

1 **Title:** Testing the unitary theory of language lateralisation using functional transcranial Doppler  
2 sonography in adults

3

4

5 **Authors:** Woodhead ZVJ\*<sup>1</sup>, Bradshaw AR<sup>1</sup>, Wilson AC<sup>1</sup>, Thompson PA<sup>1</sup>, Bishop DVM<sup>1</sup>

6 \* Corresponding author: [zoe.woodhead@psy.ox.ac.uk](mailto:zoe.woodhead@psy.ox.ac.uk)

7 <sup>1</sup> Department of Experimental Psychology, University of Oxford, Oxford, United Kingdom

8

9

10 **Major Subject Area:** Neuroscience

11 **Abstract**

12 Cerebral lateralisation for language can vary from task to task, but it is unclear if this reflects error  
13 of measurement or independent lateralisation of different language systems. We used functional  
14 transcranial Doppler sonography to assess language lateralisation in 37 adults (7 left-handers) on  
15 six tasks, each given on two occasions. Tasks taxed different aspects of language function. A  
16 preregistered structural equation analysis was used to compare models of means and covariances.  
17 For most people, a single lateralised factor explained most of the covariance between tasks. A  
18 minority, however, showed dissociation of asymmetry, giving a second factor. This was mostly  
19 derived from a receptive task, which was highly reliable but not lateralised. The results suggest  
20 that variation in strength of language lateralisation reflects true individual differences and not just  
21 error of measurement. Inclusion of several tasks in a laterality battery makes it easier to detect  
22 cases of atypical asymmetry.

## 23 **Introduction**

24 Hemispheric dominance for language is often assumed to be unidimensional and consistent across  
25 language domains, but this assumption can be questioned (Bishop, 2013; Bradshaw, Thompson,  
26 Wilson, Bishop, & Woodhead, 2017). Discrepant laterality across different language tasks (e.g.  
27 Gaillard et al., 2004; Stroobant, Buijs, & Vingerhoets, 2009; Tailby, Abbott, & Jackson, 2017) could  
28 be simply due to measurement error (Ramsey, Sommer, Rutten, & Kahn, 2001); alternatively, task  
29 differences may represent meaningful individual variation in the hemispheric organization of  
30 different language networks. It has been difficult to distinguish these possibilities, because, while  
31 we have ample evidence that the left hemisphere is heavily implicated in language function at the  
32 group level, relatively little is known about the reliability of lateralization in individuals. It is  
33 evident that a standard model based on average brain activation may give a misleading impression  
34 of uniformity (Seghier & Price, 2018). Furthermore, there is evidence that there may be subgroups  
35 of people with distinct laterality profiles, related to handedness (Mazoyer et al., 2014). Such  
36 variability in cerebral lateralisation may have functional significance, for example in terms of  
37 impaired language abilities (Bishop, 2013). In clinical neurosurgical contexts, it is important to  
38 know whether a single indicator of an individual's language laterality is sufficient, or whether a  
39 battery of measures is needed to capture laterality in multiple language domains (Gaillard et al.,  
40 2004; Stroobant et al., 2009; Tailby et al., 2017). Before we can make headway in answering such  
41 questions, we need to have reliable measures.

42 Here we report a study using functional transcranial Doppler sonography (fTCD; Knecht et al.,  
43 1998) to measure speed of blood flow in left and right middle cerebral arteries (a proxy for neural  
44 activity in language-related areas of the brain) during six different language tasks (tasks A-F). The  
45 fTCD data were used to derive laterality indices (LIs), which quantify the balance of activation in  
46 left and right hemispheres. All participants were tested on the whole battery in two separate  
47 sessions on different days in order to estimate the reliability of the LIs and the extent to which  
48 lateralization of different tasks could be explained in terms of a common factor.

### 49 **Laterality at the level of the population and the individual**

50 The question of whether language lateralisation is a unitary function has two distinct  
51 interpretations: (a) whether there are differences in extent of lateralisation across different  
52 language functions or (b) whether there are individual differences in how the strength of  
53 lateralisation varies across language functions. We first review existing literature on these

54 questions and then present simulated data to show how predictions made by the two accounts  
55 are independent and additive, but can be tested within a common framework (structural equation  
56 modelling, SEM).

### 57 **Task-related variation in extent of language lateralisation**

58 Most theories of language lateralisation have focused on how language functions are lateralised in  
59 the brain in typical humans. Such theories are not concerned with individual differences, but make  
60 theoretical statements about the properties of language that are associated with lateralised  
61 activity. An influential example of such a theory is Hickok and Poeppel's dual route model of  
62 speech processing (Hickok & Poeppel, 2007). This contrasts a dorsal stream from superior  
63 temporal to premotor cortices via the arcuate fasciculus, which is associated with sensorimotor  
64 integration of auditory speech sounds and articulatory motor actions; and a ventral stream from  
65 temporal cortex to anterior inferior frontal gyrus, which is involved in access to conceptual  
66 memory and mapping of sound to meaning (Rauschecker, 2018). Hickok and Poeppel proposed  
67 that the dorsal stream is left lateralized, whereas the ventral stream is bilateral. This kind of  
68 theory makes predictions about task-related differences that can be assessed by comparing mean  
69 LIs in a sample. Thus, the prediction from the dual route model is that mean LIs for tasks involving  
70 the dorsal stream will show left-lateralisation, whereas LIs from tasks primarily involving the  
71 ventral stream will not be lateralised.

72 Hickok and Poeppel's model contrasts with other theoretical accounts. For instance, Dhanjal et al  
73 proposed that left lateralization was a characteristic of tasks involving lexical retrieval (Dhanjal,  
74 Handunnetthi, Patel, & Wise, 2008). Evidence came from an fMRI study investigating propositional  
75 speech (e.g. sentence generation) and non-propositional speech (e.g. reciting memorized speech):  
76 articulatory jaw and tongue movements and non-propositional speech co-activated bilateral  
77 dorsal areas, including the superior temporal planes, motor and premotor cortices. Only the lexical  
78 retrieval component of propositional speech resulted in left lateralized activity (in the inferior  
79 frontal gyrus and premotor cortex).

80 Yet other accounts have focused on the complexity of the speech stimulus (Peelle, 2012), or  
81 argued that lateralization is specifically linked to aspects of complex syntactic processing (Bozic,  
82 Tyler, Ives, Randall, & Marslen-Wilson, 2010; Friederici, 2011).

### 83 **Individual differences in cerebral lateralisation**

84 Discussions about the nature of language lateralization are complicated by individual differences;  
85 although most people show the typical pattern of language laterality, some individuals show the  
86 reverse pattern – right-hemisphere language. In a large-scale comparison of left- and right-  
87 handers, Mazoyer et al (2014) reported that strong right-hemisphere bias for a sentence  
88 generation task was seen exclusively in left-handers, though milder departures from left  
89 hemisphere dominance were seen in right- as well as left-handers. A subset of people with  
90 bilateral language has also been described for many years (Milner, Branch, & Rasmussen, 1966),  
91 but this category is ambiguous. These could be people who engage both hemispheres equally  
92 during language tasks, or people who are strongly lateralized for different tasks, but in different  
93 directions. This latter scenario would provide strong evidence against a unitary hypothesis, by  
94 demonstrating that a person’s language laterality could not be predicted by a single dimension.

95 Individual differences in cerebral lateralisation have previously been observed in the comparison  
96 between left lateralised verbal functions versus right lateralised nonverbal functions. This might  
97 suggest complementarity of the two functions within the brain; however, where individual  
98 differences in these biases have been assessed, several studies have found them to be dissociated  
99 (Badzakova-Trajkov, Corballis, & Häberling, 2016; Groen, Whitehouse, Badcock, & Bishop, 2012;  
100 Rosch, Bishop, & Badcock, 2012; Whitehouse & Bishop, 2009; Zago et al., 2015; cf: Cai, Van der  
101 Haegen, & Brysbaert, 2013; Vingerhoets et al., 2013). Again, handedness has been noted as an  
102 important factor, with right-handers showing less evidence of complementarity of verbal and  
103 visuospatial functions than left-handers (Zago et al., 2015). Here, we consider whether similar  
104 dissociations might be found *within* the domain of language. Although previous investigators have  
105 considered association or dissociation in average patterns of activation for different tasks (Hesling,  
106 Labache, Jobard, & Leroux, 2018; Pinel & Dehaene, 2010), there has been little previous research  
107 documenting individual differences in task-related variation. Inconsistent LIs from task to task  
108 could simply reflect noisy measurement, making dissociations hard to interpret. Thus, in order to  
109 throw light on individual differences in language laterality, we need to include repeated measures,  
110 so that reliability of LIs from different tasks can be assessed.

### 111 **Simulated data to illustrate predictions**

112 It is possible to integrate models of task variation in lateralisation with a model of individual  
113 differences in the kind of framework shown in Figure 1. For simplicity, this shows simulated data  
114 on just two tasks, A and B, to contrast predictions from different models of the structure of  
115 language lateralisation. The Population Bias model is the simplest: it shows a population bias to

116 left-sided language laterality (i.e. positive LI values) that does not depend on the task. There are  
117 no consistent individual differences: any variation in laterality is just caused by random error. This  
118 is not a very plausible model, but provides a useful starting point from which to build more  
119 complex scenarios. Formally, the function for predicting an individual's LI is as follows:

$$120 \quad L_{ij} = a + e_{ij}$$

121 where  $i$  indexes the task, and  $j$  the individual,  $a$  is an intercept term corresponding to population  
122 bias, and  $e$  is random error.

123 In the Population Bias model, the mean LIs for different language tasks (shown by the horizontal  
124 and vertical red dotted lines) are all the same and equal to  $a$  (in this case set to 1). Note that  
125 because there are no stable individual differences, the correlations between LIs for the same task  
126 measured on different occasions (left hand panel), and between different tasks measured on the  
127 same occasion (right hand panel) are zero.

128 The second model is the Task Effect model. This incorporates consistent task-specific variation,  
129 without any stable individual differences. Formally,

$$130 \quad L_{ij} = a + t_i + e_{ij}$$

131 where  $t_i$  is a task-specific term. The only difference from the Population Bias model is that the  
132 means differ for different tasks – i.e. tasks A and B have mean LIs of 1 and 2 respectively. Again,  
133 variation in individuals' LI scores is due to random error ( $e$ ), rather than any systematic individual  
134 differences, as evidenced by zero test-retest correlations.

135 The next model is a Person Effect model. This includes stable individual differences: a person's  
136 score on any test occasion depends on an intrinsic lateral bias, which is constant from task to task  
137 but varies from person to person, i.e.

$$138 \quad L_{ij} = a + t_i + p_j + e_{ij}$$

139 where  $p_j$  is the person-specific term. This model predicts significant correlations between the  
140 same task tested on different occasions, and different tasks tested on the same occasion. An  
141 important point is that these correlations depend solely on the relative contribution of individual  
142 difference ( $p$ ) vs random noise ( $e$ ) to the LI. It does not matter whether there are also task-related  
143 effects ( $t$ ) on the LI. Thus, in the example, we have one task that is lateralised (mean LI of 2) and

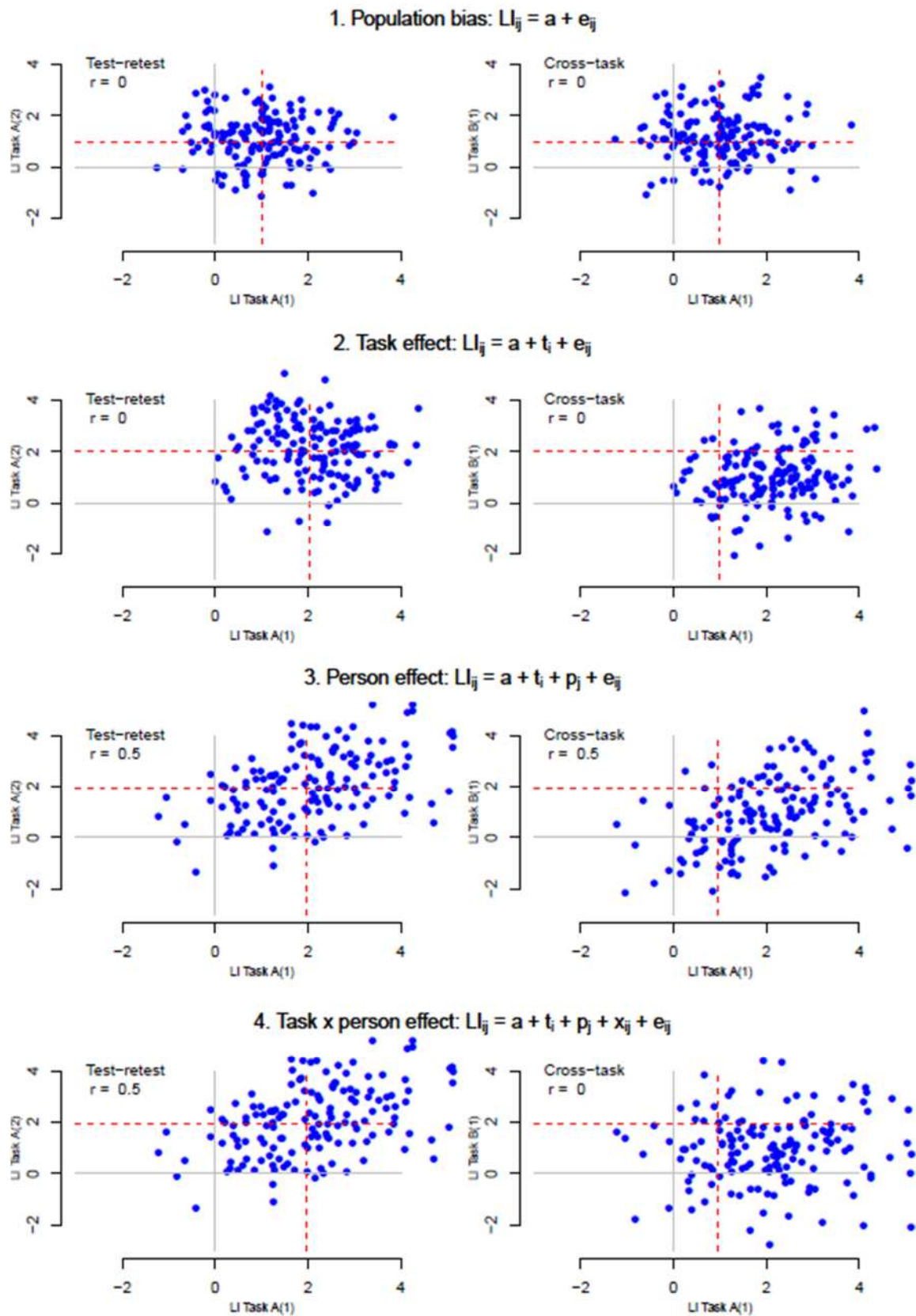
144 one that is not (mean LI of 0), yet on this model, the test-retest correlation for either task will be  
145 the same, and equivalent to the cross-task correlation.

146 The final model incorporates a Task by Person Effect: i.e., there are stable individual differences  
147 that show up as significant test-retest reliability on any one task, but the rank ordering of  
148 lateralisation varies from task to task, so cross-task correlations are low. Formally:

$$149 \quad L_{ij} = a + t_i + p_j + x_{ij} + e_{ij}$$

150 where  $x_{ij}$  reflects a contribution that is specific to the task and the individual. The depicted  
151 scenario in Figure 1 is an extreme one, with no relationship between a person's laterality on tasks  
152 A and B; in practice, there could be significant cross-task correlations, but if the within-task  
153 correlations are higher than cross-task correlations, then this would be evidence that individual  
154 differences in laterality are to some extent task-specific.

155 A key point illustrated by these simulations is that testing the multivariate model of language  
156 laterality at the population level requires different evidence – i.e. testing between means – than a  
157 multivariate model of individual differences, which requires us to consider correlations within and  
158 between tasks. Furthermore, predictions from these two types of model are independent,  
159 because correlations are not influenced by mean values. We can use structural equation  
160 modelling (SEM) to evaluate the relative fit of these four models to data on language lateralisation  
161 for participants who have LIs assessed on a range of tasks on two occasions.



162

163 **Figure 1.** Simulated data of different theoretical models of variance across sessions (1 and 2)  
164 and tasks (A and B) in language lateralization. Red dotted lines show the mean lateralization index (LI)  
165 for the task / session.

166



167 **Hypotheses**

168 We preregistered a set of hypotheses that were tested through SEM model comparison, as  
169 described in the Methods below.

170 We first tested two hypotheses concerning the group mean LI values. First, we tested the dorsal  
171 stream hypothesis (Hickok and Poeppel, 2007), which predicts that strength of lateralization  
172 depends on the extent to which tasks map on to the dorsal versus ventral speech processing  
173 streams (dorsal = stronger left lateralization). Second, following Dhanjal et al (2008), we tested  
174 the lexical retrieval hypothesis, which maintains that lateralization depends on the extent to which  
175 tasks require lexical retrieval (more lexical retrieval = stronger left lateralization).

176 A second set of hypotheses concerned individual differences in LI value. We predicted that a Task  
177 by Person Effect model, whereby covariances between tasks were modelled by two latent factors,  
178 would give a better fit to the data than a Person Effect model, where covariances were modelled  
179 by only one factor.

## 180 **Methods**

### 181 **Preregistration**

182 This project was preregistered on Open Science Framework prior to data collection  
183 (<https://osf.io/tkpm2/>). A number of changes were made to the analysis plan after collection of  
184 the data – an updated protocol is documented here: <https://osf.io/bjsv8/>. Departures to the  
185 original protocol are explained in the **Departures from pre-registered methods** section below.

### 186 **Design**

187 A test-retest, within-subject design was used. Lateralisation of brain activity was measured using  
188 Functional Transcranial Doppler Sonography (fTCD) during six language tasks: (A) List Generation,  
189 (B) Phonological Decision, (C) Semantic Decision, (D) Sentence Generation, (E) Sentence  
190 Comprehension, and (F) Syntactic Decision. Participants were tested on two sessions spaced by  
191 between 3 days and 6 weeks. Hence, each participant provided data from six tasks tested twice  
192 (A1-F1, A2-F2).

### 193 **Participants**

194 A sample size of  $n=30$  was determined by simulations of data from six tasks administered on two  
195 occasions, to determine the smallest sample size that would reliably distinguish data generated  
196 from a two factor vs single factor model, and give acceptable fit indices (see laterality\_simulations  
197 files, <https://osf.io/tkpm2/>). The simulations were based on the models of covariances, as the  
198 factor structure of the measures is our primary interest, and this gave a more conservative power  
199 estimate. We note that the sample size is small relative to those usually recruited for SEM  
200 analyses. However, because all measures were taken twice, with no practice effects expected (on  
201 the basis of previous studies with this method), there are several estimates of most parameters.  
202 For instance, the correlation between LIs for tasks A and B is estimated from A1B1, A1B2 and  
203 A2B2. Thus the repeated measures give low degrees of freedom relative to the number of  
204 measures.

205 In our original study pre-registration we did not plan to select participants according to  
206 handedness. However, both prior literature and our own preliminary data indicated that it would  
207 be advisable to treat right- and left-handers separately, as the pattern of associations between  
208 language tasks appeared to differ according to handedness, so combining handedness groups  
209 could give a misleading picture. We became concerned that results from our pre-registered

210 analysis on 30 participants (7 left-handers) were potentially misleading, as the factor structure  
211 that emerged seemed driven by a few left-handers. We therefore tested additional participants to  
212 give a total sample of 30 right-handers and seven left-handers, and we report analysis based on  
213 this larger sample as exploratory results.

214 All participants gave written informed consent. Procedures were approved by the University of  
215 Oxford's Medical Sciences Interdivisional Research Ethics Committee (approval number  
216 R40410/RE004). Subjects were recruited using the Oxford Psychology Research Participant  
217 Recruitment Scheme (<https://opr.sona-systems.com>) and by poster advertisements. The inclusion  
218 criteria were: aged 18-45 years; English native language speakers; and with normal or corrected to  
219 normal hearing and vision. Exclusion criteria were: a history of significant neurological disease or  
220 head injury; or a history of developmental language disorder.

221 It was not possible to record a Doppler signal via the temporal window in three participants. In  
222 these cases the participant was reimbursed but not tested further, and another participant was  
223 recruited in their place. One participant had excessive motion artifacts in his first session, so  
224 another participant was recruited in his place. The initial group of 30 participants (17 female and 7  
225 left-handed) had a mean age of 26.0 years (SD = 7.2 years; range: 19.2 to 45.1 years). The final  
226 group, including seven additional right-handers (2 females) had mean age 25.9 years (SD = 6.8  
227 years) with the same age range.

## 228 **Procedure**

229 The order of the six language tasks was counterbalanced between subject and session. At each  
230 session, fifteen trials of each task type were conducted with breaks in between tasks.

## 231 **Language tasks**

232 The six tasks were designed to be matched in trial structure, as far as feasible, so that differences  
233 in laterality should reflect as far as possible the linguistic task demands. The first five tasks had a  
234 visual stimulus on each trial presented against a grey background, to keep the visual demands as  
235 similar as possible; the sixth task involved presentation of written words. All stimulus materials  
236 are available on Open Science Framework (<https://osf.io/8s7vn/>).

237 The rest period prior to stimulus presentation was used for baseline correction to equate the left  
238 and right channels. Trials were 33 seconds long, and followed the structure shown in Figure 2.  
239 Trials started with the word 'CLEAR' on screen for 3 seconds, indicating that participants must

240 clear their mind in preparation for the next trial. The language task followed, lasting for 20  
 241 seconds. Procedures for each task type are detailed below, and examples of stimuli are shown in  
 242 Figure 3. Note that for tasks B, C, E and F, participants made responses to a series of stimuli on  
 243 each trial to ensure the participant was engaged in language processing throughout the activation  
 244 interval. Rapid presentation of multiple stimuli in a trial has been shown by Payne et al (Payne,  
 245 Gutierrez-Sigut, Subik, Woll, & MacSweeney, 2015) to maximise lateralised activation in fTCD.  
 246 After the task, 'REST' appeared on screen for 10 seconds, during which participants were required  
 247 to clear their minds.

248

	0	3	6	17	23	33 s
A. List Generation	Clear Mind	Stimulus	List Generation	Report	Rest	
B. Phonological Decision	Clear Mind	Phonological Decision x 6			Rest	
C. Semantic Decision	Clear Mind	Semantic Decision x 6			Rest	
D. Sentence Generation	Clear Mind	Stimulus	Sentence Generation	Report	Rest	
E. Sentence Comprehension	Clear Mind	Sentence Decision x 6			Rest	
F. Syntactic Decision	Clear Mind	Syntactic Decision x 3			Rest	

249

250 **Figure 2.** Timings within a single trial for all six task types.

251



252

253 **Figure 3.** Example stimuli for the language tasks. From left to right: picture stimulus for List  
 254 Generation task (A; recite months of the year); a matching picture pair ('book' / 'hook') for the  
 255 Phonological Decision task (B); a matching picture pair for the Semantic Decision task (C); picture  
 256 stimuli for the Sentence Generation task (D); and a picture pair for the Sentence Comprehension  
 257 task (E; 'The dog chases the girl who is jumping').

258

259 *A. List Generation*

260 This task was based on the reference task used by Mazoyer et al (2014). Participants were asked to  
261 recite an automatic sequence of words (non-propositional speech) in response to a picture. In  
262 each trial, a line drawing was displayed on a grey background for 3 seconds. Participants were  
263 trained to produce different sequences for different pictures: reciting the numbers from 1-10, the  
264 letters from A-J, the days of the week or the months of the year. A fixation cross was then  
265 presented in the center of the screen for 11 seconds, during which the participant recited the  
266 words covertly (silently) in their head. Following this, a 'REPORT' prompt was shown for 6 seconds,  
267 indicating that participants should say the sequence aloud. The list generation task involves  
268 generation of phonological output, and so should index the dorsal stream, but because it involves  
269 repeated, overlearned material, it does not implicate the ventral stream; nor does it place  
270 demands on lexical retrieval. Thus the two specific theories of interest make contrasting  
271 predictions about this task.

#### 272 *B. Phonological Decision*

273 Participants were required to make a rhyme judgement on pairs of words represented by pictures.  
274 The pictures were easily nameable line drawings of single syllable words, mostly taken from the  
275 International Picture Naming Project (IPNP) database  
276 (<https://crl.ucsd.edu/experiments/ipnp/index.html>, Szekely et al., 2004). The pictures were  
277 arranged into 45 rhyming and 45 non-rhyming pairs (based on pairings devised by Bishop &  
278 Robson, 1989). Rhyming and non-rhyming pairs did not differ significantly on orthographic  
279 similarity (assessed using MatchCalc software,  
280 <http://www.pc.rhul.ac.uk/staff/c.davis/Utilities/MatchCalc/>). For each trial, a series of 6 picture  
281 pairs was presented, each for 3.33 seconds (totaling 20 seconds). For each pair, the participant  
282 decided whether the words represented by the pictures rhymed or not, and responded by button  
283 press.

284 This task involves implicit generation of lexical items and their phonology, but does not require  
285 access to conceptual meaning. Both the dorsal-ventral stream theory and lexical retrieval theory  
286 predict it should be strongly lateralized.

#### 287 *C. Semantic Decision*

288 This task involved a semantic category judgement on objects represented in a pair of pictures.  
289 The design of this task closely matched that of the phonological decision task. The pictures were  
290 mostly taken from the IPNP database, as described above. The stimuli were matched for word

291 familiarity, orthographic neighbourhood, imageability, number of phonemes and frequency. Six  
292 picture pairs were presented, each for 3.33 seconds. For each pair, the participant decided  
293 whether the objects were from the same semantic category or not (e.g. both types of food) and  
294 responded by button press. For this task, it is necessary to access conceptual meaning, but  
295 generation of word names is not implicated. This, then, can be regarded as indexing the ventral  
296 stream. Both the dorsal-ventral stream theory and the lexical retrieval theory predict weak  
297 lateralization for this task.

#### 298 *D. Sentence Generation*

299 This task required participants to generate spoken sentences in response to line drawings,  
300 following methods described by Mazoyer and colleagues (Mazoyer et al., 2014), but using pictures  
301 that were more culturally appropriate for UK participants.

302 For each trial, a black line drawing was displayed on a grey background for 3 seconds. This was  
303 followed by a fixation cross for 11 seconds, during which the participant was required to covertly  
304 generate a sentence. Participants were trained in advance to generate sentences beginning with a  
305 subject (e.g. “the boy”), followed by a description of the subject (“with marbles”), a verb (“plays”) and  
306 ending with a detail about the action (“on the floor”). A “REPORT” prompt was then presented  
307 for six seconds, and participants were required to say their sentence aloud.

308 This task implicates both dorsal and ventral streams, and so might be expected to show weaker  
309 lateralization than purely dorsal tasks. In contrast, the lexical retrieval theory predicts strong  
310 lateralization.

#### 311 *E. Sentence Comprehension*

312 This task required participants to decide which of two pictures corresponded to a spoken  
313 sentence. Each trial comprised six picture pairs, each presented for 3.33 seconds, along with a  
314 spoken sentence that matched one of the two pictures. The sentences were spoken at a rapid  
315 pace and included some involving complex grammar with long-distance dependencies, such as  
316 ‘the shoe on the pencil is blue’, or ‘the cow that is brown is chasing the cat’. Participants indicated  
317 which of the two pictures matched the sentence by button press.

318 This task would appear to stress the ventral more than the dorsal stream, and so be relatively  
319 weakly lateralized. The task is hard to categorise in terms of lexical retrieval: it is necessary to hold

320 word meanings in memory while working out the meaning, though overt word generation is not  
321 required.

### 322 *F. Syntactic Decision*

323 This task was designed to isolate syntactic processing with minimal involvement of semantics. This  
324 task uses ‘Jabberwocky’ stimuli, based on a study by Fedorenko and colleagues (Fedorenko, Hsieh,  
325 Nieto-Castañón, Whitfield-Gabrieli, & Kanwisher, 2010), where content words of sentences are  
326 replaced by plausible non-words. Half of the stimuli were ‘sentences’, where function words, word  
327 order and morphological cues were preserved to make the stimuli recognisable as syntactically-  
328 valid sentences (e.g. ‘The tarben yipped a lev near the kruss’). The other half had a pseudorandom  
329 word order and were not perceived as sentences (e.g. ‘Kivs his porla her tal ghep in with’).

330 Each trial contained three Jabberwocky stimuli of 8 words. Words were presented sequentially at  
331 the same time as an audio recording of the spoken word. As all spoken words were recorded  
332 separately, there were no prosodic cues to whether the stimulus is a ‘sentence’ or not. Each word  
333 was presented for 0.7 seconds, and the sequence was followed by a question mark for 1 second  
334 (making a total of 6.7 seconds for each Jabberwocky stimulus). The participant was required to  
335 respond by button press following the ‘?’ prompt.

336 In terms of the dorsal-ventral stream account, this task is predicted not to show lateralization, as it  
337 is a purely receptive task. This was the only task involving nonwords, and should not be lateralized  
338 according to a lexical retrieval account.

### 339 **Behavioural Analysis**

340 For tasks A and D, the average number of words generated for each trial was calculated. For tasks  
341 B, C, E and F, percentage accuracy and average reaction time for correct trials (excluding trials  
342 where reaction time was greater than 2 standard deviations away from the mean) were  
343 calculated. The number of events where no response was received was also recorded for each task  
344 – these events were scored as incorrect.

### 345 **fTCD Analysis**

346 Our analysis of fTCD data departed from the method we preregistered in three respects; sections  
347 describing the altered methods are shown in italics, with a description and explanation of the  
348 change shown in the section ‘Departures from pre-registered methods’.

349 The dependent measures derived from the fTCD analysis were the Laterality Indices (LI) from tasks  
350 A to F at sessions 1 and 2. fTCD uses ultrasound probes positioned bilaterally over the temporal  
351 windows to measure cerebral blood flow velocity (CBFV) in the left and right middle cerebral  
352 arteries (MCA). The probes emit ultrasound pulses and detect reflected ultrasound signal. The  
353 frequency of the reflected ultrasound signal depends on the speed of the blood moving in the  
354 MCA, due to Doppler shift. Hence the difference in frequency of the emitted and reflected  
355 ultrasound signals can be used to determine the speed of blood flow. The data is recorded as CBFV  
356 (cm/s) in the left and right hemispheres.

357 The fTCD data were analysed using a custom script in R Studio (RStudio Team, 2015). The script  
358 can be found on OSF (<https://osf.io/tkpm2/>). The CBFV data was first down-sampled from 100 Hz  
359 to 25 Hz by taking every 4<sup>th</sup> datapoint. The data was segmented into epochs of 33 seconds,  
360 beginning 7 seconds before the presentation of the 'CLEAR' stimulus at the start of the trial (-7  
361 seconds peri-stimulus time). Spiking or dropout datapoints were identified as being outside of the  
362 0.0001 - 0.9999 quantiles of the CBFV data. If only a single artifact datapoint was identified within  
363 an epoch, it was replaced with the mean for that epoch. If more than one datapoint was  
364 identified, the epoch was rejected. The CBFV was then normalized (by dividing by the mean and  
365 multiplying by 100) such that the values for CBFV become independent to the angle of insonation  
366 and the diameter of the MCA. Heart cycle integration was used to normalize the data relative to  
367 rhythmic modulations in CBFV. *Each epoch was baseline corrected using the interval from -5 to 2*  
368 *seconds peri-stimulus time*. Finally, artifacts were identified as values below 60% and above 140%  
369 of the mean normalised CBFV – any epochs containing such artifacts were rejected.

370 If a participant in one session had fewer than 12 acceptable epochs for any task (i.e. more than 3  
371 of the 15 epochs were rejected), the data for that task were excluded. If a participant had more  
372 than one task excluded, all data for that participant were excluded.

373 The CBFV from left and right sensors was averaged over all epochs at each timepoint, and the  
374 mean difference (left minus right) within the period of interest was taken as the laterality index  
375 (LI). The period of interest for tasks B, C, E and F was from 6 to 23 seconds peri-stimulus time. For  
376 tasks A and D, the period of interest ended at 17 seconds to avoid activity related to overt speech  
377 production following the 'REPORT' prompt.

378 The LI value at each trial was also recorded, and used to calculate a standard error, which  
379 indicated how variable the lateralization was over trials. Outlier standard error values were



380 identified using Hoaglin and Iglewicz's procedure (Hoaglin & Iglewicz, 1987). The standard error  
381 values for every LI measurement (across all subjects, tasks and sessions; 360 values in total) were  
382 concatenated. The difference between the first and third quartiles of the data was calculated (Q3-  
383 Q1). In this dataset, outliers were defined as having standard error value more than 2.2 times this  
384 difference above the third quartile (Q3); e.g., the threshold limit =  $Q3 + 2.2*(Q3-Q1)$ . Hence, if the  
385 LI value showed exceptionally high variability across trials, it was deemed to be unreliable and  
386 therefore omitted from the final analysis.

## 387 **Departures from pre-registered methods**

388 **1. Baseline interval.** The baseline interval was 2 seconds longer than that planned in the  
389 preregistered protocol (-5 to 0 seconds), i.e. extending into the 'Clear mind' period. As shown in  
390 Supplementary Materials (<https://osf.io/g8mkv/>), this baseline gives more stable estimates of LI  
391 than the original interval.

392 **2. Definition of laterality index.** In our pre-registered protocol, we planned to use a peak-based  
393 method of measuring the Laterality Index (LI) developed by Deppe et al (Deppe, Knecht,  
394 Henningsen, & Ringelstein, 1997), which has been standard in fTCD studies of cerebral  
395 lateralization. This involves finding the absolute peak in the difference wave within the period of  
396 interest and averaging the value of the difference over a 2 second time window centered on this  
397 peak. The major limitation of this approach is that it creates a non-normal distribution of LI values,  
398 which contributed to poor model fit in our SEM analyses, which assume normality. The mean-  
399 based method that we report here gives LI values that are highly correlated with the traditional  
400 peak-based LI (Spearman  $r = 0.97$ ), but with a normal distribution (see Supplementary Materials,  
401 <https://osf.io/g8mkv/>, for further details).

402 **3. Outlier detection.** In our pre-registered document, there was an error in our description of this  
403 process; we mistakenly stated we would remove outliers based on LI scores, rather than the  
404 standard error of the LI scores. Removing LI outliers would not be sensible in the context of this  
405 study, where the focus is on individual differences: it would, for instance, lead us to exclude those  
406 with atypical right-sided language laterality, who are of particular interest for our hypothesis. Our  
407 goal in outlier removal was to exclude participants with noisy data, and the LI standard error is the  
408 appropriate measure to use to achieve this goal.

409 **4. SEM modelling.** In addition to testing the models specified in the pre-registration document,  
410 we also tested model fit of the best-fitting model using a leave-one-out procedure, which allowed

411 us to check whether the parameter estimates were unduly influenced by specific data-points. As  
412 described in Supplementary Materials (<https://osf.io/g8mkv/>), our decision to test further right-  
413 handers was prompted by discovering that there was undue influence from one left-hander, with  
414 the factor solution changing when her data were omitted. Accordingly, we present here  
415 additional analyses with 30 right-handers only, and with the full sample of 37 participants. We also  
416 computed the factor scores from the final model and plotted these to aid interpretation of the  
417 factor structure. The SEM bifactor model requires one variable to have fixed paths of 1 and 0  
418 respectively to the two factors. The fit of the model does not depend on which measure is used for  
419 this purpose, but the specific path estimates will vary. Given that List Generation task was the only  
420 task with poor test-retest reliability, we present here results using Sentence Generation for the  
421 fixed paths. This follows recommendations that the strongest indicator for a specific factor should  
422 be used for the fixed paths (Lewis, 2017).

### 423 **Structural Equation Modelling**

424 Structural Equation Modelling (SEM), as implemented in OpenMx (<https://openmx.ssri.psu.edu/>),  
425 was used to test our hypotheses. We distinguish between two sets of hypotheses: models of task  
426 effects, which concerned predictions about means, and models of person effects, which  
427 concerned covariances. As noted above, these are independent from one another. The models  
428 used to test each hypothesis are described below, and can be seen in Figure 4.

429 We will briefly describe this approach, as it not widely used in laterality research. The aim is to test  
430 how well a prespecified model fits an observed dataset. Typically SEM is used to model  
431 covariances, but it can also be used with means. Boxes denote observed variables, two-headed  
432 arrows show variances and covariances. A triangular symbol denotes a mean value, typically set to  
433 one, with the path from the box to the triangle corresponding to the mean value for that variable.  
434 Means can be set to be equivalent by giving their paths the same label. We use capital letters for  
435 paths to means. For instance, in the Population Bias model (Figure 4), all paths to the mean are set  
436 to be the same, whereas in the Task Effect model (Figure 4), the means differ from task to task,  
437 but within a task are the same from test session 1 to test session 2.

438 An oval symbol corresponds to a latent variable linking two observed variables: covariance  
439 between two observed variables is computed as the sum of the product of the paths to those  
440 variables that are linked by an oval. Paths to latent variables are shown as lower case letters. The  
441 difference between modeling of means and covariances can be appreciated by comparing the Task

442 Effect model and the Person Effect model in Figure 4. These look similar, but the former depicts  
443 the situation where the means for a task are constant across sessions, but covariances are not  
444 considered. Thus even if means are stable, tasks may be unreliable in the sense that individual  
445 differences are just due to noise, and the rank order of LIs of individuals is unstable. In contrast,  
446 the Person Effect model takes into account covariances, and is a test of the reliability of the  
447 measures, assessing how far individuals are consistent in their LI across occasions.

448 We report goodness of fit for each model relative to a 'saturated' model where all variables are  
449 unconstrained, using the Comparative Fit Index (CFI): a high CFI indicates good model fit, and it is  
450 generally recommended that CFI needs to exceed .95 for the model to be regarded as a good fit to  
451 the data. We also report the Root Mean Square Error of Approximation (RMSEA), which is a  
452 measure of badness of fit, and should ideally be below .08 (Kline, 2011).

453 Comparison of model fit to determine the most appropriate model is achieved using likelihood  
454 ratio testing. Such comparisons are valid when we have nested models. For each hypothesis, we  
455 compare two nested models computing the difference in  $-2 \log$  likelihoods, and evaluated in terms  
456 of the difference in degrees of freedom between the two models. The difference in log likelihoods  
457 follow a  $\chi^2$  distribution, so a  $\chi^2$  test can be used to evaluate whether there is a statistical  
458 difference between the models. If a significant difference is found, then one model will be a better  
459 fit to the data.

460 In general, when comparing a model against another more complex model, good model fit  
461 corresponds to a non-significant p-value, which indicates that the more parsimonious model fits as  
462 well as the more complex model, despite fewer degrees of freedom. Models that estimate many  
463 parameters (and so have fewer degrees of freedom) will tend to fit the data better, and so relative  
464 fit of models is considered using indices that take this into account. Several indices that penalize  
465 the likelihood ratio test are available, for example, Akaike's Information Criterion (AIC) or Bayesian  
466 Information Criterion (BIC). Both these indices provide a value for each nested model and the  
467 lowest value among all the models is the preferred model.

#### 468 *Step 1: Testing Stability of LI Values*

469 We began with a Fully Saturated model that modeled means and variances as totally independent,  
470 as shown in Figure 4 (top left). No correlations between LI values were modelled at this stage: the  
471 triangular symbol denotes that the paths reflect the mean for each observed variable. As an initial  
472 sanity check, we computed a Task Effect model where the LI value means and variances for each

473 task (A-F) were fixed to be the same at each testing session (i.e. the means and variances for A1 =  
474 A2, B1 = B2, etc.). We predicted that the latter model would not deteriorate compared to the Fully  
475 Saturated model, indicating that we would not need to specify separate means for different test  
476 occasions.

#### 477 *Step 2: Testing Models of Means*

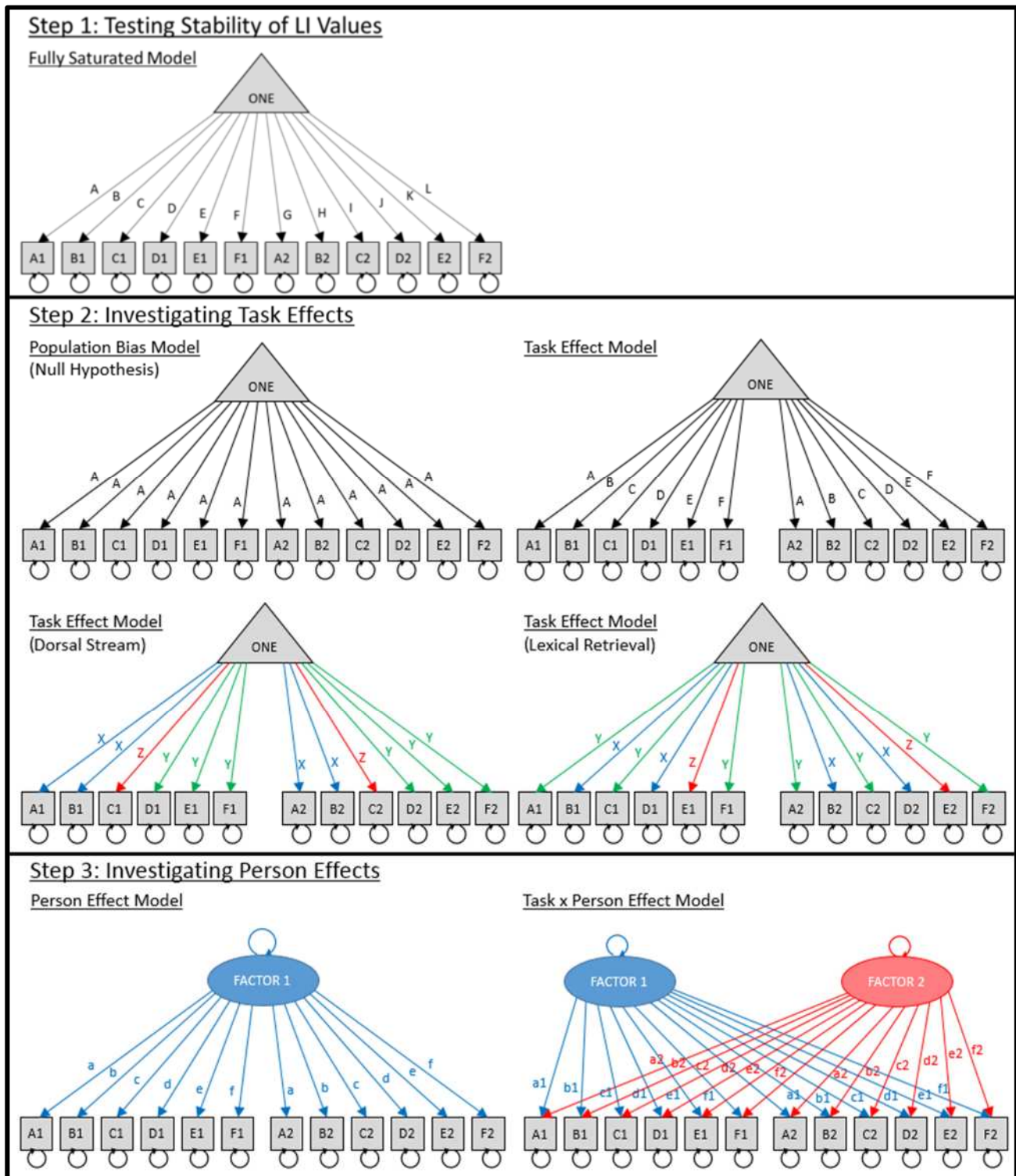
478 Our first hypothesis proposed that a significant task effect on LI value would be observed; i.e., that  
479 the mean LI values would vary between the six different tasks (tasks A-F). This was assessed by  
480 comparing the two models shown in row 2 of Figure 4: the Population Bias model and the Task  
481 Effect model.

482 The Task Effect model was then used as a baseline comparison model to test two more specific  
483 sub-hypotheses regarding which tasks would show the strongest lateralisation. In each case we  
484 divided tasks into three subsets, and fixed the means and variances for the tasks within each  
485 subset to be the same. We adopted this approach to test the Dorsal Stream hypothesis and the  
486 Lexical Retrieval hypothesis.

#### 487 *Step 3: Testing Models of Covariances*

488 Two models of covariance were compared (Figure 4, bottom). First, a person effect model was  
489 computed where covariance was predicted by a single factor, i.e. was similar across all language  
490 tasks. This was compared with a person by task effect model, with two covariance factors. The  
491 Person Effect (single factor) model is nested within the Task x Person Effect (bifactor) model, and  
492 so their relative fit can be assessed by subtraction of negative log likelihoods.

493 All analyses were conducted in R (R Core Team. & R Development Core Team, 2013). Data and  
494 analysis script are available on Open Science Framework.



495

496 **Figure 4:**

497 Step 1 (top): Simple model of means and variances. In the ‘Fully Saturated’ model the means for all  
 498 tasks could vary independently (tasks A-F, tested at sessions 1 and 2). This was compared to the  
 499 ‘Task Effect’ model, where the means for each task were fixed to be the same for each session.  
 500 The triangle symbol denotes that this is a model of means: covariances between values are not  
 501 included in the model.

502 Step 2 (middle): To test hypotheses relating to the LI means, the ‘Population Bias’ model (with  
503 means for all tasks set to be the same) was compared to the ‘Task Effect’ model (where means  
504 varied by task).

505 Furthermore, to test the ‘Dorsal Steam’ hypothesis, a model with means for subsets of dorsal (A,  
506 B), ventral (C) and mixed tasks (D, E, F) were fixed (labelled as X, Z and Y). For the ‘Lexical Retrieval’  
507 hypothesis, a model with means for subsets of tasks with lexical retrieval (B, D) and tasks without  
508 (A, C, F) were fixed (labelled as X and Y respectively).

509 Step 3 (bottom): The oval symbol denotes a common factor that determines the covariance  
510 between observed variables. To test the hypothesis relating to LI covariances, a single factor  
511 ‘Person Effect’ model, was compared to a two factor ‘Task x Person Effect’ model. To achieve  
512 model identification, one of the paths from Factor 1 to a task had to be fixed to 1, and the path  
513 from Factor 2 to that task was fixed to zero. *In our preregistration this fixed path was planned to*  
514 *be task A, but due to the low reliability of that task, it was changed in the final analysis to be task*  
515 *D.* The covariance between Factor 1 and Factor 2 was also set to zero. Note that the means were  
516 also modelled as shown in the task effect model, but this was omitted from the model diagrams  
517 here for simplicity.

## 518 **Results**

519 All data are available on OSF (<https://osf.io/s9kx6/>). Results from the pre-registered analysis  
520 protocol (i.e., using the first 30 participants only) are shown in Supplementary Materials  
521 (<https://osf.io/g8mkv/>). As noted above, the factor solution from this sample was unstable and  
522 unduly influenced by one left-hander. We report here the results based on the final sample of 30  
523 right-handers and 7 left-handers, which gives a stable solution, and we include exploratory  
524 analyses relating the findings to handedness. The LI values reported here are based on the mean  
525 difference between left and right CBFV, as this gives normally distributed variables, but the results  
526 are highly similar when the non-normal peak-based LIs are used instead. The analysis script  
527 provided on OSF (<https://osf.io/g8zka/>) facilitates comparisons between different analytic  
528 pathways.

### 529 **Behavioural results**

530 We did not have specific predictions for the behavioural results, but present them here for  
531 completeness. For List Generation (A) and Sentence Generation (D), the number of words spoken  
532 per trial was recorded. The number of words spoken in both tasks and sessions were very similar:  
533 for task A, session 1, mean = 9.5, SD = 0.42, session 2, mean = 9.6, SD = 0.29; for task D, session 1,  
534 mean = 9.2, SD = 1.21, session 2, mean = 9.4, SD = 1.24. A repeated measures ANOVA showed no  
535 significant effects of task ( $F(1,36) = 1.22, p = 0.278$ ) on the number of words spoken, but there was  
536 a significant effect of session ( $F(1,36) = 5.73, p = 0.022$ ). Trials where participants failed to  
537 respond, or responded too early were excluded from analysis: these constituted less than 0.1% of  
538 trials.

539 For decision making tasks (B, C, E and F), the accuracy and RT of each response, and the number of  
540 omitted responses, were recorded (Table 1). Note that for task F participants were required to  
541 wait until the end of the word sequence before responding, and had only a second to respond;  
542 this accounts for the fast reaction times and relatively high number of omitted responses in task F.

543 The Phonological Decision and Sentence Comprehension tasks (tasks B and E) showed evidence of  
544 practice effects, as both accuracy and reaction times improved, and the number of omitted  
545 responses fell from Session 1 to Session 2.

546

### 547 **Table 1**

548 Behavioural data for tasks B, C, E and F. The table shows mean percentage accuracy and reaction  
 549 times (with SD), and results of t-tests comparing Session 1 with Session 2 for each measure. The  
 550 number of omitted responses is reported as a percentage of all events. B = Phonological Decision;  
 551 C = Semantic Decision; E = Sentence Comprehension; F = Syntactic Decision.  
 552

Measure	Session	Task B	Task C	Task E	Task F
Accuracy (%)	1	91.3 (5.55)	95.9 (3.08)	92.5 (4.81)	89.6 (8.31)
	2	93.3 (4.28)	95.0 (3.06)	94.2 (3.79)	89.4 (8.28)
	1 vs 2	t=-3.27, p=.002	t=1.61, p=.115	t=-2.70, p=.011	t=-0.07, p=.944
Reaction times (s)	1	1.66 (0.22)	1.14 (0.2)	2.17 (0.12)	0.33 (0.08)
	2	1.49 (0.21)	1.05 (0.2)	2.11 (0.15)	0.33 (0.07)
	1 vs 2	t=8.73, p<.001	t=4.77, p<.001	t=3.27, p=.002	t=0.64, p=.528
Omitted responses (%)	1	2.34	0.84	2.79	4.20
	2	0.78	0.60	1.62	4.44

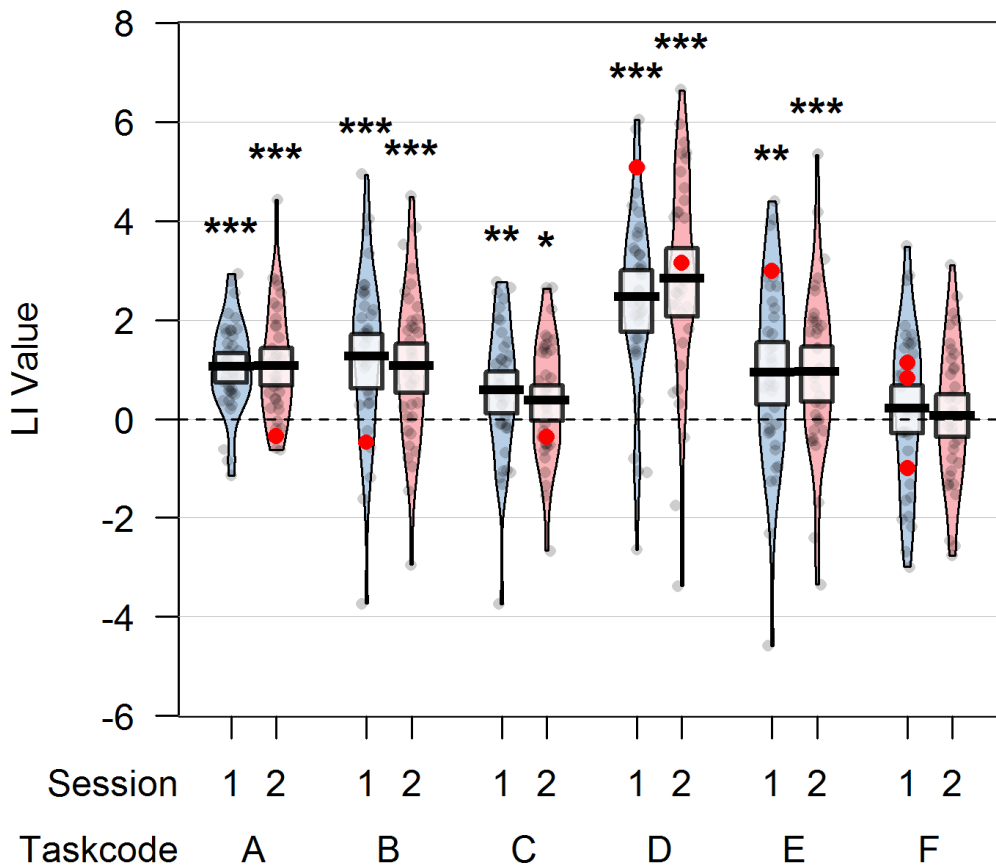
553

#### 554 Lateralisation results

555 Three outlier LI values were excluded where the standard error across trials was above the upper  
 556 cut-off. Six LI values were excluded because a subject had less than twelve useable trials for a  
 557 given task in a given session. The remaining data for these participants were retained in the  
 558 analysis. Excluded datapoints are shown as red dots in Figure 5.

559 Figure 5 shows the distribution of LIs as a pirate plot (Phillips, 2017). Task D (Sentence Generation)  
 560 showed the strongest left lateralisation. Shapiro-Wilks normality tests showed that LI values for all  
 561 12 conditions were normally distributed. One sample t-tests (testing for mean > 0) showed that all  
 562 conditions were significantly left lateralised, except task F (Syntactic Decision; Session 1: t (33) =  
 563 0.77, p = 0.224; Session 2: t (36) = 0.33, p = 0.373).





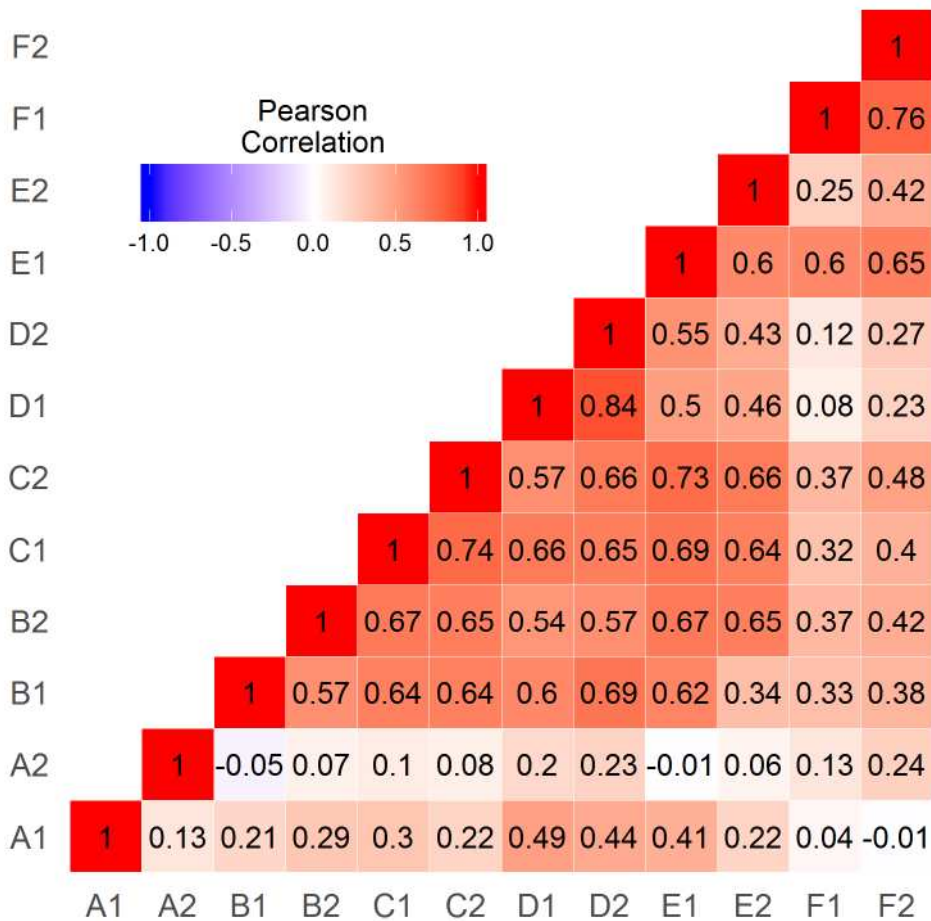
564

565 **Figure 5**

566 Pirate plot of LI values for all tasks (A-F) and sessions (blue = Session1, pink = Session2). Excluded  
567 data-points are shown in red. Asterisks show results of Wilcoxon tests comparing the LI values of  
568 the group (omitting excluded data-points) to zero (\*  $p < .05$ ; \*\*  $p < .01$ ; \*\*\*  $p < .001$ ).

569

570 Figure 6 shows a correlation matrix of LI values for all tasks and sessions. Test-retest correlations  
571 varied between tasks. Task A (List Generation) had poor test-retest reliability (Pearson's  $r = 0.13$ ),  
572 and low correlations with other tasks. Test-retest reliability for other tasks ranged from  $r = 0.57$  to  
573  $0.84$ . Tasks B, C, D and E were strongly intercorrelated. Task F (Syntactic Decision) had high test-  
574 retest reliability ( $r = 0.76$ ) but relatively low correlations with other tasks.



575

576 **Figure 6**

577 Correlation matrix for LIs from the six language tasks given on two occasions.

578

### 579 **Structural Equation Modelling**

580 The LI data were entered into the SEM analysis to test hypotheses about the group mean LI values  
 581 and covariances in LI values across subjects. Table 2 summarises the SEM results.

#### 582 *Step 1: Testing Stability of LI Values*

583 As shown in Table 2, the fit of all the means-only models was very poor. This is to be expected, as  
 584 these models ignore covariances, and, as indicated in Figure 6, there are substantial correlations  
 585 both between and within tasks. Our interest at this point, however, is in the relative fit of different  
 586 models of means, rather than overall model fit. The Fully Saturated model (with free means and  
 587 variances) was compared to the Task Effect model, which fixed the means and variances for each  
 588 task to be stable over sessions (i.e. A1 = A2, B1 = B2, etc.). The Task Effect model fit did not

589 deteriorate significantly from that of the Fully Saturated model, supporting the hypothesis that LI  
 590 means for each task were stable across sessions.

591

592 **Table 2**

593 Model fit statistics from structural equation models and model comparisons. -2LogL = -2 log  
 594 likelihoods; df = degrees of freedom; BIC = Bayesian Information Criterion; CFI = Comparative Fit  
 595 Index; RMSEA = Root Mean Square Error of Approximation.

Model	Description	-2LogL	df	BIC	CFI	RMSEA	Chi Square test	
							Compared to	p
<b>Fully Saturated Model</b>	Free means and variances	1574.4	411	90.4	NA	NA	-	NA
<b>Task Effect Model</b>	Stable means and variances	1580.8	423	53.4	0.022	0.292	Fully Saturated Model	0.896
<b>Population Bias Model</b>	Equal means and variances	1715.5	433	151.9	-0.474	0.337	Task Effect Model	<0.001
<b>Dorsal Stream Model</b>	Means for tasks AB > DEF > C	1664.8	429	115.7	-0.288	0.323	Task Effect Model	<0.001
<b>Lexical Retrieval Model</b>	Means for tasks BD > ACF	1631.6	429	82.5	-0.156	0.306	Task Effect Model	<0.001
<b>Person Effect Model</b>	Covariances have one factor structure	1378.4	417	-127.4	0.805	0.136	Task Effect Model	<0.001
<b>Task x Person Effect Model</b>	Covariances have bifactor structure	1337.8	412	-149.9	0.947	0.073	Person Effect Model	<0.001

596

597

598 *Step 2: Testing Models of Means*

599 To demonstrate whether LI means differed between tasks, the Task Effect model (with different  
 600 means for each task) was compared to the Population Bias model (with means fixed to be the  
 601 same for all tasks). This may be seen as a null hypothesis that treats all tasks as equivalent  
 602 measures of laterality. The Population Bias model gave significantly worse fit (see Table 2),  
 603 supporting the hypothesis that LI means differed between tasks.

604 Two further models were compared to the Task Effect model. The Dorsal Stream model  
 605 categorised the language tasks according to the involvement of the dorsal or ventral stream. Tasks  
 606 A and B were categorised as involving strong dorsal stream activity, task C as strong ventral stream  
 607 activity, and tasks D, E and F as intermediate (hence, means for AB > DEF > C). This model gave  
 608 significantly poorer fit than the Task Effect model – as is evident from Figure 5, which shows  
 609 relatively weak lateralisation for tasks A and B compared to task D. The Lexical Retrieval model did  
 610 not fare any better. This categorised tasks B and D as involving strong lexical retrieval, whereas  
 611 tasks A, C and F did not involve lexical retrieval, and task E was difficult to classify and so was  
 612 considered as independent of the other measures (BD > ACF). Again, this model gave a worse fit

613 than the Task Effect model, indicating that, while laterality varied between tasks, it did not fit the  
614 either of the predicted patterns. Note, however, that the pre-registered tests specified for both  
615 theories have some limitations, as discussed further below.

616 *Step 3: Testing Models of Covariances*

617 At Step 3 we tested whether the covariances between tasks had a single factor structure (Person  
618 Effect model) or a bifactor structure (Task by Person Effect model). Not surprisingly, given the  
619 strong correlations in Figure 6, both within and across tasks, the Person Effect model gave  
620 substantially better fit than the Task Effect model (see Table 2); nevertheless, the overall fit of this  
621 model was poor. The Task by Person Effect model gave a significantly improved fit. A plot of the  
622 two factors is shown in Figure 7: note that, although the model fit is not affected by task selection,  
623 the factor scores depend on which task has fixed paths to the factors. The paths for the case when  
624 Sentence Generation is fixed are shown in Table 3. It can be seen that List Generation has only a  
625 weak loading on Factor 1, whereas Phonological Decision, Semantic Decision and Sentence  
626 Comprehension have moderate loadings on both factors. Syntactic Decision has a strong loading  
627 on Factor 2 but does not load on Factor 1, reflecting the weak correlation of this task with  
628 Sentence Generation.

629

630 **Table 3**

631 Path weightings (and 95% confidence intervals) from each latent factor (Factor 1 and Factor 2) to  
632 each task (A to F) from the winning bifactor model.

633

Task	Factor 1		Factor 2	
	Path	95% CI	Path	95% CI
A: List Generation	0.18	0.05 to 0.31	-0.02	-0.27 to 0.24
B: Phonological Decision	0.61	0.40 to 0.81	0.55	0.21 to 0.89
C: Semantic Decision	0.53	0.36 to 0.69	0.52	0.23 to 0.81
D: Sentence Generation	1.00	Fixed	0.00	Fixed
E: Sentence Comprehension	0.56	0.30 to 0.82	0.95	0.54 to 1.37
F: Syntactic Decision	0.13	-0.13 to 0.40	1.16	0.75 to 1.56

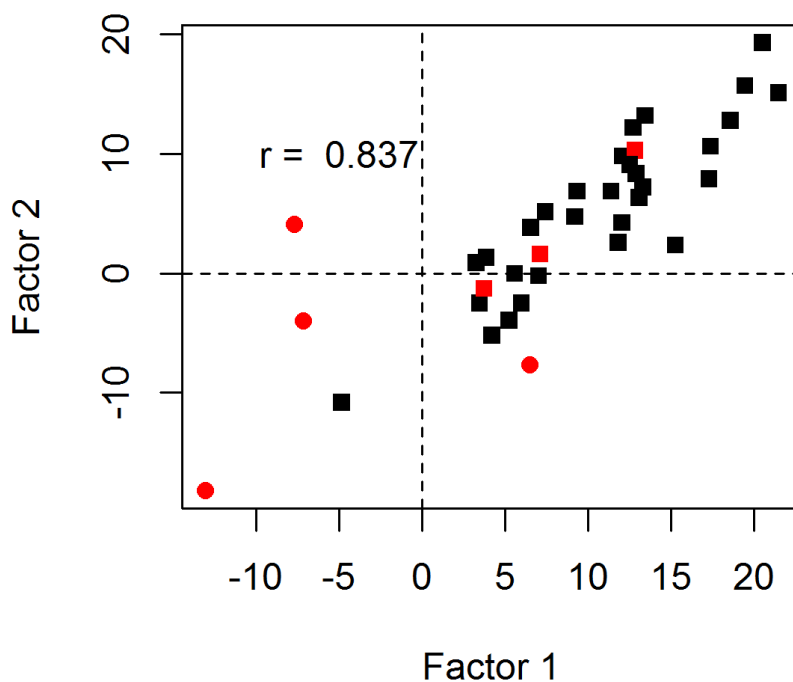
634

635 In our original analysis with 30 participants, a similar factor structure was observed, but there was  
636 a concern that this depended solely on a single left-handed participant (see Supplementary  
637 Material, <https://osf.io/g8mkv/>). With the larger sample of 37 participants, the bifactor (Task by

638 Person Effect) model was superior in all runs of a leave-one-out analysis. The bifactor model was  
639 also the best-fitting model when only the 30 right-handers were included in the analysis.

640 Nevertheless, it is clear from Figure 7 that the two factors were highly intercorrelated, and the  
641 impression is that the bifactor solution is heavily affected by some influential cases. Cook's  
642 distance identified four bivariate outliers, marked with circles in Figure 7: all four outliers were  
643 left-handers. When the analysis was re-run omitting these cases, the single factor model gave a  
644 better model fit when all N=33 subjects were included (single factor BIC=-142.7, bifactor BIC=-  
645 138.6), and in all but one run of the leave-one-out analysis.

646 We can conclude from this analysis that, although univariate normality was satisfactory, our data  
647 did not meet conditions of multivariate normality; this leads to the conclusion that the sample is  
648 not homogeneous, but contains a mixture of laterality patterns. We discuss the implications of this  
649 finding below.



650

651 **Figure 7**

652 Correlation between two factors from the bifactor (Task by Person Effect) model, with left-  
653 handers shown in red, and bivariate outliers as circles.

654

## 655 **Discussion**

656 The question of whether cerebral lateralisation is a unitary function may be interpreted at two  
657 levels: at the population level, we may ask whether all language tasks show a similar degree of  
658 lateralisation, and at the individual level, whether people show consistent differences in laterality  
659 profiles across tasks.

660 Although we used formal modelling to address these questions, a good insight into the answers  
661 can be obtained by viewing figures 5 and 6. Figure 5 shows clear differences from task to task in  
662 strength of cerebral lateralisation, whereas Figure 6 shows moderate-to-good test-retest reliability  
663 for all but one task, coupled with significant cross-task correlations.

664 The SEM analyses allowed us to explore these patterns further. Regarding means, as expected, a  
665 null hypothesis of no difference between tasks could be convincingly rejected. However, the  
666 specific patterns that we predicted should be seen on the basis of two existing models – the Dorsal  
667 Stream model and the Lexical Retrieval model – did not give a good fit. It could be argued that the  
668 data are, in fact, consistent with the Dorsal Stream model, insofar as the three tasks that involved  
669 implicit or explicit generation of speech – List Generation, Phonological Decision and Sentence  
670 Generation – were the ones that showed the strongest lateralisation (see Figure 5). The poor fit of  
671 the Dorsal Stream model was in part due to the fact that Sentence Generation was judged to  
672 implicate both streams, and was not therefore predicted to be as strongly lateralised as tasks with  
673 weaker semantic demands. However, it clearly makes demands on the phonological-articulatory  
674 system, and with hindsight it could be argued that in terms of articulatory complexity it was more  
675 demanding than the other tasks. A key question is whether blood flow measured using fTCD  
676 reflects the average of activity in a lateralised dorsal stream and a bilateral ventral stream, or  
677 whether the absolute dorsal stream activity is the main factor affecting the LI. In future we plan  
678 studies to address this question using fMRI.

679 More generally, based on the pattern of results observed in this study, it appears that whole-  
680 hemisphere lateralisation as measured by fTCD is driven most strongly by generation of  
681 meaningful, connected speech (e.g. Sentence Generation). Lateralization for this task was stronger  
682 than for automatic, non-propositional speech (List Generation) or implicit sub-vocalisation  
683 (Phonological Decision). By contrast, lateralisation was non-significant for the Syntactic Decision  
684 task.

685 We would, however, emphasise the need for caution in treating any one task as an indicator of a  
686 particular language function: it is evident that even minor modifications to task demands may  
687 affect laterality, particularly when sample size is relatively small. For instance, in a related study  
688 with a different sample of people, we recently found that List Generation was not lateralised  
689 (Woodhead, Rutherford, & Bishop, 2018). In that study we interleaved a simple number  
690 generation (counting) task with trials of Sentence Generation, whereas in the current study, List  
691 Generation was administered in a separate block, with the type of list (numbers, days of the week,  
692 months of the year) varied to engage the participants' attention throughout the block. Although  
693 the counting task used by Woodhead et al (2018) was not significantly lateralised, it had good  
694 split-half reliability and was significantly correlated with Sentence Generation, whereas the List  
695 Generation task used in the current study was the only task to show poor test-retest reliability and  
696 relatively weak correlations with other tasks. Furthermore, our Semantic Decision task was  
697 designed to tap into similar semantic processes as the Pyramids and Palm Trees test (Howard &  
698 Patterson, 1992), but resulted in weaker LIs than seen in a study by Bruckert (Bruckert, 2016) using  
699 the Pyramids and Palm Trees task. It could be that the two-alternative forced choice task used in  
700 that study was more demanding than our match/no-match decision, but this kind of difference  
701 cautions us about relying on a single test to indicate a type of linguistic processing.

702 One convincing point to emerge from the analysis of mean data is that most language tasks (B, C,  
703 D, E and F) showed stable lateralisation measured in different sessions, but they differed in terms  
704 of the strength of left-lateralisation.

705 We turn next to the findings concerning covariances. It has been argued that fTCD is not useful for  
706 studying cerebral lateralisation because it is unreliable (Cai et al., 2013), but our data support  
707 those of Stroobant and Vingerhoets (Stroobant & Vingerhoets, 2001) in demonstrating that there  
708 is significant individual variation in language laterality between people that cannot just be  
709 attributed to noise. Furthermore, by moving from a definition of laterality based on a peak in the  
710 L-R difference wave to a definition based on mean L-R difference within a period of interest, we  
711 avoid the problem that can arise when laterality is forced into a non-normal distribution (see also  
712 Woodhead et al., 2018). As shown in Figure 5 and our tests of normality, when mean L-R  
713 difference is used, the distribution of LI values is normal.

714 The SEM also tested whether a single factor could explain individual differences in language  
715 lateralisation. At first glance, the results suggested this was not the case: the bifactor (Task by  
716 Person Effect) model showed superior fit over a single factor (Person Effect) model. This was the

717 conclusion suggested by our initial pre-registered analysis, based just on a sample of 30  
718 individuals. A leave-one-out analysis, however, made us cautious about accepting that result at  
719 face value, because the factor structure changed when a single left-hander with strongly  
720 complementary laterality on two tasks was excluded. For this reason we collected more data,  
721 adding seven right-handers to the sample. With this larger sample, we again found superiority for  
722 a bifactor solution, regardless of whether we included only right-handers or the full sample  
723 including left-handers. Yet there remained misgivings about the generalisability of the result, not  
724 least because the two factors were highly correlated (Pearson's  $r = 0.84$ ). A scatterplot of the two  
725 factors revealed a number of bivariate outliers and, as with our initial analysis, the pattern of  
726 results relied on which participants were included. Of course, it is not surprising that removing  
727 participants with the strongest dissociation between factors changes the factor structure: the  
728 point we wish to make is not that the results can alter in this way, but rather that the pattern of  
729 our SEM findings appears driven by heterogeneity within the sample, reflected in the presence of  
730 bivariate outliers.

731 The answer to the question of whether laterality is a unitary function is that, clearly, there are  
732 some individuals in whom laterality is different for different aspects of language. It is not,  
733 however, the case that there are two factors that act independently in the general population.  
734 Rather, the majority of people appear to have language laterality driven by a single process  
735 affecting all types of task, with a minority showing fractionation of language asymmetry.

736 The pattern of results is consistent with accounts of laterality that postulate qualitative rather  
737 than just quantitative differences between individuals. Theoretical accounts have mostly focused  
738 on a single dimension, arguing for laterality subgroups on the basis of non-normal distributions of  
739 scores (e.g. Mazoyer et al., 2014). Our results suggest that atypical laterality may be easier to  
740 identify when more than one language measure is considered, as detection of bivariate outliers  
741 can be effective with smaller samples than those required for detecting mixtures of distributions.

742 An association between atypical laterality and left-handedness has been established for many  
743 years, ever since early observations were made of superior recovery from aphasia after gun-shot  
744 wounds in left-handers (Subirana, 1958). However, most of the emphasis has been on atypical  
745 laterality in the sense of having language mediated by the right hemisphere. Although the number  
746 of left-handers in our sample is too small for numeric analysis, the fact that three of the four  
747 bivariate outliers were left-handers is a striking departure from chance (Fisher exact probability =



748 0.016) and compatible with the idea that language lateralisation is more likely to be multifactorial  
749 in left-handers than right-handers.

750 Further studies are needed to establish the key characteristics of tasks that index the two factors  
751 seen in some people, but we offer here some speculations. The main contributor to the second  
752 factor was the Syntactic Decision task, which differed from the other tasks in several regards. It  
753 used unfamiliar, nonword stimuli, and required the listener to identify syntactic errors. It was one  
754 of two receptive language tasks that involved processing of auditory language: the other was  
755 sentence comprehension, which had moderately strong loadings on the second factor. Perhaps  
756 the most surprising finding from this study is the fact that the one task that loaded on to the  
757 second factor (Syntactic Decision) was not lateralised, yet showed high test-retest reliability  
758 ( $R=0.67$ ). We had anticipated that a lack of lateralisation on a task might be a consequence of  
759 noisy data giving poor test reliability – or alternatively a lack of individual variation if both  
760 hemispheres contributed equally in most people. Our data suggest that individuals do vary in the  
761 hemisphere used when doing the syntactic judgement task, and that this bias is reliable, but that it  
762 is not systematic across the population. This is perhaps the best evidence to date that strength as  
763 well as direction of lateralisation for a task is a stable trait.

#### 764 **Limitations**

765 As noted above, the principal limitation of fTCD is that it does not allow one to localise lateralised  
766 activity within a hemisphere. In future work, we plan to extend this line of investigation to  
767 consider whether similar patterns of lateralisation can be seen using comparable tasks with fMRI.  
768 The benefit of fTCD is that it is relatively inexpensive and quick to administer, and so enables us to  
769 gather data that can be used as a basis for developing a more hypothesis-driven approach that can  
770 then be extended and validated with fMRI.

771 A further limitation is that we lacked statistical power or range of measures that would be needed  
772 to evaluate more complex models. The bifactor model that gave the best fit in our study must be  
773 interpreted with caution. It will need to be replicated in larger samples and shown to generalise to  
774 new tasks - it remains a possibility that using a different set of tasks would reveal different or  
775 further fractionation of language lateralisation. Furthermore, although we have shown a bifactor  
776 model is a better fit than a single factor model, it is possible that more than two factors are  
777 needed to explain the full range of patterns of language lateralisation.

#### 778 **Summary**

779 In summary, these results indicate that there are meaningful differences in language lateralisation  
780 between tasks, and meaningful individual variability in lateralisation that is not simply due to  
781 measurement error. Even when a language-related task is not left-lateralised, there are stable  
782 individual differences in the contribution of the two hemispheres. Structural equation modelling  
783 of individual variability indicated that although a two-factor model gave a better fit than a single  
784 factor model, the effect was driven by a small subset of participants with discrepant laterality, and  
785 a single factor could account for variation in the majority of participants. Overall, our findings  
786 suggest there are qualitative as well as quantitative differences between people in laterality across  
787 tasks, and that consideration of asymmetry profiles on several tasks together can help identify  
788 cases of atypical laterality.

### 789 **Competing Interests**

790 None to declare.

791

### 792 **References**

- 793 Badzakova-Trajkov, G., Corballis, M. C., & Häberling, I. S. (2016). Complementarity or  
794 independence of hemispheric specializations? A brief review. *Neuropsychologia*, *93*, 386–393.  
795 <https://doi.org/10.1016/j.neuropsychologia.2015.12.018>
- 796 Bishop, D. V. M. (2013). Cerebral asymmetry and language development: Cause, correlate, or  
797 consequence? *Science*, *340*(6138). <https://doi.org/10.1126/science.1230531>
- 798 Bishop, D. V. M., & Robson, J. (1989). Unimpaired Short-term Memory and Rhyme Judgement in  
799 Congenitally Speechless Individuals: Implications for the Notion of “Articulatory Coding.” *The*  
800 *Quarterly Journal of Experimental Psychology Section A*, *41*(1), 123–140.  
801 <https://doi.org/10.1080/14640748908402356>
- 802 Bozic, M., Tyler, L. K., Ives, D. T., Randall, B., & Marslen-Wilson, W. D. (2010). Bihemispheric  
803 foundations for human speech comprehension. *Proceedings of the National Academy of*  
804 *Sciences*, *107*(40), 17439–17444. <https://doi.org/10.1073/pnas.1000531107>
- 805 Bradshaw, A. R., Thompson, P. A., Wilson, A. C., Bishop, D. V. M., & Woodhead, Z. V. J. (2017).  
806 Measuring language lateralisation with different language tasks: a systematic review. *PeerJ*, *5*,  
807 e3929. <https://doi.org/10.7717/peerj.3929>

- 808 Bruckert, L. (2016). *Is language laterality related to language abilities ?* University of Oxford.
- 809 Cai, Q., Van der Haegen, L., & Brysbaert, M. (2013). Complementary hemispheric specialization for  
810 language production and visuospatial attention. *Proceedings of the National Academy of*  
811 *Sciences of the United States of America*, 110(4), E322-30.  
812 <https://doi.org/10.1073/pnas.1212956110>
- 813 Deppe, M., Knecht, S., Henningsen, H., & Ringelstein, E. B. (1997). Average: A windows program  
814 for automated analysis of event related cerebral blood flow. *Journal of Neuroscience*  
815 *Methods*, 75(2), 147–154. [https://doi.org/10.1016/S0165-0270\(97\)00067-8](https://doi.org/10.1016/S0165-0270(97)00067-8)
- 816 Dhanjal, N. S., Handunnetthi, L., Patel, M. C., & Wise, R. J. S. (2008). Perceptual systems controlling  
817 speech production. *Journal of Neuroscience*, 28(40), 9969–9975.  
818 <https://doi.org/10.1523/JNEUROSCI.2607-08.2008>
- 819 Fedorenko, E., Hsieh, P.-J., Nieto-Castañón, A., Whitfield-Gabrieli, S., & Kanwisher, N. (2010). New  
820 method for fMRI investigations of language: defining ROIs functionally in individual subjects.  
821 *Journal of Neurophysiology*, 104, 1177–1194. <https://doi.org/10.1152/jn.00032.2010>
- 822 Friederici, A. D. (2011). The brain basis of language processing: From structure to function.  
823 *Physiological Reviews*, 91(4), 1357–92. <https://doi.org/10.1152/physrev.00006.2011>
- 824 Gaillard, W. D., Balsamo, L., Xu, B., McKinney, C., Papero, P. H., Weinstein, S., ... Theodore, W. H.  
825 (2004). fMRI language task panel improves determination of language dominance. *Neurology*,  
826 63(8), 1403–1408. <https://doi.org/10.1212/01.WNL.0000141852.65175.A7>
- 827 Groen, M. A., Whitehouse, A. J. O., Badcock, N. A., & Bishop, D. V. M. (2012). Does cerebral  
828 lateralization develop? A study using functional transcranial Doppler ultrasound assessing  
829 lateralization for language production and visuospatial memory. *Brain and Behavior*, 2(3),  
830 256–269. <https://doi.org/10.1002/brb3.56>
- 831 Hesling, I., Labache, L., Jobard, G., & Leroux, G. (2018). Brain Areas Commonly Activated and  
832 Asymmetrical in Pro- Duction , Listening and Reading Tasks At the Word Level : an Fmri Study  
833 of 144 Right-Handers, (August). <https://doi.org/10.1101/382960>
- 834 Hickok, G., & Poeppel, D. (2007). The cortical organization of speech processing. *Nature Reviews*  
835 *Neuroscience*, 8(5), 393–402. <https://doi.org/10.1038/nrn2113>
- 836 Hoaglin, D. C., & Iglewicz, B. (1987). Fine tuning some resistant rules for outlier labeling. *Journal of*

- 837 *American Statistical Association*, 82(400), 1147–1149.
- 838 Howard, D., & Patterson, K. E. (1992). The pyramids and palm trees test: A test for semantic access  
839 from words and pictures. *Bury St Edmunds, UK: Thames Valley Test Company Ltd.*, 16.  
840 Retrieved from  
841 [https://books.google.be/books/about/The\\_Pyramids\\_and\\_Palm\\_Trees\\_Test.html?id=dykON](https://books.google.be/books/about/The_Pyramids_and_Palm_Trees_Test.html?id=dykONQAACAAJ&redir_esc=y)  
842 [QAACAAJ&redir\\_esc=y](https://books.google.be/books/about/The_Pyramids_and_Palm_Trees_Test.html?id=dykONQAACAAJ&redir_esc=y)
- 843 Kline, R. B. (2011). *Principles and practice of structural equation modeling. Structural Equation*  
844 *Modeling* (Vol. 156). <https://doi.org/10.1038/156278a0>
- 845 Knecht, S., Deppe, M., Ebner, A., Henningsen, H., Huber, T., Jokeit, H., & Ringelstein, E. B. (1998).  
846 Noninvasive Determination of Language Lateralization by Functional Transcranial Doppler  
847 Sonography : A Comparison With the Wada Test. *Stroke*, 29(1), 82–86.  
848 <https://doi.org/10.1161/01.STR.29.1.82>
- 849 Lewis, T. F. (2017). Evidence regarding the internal structure: Confirmatory factor analysis.  
850 *Measurement and Evaluation in Counseling and Development*, 50(4), 239–247.  
851 <https://doi.org/10.1080/07481756.2017.1336929>
- 852 Mazoyer, B., Zago, L., Jobard, G. G., Crivello, F., Joliot, M., Perchey, G., ... Tzourio-Mazoyer, N.  
853 (2014). Gaussian mixture modeling of hemispheric lateralization for language in a large  
854 sample of healthy individuals balanced for handedness. *PLoS ONE*, 9(6), 9–14.  
855 <https://doi.org/10.1371/journal.pone.0101165>
- 856 Milner, B., Branch, C., & Rasmussen, T. (1966). Evidence for bilateral speech representation in  
857 some non-right handers. *Transactions of the American Neurological Association*, 91, 306–308.
- 858 Payne, H., Gutierrez-Sigut, E., Subik, J., Woll, B., & MacSweeney, M. (2015). Stimulus rate increases  
859 lateralisation in linguistic and non-linguistic tasks measured by functional transcranial  
860 Doppler sonography. *Neuropsychologia*, 72, 59–69.  
861 <https://doi.org/10.1016/j.neuropsychologia.2015.04.019>
- 862 Peelle, J. E. (2012). The hemispheric lateralization of speech processing depends on what “speech”  
863 is: a hierarchical perspective. *Frontiers in Human Neuroscience*, 6(November), 309.  
864 <https://doi.org/10.3389/fnhum.2012.00309>
- 865 Pinel, P., & Dehaene, S. (2010). Beyond hemispheric dominance: brain regions underlying the joint

- 866 lateralization of language and arithmetic to the left hemisphere. *Journal of Cognitive*  
867 *Neuroscience*, 22, 48–66. <https://doi.org/10.1162/jocn.2009.21184>
- 868 R Core Team., & R Development Core Team. (2013). R: A language and environment for statistical  
869 computing. R Foundation for Statistical Computing. Vienna, Austria. *R Foundation for*  
870 *Statistical Computing, Vienna, Austria.*
- 871 Ramsey, N. F., Sommer, I. E. C., Rutten, G. J., & Kahn, R. S. (2001). Combined analysis of language  
872 tasks in fMRI improves assessment of hemispheric dominance for language functions in  
873 individual subjects. *NeuroImage*, 13(4), 719–33. <https://doi.org/10.1006/nimg.2000.0722>
- 874 Rauschecker, J. P. (2018). Where, When, and How: Are they all sensorimotor? Towards a unified  
875 view of the dorsal pathway in vision and audition. *Cortex*. Elsevier Ltd.  
876 <https://doi.org/10.1016/j.cortex.2017.10.020>
- 877 Rosch, R. E., Bishop, D. V. M., & Badcock, N. A. (2012). Lateralised visual attention is unrelated to  
878 language lateralisation, and not influenced by task difficulty - A functional transcranial  
879 Doppler study. *Neuropsychologia*, 50(5), 810–815.  
880 <https://doi.org/10.1016/j.neuropsychologia.2012.01.015>
- 881 Seghier, M. L., & Price, C. J. (2018). Interpreting and Utilising Intersubject Variability in Brain  
882 Function. *Trends in Cognitive Sciences*, 22(6), 517–530.  
883 <https://doi.org/10.1016/j.tics.2018.03.003>
- 884 Stroobant, N., Buijs, D., & Vingerhoets, G. (2009). Variation in brain lateralization during various  
885 language tasks: A functional transcranial Doppler study. *Behavioural Brain Research*, 199(2),  
886 190–196. <https://doi.org/10.1016/j.bbr.2008.11.040>
- 887 Stroobant, N., & Vingerhoets, G. (2001). Test-retest reliability of functional transcranial Doppler  
888 ultrasonography. *Ultrasound in Medicine & Biology*, 27(4), 509–514.  
889 [https://doi.org/10.1016/S0301-5629\(00\)00325-2](https://doi.org/10.1016/S0301-5629(00)00325-2)
- 890 Subirana, A. (1958). The prognosis in aphasia in relation to cerebral dominance and handedness.  
891 *Brain*, 81(3), 415–425. <https://doi.org/10.1093/brain/81.3.415>
- 892 Szekely, A., Jacobsen, T., D’Amico, S., Devescovi, A., Andonova, E., Herron, D., ... Bates, E. (2004). A  
893 new on-line resource for psycholinguistic studies. *Journal of Memory and Language*, 51(2),  
894 247–250. <https://doi.org/10.1016/j.jml.2004.03.002>

- 895 Tailby, C., Abbott, D. F., & Jackson, G. D. (2017). The diminishing dominance of the dominant  
896 hemisphere: Language fMRI in focal epilepsy. *NeuroImage: Clinical*, *14*, 141–150.  
897 <https://doi.org/10.1016/j.nicl.2017.01.011>
- 898 Vingerhoets, G., Alderweireldt, A.-S. S., Vandemaele, P., Cai, Q., Van der Haegen, L., Brysbaert, M.,  
899 & Achten, E. (2013). Praxis and language are linked: Evidence from co-lateralization in  
900 individuals, with atypical language dominance. *Cortex*, *49*(1), 172–183.  
901 <https://doi.org/10.1016/j.cortex.2011.11.003>
- 902 Whitehouse, A. J. O., & Bishop, D. V. M. (2009). Hemispheric division of function is the result of  
903 independent probabilistic biases. *Neuropsychologia*, *47*(8–9), 1938–1943.  
904 <https://doi.org/10.1016/j.neuropsychologia.2009.03.005>
- 905 Woodhead, Z. V. J., Rutherford, H. A., & Bishop, D. V. M. (2018). Measurement of language  
906 laterality using functional transcranial Doppler ultrasound: a comparison of different tasks.  
907 *Wellcome Open Research*, *3*(0), 104. <https://doi.org/10.12688/wellcomeopenres.14720.1>
- 908 Zago, L., Petit, L., Mellet, E., Jobard, G., Crivello, F., Joliot, M., ... Tzourio-Mazoyer, N. (2015). The  
909 association between hemispheric specialization for language production and for spatial  
910 attention depends on left-hand preference strength. *Neuropsychologia*, *93*, 394–406.  
911 <https://doi.org/10.1016/j.neuropsychologia.2015.11.018>
- 912