conclusion of Skosnik *et al.*³⁵ that grammaticality and recognition judgements rely on different
39 networks. Since both words and phantoms are statistically congruent, the recognition of words
40 must rely on additional support from memory representations derived from the learning phase,
41 which recruits additional neural networks, different from those engaged in phantom recognition.
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

To conclude, we found that the neural substrates underlying online statistical learning processes and offline recognition of the learned patterns rely on different neural substrates, indicating that the neurocognitive mechanisms that support the initial formation and subsequent maintenance of structural representations are distinct. Also, we dissociated, at the neural level, recognition of discrete constituents from the sensory input vs. recognition of mere statistical structure that is used to build constituents, enabling recognition of novel constituents that were not experienced before. Mechanisms for statistical learning have been shaped by fitness needs in the environment of evolutionary adaptiveness. We suggest that statistical mechanisms for detecting breaches in statistical structure are more essential to fitness than those that detect structural units; they are more evolutionary ancient and prevail over those that allow us to recognize structural units.

Acknowledgements

The research was supported by the Spanish Ministry of Economy and Competitiveness (MINECO) through the “Severo Ochoa” Programme for Centres/Units of Excellence in R&D (SEV-2015-490), and project grants RTI2018-098317-B-I00 awarded to MO, by the Basque Government through project grant PI-2017-25 awarded to DS, and by European Commission as Marie Skłodowska-Curie fellowship DLV-792331 to LP.

Competing interests

The authors declare no competing interests.

References

1. Kikuchi, Y., Sedley, W., Griffiths, T. D. & Petkov, C. I. Evolutionarily conserved neural signatures involved in sequencing predictions and their relevance for language. *Curr. Opin. Behav. Sci.* **21**, 145–153 (2018).
2. Milne, A. E., Petkov, C. I. & Wilson, B. Auditory and Visual Sequence Learning in Humans and Monkeys using an Artificial Grammar Learning Paradigm. *Neuroscience* **389**, 104–117 (2018).
3. Wilson, B. *et al.* Auditory Artificial Grammar Learning in Macaque and Marmoset Monkeys. *J. Neurosci.* **33**, 18825–18835 (2013).
4. Wilson, B. *et al.* Auditory sequence processing reveals evolutionarily conserved regions of frontal cortex in macaques and humans. *Nat. Commun.* **6**, 8901 (2015).
5. Saffran, J. *et al.* Grammatical pattern learning by human infants and cotton-top tamarin monkeys. *Cognition* **107**, 479–500 (2008).

- 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
6. Spierings, M. J. & Ten Cate, C. Budgerigars and zebra finches differ in how they generalize in an artificial grammar learning experiment. *Proc. Natl. Acad. Sci. U. S. A.* **113**, E3977-84 (2016).
7. van Heijningen, C. A. A., Chen, J., van Laatum, I., van der Hulst, B. & ten Cate, C. Rule learning by zebra finches in an artificial grammar learning task: which rule? *Anim. Cogn.* **16**, 165–175 (2013).
8. Wilson, B., Smith, K. & Petkov, C. I. Mixed-complexity artificial grammar learning in humans and macaque monkeys: evaluating learning strategies. *Eur. J. Neurosci.* **41**, 568–578 (2015).
9. Thiessen, E. D., Kronstein, A. T. & Hufnagle, D. G. The extraction and integration framework: A two-process account of statistical learning. *Psychol. Bull.* **139**, 792–814 (2013).
10. Deutsch, D. Grouping Mechanisms in Music. in *The Psychology of Music* 183–248 (Elsevier, 2013). doi:10.1016/B978-0-12-381460-9.00006-7
11. Povel, D.-J. & Essens, P. Perception of Temporal Patterns. *Music Percept. An Interdiscip. J.* **2**, 411–440 (1985).
12. Ravnani, A., Delgado, T. & Kirby, S. Musical evolution in the lab exhibits rhythmic universals. *Nat. Hum. Behav.* **1**, 0007 (2017).
13. Baldwin, D. A. & Baird, J. A. Discerning intentions in dynamic human action. *Trends Cogn. Sci.* **5**, 171–178 (2001).
14. Hard, B. M., Meyer, M. & Baldwin, D. Attention reorganizes as structure is detected in dynamic action. *Mem. Cognit.* **47**, 17–32 (2019).
15. Hard, B. M., Recchia, G. & Tversky, B. The shape of action. *J. Exp. Psychol. Gen.* **140**, 586–604 (2011).
16. Baldwin, D., Andersson, A., Saffran, J. & Meyer, M. Segmenting dynamic human action via statistical structure. *Cognition* **106**, 1382–1407 (2008).
17. Gambell, Timothy, Yang, C. *Word segmentation: Quick but not dirty.* (2006).
18. Frank, M. C., Goldwater, S., Griffiths, T. L. & Tenenbaum, J. B. Modeling human performance in statistical word segmentation. *Cognition* **117**, 107–125 (2010).
19. Perruchet, P. & Vinter, A. PARSER: A Model for Word Segmentation. *J. Mem. Lang.* **39**, 246–263 (1998).
20. Elman, J. L. Finding structure in time. *Cogn. Sci.* **14**, 179–211 (1990).
21. Endress, A. D. & Mehler, J. The surprising power of statistical learning: When fragment knowledge leads to false memories of unheard words. *J. Mem. Lang.* **60**, 351–367 (2009).
22. Endress, A. D. & Langus, A. Transitional probabilities count more than frequency, but

- 1
2
3 might not be used for memorization. *Cogn. Psychol.* **92**, 37–64 (2017).
4
5 23. Knowlton, B. J. & Squire, L. R. Artificial grammar learning depends on implicit acquisition
6 of both abstract and exemplar-specific information. *J. Exp. Psychol. Learn. Mem. Cogn.*
7 **22**, 169–181 (1996).
8
9 24. Erickson, L. C. & Thiessen, E. D. Statistical learning of language: Theory, validity, and
10 predictions of a statistical learning account of language acquisition. *Dev. Rev.* **37**, 66–108
11 (2015).
12
13 25. Aslin, R. N. & Newport, E. L. Statistical learning: From acquiring specific items to forming
14 general rules. *Curr. Dir. Psychol. Sci.* **21**, 170–176 (2012).
15
16 26. Perruchet, P. & Poulin-Charronnat, B. Beyond transitional probability computations:
17 Extracting word-like units when only statistical information is available. *J. Mem. Lang.* **66**,
18 807–818 (2012).
19
20 27. Saksida, A., Langus, A. & Nespors, M. Co-occurrence statistics as a language-dependent
21 cue for speech segmentation. *Dev. Sci.* **20**, e12390 (2017).
22
23 28. Karuza, E. A. *et al.* The neural correlates of statistical learning in a word segmentation
24 task: An fMRI study. *Brain Lang.* **127**, 46–54 (2013).
25
26 29. McNealy, K., Mazziotta, J. C. & Dapretto, M. Cracking the Language Code: Neural
27 Mechanisms Underlying Speech Parsing. *J. Neurosci.* **26**, 7629–7639 (2006).
28
29 30. Cunillera, T. *et al.* Time course and functional neuroanatomy of speech segmentation in
30 adults. *Neuroimage* **48**, 541–553 (2009).
31
32 31. Blakemore, S. J., Rees, G. & Frith, C. D. How do we predict the consequences of our
33 actions? A functional imaging study. *Neuropsychologia* **36**, 521–9 (1998).
34
35 32. Abia, D. & Okanoya, K. Statistical segmentation of tone sequences activates the left
36 inferior frontal cortex: A near-infrared spectroscopy study. *Neuropsychologia* **46**, 2787–
37 2795 (2008).
38
39 33. Petersson, K.-M., Folia, V. & Hagoort, P. What artificial grammar learning reveals about
40 the neurobiology of syntax. *Brain Lang.* **120**, 83–95 (2012).
41
42 34. Bischoff-Grethe, A., Proper, S. M., Mao, H., Daniels, K. A. & Berns, G. S. Conscious and
43 unconscious processing of nonverbal predictability in Wernicke's area. *J. Neurosci.* **20**,
44 1975–81 (2000).
45
46 35. Skosnik, P. D. *et al.* Neural correlates of artificial grammar learning. *Neuroimage* **17**,
47 1306–14 (2002).
48
49 36. Furl, N. *et al.* Neural prediction of higher-order auditory sequence statistics. *Neuroimage*
50 **54**, 2267–2277 (2011).
51
52 37. Nastase, S., Iacovella, V. & Hasson, U. Uncertainty in visual and auditory series is coded
53 by modality-general and modality-specific neural systems. *Hum. Brain Mapp.* **35**, 1111–
54
55
56
57
58
59
60

- 1
2
3 28 (2014).
4
5 38. Corbetta, M., Patel, G. & Shulman, G. L. The Reorienting System of the Human Brain:
6 From Environment to Theory of Mind. *Neuron* **58**, 306–324 (2008).
7
8 39. Vilberg, K. L. & Rugg, M. D. Memory retrieval and the parietal cortex: A review of
9 evidence from a dual-process perspective. *Neuropsychologia* **46**, 1787–1799 (2008).
10
11 40. Rugg, M. D., Johnson, J. D., Park, H. & Uncapher, M. R. Chapter 21 Encoding-retrieval
12 overlap in human episodic memory: A functional neuroimaging perspective. in *Progress*
13 *in brain research* **169**, 339–352 (2008).
14
15 41. Ordin, M., Polyanskaya, L., Laka, I. & Nespors, M. Cross-linguistic differences in the use of
16 durational cues for the segmentation of a novel language. *Mem. Cognit.* **45**, 863–876
17 (2017).
18
19 42. Ordin, M., Polyanskaya, L., & Soto, D. Metacognitive processing in language learning
20 tasks is affected by bilingualism. *J. Exp. Psychol. Learn. Mem. Cogn.*
21
22 43. Gonzalez-Castillo, J. *et al.* Imaging the spontaneous flow of thought: Distinct periods of
23 cognition contribute to dynamic functional connectivity during rest. *Neuroimage* **202**,
24 116129 (2019).
25
26 44. Smith, S. M. Fast robust automated brain extraction. *Hum. Brain Mapp.* **17**, 143–155
27 (2002).
28
29 45. Greve, D. N. & Fischl, B. Accurate and robust brain image alignment using boundary-
30 based registration. *Neuroimage* **48**, 63–72 (2009).
31
32 46. Rosenthal, C. R., Andrews, S. K., Antoniadou, C. A., Kennard, C. & Soto, D. Learning
33 and Recognition of a Non-conscious Sequence of Events in Human Primary Visual
34 Cortex. *Curr. Biol.* **26**, 834–41 (2016).
35
36 47. Woolrich, M. W., Ripley, B. D., Brady, M. & Smith, S. M. Temporal Autocorrelation in
37 Univariate Linear Modeling of fMRI Data. *Neuroimage* **14**, 1370–1386 (2001).
38
39 48. Jenkinson, M., Bannister, P., Brady, M. & Smith, S. Improved optimization for the robust
40 and accurate linear registration and motion correction of brain images. *Neuroimage* **17**,
41 825–41 (2002).
42
43 49. Woolrich, M. W., Behrens, T. E. J., Beckmann, C. F., Jenkinson, M. & Smith, S. M.
44 Multilevel linear modelling for fMRI group analysis using Bayesian inference.
45 *Neuroimage* **21**, 1732–1747 (2004).
46
47 50. Beckmann, C. F., Jenkinson, M. & Smith, S. M. General multilevel linear modeling for
48 group analysis in fMRI. *Neuroimage* **20**, 1052–1063 (2003).
49
50 51. Poldrack, R. A. *et al.* The Neural Correlates of Motor Skill Automaticity. *J. Neurosci.* **25**,
51 5356–5364 (2005).
52
53 52. Rosenthal, C. R., Roche-Kelly, E. E., Husain, M. & Kennard, C. Response-dependent
54
55
56
57
58
59
60

- 1
2
3 contributions of human primary motor cortex and angular gyrus to manual and perceptual
4 sequence learning. *J. Neurosci.* **29**, 15115–25 (2009).
5
- 6 53. de Bourbon-Teles, J. *et al.* Thalamic Control of Human Attention Driven by Memory and
7 Learning. *Curr. Biol.* **24**, 993–999 (2014).
8
- 9 54. Endress, A. D., Nespors, M. & Mehler, J. Perceptual and memory constraints on language
10 acquisition. *Trends Cogn. Sci.* **13**, 348–353 (2009).
11
- 12 55. Prince, S. E., Daselaar, S. M. & Cabeza, R. Neural Correlates of Relational Memory:
13 Successful Encoding and Retrieval of Semantic and Perceptual Associations. *J.*
14 *Neurosci.* **25**, 1203–1210 (2005).
15
- 16 56. Turk-Browne, N. B., Scholl, B. J., Chun, M. M. & Johnson, M. K. Neural Evidence of
17 Statistical Learning: Efficient Detection of Visual Regularities Without Awareness. *J.*
18 *Cogn. Neurosci.* **21**, 1934–1945 (2009).
19
- 20 57. Gennari, S. P., Millman, R. E., Hymers, M. & Mattys, S. L. Anterior paracingulate and
21 cingulate cortex mediates the effects of cognitive load on speech sound discrimination.
22 *Neuroimage* **178**, 735–743 (2018).
23
- 24 58. Ciaramelli, E., Grady, C. L. & Moscovitch, M. Top-down and bottom-up attention to
25 memory: A hypothesis (AtoM) on the role of the posterior parietal cortex in memory
26 retrieval. *Neuropsychologia* **46**, 1828–1851 (2008).
27
- 28 59. Gehring, W. J. & Fencsik, D. E. Functions of the Medial Frontal Cortex in the Processing
29 of Conflict and Errors. *J. Neurosci.* **21**, 9430–9437 (2001).
30
- 31 60. Scargle, J. D., Norris, J. P., Jackson, B. & Chiang, J. STUDIES IN ASTRONOMICAL
32 TIME SERIES ANALYSIS. VI. BAYESIAN BLOCK REPRESENTATIONS. *Astrophys. J.*
33 **764**, 167 (2013).
34
- 35 61. Alloy, L. B. & Tabachnik, N. Assessment of covariation by humans and animals: the joint
36 influence of prior expectations and current situational information. *Psychol. Rev.* **91**, 112–
37 49 (1984).
38
- 39 62. Jessen, F. *et al.* Activation of human language processing brain regions after the
40 presentation of random letter strings demonstrated with event-related functional magnetic
41 resonance imaging. *Neurosci. Lett.* **270**, 13–6 (1999).
42
- 43 63. Mechelli, A. *et al.* Structural plasticity in the bilingual brain. *Nature* **431**, 757–757 (2004).
44
- 45 64. Catani, M., Jones, D. K. & ffytche, D. H. Perisylvian language networks of the human
46 brain. *Ann. Neurol.* **57**, 8–16 (2005).
47
- 48 65. Frey, S., Campbell, J. S. W., Pike, G. B. & Petrides, M. Dissociating the Human
49 Language Pathways with High Angular Resolution Diffusion Fiber Tractography. *J.*
50 *Neurosci.* **28**, 11435–11444 (2008).
51
- 52 66. Kelly, C. *et al.* Broca's region: linking human brain functional connectivity data and non-
53 human primate tracing anatomy studies. *Eur. J. Neurosci.* **32**, 383–398 (2010).
54
55
56
57
58
59
60

- 1
2
3 67. Seghier, M. L. The Angular Gyrus. *Neurosci.* **19**, 43–61 (2013).
4
5 68. Paraskevopoulos, E., Chalas, N. & Bamidis, P. Functional connectivity of the cortical
6 network supporting statistical learning in musicians and non-musicians: an MEG study.
7 *Sci. Rep.* **7**, 16268 (2017).
8
9
10
11
12
13
14
15

16 FIGURE CAPTIONS:

17 **Figure 1.** Example of a learning phase run (above), where 25.92-second pseudorandom blocks
18 are interspersed with 181.44-second structured blocks; and the structure of the recognition trial
19 (below).
20

21 **Figure 2.** Percentage of endorsed trials for words, non-words and phantoms. Error bars show
22 95% CI.
23

24 **Figure 3.** Mean confidence rating (weighted average) assigned to the trails with accepted and
25 rejected words, phantoms and non-words. Error bars show 95%CI.
26

27 **Figure 4.** BOLD responses during the learning phase ($Z > 2.3$, $P = 0.05$ corrected). Brain
28 regions showing BOLD response increases in structured vs pseudorandom chunks in the
29 course of training.
30

31 **Figure 5.** BOLD responses during the recognition phase ($Z > 2.3$, $P = 0.05$ corrected). **(A)** Brain
32 regions showing increased response for words accepted relative to nonword rejected ($w_acc >$
33 $nonw_rej$). Correctly identified words elicit larger activation in the LIFG compared to correctly
34 identified non-words. **(B)** Brain regions showing increased activity on trials with words accepted
35 compared to phantoms rejected ($w_acc > ph_rej$). Endorsed words, relative to rejected
36 phantoms, elicit larger activation in the angular gyrus and intra-parietal sulcus, the anterior
37 division of the cingulate gyrus, and the posterior division of the right STG.
38
39

40 **Figure 6.** Illustration of the brain activity maps that showed increased activity during learning
41 compared to the recognition phase ($Z > 2.3$, $P = 0.05$ corrected).
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1. The list of words, phantoms and non-words used in the experiment

| | Words | Phantoms | Non-words |
|----|--------------|-----------------|------------------|
| 1 | | | |
| 2 | | | |
| 3 | | | |
| 4 | ROSENU | PASENU | ROTIMO |
| 5 | ROKAFA | LEKAFA | SEPAKO |
| 6 | PASETI | ROSETI | FALUSA |
| 7 | LEKATI | ROKATI | FOLERI |
| 8 | PAMONU | PAMOFA | TAMUPE |
| 9 | LEMOFA | LEMONU | NIKANU |
| 10 | PERIKO | MURIKO | NURIFE |
| 11 | MURIFO | LUTASA | FOLUKA |
| 12 | PETASA | PERIFO | NIMUKO |
| 13 | LUTAFO | PETAFO | MOPARO |
| 14 | MUNIKO | MUNISA | LESATI |
| 15 | LUNISA | LUNIKO | TASEFA |
| 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 | | | |

Unedited manuscript

Table 2. Location of peaks related to the increase in the activation difference between structured and random chunks

| Cluster | Extent (voxels) | Anatomical region | Z max | x | y | z |
|---------|-----------------|---|-------|-----|-----|----|
| 1 | 1008 | Anterior division of cingulate gyrus, paracingulate gyrus, superior frontal gyrus | 3.46 | 14 | 38 | 24 |
| 2 | 1002 | Superior temporal gyrus, left (slightly extending to middle temporal gyrus) | 3.97 | -60 | -14 | -2 |
| 3 | 777 | Superior temporal gyrus, right | 3.57 | 64 | -8 | -2 |

Unedited manuscript

Table 3. Location of peak activation differences (with NMI coordinates of the activation peaks) for the relevant contrasts in the recognition phase

| Contrast | Cluster | Extent (voxels) | Anatomical region | Z max | X | y | z |
|------------------|---------|-----------------|---|-------|-----|-----|----|
| w_acc > nonw_rej | 1 | 6207 | LIFG (par opercularis extending to par triangularis), i.e., BA44 extending to BA45. | 3.76 | -38 | 18 | 22 |
| w_acc > ph_rej | 1 | 297 | Anterior division of the cingulate gyrus, paracingulate gyrus | 3.49 | -2 | 26 | 36 |
| | 2 | 292 | Posterior divisions of the right superior temporal and middle temporal gyri | 3.7 | 64 | -26 | -2 |
| | 3 | 227 | Anterior division of the left intra-parietal sulcus, angular gyrus | 3.73 | -34 | -50 | 34 |

Table 4. Location of peak activation differences (with NMI coordinates of the activation peaks) for the *learning > recognition* contrast

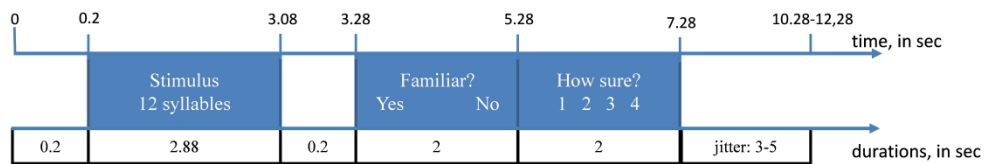
| Cluster | Extent (voxels) | Anatomical region | Z max | x | y | z |
|---------|-----------------|---|-------|-----|-----|----|
| 1 | 688 | The superior temporal gyrus, slightly extending to the middle temporal gyrus, left. | 3.88 | -60 | -14 | -2 |
| 2 | 699 | Anterior division of the cingulate gyrus, the paracingulate gyrus, right, the superior frontal gyrus. | 3.51 | 14 | 40 | 24 |

Unedited manuscript

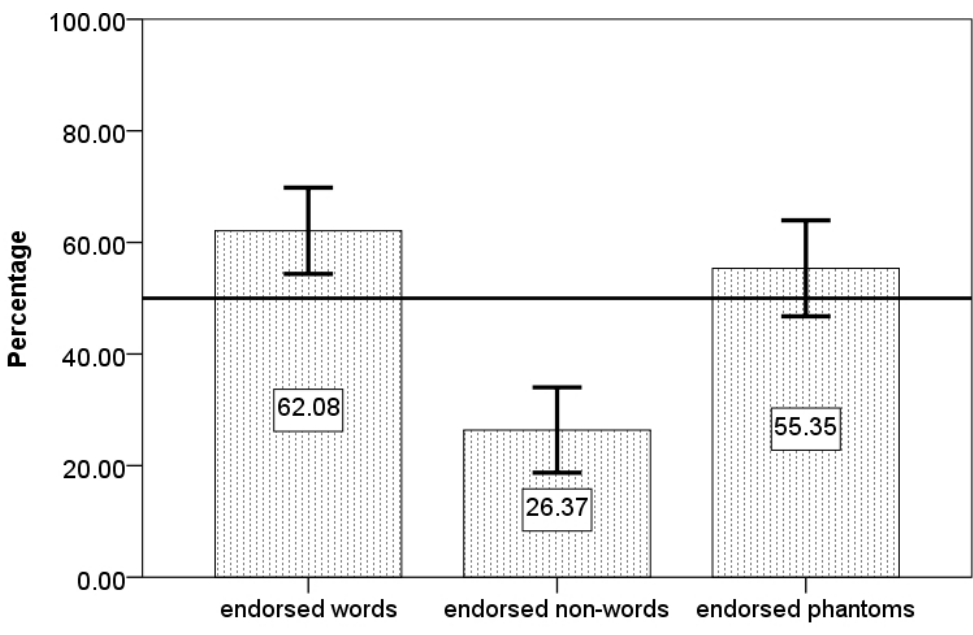
Example of a learning run



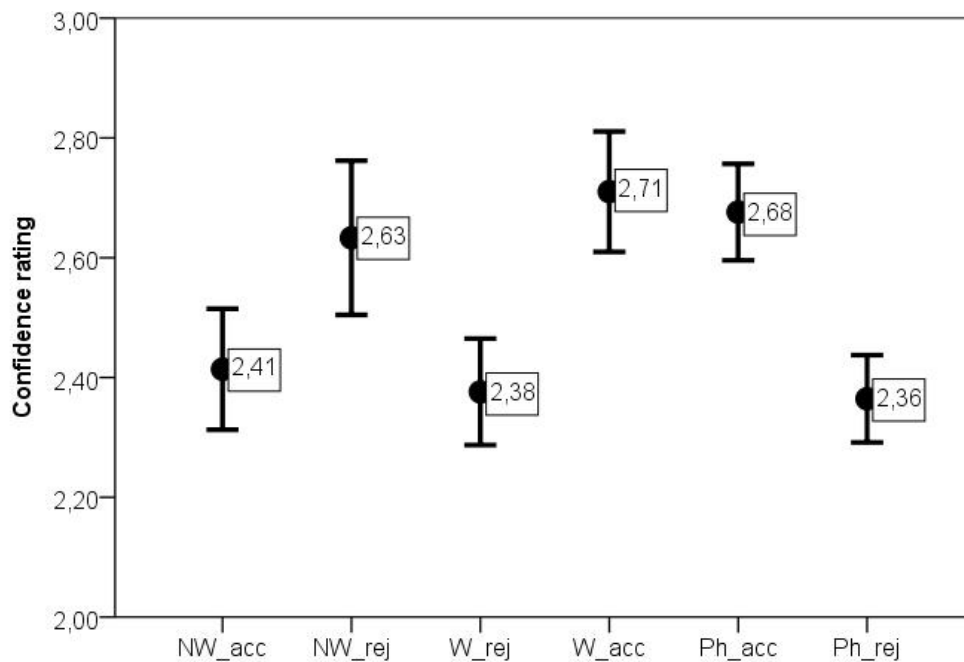
Structure of the recognition trial



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

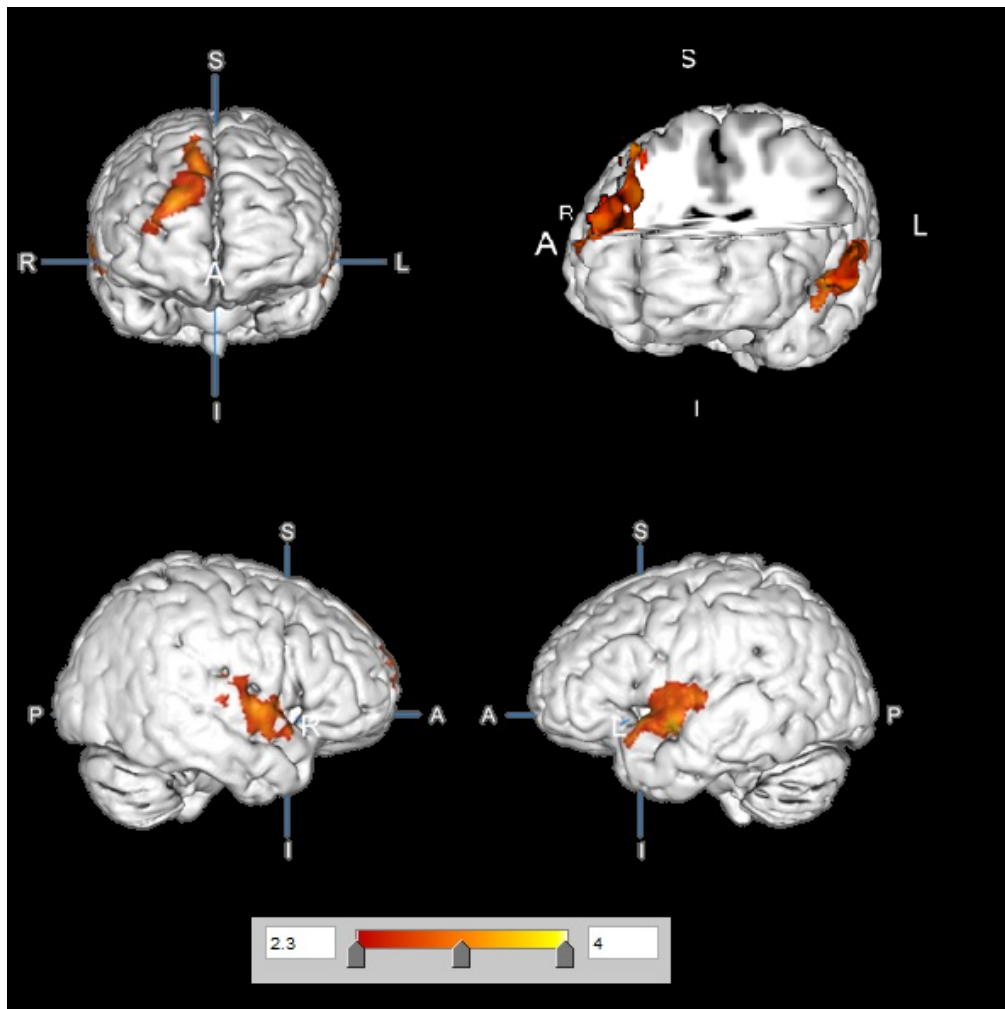


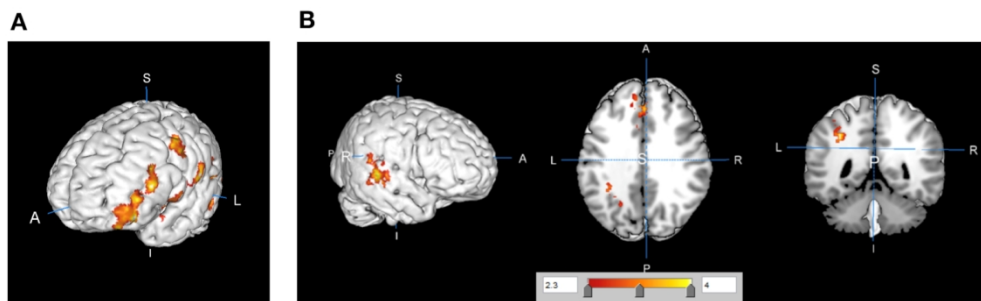
271x169mm (72 x 72 DPI)



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

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





184x58mm (220 x 220 DPI)

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

