

**Involvement of prefrontal cortex in scalar implicatures:
evidence from magnetoencephalography**

Journal:	<i>Language, Cognition and Neuroscience</i>
Manuscript ID:	PLCP-2014-OP-9658.R2
Manuscript Type:	Original Paper
Date Submitted by the Author:	23-Feb-2015
Complete List of Authors:	Politzer-Ahles, Stephen; New York University, Abu Dhabi, NYUAD Institute Gwilliams, Laura; New York University, Psycholinguistics
Keywords:	scalar implicature, pragmatics, magnetoencephalography, prefrontal cortex

SCHOLARONE™
Manuscripts

1
2
3 **Involvement of prefrontal cortex in scalar implicatures: evidence from**
4 **magnetoencephalography**
5
6
7

8 Stephen Politzer-Ahles*, Laura Gwilliams
9
10 *NYU Abu Dhabi Institute, New York University, Abu Dhabi, United Arab Emirates*

11
12
13 **Address correspondence to*
14 Stephen Politzer-Ahles
15 NYUAD
16 PO Box 129188
17 Abu Dhabi
18 UAE
19 Tel.: +971-56-689-9497
20 E-mail: spa268@nyu.edu
21
22

23 *Additional contacts*
24 Laura Gwilliams
25 NYUAD
26 PO Box 129188
27 Abu Dhabi
28 UAE
29 Tel.: +971 563190446
30 E-mail: laura.gwilliams@nyu.edu
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 Involvement of prefrontal cortex in scalar implicatures: evidence from magnetoencephalography

8
9
10 While many recent studies have investigated scalar inferences (i.e., the
11 interpretation of *some* as meaning *some but not all*) using offline as well as online
12 behavioural methods, little is known about the neural correlates of scalar
13 inference realisation. The present study used magnetoencephalography, which
14 has high temporal resolution, to measure neural activity while participants heard
15 stories that included the scalar inference trigger *some* in contexts that either
16 provide strong cues for a scalar inference or provide weaker cues. The middle
17 portion of the lateral prefrontal cortex (Brodmann area 46) showed an increased
18 response to *some* in contexts with fewer cues to the inference, suggesting that this
19 condition elicited greater effort. While the results are not predicted by traditional
20 all-or-nothing accounts of scalar inferencing that assume the process is always
21 automatic or always effortful, they are consistent with more recent gradient
22 accounts which predict that the speed and effort of scalar inferences is strongly
23 modulated by numerous contextual factors.
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

Keywords: scalar implicature; pragmatics; magnetoencephalography; prefrontal
cortex

Introduction

Comprehending a natural utterance generally involves inferring messages that were not explicitly expressed. For example, the literal meaning of (1) is (1a), but in many contexts a competent listener will interpret (1) as meaning (1b).

- (1) Some of the flowers have bloomed.
 - (a) At least one of the flowers has bloomed.
 - (b) Some, *but not all*, of the flowers have bloomed.

Meaning (1b), the *pragmatic interpretation*, is presumably realised by generating a set of stronger alternative sentences that were not spoken (e.g., "All of the flowers have bloomed") and inferring that a competent speaker choosing not to utter a stronger

1
2 sentence must have intended to convey that the stronger sentence is not true, e.g., that
3
4 not all of the flowers had bloomed (Chemla & Singh, 2014; Horn, 1972; Katsos &
5 Cummins, 2010; Noveck & Sperber, 2007; Sauerland, 2012). This process is known as
6 a *scalar inference*. The pragmatic interpretation (1b) can be cancelled without yielding a
7 self-contradictory sentence, as shown in (2); furthermore, in certain contexts the
8 pragmatic interpretation may not arise at all, such as in semantic contexts where the
9 alternative sentence is weaker rather than stronger (3) and epistemic contexts where the
10 speaker does not have enough knowledge of the situation to license the scalar inference
11 (4). Meaning (1a), the lexical interpretation, does not have these properties.
12
13

- 14 (2) Some of the flowers have bloomed; in fact, all of them have.
15
16 (3) If some of the flowers bloom, the garden will look beautiful.
17
18 (4) I took a quick glance at the garden and saw that some of the flowers had
19 bloomed.

20
21 A major question of interest in psycholinguistics is how scalar inferences are
22 cognitively realised. In recent years there has been substantial debate over whether the
23 mechanisms for scalar inferencing are pragmatic (per Horn, 1972) or grammatical (per
24 Chierchia, Fox, & Spector, 2012), and whether the pragmatic meaning is realised
25 immediately and automatically (per Chierchia, 2004, and Levinson, 2000) or with a
26 processing cost (per e.g. Sperber & Wilson, 1995). The debate over whether inference
27 realisation is pragmatic or grammatical is mainly a question of listeners' competence
28 and the unit of analysis is often overt judgments of whether or not enriched meanings
29 are realised in various types of complex sentences (see e.g. Chemla & Singh, 2014;
30 Sauerland, 2012), although the answer to this debate could also have implications for
31 processing. On the other hand, the debate over the speed and context-dependence of
32 scalar inferencing (see e.g. Katsos & Cummins, 2010; Noveck & Reboul, 2008) is
33
34

1 concerned mainly with psycholinguistic measures, but the specific grammatical or
2 pragmatic mechanisms underlying putative processing costs (or lack thereof) are not as
3 clearly spelled out. Both the question of pragmatic vs. grammatical processing and that
4 of cost vs. automaticity have been addressed in numerous behavioural experiments
5 using methods including online forced-choice paradigms (Bott, Bailey, & Grodner,
6 2012; Bott & Noveck, 2004; Chevallier et al., 2008; Feeney et al., 2004), offline forced-
7 choice paradigms (Chemla & Spector, 2011; Degen & Tanenhaus, 2014; Geurts &
8 Poussoulous, 2009), self-paced reading (Bergen & Grodner, 2012; Breheny, Katsos, &
9 Williams, 2006; Hartshorne & Snedeker, submitted; Lewis, 2013; Politzer-Ahles &
10 Fiorentino, 2013), dual-task (De Neys & Schaeken, 2007; Dieussaert, Verkerk, Gillard,
11 & Shaeken, 2011; Marty & Chemla, 2013; Marty, Chemla, & Spector, 2013), and visual
12 world eye-tracking (Breheny, Ferguson, & Katsos, 2012, 2013; Degen & Tanenhaus,
13 2014; Grodner, Klein, Carbay, & Tanenhaus, 2010; Huang & Snedeker, 2009).

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
Less attention has been given, however, to the neural correlates of scalar
implicature. Thus far, only one experiment has investigated the neural substrates of
scalar inferencing using a method with high spatial resolution. Shetreet, Chierchia, and
Gaab (2014) used functional magnetic resonance imaging (fMRI) to examine brain
regions activated in *some*-sentences (e.g., "Some mice have grapes") and in *every*-
sentences (e.g., "Every penguin is on the bus") during a picture-sentence verification
task. They found that left inferior frontal gyrus (LIFG; Brodmann area 47) was
activated more for *some*-sentences, which evoked scalar inferences, than for *every*-
sentences, which did not. The authors argue that LIFG may thus be the locus of
semantic aspects of inference computation, such as the generation of relevant
alternatives and enrichment of the quantifier's interpretation via negation of the
alternatives (see Chierchia, 2004; Chierchia et al., 2012).

This finding is somewhat complicated, however, by several aspects of the study's methods and design. While fMRI has excellent spatial resolution, it has poor temporal resolution, making it difficult to know at what point in the sentence the LIFG effect was elicited (an issue of importance, since several accounts of scalar inferencing argue that it occurs immediately when the scalar expression, e.g. *some*, is encountered). On the other hand, methods with high temporal resolution, such as magnetoencephalography and electroencephalography, would allow researchers to disentangle inference-related neural activity elicited at the moment a scalar expression is presented, from other processes (such as verification of the upper-bounded meaning relative to the context) that may be elicited later as the sentence unfolds. Furthermore, the critical comparison in the experiment was between different words, *some* and *every*. The effects observed could reflect downstream verification strategies rather than inferencing itself. Crucially, the denotations of *some* and *every*, and thus the verification strategies they may induce, are different. Verifying *every* requires only checking for one counterexample; the same is true for verifying the semantic interpretation of *some*, as it requires only finding one instance of e.g. a mouse with grapes. Verifying *some* [*but not all*], on the other hand, requires identifying two different subsets (e.g., mice that have grapes and mice that do not). Thus, the effects could be due to processing different denotations rather than to the actual process of realising the enriched meaning.¹ To rule out such differences, it would be valuable to test scalar inferencing in a paradigm that 1) does not require an explicit verification task; 2) compares *some* in an inference-triggering context to *some* in a non-inference-triggering context, rather than comparing *some* to a different word; and 3) includes controls that replicate the differences in denotation between lexical *some* ("at least one") and pragmatic *some* ("at least one, but not all") but do not involve scalar inferences. Observing similar frontal activation in

such a design would add converging evidence that this region is involved in some aspect of scalar inferencing.

The present study

The present study used magnetoencephalography (MEG), which has the requisite combination of good spatial and temporal resolution, to investigate the neural substrates of realising scalar inferences. While several previous studies have used high-temporal-resolution techniques to examine inferencing (Chevallier, Bonnefond, Van der Henst, & Noveck, 2010; Hartshorne, Liem Azar, Snedeker, & Kim, in press; Hunt, Politzer-Ahles, Gibson, Minai, & Fiorentino, 2013; Nieuwland, Ditman, & Kuperberg, 2010; Noveck & Posada, 2003; Politzer-Ahles, Fiorentino, Jiang, & Zhou, 2013; Sikos, Tomlinson, Traut, & Grodner, 2013; Zhao, Liu, Chen, & Chen, 2015), they have all used electroencephalography (EEG), which has poorer spatial resolution; furthermore, other than Hartshorne et al. (in press) and Sikos et al. (2013), these studies have all used violation paradigms and/or examined words downstream of the quantifier. The present study is the first study with high spatial resolution to examine successful scalar inferencing while also controlling for the lexical issues described in the previous section.

We adopted a paradigm that has been widely used in reading time studies (Bergen & Grodner, 2012; Breheny et al., 2006; Hartshorne & Snedeker, submitted; Lewis, 2013; Politzer-Ahles & Fiorentino, 2013) and one EEG study (Hartshorne et al., in press), in order to contrast *some* in contexts that strongly support a scalar inference and in contexts that do not. An example is shown in (5). The first two sentences establish a context that introduces a salient set of referents and asks a question about either *all* of them (creating an upper-bounded context, 5a,c) or about *any* of them

(creating a lower-bounded context, 5b,d). *Some* is more likely to be interpreted pragmatically in an upper-bound context, where there is an explicit stronger alternative, than in a lower-bound context. Thus, brain regions associated with making scalar inferences may be expected to show greater activation in response to the word *some* in the upper-bound items.

(5)

(a) **Upper-bound *some*:** Mary was preparing to throw a party for John's relatives. She asked John whether all of them were staying in his apartment. John said that some of them were.

(b) **Lower-bound *some*:** Mary was preparing to throw a party for John's relatives. She asked John whether any of them were staying in his apartment. John said that some of them were.

(c) **Upper-bound *only some*:** Mary was preparing to throw a party for John's relatives. She asked John whether all of them were staying in his apartment. John said that *only* some of them were.

(d) **Lower-bound *only some*:** Mary was preparing to throw a party for John's relatives. She asked John whether any of them were staying in his apartment. John said that *only* some of them were.

While a lower-bound context like (5b) is often treated as an instance in which the scalar inference is not licensed (see, e.g., Breheny et al., 2006), and this context manipulation has been shown to modulate downstream reading times for words whose interpretation depends on the inference (Breheny et al., 2006; Lewis, 2013; Politzer-Ahles & Fiorentino, 2013), offline ratings we conducted (see Methods) showed that participants still incorporated the implicature into their ultimate interpretations of these sentences in this set of stimuli. Thus, we refer to these conditions as contexts that

provide stronger (5a, upper-bound) or weaker (5b, lower-bound) cues supporting the eventual scalar inference, rather than as conditions that do or do not license the inference.

We also included a pair of control conditions (5c,d) where the *not all* interpretation is made explicit via the addition of *only*. Unlike bare *some*, *only some* is semantically specified to mean *some but not all* and this meaning cannot be cancelled (Minai & Fiorentino, 2010). These conditions have the same context manipulation as the critical conditions (5a,b) but do not involve scalar inferencing. Thus, brain regions involved in scalar inferencing can be limited to those showing an interaction such that there is a context effect in the *some* conditions (5a,b) and not the *only some* conditions (5c,d).²

The direction of the context effect to be expected is an empirical question. Under accounts that assume a processing cost for scalar inferencing across the board (e.g. Noveck & Sperber, 2007), one would expect to see greater activation in the more strongly inference-supporting (upper-bounded) context. This may reflect extra effort involved in perspective-taking to make primary implicatures, under a pragmatic account, or extra effort involved in applying semantic operations to parse the strengthened interpretation, under a semantic account; some of these putative operations are likely to be required under both pragmatic and semantic accounts (see Chemla & Singh, 2014). On the other hand, under accounts that assume scalar inferences are made by default (e.g. Levinson, 2000), no difference in neural activation would be predicted between the two conditions (given that scalar inferences were realised at similar rates in both conditions). Results from reading time experiments using this paradigm are mixed; some find slowdowns for *some* in the more strongly inference-supporting context (Bergen & Grodner, 2012; Breheny et al., 2006; Hartshorne & Snedeker, submitted,

1
2 Experiment 1), consistent with accounts assuming extra processing effort for
3 inferencing, whereas others find no slowdown for *some* (Hartshorne & Snedeker,
4 submitted, Experiment 2; Lewis, 2013; Politzer-Ahles & Fiorentino, 2013), and none
5 find slowdowns for *some* in the less strongly inference-supporting (lower-bounded)
6 context. In the only electrophysiological study using this paradigm, Hartshorne and
7 colleagues (in press) found no difference in scalp EEG for *some* (although they did, like
8 all of the reading-time studies, find effects downstream confirming that scalar
9 inferences were more available in the upper-bounded than lower-bounded contexts). In
10 the present study we focus on source-level analysis, but in order to compare the results
11 to those of Hartshorne and colleagues (in press) we also report sensor-level findings.
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

Methods

Participants

Eleven native English speakers (9 female, mean age = 23, SD = 5.39) took part in the experiment; four additional volunteers participated but were not included in the data analysis because of excessive artifacts in their recordings. All had normal or corrected-to-normal vision and hearing and were right-handed as assessed by the Edinburgh Handedness Questionnaire (Oldfield, 1971). All participants provided their informed consent and were paid for their participation, and experimental procedures were approved by the Institutional Review Board of New York University Abu Dhabi.

Materials

One hundred and twenty-eight vignettes were created for the critical trials, according to the template shown in (5). The first sentence of each vignette introduced a set of

referents (e.g., John's relatives) to the discourse, and the second established an upper- or lower-bounded question under discussion using either the quantifier *all* or the quantifier *any*. The only difference between contexts is the use of *all* or *any* in the second sentence. Finally, in the implicit upper bound (*some*) conditions (5a,b), the third sentence responded to the question using the phrase *some of them*, which relies on a scalar inference to establish the "some but not all" interpretation. In the explicit upper bound (*only some*) conditions (5c,d), the response sentence uses the phrase *only some of them*, the "not all" interpretation of which is semantically specified without any scalar inference, and which does not rely on contextual support from the discourse cues. This yielded a 2 (CONTEXT: upper- vs. lower-bounded) × 2 (QUANTIFIER: *some* vs. *only some*) design.

An additional 128 vignettes were created to serve as fillers. Sixty-four introduced upper-bounded (*all*) questions under discussion but included answers with the quantifier *all* in the position where *some* or *only some* occurred in the critical items. The other 64 introduced lower-bounded (*any*) questions under discussion but included answers with the quantifier *none*. These fillers served to keep the critical quantifiers from being wholly predictable, and to reinforce the contrast between *some* and *all* in the upper-bounded conditions and the contrast between *some* and *none* in the lower-bounded conditions.

The stimuli were read aloud by a female native speaker of American English in an anechoic chamber at the University of Kansas, who was instructed to avoid placing contrastive stress on the quantifiers. For the critical items, the context sentences of the vignettes were read separately from the critical sentences.³ Filler items were read as complete vignettes. The recordings were digitised at 44100 Hz and later segmented and intensity-normalised using Praat (Boersma & Weenik, 2014), and the four combinations

1
2
3 of context and critical sentences were then spliced together to create the items. The
4
5 onset latencies of the quantifiers were measured by hand.
6
7
8

9 ***Offline questionnaire***
10

11 To evaluate the likelihood of making an upper-bounded interpretation of *some* in the
12 stimuli used in this experiment, a separate offline judgment task was conducted (after
13 the MEG experiment, and with different participants) in which participants answered,
14 for each vignette, whether or not an "all" interpretation was possible. For example,
15 participants heard vignette (6) and were then asked whether or not it was possible that
16 all of John's relatives were staying in his apartment.
17
18

- 19 (6) Mary was preparing to throw a party for John's relatives. She asked John
20 whether all of them were staying in his apartment. John said that some of them
21 were.
22
23

24 The question probed the "all" interpretation, rather than the "not all" interpretation (as
25 was done by Degen, *in press*), because the presence of an "all" interpretation guarantees
26 that a scalar inference *was not* computed, whereas the presence of a "not all"
27 interpretation does not guarantee that an inference *was* computed—the lower-bounded,
28 semantic interpretation of *some* may still be consistent with "not all".
29
30

31 The 128 critical items and 128 filler items from the MEG experiment were
32 divided into eight sub-experiments, with 16 critical and 16 filler items each. Within
33 each sub-experiment, the 16 critical items were organised into four versions in a Latin
34 square design, such that each version had four critical items per condition. This resulted
35 in 32 lists for the offline experiment. The lists were administered over the Internet using
36 Qualtrics, and each was completed by three participants, recruited via Amazon
37
38

1
2
3 Mechanical Turk. A given participant was allowed to complete multiple sub-
4 experiments, with the constraint that they could not complete more than one Latin-
5 square list within the same sub-experiment (which were clearly indicated as such). A
6 total of 28 unique participants completed the task, with each participant completing
7 between one and eight sub-experiments. Completion of each took about ten minutes.
8
9
10
11
12
13
14
15

MEG procedure

16
17 The vignettes were presented binaurally to participants over tube earphones (Aearo
18 Technologies), using the Presentation stimulus delivery software (Neurobehavioral
19 Systems). The participants' task was to listen to each vignette for comprehension. Each
20 trial began with the presentation of a fixation cross, which remained on the screen while
21 the vignette was playing. The vignette began playing 250-750 ms after the appearance
22 of the cross, and the cross remained on the screen for an additional 250-750 ms after the
23 sentence offset. Thirty-three percent of trials were followed by a comprehension
24 question presented visually on the screen. The comprehension questions probed various
25 portions of the vignettes, but on critical trials the comprehension questions never asked
26 questions that depended on the interpretation of the critical quantifier (e.g., questions
27 such as "How many of John's relatives were staying?" were not asked). Each question
28 was presented along with two possible answers, and the participant indicated her choice
29 using a button box placed below her left index and middle fingers. The next trial began
30 after the participant made her response (in the case of trials with comprehension
31 questions) or after the participant pressed either button (in the case of trials with no
32 comprehension question).

33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54 The stimuli were organised into four lists in a Latin square design. The item
55 order was randomised at runtime for each participant, with the restriction that the first
56
57
58
59
60

1
2 five trials were always fillers. The recording took about 50 minutes to complete.
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

Prior to the recording, the participant's head shape was digitised using a dual source handheld FastSCAN laser scanner (Polhemus, VT, USA). Three fiducial points and the position of five marker coils placed around the participant's face were also digitised. This was in order to localise the position of the participant's head in relation to the MEG sensors to allow for source reconstruction.

Data acquisition and preprocessing

MEG recordings used a whole-head 208 channel axial gradiometer system (Kanazawa Institute of Technology, Kanazawa, Japan) with a sampling frequency of 1000 Hz; data were filtered online with a low-pass of 200 Hz and a high-pass of 0.1 Hz. Data were epoched from -200 to 1200 ms, relative to the acoustic onset of the quantifier *some*.⁴ Artifact rejection consisted of manual inspection of the data to remove blinks, and amplitude cut-offs at 3000 fT. This resulted in removal of 15% of the trials, leaving an average of 26 trials per condition for each subject. Participants with fewer than 15 trials in any condition were excluded from further analysis. Preprocessed data were averaged for each condition and subject, using the prestimulus interval for baseline correction.

L2 minimum norm estimates of source activity were calculated using BESA 6.0 (MEGIS Software GmbH). Unsigned current estimates were created for each condition and low-pass filtered at 40 Hz. The source space was parcellated into Brodmann areas by converting sources to Talairach space using the Talairach Daemon (www.talairach.org) and matching each point to the nearest labelled Brodmann area in Talairach space. For each condition and each participant, the current estimate for each region was found by averaging the current estimates for all source vertices within that region. Source-space analysis was restricted to the following regions in both

hemispheres: ventromedial prefrontal cortex (BA 11) and temporal pole (BA 38), which have been implicated in basic phrasal composition in MEG (e.g., Bemis & Pylkkänen, 2011); temporo-parietal junction (BA 39), which has been implicated in speech act processing and theory of mind (e.g. Egorova, Shtyrov, & Pulvermüller, 2013); temporal regions implicated in lexical semantics and auditory processing (BAs 21, 22, and 42); and frontal regions (inferior frontal gyrus: BAs 44-45⁵; middle prefrontal cortex: BA 46; ventral orbitofrontal cortex: BA 47) implicated in aspects of structural and grammatical processing and recently implicated in scalar inference generation (Shetreet et al., 2014). These regions were chosen based on *a priori* predictions that they may be involved in some aspect of pragmatic processing or semantic composition.

Statistical analysis

For the sensor-level statistical analysis, we conducted non-parametric spatiotemporal clustering (Maris & Oostenveld, 2007). While the experiment used a 2×2 design (see Materials), the 2×2 test was computationally implemented using a series of t-tests on selected subsets of the data, which is conceptually comparable to running factorial analyses of variance when the factors only involve two levels (see, e.g., <http://mailman.science.ru.nl/pipermail/fieldtrip/2011-January/003447.html>). To test the main effect of CONTEXT, two datasets were created by averaging, respectively, the upper-bounded *some* and *only some* datasets, and the lower-bounded *some* and *only some* datasets; these two datasets were then compared using paired *t*-tests. The main effect of QUANTIFIER was tested in a similar way. The interaction between CONTEXT and QUANTIFIER (i.e., the difference of differences) was tested by creating two datasets representing the two context effects (subtracting the lower-bounded *some* from the upper-bounded *some* dataset, and the lower-bounded *only some* from the upper-bounded *only some* dataset), and comparing these context effects using paired *t*-tests.

Source-level statistics were performed by conducting nonparametric temporal clustering (Maris & Oostenveld, 2007) on the estimates from 200-1000 ms post stimulus onset within each region, and correcting for false discovery rate (Benjamini & Yekutieli, 2001) across regions. F-tests were used to measure the main effects of CONTEXT and QUANTIFIER, and a custom statistic was used to measure the crucial interaction. Nonparametric temporal clustering allows for the use of user-defined test statistics, and in this case we constructed a test statistic that quantifies the presence of a CONTEXT \times QUANTIFIER interaction such that there is a large effect of CONTEXT in either direction for *some* sentences, but no effect of CONTEXT for *only some* sentences; see equation (1). This statistic is large when such an interaction is present, and small otherwise.

$$|t_{(\text{upper-bound } \textit{some})} - t_{(\text{lower-bound } \textit{some})}| - |t_{(\text{upper-bound } \textit{only some})} - t_{(\text{lower-bound } \textit{only some})}| \quad (1)$$

For the largest cluster (any series of 10 or more contiguous timepoints with *p*-values⁶ of .3 or less) in each region, a cluster test statistic was computed by summing the sample test statistics (the F-tests of the main effects, or the statistic in (1) for the interaction) over all samples in the cluster, and then compared to 10,000 permutation test statistics each calculated by permuting the condition labels of the data and recalculating the test statistic. The proportion of permutation test statistics greater than the original cluster test statistic was the *p*-value for that region. False discovery rate adjustment was then applied to the *p*-values of the largest clusters of the regions.

Results

Offline questionnaire

On the filler items, participants correctly answered "yes" (i.e., indicated that an "all"

interpretation was possible) to 98.70% ($SD=3.39\%$) of *all*-quantifier fillers and only answered "yes" to 0.52% ($SD=3.33\%$) of *no*-quantifier fillers, indicating that participants were attentive to the final sentences of the vignettes. One of the critical items was incorrectly coded in the web version of the experiment and was removed from subsequent analyses. The mean percentages of "yes" responses (indicating upper-bounded or non-pragmatic readings) for the rest of the critical items are shown in the left portion of Figure 1 below. For statistical analysis, responses were treated as a categorical variable and were modelled with generalised linear mixed models, with fixed effects of QUANTIFIER, CONTEXT, and the QUANTIFIER \times CONTEXT interaction, and maximal random effects structures for SUBJECT, ITEM, and Latin square LIST.

--Figure 1 about here--

A slight numerical effect of Quantifier is evident, with *some* vignettes more likely than *only some* vignettes to be consistent with an "all" interpretation. This effect failed to reach significance, however, in the mixed model with maximal random effects ($\chi^2(1) = 0.08, p = .776$), although it did reach significance in a model with only random intercepts ($\chi^2(1) = 38.52, p < .001$) and in repeated-measures ANOVA ($F(1,27) = 5.19, p = .019$). The effect of Context was not significant in any analysis (maximal mixed model: $\chi^2(1) = 0.07, p = .787$; intercept-only mixed model: $\chi^2(1) = 1.03, p = .31$; ANOVA: $F(1,27) = 2.69, p = .113$), nor did the interaction (maximal mixed model: $\chi^2(1) = 0.19, p = .663$; intercept-only mixed model: $\chi^2(1) = 0.32, p = .571$; ANOVA: $F(1,27) = 0.44, p = .513$)

As is evident from Figure 1, the variance between participants was much greater for *some*, which is ambiguous, than for *only some*, which is unambiguous. This is due to a split in the *some* responses between pragmatic responders, who tend to interpret *some* as upper-bounded, and semantic responders, who do not; such group differences are

common in the experimental literature on scalar inferences (Noveck & Posada, 2003; Bott & Noveck, 2004; Hunt et al., 2013). The higher average rate of "all" interpretations for *some* is apparently driven by the "tail" of a few semantic responders. Histograms of the subject means for *some* and *only some* items (averaged across upper-bound and lower-bound contexts) are shown in the right-hand portion of Figure 1.

15 *Behavioural*

All participants kept in the electrophysiological analysis responded with greater than 80% accuracy on critical trials. Comprehension accuracy was 93.3% for SOME_all items (items with the quantifier *some*, in the upper-bounded *all* context), 86.2% for SOME_any items, 94.1% in ONLYSOME_all items, and 94.1% in ONLYSOME_any items. Accuracy was measured using generalised linear mixed models comprising fixed effects of CONTEXT, QUANTIFIER, the CONTEXT × QUANTIFIER interaction, and TRIAL NUMBER, and random effects (including both random intercepts and maximal random slope structure, per Barr, Levy, Scheepers & Tily, 2013) of PARTICIPANT, ITEM, and LIST. Model comparisons using log-likelihood tests showed no significant fixed effects (χ^2 s < 2.59, p s > .1).

42 *Sensor space*

Event-related fields evoked by *some* are shown in Supplementary Figure 1. The spatiotemporal clustering analysis in sensor space revealed no significant clusters for either the main effect of CONTEXT, the main effect of QUANTIFIER, or the CONTEXT × QUANTIFIER interaction.

55 *Source space*

One of the regions tested, left BA 46 (corresponding to the middle prefrontal cortex,

1
2
3 MPFC; $p_{\text{unc.}} = .002$, $p_{\text{FDR}} = .04$), showed an interaction in the 571-711 ms time window
4
5 that reached significance after FDR adjustment. Another region, left BAs 44-45
6
7 (corresponding to the inferior frontal gyrus; $p_{\text{unc.}} = .01$, $p_{\text{FDR}} = .1$) showed a very
8
9 marginal interaction in the 561-669 ms time window. The current estimates for these
10
11 regions are displayed in Figure 2. Current estimates for the rest of the regions tested are
12
13 shown in Supplementary Figure 2 in the online version of the article, and whole-brain
14
15 maps of the current estimates are shown in Supplementary Figure 3. The results of the
16
17 statistical analysis are shown in Table 1.
18
19

20 --Figure 2 about here--
21
22
23 --Table 1 about here--
24

25 In each of the significant regions, the FDR-significant interaction was due to
26 greater activity evoked by bare *some* in the lower-bounded context with weaker cues for
27 the inference than upper-bounded context with stronger cues (BA 46: $t(10) = 3.41$, 95%
28 CI = 1.02–4.87; BAs 44-45: $t(10) = 4.44$, 95% CI = 1.65–4.98), and no substantial
29 context effect evoked by *only some* (BA 46: $t(10) = -0.15$, 95% CI = -1.34–1.17; BAs
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 44-45: $t(10) = -0.66$, 95% CI = -1.55–0.84).

No regions showed FDR-significant main effects of either CONTEXT ($p_{\text{unc.}} > .018$, $p_{\text{FDR}} > .38$) or QUANTIFIER ($p_{\text{unc.}} > .125$, $p_{\text{FDR}} > .999$).

Discussion

The present study tested the neural-level computation of scalar inferences using the high temporal and spatial resolution of magnetoencephalography, and using a paradigm which allows for a direct comparison of the same word (*some*) occurring either in contexts that provide strong cues for incorporating a scalar implicature into its meaning, or contexts that do not. At the sensor level we observed no effect of scalar inferencing, consistent with Hartshorne and colleagues (in press), who only observed EEG effects of

1
2 inferencing well after the triggering expression. On the other hand, at the source level,
3
4 we found that the MPFC (BA 46) was selectively sensitive to scalar inferencing:
5
6 specifically, greater activation was evoked by *some* in context that provided *less* support
7
8 for the inference, compared to the contexts that provided more support. The location of
9
10 this activation is consistent with, although somewhat dorsal of, the region implicated in
11
12 scalar inferencing by the fMRI study of Shetreet and colleagues (2014); nonetheless, the
13
14 nature of the effect may be different.
15
16

17
18
19 ***Neural correlates of scalar inferencing***
20

21 The present study observed activation related to scalar inferences in the left MPFC (BA
22 46), as well as a similar trend that did not reach significance in left inferior frontal gyrus
23 (BAs 44–45). Of these two regions, the latter is spatially quite close to the activation that
24 Shetreet and colleagues (2014) observed in BA 47 (MNI coordinates: -36, 17, -22) for
25 the comparison between *some* sentences and *every* sentences (i.e., their "implicature
26 generation") comparison—the region the authors implicated in the realisation of scalar
27 inferences. On the other hand, the most significant region identified in the present study,
28 BA 46, is closer to activation Shetreet and colleagues observed in BA 10 (middle frontal
29 gyrus; MNI coordinates: -27, 44, -2) for the comparison between felicitous *some* (i.e.,
30 *some* describing pictures in which some but not all of the elements have the property
31 described) versus infelicitous *some* (i.e., *some* describing pictures in which all of the
32 elements have the property described), that is, the "implicature mismatch" comparison.
33 They argued that this region is involved in managing the conflict between the inference-
34 based meaning and the pragmatically mismatching context. In short, the present study
35 observed activation that partially overlaps with the LIFG activation in Shetreet et al.
36 (2014) but also extends to slightly more dorsal and anterior regions, and may overlap
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 with the region implicated in violation sentences as well as the region implicated in
3 basic inference realisation.
4

5
6 The present results offer converging evidence that these regions are indeed
7 involved in comprehending scalar inferences, rather than, for instance, verifying
8 different denotations or performing different kinds of quantification. These results also
9 extend the previous finding by showing that this activation occurs rapidly, within 600
10 ms after the realisation of the quantifier (adding to previous evidence that scalar
11 implicatures, rather than being postponed until the end of a sentence or proposition as
12 traditional Gricean accounts assume, are realised immediately and *in situ* during
13 incremental sentence processing; see Nieuwland et al., 2010; Politzer-Ahles et al., 2013,
14 Sikos et al., 2013; see Geurts, 2010, for discussion of how immediate incremental
15 computation of implicatures can be reconciled with a Gricean view; see also Chemla &
16 Spector, 2011, Chemla & Singh, 2014, Geurts & Pousoulous, 2009, among others, for
17 further discussion of local vs. global computation of scalar inferences). The fact that this
18 activation was strongest in MPFC is interesting because the lateral prefrontal cortex has
19 been suggested to play a role in working memory (e.g. Barbey, Koenigs, & Grafman,
20 2013, among others), and there is also some evidence that scalar inference realisation
21 may be modulated by working memory (Dieussaert et al., 2011; Feeney et al., 2004;
22 Marty & Chemla, 2013; Politzer-Ahles, Fiorentino, Durbin, & Li, 2014). This region,
23 however, is also implicated in many other aspects of cognition, such as planning and
24 cognitive control (Caplan, 2006), comprehension of syntactic structure (Hashimoto &
25 Sakai, 2002), and lexical retrieval (Chee, Hon, Caplan, Lee, & Goh, 2002); therefore,
26 the relationship between working memory and scalar inferencing at the neural level
27 requires further investigation before strong conclusions regarding this topic can be
28 made.
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 While the present experiment supports the conclusions of Shetreet and
3 colleagues (2014) that frontal areas are involved in scalar inferencing, there are some
4 important differences between the findings. Perhaps most noticeably, the present study
5 observed the greatest activation for a condition that provides fewer cues for an inference
6 (*some* in a context with a lower-bounded question under discussion), whereas Shetreet
7 and colleagues observed the greatest activation in an inference condition (*some*
8 sentences) compared to a no-inference condition (*every* sentences). The numerous
9 differences between the studies in terms of design, critical comparisons, and
10 methodologies used make it difficult to determine why the activation pattern is
11 reversed—for example, fMRI activation observed in these brain regions in the previous
12 study may reflect processes occurring at very different times than what was measured in
13 the present study. The following section will discuss in more detail the potential
14 implications of the direction of the effect observed in the present study. Another
15 important difference between the present study and that of Shetreet and colleagues is
16 that the previous study attributed two different loci of activation (inferior frontal and
17 prefrontal) to two different cognitive processes (realising inferences, or "implicature
18 generation", and comprehending infelicitous inferences, or "implicature mismatch"),
19 whereas the present experiment observed activation overlapping with both of these
20 areas in the same contrast. Again, the differences in contrasts and methods between
21 these two studies make it premature to draw strong conclusions; in particular, more
22 research into the neural correlates of implicature violations using techniques like MEG
23 is needed, as the present study did not test violations and previous electrophysiological
24 studies using violation paradigms (e.g., Hunt et al., 2013; Nieuwland et al., 2010;
25 Noveck & Posada, 2003; Panizza et al., 2014; Politzer-Ahles et al., 2013; Zhao et al.,
26 2015) did not perform source localisation. As the spatial resolution of fMRI is superior
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 to that of MEG, and the different regions identified in Shetreet et al. (2014) are close
3 together, care should be taken in interpreting apparent differences between the results of
4 the present study and of that study.
5
6

7 A final point to take note of is that in the present study, effects of scalar
8 inferencing were only observed in the source analysis, not in the sensor analysis. This is
9 consistent with an EEG study by Hartshorne and colleagues (in press), using a similar
10 design as the present study's, which did not observe an effect of scalar inferencing at the
11 quantifier, but confirmed through analysis of downstream words that the context did
12 indeed modulate comprehension of *some*. Thus, it appears that the neural correlates of
13 making scalar inferences that are observed at the source level may be difficult to detect
14 at the scalp level in EEG or MEG.
15
16

17 ***Implications for processing models***

18

19 Before addressing what the brain data can tell us about how scalar inferences are
20 realised, an important question is whether the participants in the experiment even did
21 realise them at all. While previous reading time evidence suggests that participants
22 differentially realise or fail to realise scalar inferences as a function of upper- versus
23 lower-bounded context with similar stimulus manipulations (Lewis, 2013; Politzer-
24 Ahles & Fiorentino, 2013) and intuition suggests that the inference is less available in
25 the nonsupporting lower-bounded context (Katsos & Cummins, 2010), our post-hoc
26 offline questionnaire found that these participants were about equally likely to make
27 scalar inferences in both upper- and lower-bound contexts with the present set of
28 materials . This could be because *some* is such a strong cue for the upper-bounded
29 interpretation, especially in an experiment where so many similar vignettes are
30 presented together that this interpretation is generated even when it is not relevant. That
31 indeed is the prediction of defaultist processing models (Levinson, 2000), but it could
32
33

also be modelled in Bayesian frameworks (Goodman & Stuhlmüller, 2013; Grodner & Russell, 2013), where in the lower-bounded context that makes the inference less relevant to the question under discussion, the probability of the speaker uttering *some* given the intended message "more than none" is still low. While our offline data have limitations (they are from a separate group of participants as the MEG experiment, and rather than using online measures of scalar inference realisation they only probe the eventual outcome of the inference), they do suggest that the context manipulation in the present experiment may have affected the ease of realising inferences, rather than the ultimate interpretation. In the future, experiments testing this question using MEG data from a downstream anaphorical expression (e.g., following the design of Hartshorne and colleagues, *in press*) would be valuable in further investigating this issue.⁷

As for the measures at the quantifier itself, traditional processing accounts of scalar inferencing make distinct predictions about the direction of context effects in behavioural measures and, by extension, neural activity. Context-driven processing accounts of scalar implicature, including that of Relevance theory (e.g. Noveck & Sperber, 2007) predict that realising scalar inferences always engenders a processing cost, and therefore there should be greater cognitive effort (which may be reflected by increased reading times, delayed eye movements, and possibly greater neural activity) when comprehending expressions that trigger an inference, compared to expressions that do not. This account does not necessarily predict that the amount of effort should be the same across contexts that do trigger inferences; it is feasible that more effort could be required in a context that provides less support for the inference. Such an account is thus consistent with the present results, although this pattern of results can be accounted for more explicitly by a constraint-based formulation of context-driven processing (see below).

Defaultist processing models (e.g., Levinson, 2000), on the other hand, suppose that scalar inferences of the sort tested here are realised automatically by default, and that realising the lower-bounded interpretation (i.e., realising that *some* may be consistent with *all* when its meaning is not enriched to include "and not all") is what takes more effort. Since our offline survey suggested that context did not affect the extent to which participants cancelled inferences and realised the lower-bounded interpretation of *some*, such models would have difficulty explaining why the lower-bounded condition in the present study engendered greater MEG activity: if inferences were realised automatically and with equal effort in both conditions, and did not require more cancellation in either condition, it is not clear what other mechanism the additional activation could be reflecting. Even if we assume that only the participants in the MEG study performed more inference cancellation in the lower-bounded than upper-bounded contexts, and that the participants in the offline questionnaire did not, there are conceptual challenges to attributing the observed MEG effect to cancellation. First of all, the prediction of greater processing cost in an inference-nonsupporting context only makes sense along with the assumption that the participant *must* actively cancel illicit inferences. In the paradigm used here, however, it is not clear why a participant would need to cancel even an illicit inference, as the inference does not conflict with anything else in the utterance, and the exchange can still be fully understood even with the unnecessary information added by the inference (see Politzer-Ahles & Fiorentino, 2013, for further discussion). Furthermore, some formal accounts of scalar implicature are not conducive to mechanisms for cancellation. For example, while early formal accounts included mechanisms for cancellation (e.g. Chierchia, 2004; Levinson, 2000), grammatical or lexicalist accounts of inferencing could also be formulated as involving competition between two readings (the enriched and non-

enriched) which are both realised before one is ultimately selected (e.g. Chemla & Singh, 2014); under such a view, it is not clear why a cancellation mechanism would be necessary, or why selecting one reading would be more computationally costly than selecting another reading.

While the psycholinguistic literature on scalar implicatures has traditionally been framed in terms of the competing context-driven and defaultist accounts, other formulations of scalar inferencing may better account for the present results. For instance, Degen and Tanenhaus (2011, 2014) have proposed a constraint-based account of inferencing in which inferences are always derived via the same mechanisms, but may appear context-based (i.e., slow and effortful) or default (i.e., rapid and easy) depending on the strength of numerous cues, such as context and prosody, that can interact to facilitate or inhibit the inference. The results of the present study could easily be explained under such an account: participants made the inferences in both contexts that strongly supported the inference and contexts that only weakly supported it (perhaps because the global experimental context introduced so many sentences of the type in (5)), but the inference was more difficult to make in the context with weak support, which provided fewer cues to facilitate the inference. While a computationally explicit account of how the facilitation occurs is beyond the scope of the present paper, one might imagine that the upper-bounded and lower-bounded items in the present experiment include many of the same cues (the word *some*, uttered with the same prosody and in the same experimental context—i.e., a context that did not include numerals and other lexical alternatives to *some*), but the upper-bounded items include an additional cue that the lower-bounded items do not: the explicit upper-bounded question under discussion (whether "all" was true) is a much stronger cue in favour of realising a scalar inference. Similarly, Huang and Snedeker (2011) propose that there

1
2
3 are both bottom-up and top-down routes to realising scalar inferences, and that
4
5 inferences may be realised rapidly and effortlessly via top-down mechanisms when, for
6 example, the information structure of the context makes the inference readily available
7 even before the triggering expression. Such an account would also predict greater
8 activation in the less supportive context, which might require a more bottom-up
9 mechanism to derive the inference, compared to the more supportive context which
10 explicitly introduces the quantifier *all* and thus allows the comprehender to verbally
11 pre-encode the alternatives (thus giving top-down support to the derivation of
12 alternative messages, which is a crucial part of making scalar inferences, see Chemla &
13 Singh, 2014).

14
15
16
17
18
19
20
21
22
23
24
25 The final alternative we will discuss is a lexical/grammatical ambiguity account,
26 whereby the enriched reading of *some* is not itself derived by an inference but is
27 underlyingly specified in either the lexical entry for *some* (which may have both the
28 lower-bounded and upper-bounded readings listed in the lexicon), or the grammatical
29 parse of the utterance (which may yield parses both with and without an
30 "Exhaustification" operator responsible for triggering the enriched meaning; see
31 Chierchia et al., 2012). Under such accounts, the role of pragmatics in scalar inference
32 is not to actually realise the enriched meaning, but rather to choose between the two
33 competing meanings. In this view, the increased activation observed in the lower-
34 bounded context may reflect greater ambiguity, as inferior frontal gyrus is in some
35 studies more strongly activated by ambiguous words than nonambiguous words, or by
36 subordinate as compared to dominant readings of ambiguous words (Bilenko, Grindrod,
37 Myers, & Blumstein, 2009; Fiebach, Vos, & Friderici, 2004; c.f. Copland et al., 2007).
38 Specifically, in the more supportive context, the salient upper-bounded question under
39 discussion strongly biases the hearer towards realising a "some but not all"
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2 interpretation; on the other hand, the less supportive context does not introduce this sort
3 of bias. Arguably, however, the lower-bounded context may bias the reader just as
4 strongly *against* realising a "some but not all" interpretation. An ambiguity-based
5 account of the present findings would require a more detailed, and perhaps quantitative,
6 account of how the questions under discussion create biases regarding the reading of
7 *some*, as well as an account of the locus of the ambiguity.
8
9

10 In short, further research is required to disentangle several competing accounts
11 with different explanations of the pattern of results observed in the present study. The
12 pattern of results is potentially consistent with gradient constraint-based processing and
13 with an ambiguity account of scalar inferences. One important conclusion that can be
14 made is that the neural findings suggest a more detailed processing account is needed
15 compared to the traditional defaultist vs. context-driven dichotomy that has been drawn
16 in most of the psycholinguistic literature until recently. For example, under the context-
17 driven view which has seen much support from psycholinguistic data, it would have
18 been easy to predict that there is an EEG/MEG component associated with making
19 scalar inferences, that such a component appears in sentences where inferences are
20 made and not sentences where inferences are not made, and that neurolinguistic
21 experiments simply need to identify when and where the component surfaces. The
22 present results suggest that the actual situation is much more complicated, and that the
23 search for scalar implicatures in the brain needs to be informed by more explicit models
24 of what specific operations (such as constraint evaluation or ambiguity resolution
25 strategies) feed into the derivation of inferences, as well as by gradient data on how
26 strongly a given sentence supports an inference rather than just all-or-nothing judgments
27 of whether it does or does not trigger an inference (as in Degen, in press). For future
28 research, it will likely not be enough to just compare inference-supporting and
29
30

1
2 inference-nonsupporting stimuli, but to also develop ways to test more specific
3 mechanisms such as those described above.
4
5
6
7
8

9 Conclusion

10

11 The present study, combining for the first time a design that is linguistically motivated
12 and controlled for all relevant variables with neurolinguistic techniques that provide
13 good temporal and spatial resolution, showed that the prefrontal cortex is involved in
14 the comprehension of scalar inferences. The specific pattern of activation shown, with
15 greater activity elicited in contexts that provide fewer cues to support a scalar inference
16 rather than contexts that provide more cues, is in line with both constraint-based and
17 lexical/grammatical ambiguity-based proposals for how scalar inferences are realised.
18
19 The data suggest that inference realisation is not a monolithic process, but rather that it
20 may involve multiple more fine-grained mechanisms or cues. The present results offer
21 the first brain-level evidence for an electrophysiological component associated with the
22 realisation of scalar inferences, while also highlighting the need for a more spelled-out
23 model of the specific cognitive operations that contribute to inferential processing.
24
25

26 Acknowledgements

27

28 This research was funded by grant G1001 from the NYUAD Institute, New York
29 University Abu Dhabi. We also thank Liina Pylkkänen, Diogo Almeida, and two
30 anonymous reviewers for feedback on previous versions of this paper; Estibaliz Blanco
31 Elorrieta for assistance in data collection; and Kelly Berkson, Natalie Pak, and Kate
32 Coughlin for assistance in stimulus development. The authors are responsible for any
33 errors in the manuscript.
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

References

- 7
- Barbey, A., Koenigs, M., & Grafman, J. (2013). Dorsolateral prefrontal contributions to human working memory. *Cortex*, 49, 1195-1205.
- 8
- Barr, D., Levy, R., Scheepers, C., & Tily, H. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*, 68, 255-278.
- 9
- Bemis, D., & Pylkkänen, L. (2011). Simple composition: a magnetoencephalography investigation into the comprehension of minimal linguistic phrases. *The Journal of Neuroscience*, 31, 2801-2814.
- 10
- Benjamini, Y., & Yekutieli, D. (2001). The control of the false discovery rate in multiple testing under dependency. *Annals of Statistics*, 29, 1165-1188.
- 11
- Bergen, L., & Grodner, D. J. (2012). Speaker knowledge influences the comprehension of pragmatic inferences. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 38, 1450.
- 12
- Bilenko, N., Grindrod, C., Myers, E., & Blumstein, S. (2009). Neural correlates of semantic competition during processing of ambiguous words. *Journal of Cognitive Neuroscience*, 21, 960-975.
- 13
- Boersma, P., & Weenink, D (2014). Praat: doing phonetics by computer [Computer program]. <http://www.praat.org/>.
- 14
- Bott, L., Bailey, T., & Grodner, D. (2012). Distinguishing speed from accuracy in scalar implicatures. *Journal of Memory and Language*, 66, 123-142.
- 15
- Bott, L., & Noveck, I. (2004). Some utterances are underinformative: The onset and time course of scalar inferences. *Journal of Memory and Language*, 51, 437-457.
- 16
- Breheny, R., Katsos, N., & Williams, J. (2006). Are generalised scalar implicatures generated by default? An on-line investigation into the role of context in generating pragmatic inferences. *Cognition*, 100, 434-463.
- 17
- Breheny, R., Ferguson, H., Katsos, N. (2012). Investigating the time-course of accessing conversational implicatures during incremental sentence interpretation. *Language and Cognitive Processes*, 28, 443-467.
- 18
- Breheny, R., Ferguson, H. J., & Katsos, N. (2013). Taking the epistemic step: toward a model of on-line access to conversational implicatures. *Cognition*, 126, 423-440.
- 19
- Caplan, D. (2006). Why is Broca's area involved in syntax? *Cortex*, 42, 46-471.
- 20
- Chee, M. W., Hon, N. H., Caplan, D., Lee, H. L., & Goh, J. (2002). Frequency of concrete words modulates prefrontal activation during semantic judgments. *Neuroimage*, 16, 259-268.
- 21

- 1
2 Chemla, E. & Singh, R. (2014). Remarks on the experimental turn in the study of scalar
3 implicature. *Language and Linguistics Compass*, 8, 373-386.
4
5 Chemla, E., & Spector, B. (2011). Experimental evidence for embedded scalar
6 implicatures. *Journal of Semantics*, 28, 359-400.
7
8 Chevallier, C., Bonnefond, M., Van der Henst, J-B, & Noveck, I (2010). Using ERPs to
9 capture inferential processes guided by prosodic cues. *Italian Journal of*
10 *Linguistics*, 22, 125-152.
11
12 Chevallier, C., Noveck, I., Nazir, T., Bott, L., Lanzetti, V., & Sperber, D. (2008).
13 Making disjunctions exclusive. *Quarterly Journal of Experimental Psychology*,
14 61, 1741-1760.
15
16 Chierchia, G. (2004). Scalar implicatures, polarity phenomena, and the
17 syntax/pragmatics interface. *Structures and Beyond*, 3, 39-103.
18
19 Chierchia, G., Fox, D., & Spector, B. (2012). Scalar implicature as a grammatical
20 phenomenon. In Maienborn, von Heusinger & Portner (Eds.), *Semantics: An*
21 *International Handbook of Natural Language Meaning*, Vol. 3 (pp. 2297-2331).
22 Berlin: Mouton de Gruyter.
23
24 Copland, D., Zubizaray, G., McMahon, K., & Eastburn, M. (2007). Neural correlates of
25 semantic priming for ambiguous words: an event-related fMRI study. *Brain*
26 *Research*, 1131, 163-172.
27
28 De Neys, W., & Schaeken, W. (2007). When people are more logical under cognitive
29 load. *Experimental Psychology*, 54, 128-133.
30
31 Degen, J. (in press). Investigating the distribution of *some* (but not *all*) implicatures
32 using corpora and web-based methods. *Semantics and Pragmatics*.
33
34 Degen, J. and Tanenhaus, M. (in press). Processing scalar implicature: a constraint-
35 based approach. *Cognitive Science*. doi: 10.1111/cogs.12171
36
37 Dieussaert, K., Verkerk, S., Gillard, E., & Schaeken, W. (2011). Some effort for some:
38 further evidence that scalar implicatures are effortful. *Quarterly Journal of*
39 *Experimental Psychology*, 64, 2352-2367.
40
41 Egorova, N., Shtyrov, Y., & Pulvermüller, F. (2013). Early and parallel processing of
42 pragmatic and semantic information in speech acts: neurophysiological
43 evidence. *Frontiers in Human Neuroscience*, 7.
44
45 Feeney, A., Scafton, S., Duckworth, A., & Handley, S. (2004). The story of some:
46 everyday pragmatic inferences by children and adults. *Canadian Journal of*
47 *Experimental Psychology*, 54, 128-133.
48
49 Fiebach, C., Vos, S., & Friderici, A. (2004). Neural correlates of syntactic ambiguity in
50 sentence comprehension for low and high span readers. *Journal of Cognitive*
51 *Neuroscience*, 16, 1562-1575.
52
53 Geurts, B. (2010). *Quantity Implicatures*. Cambridge: Cambridge University Press.
54
55
56
57
58
59
60

- 1
2
3 Geurts, B., & Pousoulous, N. (2009). Embedded implicatures?!? *Semantics and*
4 *Pragmatics*, 2, 4-1.
- 5
6 Goodman, N., & Stuhlmüller, A. (2013). Knowledge and implicature: modeling
7 language understanding as social cognition. *Topics in Cognitive Science*, 5, 173-
8 184.
- 9
10 Grodner, D., Klein, N., Carbury, K., & Tanenhaus, M. (2010). "Some," and possibly all,
11 scalar inferences are not delayed: Evidence for immediate pragmatic
12 enrichment. *Cognition*, 116, 42.
- 13
14 Grodner, D., & Russell, B. (2013). Evidence for a rational probabilistic account of
15 Gricean implicatures. *Poster presented at 26th CUNY Conference on Human*
16 *Sentence Processing*.
- 17
18 Hartshorne, J., & Snedeker, J. (submitted). The speed of inference: Evidence against
19 rapid use of context in calculation of scalar implicatures.
- 20
21 Hartshorne, J., Liem Azar, S., Snedeker, J., & Kim, A. (in press). The neural
22 computation of scalar implicature. *Language, Cognition and Neuroscience*.
- 23
24 Horn, L. (1972). *On the semantic properties of logical operators in English*. Ph.D.
25 dissertation, University of California, Los Angeles.
- 26
27 Huang, Y., & Snedeker, J. (2009). Online interpretation of scalar quantifiers: Insight
28 into the semantics-pragmatics interface. *Cognitive Psychology*, 58, 376-415.
- 29
30 Huang, Y., & Snedeker, J. (2011). Logic and conversation revisited: Evidence for a
31 division between semantic and pragmatic content in real-time language
32 comprehension. *Language and Cognitive Processes*, 26, 1161-1172.
- 33
34 Hunt, L., Politzer-Ahles, S., Gibson, L., Minai, U., & Fiorentino, R. (2013). Pragmatic
35 inferences modulate N400 during sentence comprehension: Evidence from
36 picture-sentence verification. *Neuroscience Letters*, 534, 246-251.
- 37
38 Katsos, N., & Cummins, C. (2010). Pragmatics: from theory to experiment and back
39 again. *Language and Linguistics Compass*, 4, 282-295.
- 40
41 Levinson, S. C. (2000). *Presumptive meanings: The theory of generalized*
42 *conversational implicature*. Cambridge: MIT press.
- 43
44 Lewis, S. (2013). *Pragmatic enrichment in language processing and development*. PhD
45 dissertation, University of Maryland.
- 46
47 Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and
48 MEG-data. *Journal of Neuroscience Methods*, 164, 177-190.
- 49
50 Marty, P., & Chemla, E. (2013). Scalar implicatures: working memory and a
51 comparison with only. *Frontiers in Psychology*, 4.
- 52
53 Marty, P., Chemla, E., & Spector, B. (2013). Interpreting numerals and scalar items
54 under memory load. *Lingua*, 133, 152-163.
- 55
56 Minai, U., & Fiorentino, R. (2010). The role of the focus operator only in children's
57 computation of sentence meaning. *Language Acquisition*, 17, 183-190.
- 58
59
60

- Nieuwland, M. S., Ditman, T., & Kuperberg, G. R. (2010). On the incrementality of pragmatic processing: An ERP investigation of informativeness and pragmatic abilities. *Journal of Memory and Language*, 63, 324-346.
- Noveck, I., & Reboul, A. (2008). Experimental pragmatics: a Gricean turn in the study of language. *Trends in Cognitive Sciences*, 12, 425-431.
- Noveck, I., & Posada, A. (2003). Characterizing the time course of an implicature: An evoked potentials study. *Brain and Language*, 85, 203-210.
- Noveck, I., & Sperber, D. (2007). The why and how of experimental pragmatics: The case of 'scalar inferences'. In N. Burton-Roberts (Ed.) *Advances in Pragmatics*. Basingstoke: Palgrave.
- Oldfield, R. C. (1971). The assessment and analysis of handedness: the Edinburgh inventory. *Neuropsychologia*, 9(1), 97-113.
- Panizza, D., & Onea, E. (2014). Some implicatures take their time: an ERP study on scalar implicatures with 'sentence-picture vs. picture-sentence' verification task. Talk presented at 27th CUNY Conference on Human Sentence Processing. Columbus, OH.
- Politzer-Ahles, S., & Fiorentino, R. (2013). The realization of scalar inferences: context sensitivity without processing cost. *PLoS ONE*, 8(5), e63943.
- Politzer-Ahles, S., Fiorentino, R., Durbin, J., & Li, L. (2014). The role of working memory in the online realization of scalar inferences. Poster presented at 27th CUNY Conference on Human Sentence Processing. Columbus, OH.
- Politzer-Ahles, S., Fiorentino, R., Jiang, X., & Zhou, X. (2013). Distinct neural correlates for pragmatic and semantic meaning processing: An event-related potential investigation of scalar implicature processing using picture-sentence verification. *Brain Research*, 1490, 134-152.
- Sauerland, U. (2012). The Computation of scalar implicatures: Pragmatic, lexical or grammatical? *Language and Linguistics Compass*, 6, 36-49.
- Shetreet, E., Chierchia, G., & Gaab, N. (2014). When some is not every: Dissociating scalar implicature generation and mismatch. *Human Brain Mapping*, 35, 1503-1514.
- Sikos, L., Tomlinson, S., Traut, H., & Grodner, D. (2013). Incremental computation of scalar implicatures: An ERP study. Poster presented at 26th CUNY Conference on Human Sentence Processing. Columbia, SC.
- Sperber, D., & Wilson, C. (1995). *Relevance: Communication and Cognition*. Oxford: Blackwell.
- Yekutieli, D., & Benjamini, Y. (1999). Resampling-based false discovery rate controlling multiple test procedures for correlated test statistics. *Journal of Statistical Planning and Inference*, 82, 171-196.

- 1
2
3 Zhao, M., Liu, T., Chen, G., & Chen, F. (2015). Are scalar implicatures automatically
4 processed and different for each individual? A mismatch negativity (MMN)
5 study. *Brain Research*, 1599, 137-149.
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

For Peer Review Only

Table 1. Results of the statistical analysis of current estimates. Each row gives uncorrected p -value (controlled for multiple comparisons over time, but not for false discovery rate across regions), FDR-adjusted p -value (Yekutieli & Benjamini, 1999), and time window (begin and end times in ms) of the largest interaction cluster in that region. See Methods for a description of the test statistic used to quantify the interaction for the purpose of defining clusters. The final column lists the MNI coordinates for the center of the region.

ROI	$p_{\text{unc.}}$	p_{FDR}	Cluster times	Coordinates (x, y, z)
Left SFG (BA 10)	.479	.798	[933 951]	-21, 60, 8
Left MeFG (BA 11)	.404	.808	[927 960]	-13, 38, -16
Left MTG (BA 21)	>.99	>.99	[161 182]	-62, -18, -11
Left STG (BA 22)	.235	.783	[800 839]	-63, -27, 8
Left ATL (BA 38)	.567	.872	[118 141]	-36, -11, -32
Left TPJ (BA 39)	.737	>.99	[816 835]	-51, -67, 25
Left TTG (BA 42)	.067	.335	[773 851]	-66, -24, 11
Left IFG (BA 44-45)	.01	.1	[563 668]	-55, 24, 13
Left MPFC (BA 46)	.002	.040	[572 709]	-48, 40, 18
Left vOFC (BA 47)	.212	.848	[116 143]	-37, 27, -15
Right SFG (BA 10)	.807	>.99	[493 523]	25, 60, 8
Right MeFG (BA 11)	.440	.800	[972 1000]	16, 36, -17

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

Right MTG (BA 21)	.241	.689	[808 839]	62, -22, -9
Right STG (BA 22)	.271	.602	[161 207]	63, -33, 10
Right ATL (BA 38)	.035	.233	[769 902]	35, 10, -32
Right TPJ (BA 39)	.902	>.99	[965 982]	51, -66, 28
Right TTG (BA 42)	.254	.635	[503 571]	65, -20, 11
Right IFG (BA 44-45)	>.99	>.99	[110 250]	54, -26, 6
Right MPFC (BA 46)	>.99	>.99	[110 343]	49, -40, 17
Right vOFC (BA 47)	>.99	>.99	[137 197]	40, 27, -13

Figure captions

Figure 1. Left: mean percentages of "yes" responses (indicating lower-bounded or non-pragmatic readings); error bars represent $\pm 2 \times \text{SE}$ (the standard error of the by-subject means). Right: histograms of the subject means for *some* and *only some* items (averaged across upper-bounded and lower-bounded contexts)

Figure 2. Current estimates for regions showing FDR-significant or -marginal interactions. The time window of the significant cluster is highlighted in gray, and the upper left portion of each plot shows the spatial location of the vertices comprising that region. SOME_all: *some* in the "all" (upper-bounded, inference-supporting) context; SOME_any: *some* in the "any" (lower-bounded, inference-nonsupporting) context; ONLYSOME_all: *only some* in the "all" (upper-bounded) context; ONLYSOME_any: *only some* in the "any" (lower-bounded) context.

1
2
3
4
5
6
7

¹ The authors did include baseline conjunction analyses to argue against this possibility. These analyses involved comparing activations for *some*-sentences in various conditions to activations for correct *every*-sentences (that is, those accompanying pictures in which every element has the property described), and showed that processing the lexical meaning of *some* either did not elicit greater activation than the lexical meaning of *every*, or only did so in other brain regions. This analysis may not have ruled out potential differences due to lexical meanings, however, because it included *some*-NONE items—that is, *some*-sentences accompanying pictures in which none of the elements have the property described. For example, if verification of *some but not all* followed a two-step procedure (first verify whether any Xs have the property, then check whether all the Xs have the property), then *some*-NONE items would not elicit the second step, and therefore might not have shown greater verification costs than *every*-sentences. Crucially, *some*-NONE sentences were included in this baseline analysis but not in the critical analysis, which may have caused lexical differences to emerge in only the critical analysis.

² An anonymous reviewer notes that "only some" is somewhat infelicitous in the lower-bounded context (i.e., "She asked John whether any of them were staying his apartment. John said that only some of them were"). While this may indeed introduce additional differences in activity between the two "only some" conditions, the statistical analysis conducted in the present study (see Methods) used a test statistic designed to identify brain regions that showed an effect of Context in the "some" sentences and not the "only some" sentences. Therefore, while this aspect of the control items may have caused our analysis to miss some regions that could have been involved in scalar inferencing, it would not have caused us to falsely detect any regions.

³ That is to say, the following four sentences corresponding to the possible versions of the item shown in (5) were each read during different blocks of the recording, with multiple other sentences in between:

- Mary was preparing to throw a party for John's relatives. She asked John whether all of them were staying in his apartment.
- Mary was preparing to throw a party for John's relatives. She asked John whether any of them were staying in his apartment.
- John said that some of them were.
- John said that only some of them were.

⁴ The data also showed the same pattern of results when the *only some* conditions were time-locked to the onset of *only* rather than the onset of *some*.

⁵ The small sizes of BAs 44 and 45, along with the coarse-grained nature of the decimation of the BESA cortex into sources, meant that each of these BAs contained a very small number of source vertices. Therefore, in the parcelation used for our analysis, these two BAs are combined.

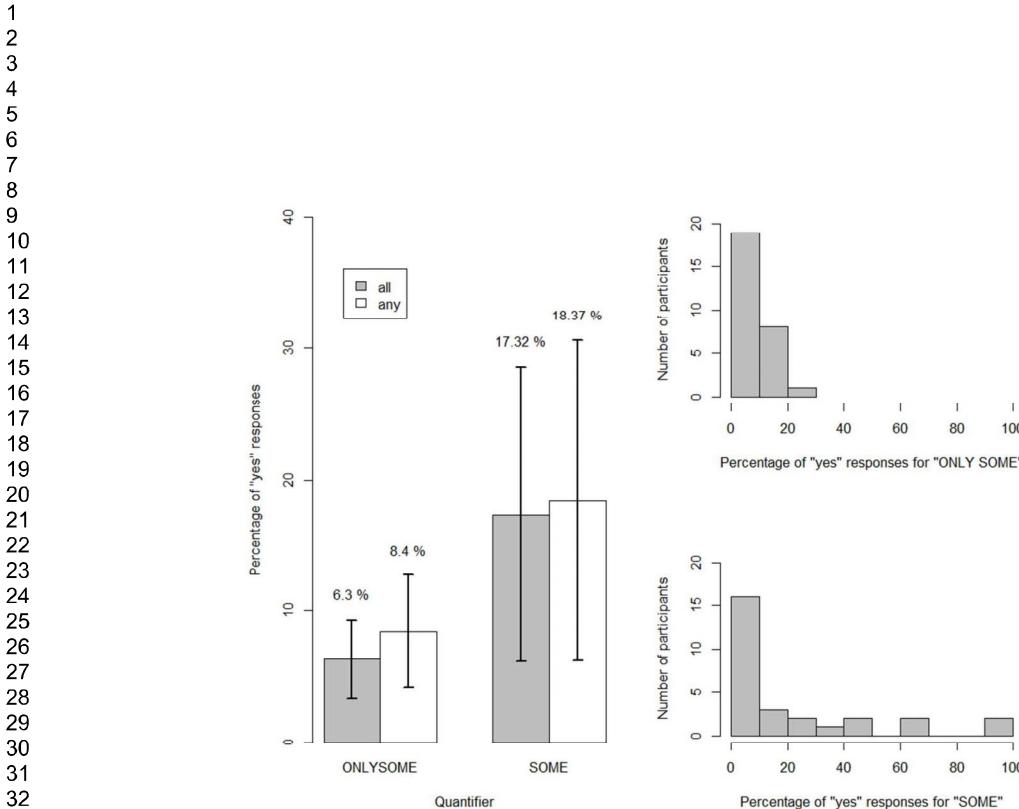
⁶ Because the test statistic used for identifying clusters was not a true *t*-statistic, but rather a difference of *t*-statistics, it does not necessarily fit a *t* distribution, and thus *p*-values based on the *t* distribution are not accurate. In other words, $p < .3$ in the present context does not necessarily mean a test-wise Type I error rate of less than .3; rather, it is a somewhat arbitrary value. This does not matter for a permutation test, however, as the threshold used for creating clusters does not affect the familywise error rate of the overall analysis (see Maris & Oostenveld, 2007).

⁷ An additional question raised by the offline data is as follows: why is it the case that both the MEG data at the quantifier and the offline data at the end of the sentence suggest that that inference is realised in both contexts, whereas reading time and EEG data at a downstream expression (e.g. Breheny et al., 2006; Bergen & Grodner, 2012; Hartshorne et al., in press; Lewis, 2013; Politzer-Ahles & Fiorentino, 2013) suggest that the inference is more likely in

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

the upper-bounded context? It is difficult to compare these findings, given that they come from different measures, different samples of participants, and different experimental designs. One possible explanation is that, even if participants realise the inference in both contexts, only in the upper-bounded context do they commit to the inference so fully as to make a strong forward prediction that "the rest" will appear later in the sentence. This explanation is ad-hoc, however, and this topic requires further investigation.

For Peer Review Only



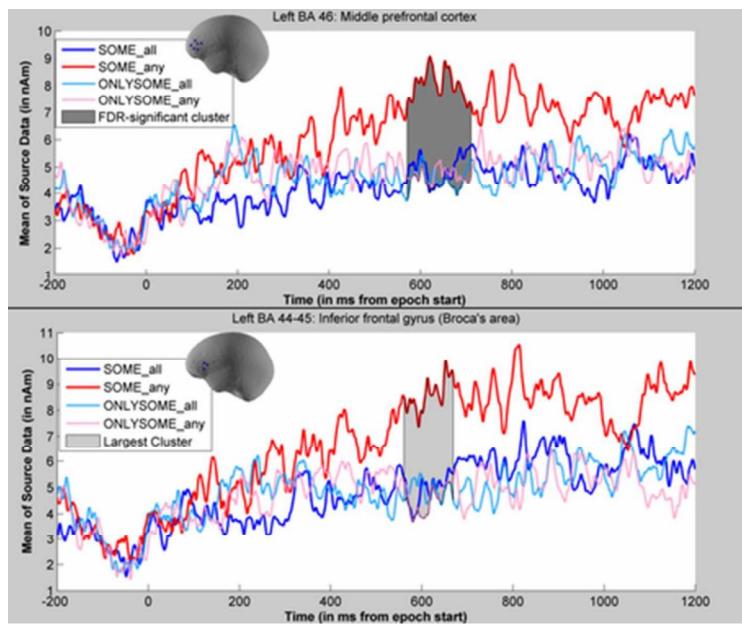


Figure 2. Current estimates for regions showing FDR-significant or -marginal interactions. The temporal window of the significant cluster is highlighted in gray, and the upper left portion of each plot shows the spatial location of the vertices comprising that region. SOME_all: some in the "all" (upper-bounded, inference-supporting) context; SOME_any: some in the "any" (lower-bounded, inference-nonsupporting) context; ONLYSOME_all: only some in the "all" (upper-bounded) context; ONLYSOME_any: only some in the "any" (lower-bounded) context.

19x15mm (600 x 600 DPI)

1
2
3
4
5
6

Supplementary figure captions

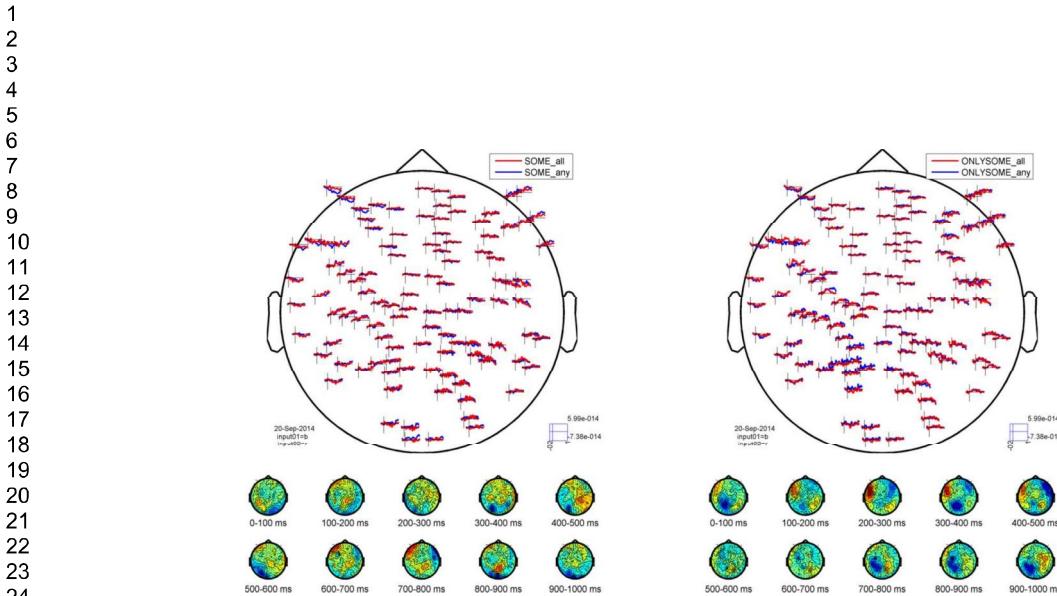
7 Supplementary figure 1. Top: Waveforms showing the event-related fields evoked by the critical
8 word. SOME_all: *some* in the "all" (upper-bounded, inference-supporting) context; SOME_any:
9 some in the "any" (lower-bounded, inference-nonsupporting) context; ONLYSOME_all: *only*
10 *some* in the "all" (upper-bounded) context; ONLYSOME_any: *only some* in the "any" (lower-
11 bounded) context. Waveforms represent a subset of the channels in the sensor space. Bottom:
12 Topographic plots of the difference waves for the context effect at *some* (SOME_all –
13 SOME_any) and at *only some* (ONLYSOME_all – ONLYSOME_any).

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

Supplementary figure 2. Current estimates for all regions tested (except BA 46 and BAs 44–45, which are shown in Figure 1). For each region, the temporal window of the largest cluster is highlighted in gray, and the upper left portion of each plot shows the spatial location of the vertices comprising that region. SOME_all: *some* in the "all" (upper-bounded, inference-supporting) context; SOME_any: *some* in the "any" (lower-bounded, inference-nonsupporting) context; ONLYSOME_all: *only some* in the "all" (upper-bounded) context; ONLYSOME_any: *only some* in the "any" (lower-bounded) context.

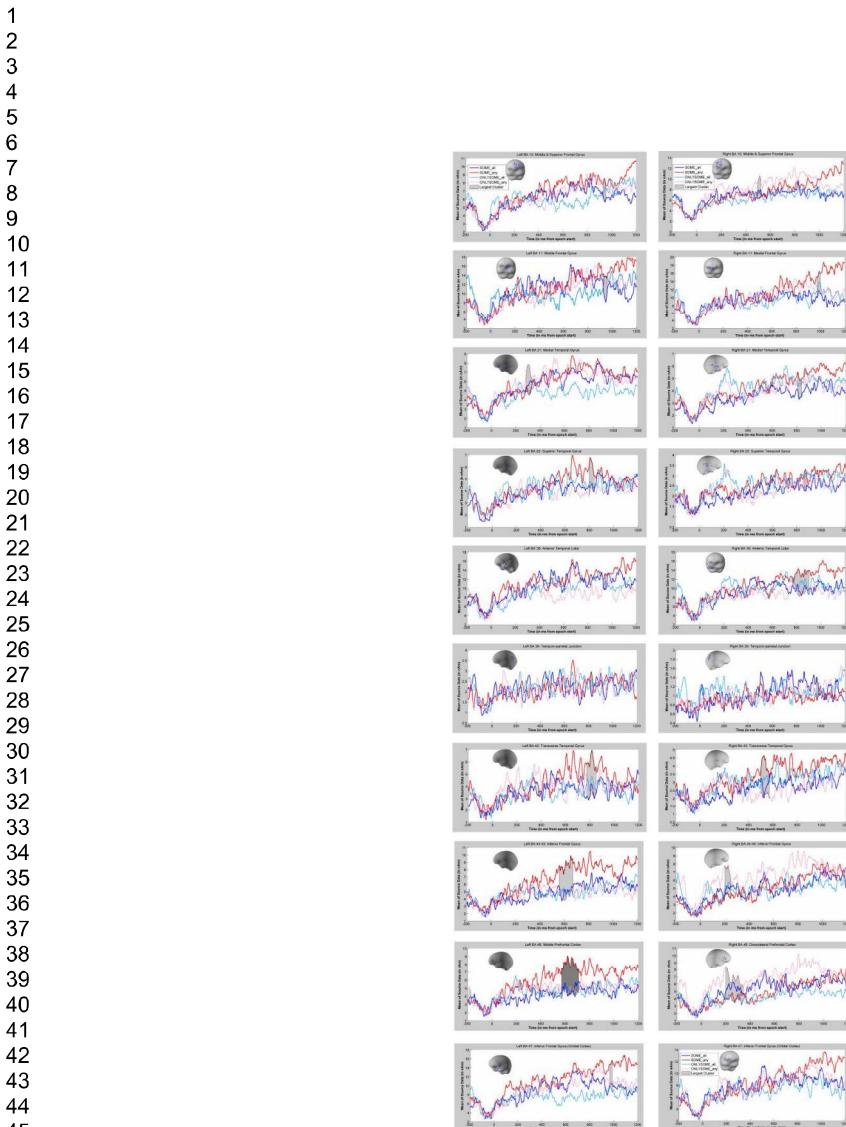
Supplementary figure 3. Maps of current estimates overlain on the whole brain. The difference between SOME_all and SOME_any is shown in the upper portion (A), and the difference between ONLYSOME_all and ONLYSOME_any in the lower portion (B). For each contrast, the plots are masked at $p = .01$ (paired t -tests). Each row shows a different view (left, right, front, back), and each column shows a different time window.

Supplementary file 1. Description of the post-hoc offline survey carried out on Mechanical Turk.



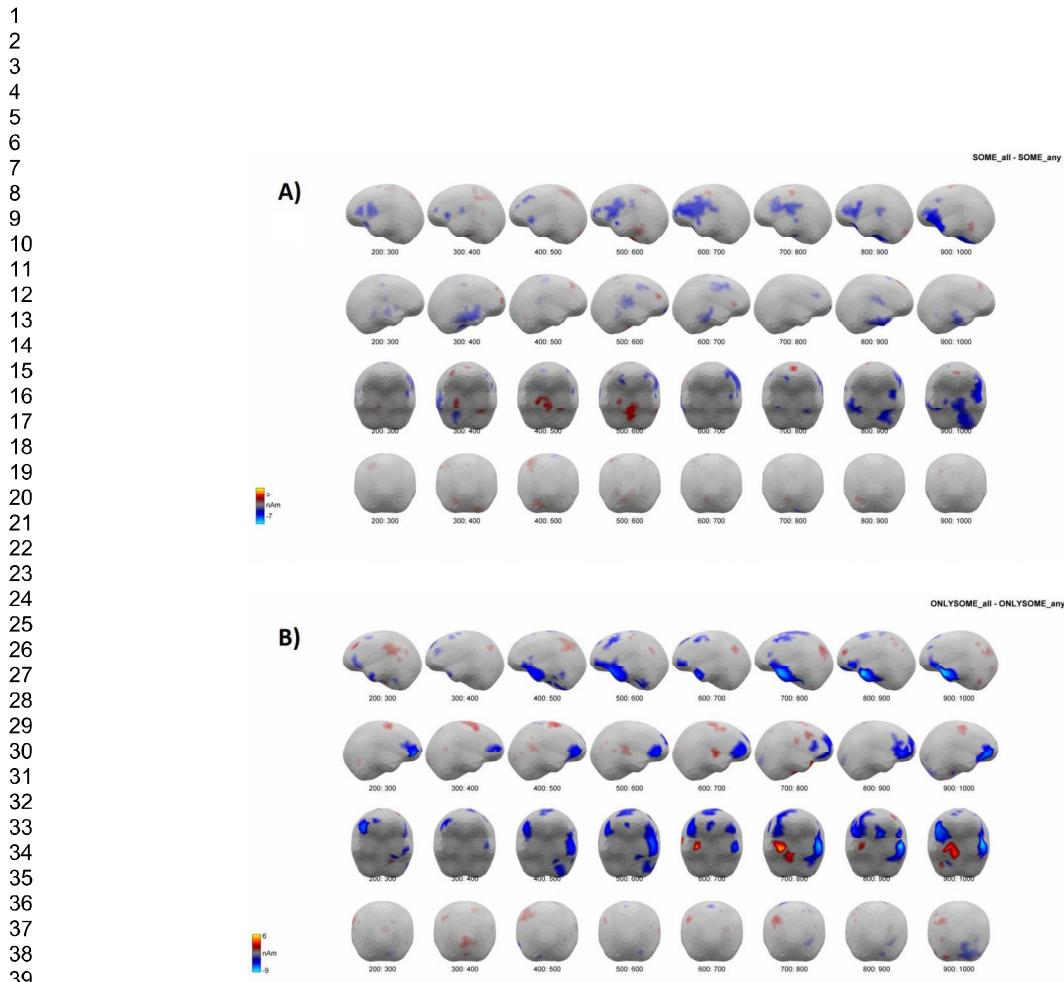
Waveforms showing the event-related fields evoked by the critical word. **SOME_all**: some in the "all" (upper-bounded, inference-supporting) context; **SOME_any**: some in the "any" (lower-bounded, inference-nonsupporting) context; **ONLYSOME_all**: only some in the "all" (upper-bounded) context; **ONLYSOME_any**: only some in the "any" (lower-bounded) context. Waveforms represent a subset of the channels in the sensor space. Bottom: Topographic plots of the difference waves for the context effect at some (**SOME_all** - **SOME_any**) and at only some (**ONLYSOME_all** - **ONLYSOME_any**).

322x174mm (116 x 116 DPI)



Supplementary Figure 2. Current estimates for all regions tested (except BA 46 and BA 44-45, which are shown in Figure 1). For each region, the temporal window of the largest cluster is highlighted in gray, and the upper left portion of each plot shows the spatial location of the vertices comprising that region.
SOME_all: some in the "all" (upper-bounded, inference-supporting) context; SOME_any: some in the "any" (lower-bounded, inference-nonsupporting) context; ONLYSOME_all: only some in the "all" (upper-bounded) context; ONLYSOME_any: only some in the "any" (lower-bounded) context,

150x369mm (300 x 300 DPI)



Supplementary Figure 3. Maps of current estimates overlain on the whole brain. The difference between SOME_all and SOME_any is shown in the upper portion (A), and the difference between ONLYSOME_all and ONLYSOME_any in the lower portion (B). For each contrast, the plots are masked at $p = .01$ (paired t-tests). Each row shows a different view (left, right, front, back), and each column shows a different time window.
402x423mm (120 x 120 DPI)