

1 **Tracking the behavioral and neural dynamics of semantic representations through negation**

2
3 Arianna Zuanazzi^{1*}, Pablo Ripollés^{1,2,3}, Wy Ming Lin⁴, Laura Gwilliams⁵, Jean-Rémi King^{1,6†},
4 David Poeppel^{1,3,7†}

5
6 1 Department of Psychology, New York University, New York, NY, USA.

7 2 Music and Audio Research Lab (MARL), New York University, New York, NY, USA.

8 3 Center for Language, Music and Emotion (CLaME), New York University, New York, NY, USA.

9 4 Hector Research Institute for Education Sciences and Psychology, University of Tübingen,
10 Tübingen, Germany.

11 5 Department of Neurological Surgery, University of California, San Francisco, CA, USA.

12 6 Ecole Normale Supérieure, PSL University, Paris, France.

13 7 Ernst Strüngmann Institute for Neuroscience, Frankfurt, Germany.

14
15 †These authors contributed equally to this work.

16 *Corresponding author. Email: az1864@nyu.edu.

17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

35 **Abstract**

36 Combinatoric linguistic operations underpin human language processes, but how meaning changes
37 over time is not well understood. We address this puzzle by exploiting the ubiquitous function of
38 negation. We track the online effects of negation (“not”) and intensifiers (“really”) on the
39 representation of scalar adjectives (e.g., “good”) in parametrically designed behavioral and
40 neurophysiological (MEG) experiments. The behavioral data show that participants first interpret
41 negated adjectives as affirmative and then modify their interpretation towards, but never exactly as,
42 the opposite meaning. Decoding analyses of neural activity further reveal that negation does not
43 *invert* the representation of adjectives (i.e., “not bad” represented as “good”) but rather *mitigates*
44 their representation, at early lexical-semantic processing stages. This putative suppression
45 mechanism of negation is supported by increased synchronization of beta-band neural activity in
46 sensorimotor areas. The analysis of negation provides a steppingstone to understand how the human
47 brain represents changes of meaning over time.

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69 **Introduction**

70 A hallmark of language processing is that we combine elements of the stored inventory - informally
71 speaking, words - and thereby flexibly generate new meanings or change current meanings. The
72 final representations derive in systematic ways from the combination of individual pieces. The
73 composed meanings can relate in relatively straightforward ways to the building blocks (e.g., “the
74 cat sat on the mat”) or stem from more subtle inferential processes (e.g., “this theory is not even
75 wrong”). A mechanistic understanding of the underlying processes requires characterization of how
76 meaning representations are constructed in real time. There has been steady progress and productive
77 debate on syntactic structure building¹⁻⁶. In contrast, how novel semantic configurations are
78 represented over time is less widely investigated. In the experimental approach pursued here, we
79 build on the existing literature on precisely controlled *minimal* linguistic environments^{7,8}. We
80 deploy a new, simple parametric experimental paradigm that capitalizes on the powerful role that
81 *negation* plays in shaping semantic representations of words. While negation is undoubtedly a
82 complex linguistic operation that can affect comprehension as a function of other linguistic factors
83 (such as discourse and pragmatics⁹⁻¹¹), our investigation specifically focuses on how negation
84 operates in phrasal structures. Combining behavioral and neurophysiological data, we show how
85 word meaning is (and is not) modulated in controlled contexts that contrast affirmative (e.g., “really
86 good”) and negated (e.g., “not good”) phrases. The results identify models and mechanisms of how
87 negation, a compelling window into semantic representation, operates in real time.

88 Negation is ubiquitous – and therefore interesting in its own right. Furthermore, it offers a
89 compelling linguistic framework to understand how the human brain builds meaning through
90 combinatoric processes. Intuitively, negated concepts (e.g., “not good”) entertain some relation
91 with the affirmative concept (e.g., “good”) as well as their counterpart (e.g., “bad”). The function
92 of negation in natural language has been a matter of longstanding debate among philosophers,
93 psychologists, logicians, and linguists¹². In spite of its intellectual history and relevance
94 (interpreting negation was, famously, a point of debate between Bertrand Russell and Ludwig
95 Wittgenstein), comparatively little research investigates the cognitive and neural mechanisms
96 underpinning negation. Previous work shows that negated phrases/sentences are processed with
97 more difficulty (slower, with more errors) than the affirmative counterparts, suggesting an
98 asymmetry between negated and affirmative representations; furthermore, state-of-the-art artificial
99 neural networks appear to be largely insensitive to the contextual impacts of negation¹³⁻²⁰. This
100 asymmetry motivates one fundamental question: *how* does negation operate?

101 Studies addressing this question suggest that negation operates as a suppression mechanism
102 by reducing the extent of available information²¹⁻²³, either in two steps^{18,24-28} or in one incremental

103 step^{12,29–31}; other studies demonstrate that negation is rapidly and dynamically integrated into
104 meaning representations^{10,32}, even unconsciously³³. Within the context of action representation
105 (e.g., “cut”, “wish”), previous research suggests that negation recruits general-purpose inhibitory
106 and cognitive control systems^{34–41}.

107 While the majority of neuroimaging studies focused on how negation affects action
108 representation, psycholinguistic research shows that scalar adjectives (e.g., “bad-good”, “close-
109 open”, “empty-full”) offer insight into how negation operates on semantic representations of single
110 words. These studies provide behavioral evidence that negation can either *eliminate* the negated
111 concept and convey the opposite meaning (“not good” = “bad”) or *mitigate* the meaning of its
112 antonym along a semantic continuum (“not good” = “less good”, “average”, or “somehow bad”;
113 ^{11,12,42–44}). Thus, the system of polar opposites generated by scalar adjectives provides an especially
114 useful testbed to investigate changes in representation of abstract concepts along a semantic scale
115 (e.g., “bad” to “good”), as a function of negation (e.g., “bad” vs. “not good”).

116 Here, we capitalize on the semantic continuum offered by scalar adjectives to investigate
117 *how* negation operates on the representation of abstract concepts (e.g., “bad” vs. “good”). First, we
118 track how negation affects semantic representations over time in a behavioral study. Next, we use
119 magnetoencephalography (MEG) and a decoding approach to track the evolution of neural
120 representations of target adjectives in affirmative and negated phrases. We test four hypotheses: (1)
121 negation does not change the representation of adjectives (e.g., “not good” = “good”), (2) negation
122 weakens the representation of adjectives (e.g., “not good” < “good”), (3) negation inverts the
123 representation of adjectives (e.g., “not good” = “bad”), and (4) negation changes the representation
124 of adjectives to another representation (e.g., “not good” = e.g., “unacceptable”). The combined
125 behavioral and neurophysiological data adjudicate among these hypotheses and identify potential
126 mechanisms that underlie how negation functions in online meaning construction.

127

128

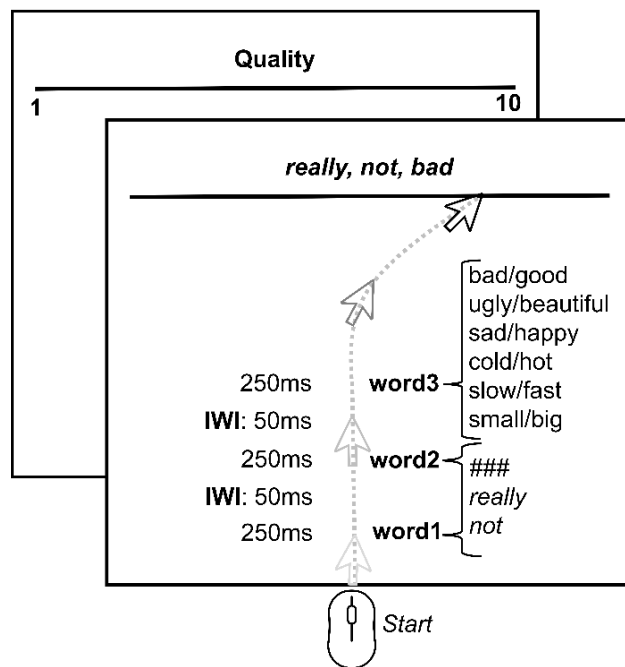
129 **Results**

130 *Experiment 1: Continuous mouse tracking reveals a two-stage representation of negated* 131 *adjectives*

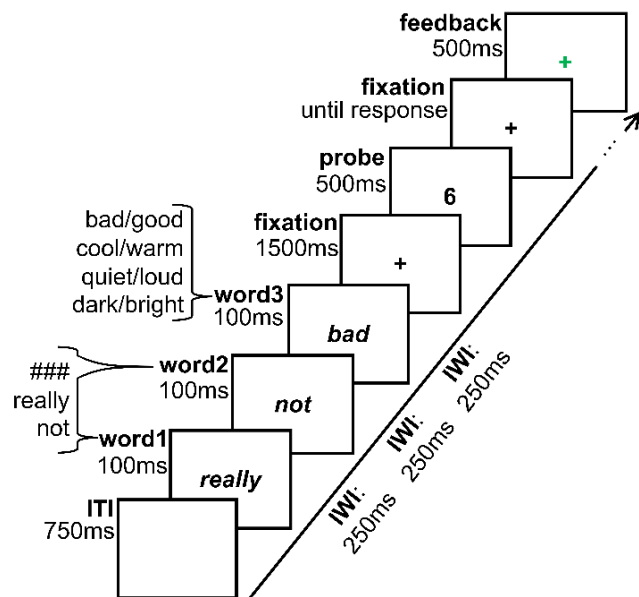
132 Experiment 1 (online behavioral experiment; N = 78) aimed to track changes in representation over
133 time of scalar adjectives in affirmative and negated phrases. Participants read two-to-three-word
134 phrases comprising one or two modifiers (“not” and “really”) and a scalar adjective (e.g., “really
135 really good”, “really not quiet”, “not ### fast”). The number and position of modifiers were
136 manipulated to allow for a characterization of negation in simple and complex phrasal contexts,

137 above and beyond single word processing. Adjectives were selected to represent opposite poles
 138 (i.e., antonyms) of the respective semantic scales: *low* pole of the scale (e.g., “bad”, “ugly”, “sad”,
 139 “cold”, “slow”, and “small”) and *high* pole of the scale (e.g., “good”, “beautiful”, “happy”, “hot”,
 140 “fast”, and “big”). A sequence of dashes was used to indicate the absence of a modifier. **Fig. 1A**
 141 and **Table S1** provide a comprehensive list of the linguistic stimuli. On every trial, participants
 142 rated the overall meaning of each phrase on a scale defined by each antonym pair (**Fig. 1A**). We
 143 analyzed reaction times and continuous mouse trajectories, which consist of the positions of the
 144 participant’s mouse cursor while rating the phrase meaning. Continuous mouse trajectories offer
 145 the opportunity to measure the unfolding of word and phrase comprehension over time, thus
 146 providing time-resolved dynamic data that reflect changes in meaning representation^{15,45,46}.

A. Behavioral experiment: mouse trajectories



B. MEG experiment: behavioral task



147 **Figure 1. Experimental procedures.**

148 (A) Behavioral procedure. Participants read affirmative or negated adjective phrases (e.g., “really really good”, “###
 149 not bad”) word by word and rated the overall meaning of each phrase on a scale. Each trial consisted of combinations
 150 of “###”, “really”, and “not” in word positions 1 and 2, followed by an adjective representing the low or high pole
 151 across six possible scalar dimensions. Before each trial, participants were informed about the scale direction, e.g., “bad”
 152 to “good”, i.e., 1 to 10. Scale direction was pseudorandomized across blocks. For each trial, we collected continuous
 153 mouse trajectories throughout the entire trial as well as reaction times. (B) MEG procedure. Participants read
 154 affirmative or negated adjective phrases and were instructed to derive the overall meaning of each adjective phrase on
 155 a scale from 0 to 8, e.g., from “really really bad” to “really really good”. After each phrase, a probe (e.g., 6) was

156 presented, and participants were required to indicate whether the probe number correctly represented the overall
157 meaning of the phrase on the scale (*yes/no* answer, using a keypad). Feedback was provided at the end of each trial
158 (green or red cross). While performing the task, participants lay supine in a magnetically shielded room while
159 continuous MEG data were recorded through a 157-channel whole-head axial gradiometer system. Panels A and B:
160 “####” = no modifier; IWV = inter-word-interval.

161

162

163 *Reaction times.* To evaluate the effect of antonyms and of negation on reaction times in behavioral
164 Experiment 1, we performed a 2 (*antonym*: low vs. high) x 2 (*negation*: negated vs. affirmative)
165 repeated-measures ANOVA. The results revealed a significant main effect of antonyms ($F(1,77) =$
166 $60.83, p < 0.001, \eta_p^2 = 0.44$) and a significant main effect of negation ($F(1,77) = 104.21, p < 0.001,$
167 $\eta_p^2 = 0.57$, **Fig.2A**). No significant crossover interaction between antonyms and negation was
168 observed ($p > 0.05$). Participants were faster for high adjectives (e.g., “good”) than for low
169 adjectives (e.g., “bad”) and for affirmative phrases (e.g., “really really good”) than for negated
170 phrases (e.g., “really not good”). These results support previous behavioral data showing that
171 negation is associated with increased processing difficulty^{15,16}. A further analysis including the
172 number of modifiers as factor (i.e., *complexity*) indicates that participants were faster for phrases
173 with two modifiers, e.g., “not really”, than phrases with one modifier, e.g., “not ####” ($F(1,77) =$
174 $16.02, p < 0.001, \eta_p^2 = 0.17$), suggesting that the placeholder “####” may induce some interference
175 to this otherwise relatively natural language task.

176

177 *Continuous mouse trajectories.* Continuous mouse trajectories across all adjective pairs and across
178 all participants are depicted in **Fig.2B** and **Fig.2C** (*low* and *high* summarize the two antonyms
179 across all scalar dimensions, see **Fig.S1** for each adjective dimension separately).

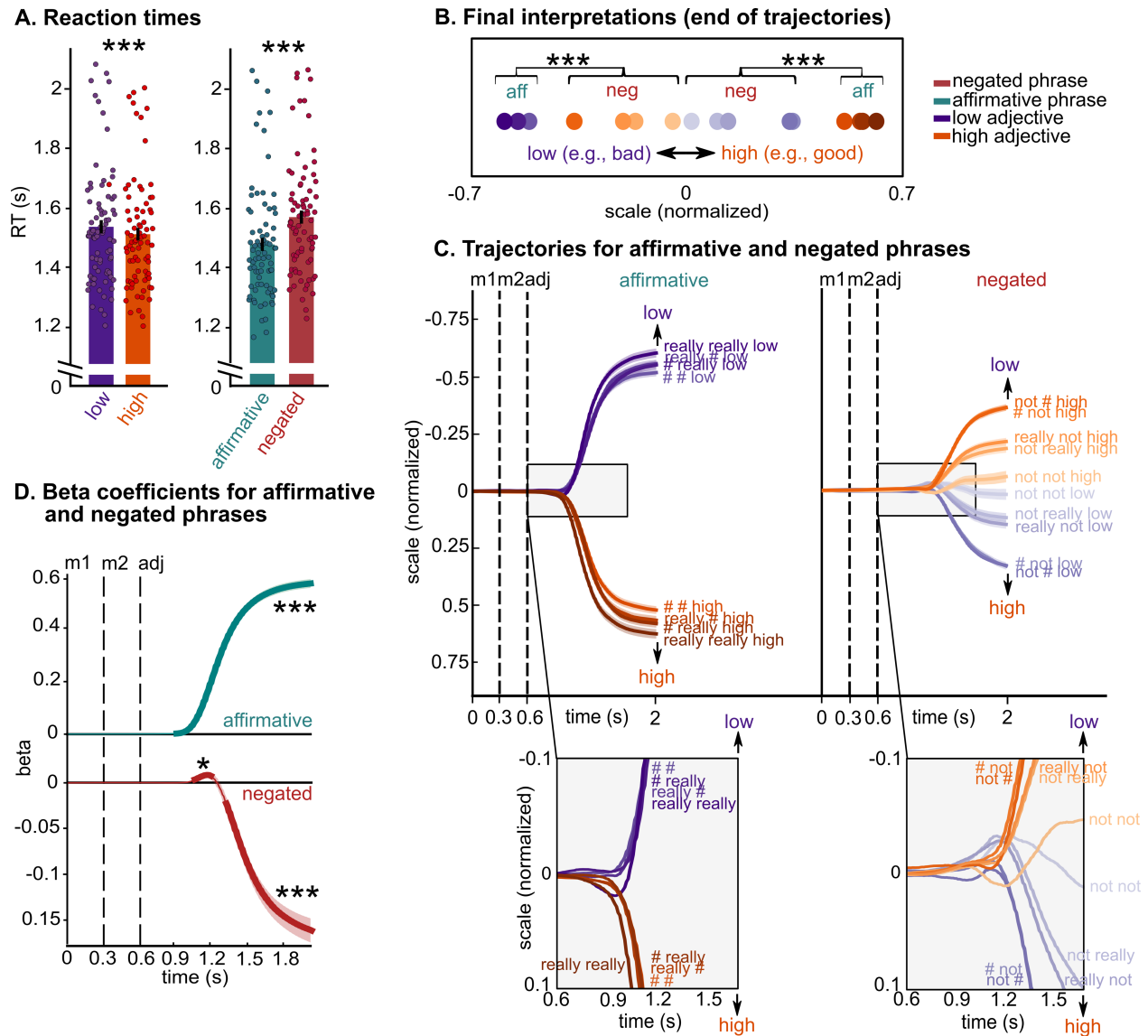
180 To quantify how the final interpretation of scalar adjectives changes as a function of
181 negation, we first performed a 2 (*antonym*: low vs. high) x 2 (*negation*: negated vs. affirmative)
182 repeated-measures ANOVA for participants’ ends of trajectories (filled circles in **Fig.2B**), which
183 reveal a significant main effect of antonyms ($F(1,77) = 338.57, p < 0.001, \eta_p^2 = 0.83$), a significant
184 main effect of negation ($F(1,77) = 65.50, p < 0.001, \eta_p^2 = 0.46$), and a significant antonyms by
185 negation interaction ($F(1,77) = 1346.07, p < 0.001, \eta_p^2 = 0.95$). Post-hoc tests show that the final
186 interpretation of negated phrases is located at a more central portion on the semantic scale than that
187 of affirmative phrases (affirmative low < negated high, and affirmative high > negated low, $p_{\text{holm}} <$
188 0.001). Furthermore, the final interpretation of negated phrases is significantly more variable
189 (measured as standard deviations) than that of affirmative phrases ($F(1,77) = 78.14, p < 0.001, \eta_p^2$
190 $= 0.50$). Taken together, these results suggest that negation shifts the final interpretation of

191 adjectives towards the antonyms, but never to a degree that overlaps with the interpretation of the
192 affirmative antonym.

193 Second, we explored the temporal dynamics of adjective representation as a function of
194 negation (i.e., from the presentation of word 1 to the final interpretation; lines in **Fig.2C**). While
195 mouse trajectories of affirmative phrases branch towards either side of the scale and remain on that
196 side until the final interpretation (lines in the left, gray, zoomed-in panel in **Fig.2C**), trajectories of
197 negated phrases first deviate towards the side of the adjective and then towards the side of the
198 antonym, to reach the final interpretation (i.e., “not low” first towards “low” and then towards
199 “high”; right, gray, zoomed-in panel in **Fig.2C**; see **Fig.S1** for each adjective dimension separately).
200 To characterize the degree of deviation towards each side of the scale, we performed regression
201 analyses with antonyms as the predictor and mouse trajectories as the dependent variable (see
202 **Methods**). The results confirm this observation, showing that (1) in affirmative phrases, betas are
203 positive (i.e., mouse trajectories moving towards the adjective) starting at 300 ms from adjective
204 onset ($p < 0.001$, green line in **Fig.2D**); and that (2) in negated phrases, betas are positive between
205 450 and 580 ms from adjective onset (i.e., mouse trajectories moving towards the adjective, $p =$
206 0.04), and only become negative (i.e., mouse trajectories moving towards the antonym, $p < 0.001$)
207 from 700 ms from adjective onset (red line in **Fig.2D**). Note that beta values of negated phrases are
208 smaller than that for affirmative phrases, again suggesting that negation does not invert the
209 interpretation of the adjective to that of the antonym.

210 Finally, we replicated this experiment in a new group of 55 online participants (**Fig.S2**).
211 The replication illustrates the robustness of the behavioral mouse tracking findings, even in the
212 absence of feedback. Taken together, these results suggest that participants initially interpreted
213 negated phrases as affirmative (e.g., “not good” interpreted along the “good” side of the scale) and
214 later as a mitigated interpretation of the opposite meaning (e.g., the antonym “bad”).

215



216 **Figure 2. Behavioral results.**

217 (A) Reaction times results for the online behavioral study (n=78). Bars represent the participants' mean ± SEM and
 218 dots represent individual participants. Participants were faster for high adjectives (e.g., “good”) than for low adjectives
 219 (e.g., “bad”) and for affirmative phrases (e.g., “really really good”) than for negated phrases (e.g., “really not good”).
 220 The results support previous behavioral data showing that negation is associated with increased processing difficulty.
 221 (B) Final interpretations (i.e., end of trajectories) of each phrase, represented by filled circles (purple = low, orange =
 222 high), averaged across adjective dimensions and participants, showing that negation never inverts the interpretation of
 223 adjectives to that of their antonyms. (C) Mouse trajectories for low (purple) and high (orange) antonyms, for each
 224 modifier (shades of orange and purple) and for affirmative (left panel) and negated (right panel) phrases. Zoomed-in
 225 panels at the bottom demonstrate that mouse trajectories of affirmative phrases branch towards the adjective’s side of
 226 the scale and remain on that side until the final interpretation; in contrast, the trajectories of negated phrases first deviate
 227 towards the side of the adjective and subsequently towards the side of the antonym. This result is confirmed by linear
 228 models fitted to the data at each timepoint in D. (D) Beta values (average over 78 participants) over time, separately
 229 for affirmative and negated phrases. Thicker lines indicate significant time windows. Panels C, D: black vertical dashed

230 lines indicate the presentation onset of each word: modifier 1, modifier 2 and adjective; each line and shading represent
231 participants' mean \pm SEM; Panels A,B,D: *** $p < 0.001$; * $p < 0.05$.

232

233 ***Experiment 2: MEG shows that negation weakens the representation of adjectives and recruits*** 234 ***response inhibition networks***

235 In this study (MEG experiment, $N = 26$), participants read adjective phrases comprising one or two
236 modifiers (“not” and “really”) and scalar adjectives across different dimensions (e.g., “really really
237 good”, “really not quiet”, “not ### dark”). Adjectives were selected to represent opposite poles
238 (i.e., the antonyms) of the respective semantic scales: *low* pole of the scale (e.g., “bad”, “cool”,
239 “quiet”, “dark”) and *high* pole of the scale (e.g., “good”, “warm”, “loud”, “bright”). A sequence of
240 dashes was used to indicate the absence of a modifier. **Fig. 1B** and **Table S2** provide the
241 comprehensive list of the linguistic stimuli. Participants were asked to indicate whether a probe
242 (e.g., 6) correctly represented the meaning of the phrase on a scale from “really really low” (0) to
243 “really really high” (8) (*yes/no* answer, **Fig.1B**). Behavioral data of Experiment 2 replicate that of
244 Experiment 1: negated phrases are processed slower and with more errors than affirmative phrases
245 (main effect of negation for RTs: $F(1,25) = 26.44$, $p < 0.001$, $\eta_p^2 = 0.51$; main effect of negation for
246 accuracy: $F(1,25) = 8.03$, $p = 0.009$, $\eta_p^2 = 0.24$).

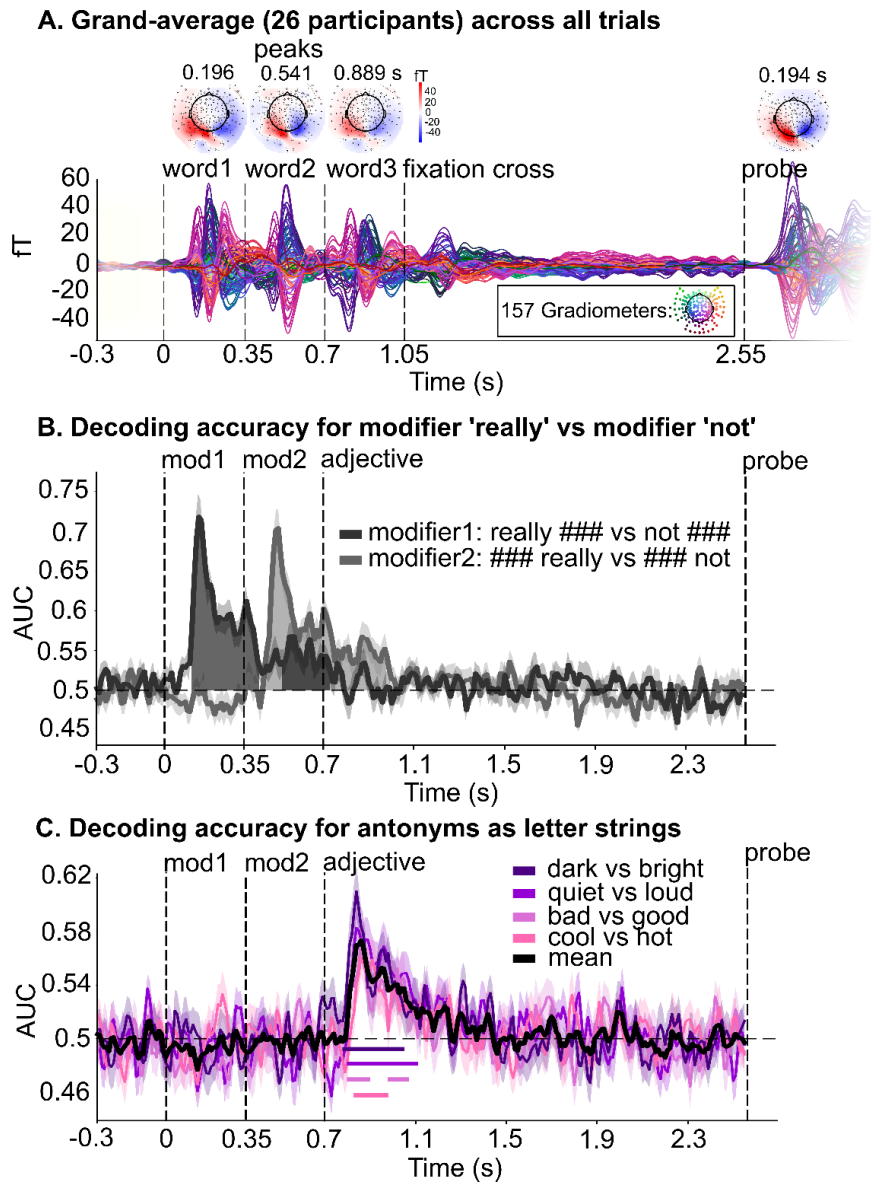
247 The MEG analysis, using largely temporal and spatial decoding approaches⁴⁷, comprises
248 four steps: (1) we first identify the temporal correlates of simple word representation (i.e., the words
249 “really” and “not” in the modifier position, and each pair of scalar adjectives in the second word
250 position, i.e., the head position); (2) we test lexical-semantic representations of adjectives over time
251 beyond the single word level, by entering *low* (“bad”, “cool”, “quiet” and “dark”) and *high* (“good”,
252 “warm”, “loud” and “bright”) antonyms in the same model. We then test the representation of the
253 negation operator over time; (3) we then ask how negation operates on the representation of
254 adjectives, by teasing apart four possible mechanisms (i.e., *No effect*, *Mitigation*, *Inversion*,
255 *Change*); (4) we explore changes in beta power as a function of negation (motivated by the literature
256 implicating beta-band neural activity).

257

258 (1) *Temporal decoding of single word processing*

259 Results show that the temporal decoding (see **Methods**) of “really” vs. “not” is significant between
260 120 and 430 ms and between 520 and 740 ms from the onset of the first modifier (dark gray areas,
261 $p < 0.001$ and $p = 0.001$) and between 90 and 640 ms from the onset of the second modifier (light
262 gray areas, $p < 0.001$, **Fig.3B**). Pairs of antonyms from different scales were similarly decodable
263 between 90 and 410 ms from adjective onset (quality: 110 to 200 ms, $p = 0.002$ and 290 to 370 ms,

264 $p = 0.018$; temperature: 140 to 280 ms, $p < 0.001$; loudness: 110 to 410 ms, $p < 0.001$; brightness:
265 90 to 350 ms, $p < 0.001$, **Fig.3C**), reflecting time windows during which the brain represents visual,
266 lexical, and semantic information (e.g., ^{7,48}).



267 **Figure 3. Evoked activity and temporal decoding of modifiers and adjectives as letter strings.**
268 **(A)** The butterfly (bottom) and topo plots (top) illustrate the event-related fields elicited by the presentation of each
269 word as well as the probe, with a primarily visual distribution of neural activity right after visual onset (i.e., letter string
270 processing). We performed multivariate decoding analyses on these preprocessed MEG data. Detector distribution of
271 MEG system in inset box. ft: femtoTesla magnetic field strength. **(B)** We estimated the ability of the decoder to
272 discriminate “really” vs. “not” in either modifier’s position, from all MEG sensors. We contrasted phrases with
273 modifiers “really ###” and “not ###”, and phrases with modifiers “### not” and “### really”. **(C)** We evaluated whether
274 the brain encodes representational differences between each pair of antonyms (e.g., “bad” vs. “good”), in each of the
275 four dimensions (quality, temperature, loudness, and brightness). The mean across adjective pairs is represented as a

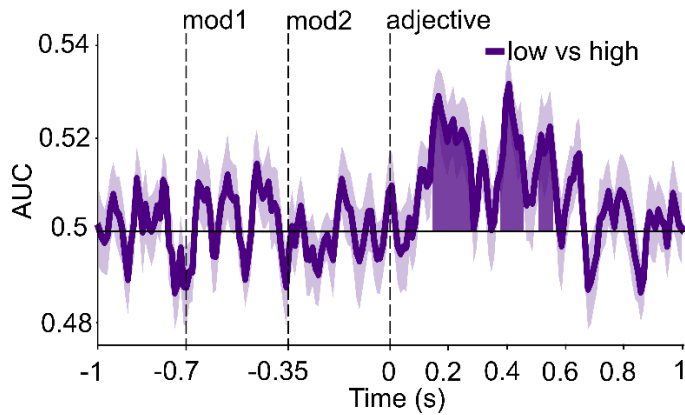
276 solid black line; significant windows are indicated by horizontal solid lines below. For panels B and C: AUC = area
277 under the receiver operating characteristic curve, chance = 0.5 (black horizontal dashed line); For all panels: black
278 vertical dashed lines indicate the presentation onset of each word: modifier 1, modifier 2, and adjective; each line and
279 shading represent participants' mean \pm SEM.

280

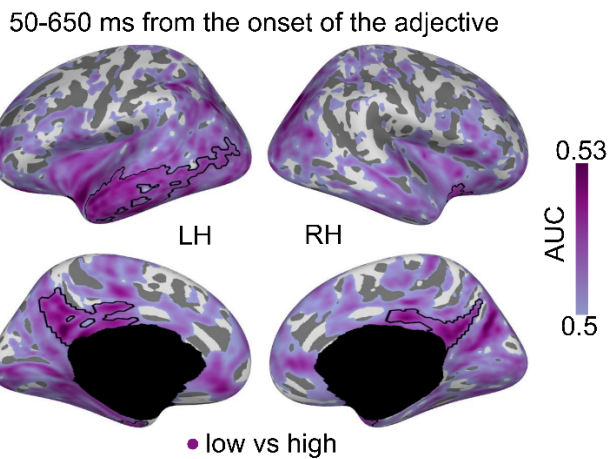
281 (2) *Temporal and spatial decoding of adjectives and negation*

282 After establishing that single words' features can be successfully decoded in sensible time windows
283 (see **Fig.3**), we moved beyond single word representation to selectively evaluate lexical-semantic
284 differences between *low* ("bad", "cool", "quiet" and "dark") and *high* ("good", "warm", "loud" and
285 "bright") adjectives, regardless of the specific scale (i.e., pooling over *quality*, *temperature*,
286 *loudness*, and *brightness*). Temporal decoding analyses (see **Methods**) reveal significant
287 decodability of *low* vs. *high* antonyms in three time windows between 140 and 560 ms from
288 adjective onset (140 to 280 ms, $p < 0.001$; 370 to 460 ms: $p = 0.009$; 500 to 560 ms: $p = 0.044$,
289 purple areas in **Fig.4A**). No significant differences in lexical-semantic representation between *low*
290 and *high* antonyms were observed in later time windows (i.e., after 560 ms from adjective onset).
291 The spatial decoding analysis illustrated in **Fig.4B** (limited to 50-650 ms from adjective onset, see
292 **Methods**) show that decoding accuracy for *low* vs. *high* antonyms is significantly above chance in
293 a widespread left-lateralized brain network, encompassing the anterior portion of the superior
294 temporal lobe, the middle, and the inferior temporal lobe (purple areas in **Fig.4B**, significant
295 clusters are indicated by a black contour: left temporal lobe cluster, $p = 0.002$). A significant cluster
296 was also found in the right temporal pole, into the insula ($p = 0.007$). Moreover, we found
297 significant clusters in the bilateral cingulate gyri (posterior and isthmus) and precunei (left
298 precuneus/cingulate cluster, $p = 0.009$; right precuneus/cingulate cluster, $p = 0.037$). Overall, these
299 regions are part of the (predominantly left-lateralized) frontotemporal brain network that underpins
300 lexical-semantic representation and composition^{7,8,48-55}.

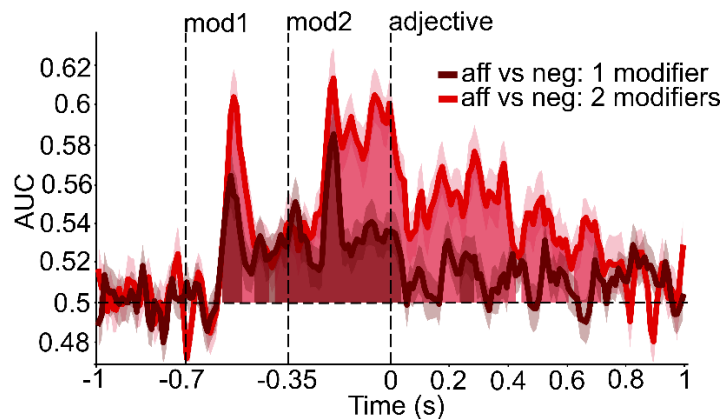
A. Temporal decoding of antonyms: word meaning



B. Spatial decoding of antonyms



C. Temporal decoding of negation as a function of complexity



301 **Figure 4. Temporal and spatial decoding of antonyms across all scales and temporal decoding of negation.**
302 (A) Decoding accuracy (purple line) of lexical-semantic differences between antonyms across all scales (i.e., pooling
303 over “bad”, “cool”, “quiet” and “dark”; and “good”, “warm”, “loud” and “bright” before fitting the estimators) over
304 time; significant time windows are indicated by purple areas; (B) Decoding accuracy (shades of purple) for antonyms
305 across all scales over brain sources (after pooling over the four dimensions), between 50 and 650 ms from adjective
306 onset. Significant spatial clusters are indicated by a black contour. (C) Decoding accuracy of negation over time, as a

307 function of the number of modifiers (1 modifier: dark red line and shading; 2 modifiers: light red line and shading).
308 Significant time windows are indicated by dark red (1 modifier) and light red (2 modifiers) areas. For all panels: AUC:
309 area under the receiver operating characteristic curve, chance = 0.5 (black horizontal dashed line); black vertical dashed
310 lines indicate the presentation onset of each word: modifier1, modifier2 and adjective; each line and shading represent
311 participants' mean \pm SEM; aff = affirmative, neg = negated; LH = left hemisphere; RH = right hemisphere.

312

313 Next, we turn to representations of negation over time. We performed a temporal decoding analysis
314 for phrases containing “not” vs. phrases not containing “not”, separately for phrases with one and
315 two modifiers (to account for phrase complexity; see **Table S2** for a list of all trials). For phrases
316 with one modifier, the decoding of negation is significantly higher than chance throughout word 1
317 (-580 to -500 ms from adjective onset, $p = 0.005$), then again throughout word 2 (-470 to 0 ms from
318 adjective onset, $p < 0.001$). After the presentation of the adjective, negation decodability is again
319 significantly above chance between 0 and 40 ms ($p = 0.034$) and between 230 and 290 ms from
320 adjective onset ($p = 0.018$; dark red line and shading in **Fig.4C**). Similarly, for phrases with two
321 modifiers, the decoding of negation is significantly higher than chance throughout word 1 (-580 to
322 -410 ms from adjective onset, $p = 0.002$), throughout word 2 (-400 to 0 ms from adjective onset, p
323 < 0.001), and for a longer time window from adjective onset compared to phrases with one modifier,
324 i.e., between 0 and 720 ms (0 to 430 ms, $p < 0.001$; 440 to 500 ms, $p = 0.030$; 500 to 610 ms, $p <$
325 0.001 ; 620 to 720 ms, $p < 0.001$; light red line and shading in **Fig.4C**). The same analysis time-
326 locked to the onset of the probe shows that negation is once again significantly decodable between
327 230 and 930 ms after the probe (**Fig.S3**).

328 Cumulatively, these results suggest that the brain encodes negation every time a “not” is
329 presented and maintains this information up to 720 ms after adjective onset. Further, they show that
330 the duration of negation maintenance is amplified by the presence of a second modifier, highlighting
331 combinatoric effects ^{2,6,56}.

332

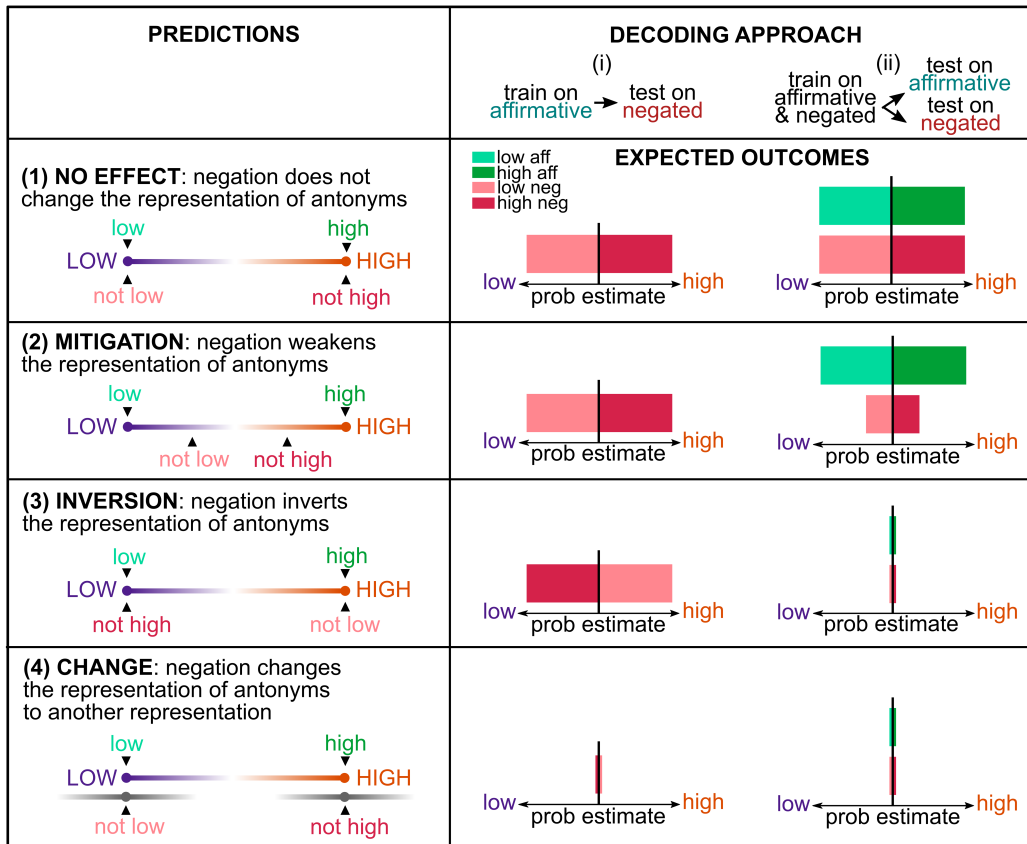
333 (3) *Effect of negation on lexical-semantic representations of antonyms over time*

334 The temporal decoding analyses performed separately for adjectives and for negation demonstrates
335 that the brain maintains the representation of the modifiers available throughout the presentation of
336 the adjective. Here we ask how negation *operates on* the representation of the antonyms at the
337 neural level, leveraging theoretical accounts of negation ^{11,12,42-44}, behavioral results of Experiment
338 1, and two complementary decoding approaches. We test four hypotheses (see *Predictions* in
339 **Fig.5A**): (1) *No effect of negation*: negation does not change the representation of adjectives (i.e.,
340 “not low” = “low”). We included this hypothesis based on the two-step theory of negation, wherein
341 the initial representation of negated adjectives would not be affected by negation ²⁷. (2) *Mitigation*:

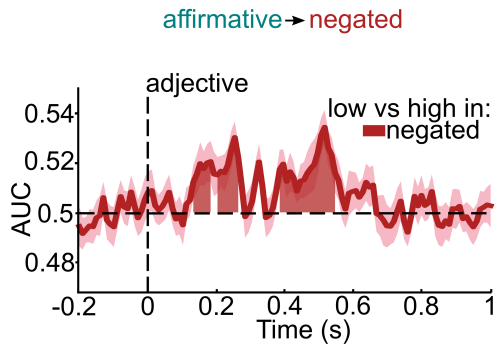
342 negation weakens the representation of adjectives (i.e., “not low” < “low”). (3) *Inversion*: negation
343 inverts the representation of adjectives (i.e., “not low” = “high”). Hypotheses (3) and (4) are derived
344 from previous linguistics and psycholinguistics accounts on comprehension of negated adjectives
345 ^{42–44}. Finally, (4) *Change*: we evaluated the possibility that negation might change the representation
346 of adjectives to another representation outside the semantic scale defined by the two antonyms (e.g.,
347 “not low” = e.g., “fair”).

348 To adjudicate between these four hypotheses, we performed two sets of decoding analyses.
349 Decoding approach (i): we computed the accuracy with which estimators trained on *low* vs. *high*
350 antonyms in affirmative phrases (e.g., “really really bad” vs. “really really good”) generalize to the
351 representation of *low* vs. *high* antonyms in negated phrases (e.g., “really not bad” vs. “really not
352 good”) at each time sample time-locked to adjective onset (see **Methods**); decoding approach (ii):
353 we trained estimators on *low* vs. *high* antonyms in affirmative and negated phrases together (in 90%
354 of the trials) and computed the accuracy of the model in predicting the representation of *low* vs.
355 *high* antonyms in affirmative and negated phrases separately (in the remaining 10% of the trials;
356 see **Methods**). Decoding approach (ii) allows for direct comparison between AUC and probability
357 estimates in affirmative and negated phrases. Expected probability estimates (i.e., the averaged
358 class probabilities for *low* and *high* classes) as a result of decoding approach (i) and (ii) are depicted
359 as light and dark, green and red bars under *Decoding approach* in **Fig.5A**.

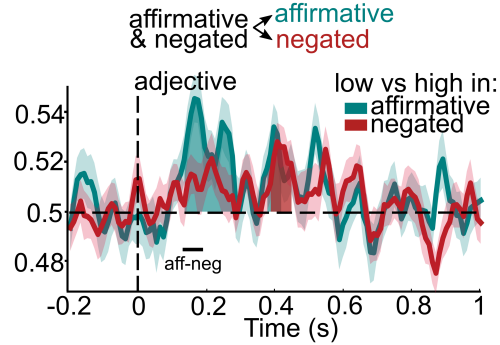
A. Expected outcomes for the effect of negation on the representation of antonyms



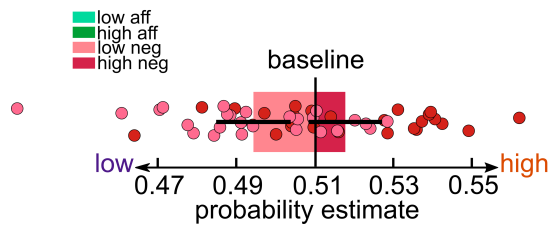
B. Results of decoding approach (i)



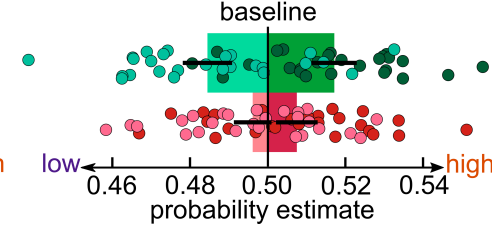
C. Results of decoding approach (ii)



D. Significant clusters (averaged)



E. Significant clusters (averaged)



360 **Figure 5. Predictions, decoding approaches, and results of the effect of negation on the representation of**
 361 **adjectives.**

362 (A) We tested four possible effects of negation on the representation of adjectives: (1) *No effect*, (2) *Mitigation*, (3)
 363 *Inversion*, (4) *Change* (left column). Note that we depicted predictions of (3) *Inversion* on the extremes of the scale,
 364 but a combination of inversion and mitigation would predict the same outcomes. We performed two sets of decoding

365 analyses (right column): (i) We trained estimators on low (purple) vs. high (orange) antonyms in affirmative phrases
366 and predicted model accuracy and probability estimates of low vs. high antonyms in negated phrases (light and dark
367 red bars). (ii) We trained estimators on low vs. high antonyms in affirmative and negated phrases together and predicted
368 model accuracy and probability estimates in affirmative (light and dark green bars) and negated phrases (light and dark
369 red bars) separately. **(B)** Decoding accuracy (red line) over time of antonyms for negated phrases, as a result of decoding
370 approach (i). Significant time windows are indicated by red areas. **(C)** Decoding accuracy of antonyms over time for
371 affirmative (green line) and negated (red line) phrases, as a result of decoding approach (ii). Significant time windows
372 for affirmative and negated phrases are indicated by green and red areas. The significant time window of the difference
373 between affirmative and negated phrases is indicated by a black horizontal solid line. **(D)** Probability estimates for low
374 (light red) and high (dark red) negated antonyms averaged across the significant time windows depicted in **B**. Bars
375 represent the participants' mean \pm SEM and dots represent individual participants. **(E)** Probability estimates for low
376 (light green) and high (dark green) affirmative adjectives and for low (light red) and high (dark red) negated adjectives,
377 averaged across the significant time window depicted as a black horizontal line in **C**. Chance level of probability
378 estimates was computed by averaging probability estimates of the respective baseline (note that the baseline differs
379 from 0.5 due to the different number of trials for each class in the training set of decoding approach (i)). Bars represent
380 the participants' mean \pm SEM and dots represent individual participants. For panels **B** and **C**: AUC: area under the
381 receiver operating characteristic curve, chance = 0.5 (black horizontal dashed line); each line and shading represent
382 participants' mean \pm SEM. Panels **B,C,D,E**: the black vertical dashed line indicates the presentation onset of the
383 adjective; green = affirmative phrases, red = negated phrases.

384
385 Temporal decoding approach (i) reveals that the estimators trained on the representation of
386 *low* vs. *high* antonyms in affirmative phrases significantly generalize to the representation of *low*
387 vs. *high* antonyms in negated phrases, in four time windows between 130 and 550 ms from adjective
388 onset (130 to 190 ms, $p = 0.039$; 200 to 270 ms: $p = 0.003$; 380 to 500 ms: $p < 0.001$; 500 to 550
389 ms: $p = 0.008$; red areas in **Fig.5B**). **Fig.5D** depicts the probability estimates averaged over the
390 significant time windows for *low* and *high* antonyms in negated phrases. These results only support
391 predictions (1) *No effect* and (2) *Mitigation*, thus invalidating predictions (3) *Inversion* and (4)
392 *Change*. **Fig.S4** illustrates a different approach that similarly leads to the exclusion of prediction
393 (3) *Inversion*.

394 Temporal decoding approach (ii) shows significant above chance decoding accuracy for
395 affirmative phrases between 130 and 280 ms ($p < 0.001$) and between 370 and 420 ms ($p = 0.035$)
396 from adjective onset. Conversely, decoding accuracy for negated phrases is significantly above
397 chance only between 380 and 450 ms after the onset of the adjective ($p = 0.004$). Strikingly, negated
398 phrases are associated with significantly lower decoding accuracy than affirmative phrases in the
399 time window between 130 and 190 ms from adjective onset ($p = 0.040$; black horizontal line in
400 **Fig.5C**). **Fig.5E** represents the probability estimates averaged over this 130-190 ms significant time

401 window for *low* and *high* antonyms, separately in affirmative and negated phrases, illustrating
402 reduced probability estimates for negated compared to affirmative phrases.

403 Overall, the generalization of representation from affirmative to negated phrases and the
404 higher decoding accuracy (and probability estimates) for affirmative than negated phrases within
405 the first 500 ms from adjective onset (i.e., within the time window of lexical-semantic processing
406 shown in **Fig.4A**) provide direct evidence in support of prediction (2) *Mitigation*, wherein negation
407 weakens the representation of adjectives. The alternative hypotheses did not survive the different
408 decoding approaches.

409

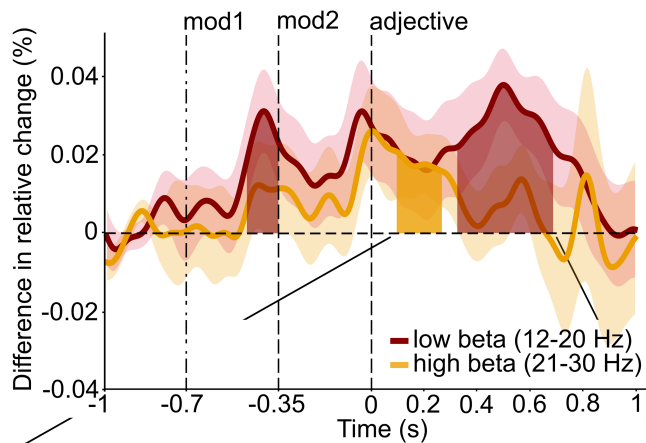
410 (4) *Changes in beta power as a function of negation*

411 We distinguished among four possible mechanisms of how negation could operate on the
412 representation of adjectives and demonstrated that negation does not invert or change the
413 representation of adjectives but rather weakens the decodability of *low* vs. *high* antonyms,
414 significantly for about 60 ms from adjective onset. The availability of negation upon the processing
415 of the adjective (**Fig.4C**) and the reduced decoding accuracy for antonyms in negated phrases
416 (**Fig.5C**) raise the question of whether negation operates through inhibitory mechanisms, as
417 suggested by previous research employing action-related verbal material^{35–37}. We therefore
418 performed time-frequency analyses, focusing on beta power (including low-beta: 12 to 20 Hz, and
419 high-beta: 20 to 30 Hz,⁵⁷, see **Methods**), which has been previously associated with inhibitory
420 control⁵⁸ (see **Fig.S5** for comprehensive time-frequency results). We reasoned that, if negation
421 operates through general-purpose inhibitory systems, we should observe higher beta power for
422 negated than affirmative phrases in sensorimotor brain regions.

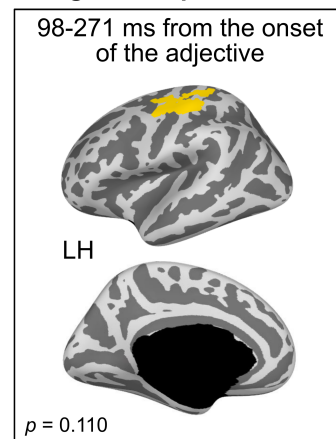
423 Our results are consistent with this hypothesis, showing significantly higher low-beta power
424 (from 229 to 350 ms from the onset of modifier1: $p = 0.036$; from 326 to 690 ms from adjective
425 onset: $p = 0.012$; red line in **Fig.6A**) and high-beta power (from 98 to 271 ms from adjective onset:
426 $p = 0.044$; yellow line in **Fig.6A**) for negated than affirmative phrases. **Fig.S6** further shows low
427 and high-beta power separately for negated and affirmative phrases, compared to phrases with no
428 modifier.

429 Our whole-brain source localization analysis shows significantly higher low-beta power for
430 negated than affirmative phrases in the left precentral, postcentral, and paracentral gyri ($p = 0.012$;
431 between 326 and 690 ms from adjective onset, red cluster in **Fig.6C**). For high-beta power, similar
432 (albeit not significant) sensorimotor spatial patterns emerge (yellow cluster in **Fig.6B**).

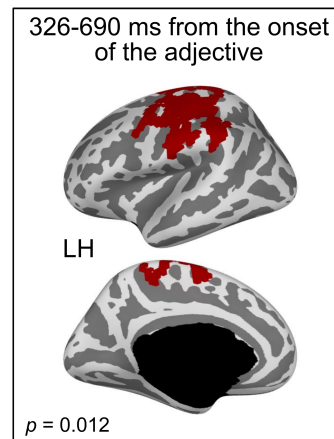
A. Low- and high-beta power of (negated - affirmative phrases)



B. High-beta spatial cluster



C. Low-beta spatial cluster



433 **Figure 6. Differences in beta power over time between negated and affirmative phrases.**

434 (A) Differences in low (12-20 Hz, red) and high (21-30 Hz, yellow) beta power over time between negated and
435 affirmative phrases. Negated phrases show higher beta power compared to affirmative phrases throughout the
436 presentation of the modifiers and for a sustained time window from adjective onset up to ~700 ms; significant time
437 windows are indicated by red (low-beta) and yellow (high-beta) areas; black vertical dashed lines indicate the
438 presentation onset of each word: modifier1, modifier2 and adjective; each line and shading represent participants' mean
439 \pm SEM. (B) Differences (however not reaching statistical significance, $\alpha = 0.05$) in high-beta power between negated
440 and affirmative phrases (restricted between 97 and 271 ms from adjective onset, yellow cluster). (C) Significant
441 differences in low-beta power between negated and affirmative phrases (restricted between 326 and 690 ms from
442 adjective onset) in the left precentral, postcentral and paracentral gyrus (red cluster). Note that no significant spatial
443 clusters were found in the right hemisphere.

444
445
446
447
448
449

450 **Discussion**

451 We tracked changes over time in lexical-semantic representations of scalar adjectives, as a function
452 of the intensifier “really” and the negation operator “not”. Neural correlates of negation have
453 typically been investigated in the context of action verbs ^{29,35–37,40,41,59–63}. Our study employs
454 minimal linguistic contexts to characterize in detail how negation operates on abstract, non-action-
455 related lexical-semantic representations. We leveraged (1) psycholinguistic findings on adjectives
456 that offer a framework wherein meaning is represented on a continuum ^{42,43}, (2) time-resolved
457 behavioral and neural data, and (3) multivariate analysis methods (decoding) which can
458 discriminate complex lexical-semantic representations from distributed neuronal patterns (e.g., ⁶²).

459 The longer RTs and decreased accuracy for negated phases shown in Experiment 1
460 (**Fig.2A**), in the replication experiment (**Fig.S2A**), and in Experiment 2, are consistent with data
461 demonstrating that negation incurs increased processing costs ^{13–18,27,32}. More significantly, mouse
462 trajectories show that participants initially interpreted negated phrases as affirmative (e.g., “not
463 good” is located on the “good” side of the scale, for ~130 ms, **Fig.2C** and **Fig.S2C**), indicating that
464 initial representations of negated scalar adjectives are closer to the representations of the adjectives
465 rather than that of their antonyms. Similarly, participants’ final interpretations of negated adjectives
466 (e.g., “not good”, “really not good”) never overlapped with the final interpretations of the
467 corresponding affirmative antonyms (e.g., “bad”, “really bad”, “really really bad”; **Fig.2B** and
468 **Fig.S2B**) highlighting how negation never inverts the meaning of an adjective to that of its antonym,
469 even when participants are making decisions on a binary semantic scale (9,37-40).

470 Continuous mouse trajectories allowed us to quantify dynamic changes in participants’
471 interpretations. MEG provided a means to directly track neural representations over time. We first
472 identified the temporal correlates of lexical-semantic processing *separately* for scalar adjectives
473 and for the negation operator. The time window of adjective representation (~140-560 ms from
474 adjective onset, **Fig.4A**) is consistent with previous studies investigating lexical-semantic
475 processing in language comprehension (130–200 ms up to ~550 ms from word onset ^{64–68}). Spatial
476 decoding results corroborate temporal results, highlighting the involvement of the left-lateralized
477 frontotemporal brain network in adjective processing (**Fig.4B**, ^{7,8,48–55}). Our data further
478 demonstrate that negation is processed in parallel to the processing of the adjective (up to ~700 ms;
479 **Fig.4C**), not serially (see ^{69,70} for related patterns in the context of negation + auxiliary verb and
480 adjective + noun). Finally, they show that the decodability of negation increases in phrases with
481 two modifiers (e.g., “really not”, “not really”, **Fig.4C**, **Fig.S3**), highlighting compositional effects
482 ⁶.

483 We then evaluated the effects of the negation operator *on* adjective representation, to
484 address the question of *how* negation operates on lexical-semantic representations of antonyms. We
485 contrasted four hypotheses (**Fig.5A**): negation (1) does not change the representation of scalar
486 adjectives (e.g., “not good” = “good”, *No effect*), (2) weakens the representation of scalar adjectives
487 (e.g., “not good” < “good”, *Mitigation*), (3) inverts the representation of scalar adjectives (e.g., “not
488 good” = “bad”, *Inversion*), or (4) changes the representation of scalar adjectives to another
489 representation (e.g., “not good” = e.g., “unacceptable”, *Change*). First, we demonstrated that, within
490 the time window of adjective encoding, the representation of affirmative adjectives generalizes to
491 that of negated adjectives (**Fig.5B** and **Fig.5D**). This finding rules out predictions (3) *Inversion* and
492 (4) *Change*. Moreover, these findings complement our behavioral data that show that negated
493 adjectives are initially interpreted by participants as affirmative. Second, we showed that the
494 representation of adjectives in affirmative and negated phrases is not identical but is weakened by
495 negation (**Fig.5C** and **Fig.5E**). This result rules out prediction (1) *No effect* and supports prediction
496 (2) *Mitigation*, wherein negation weakens the representation of adjectives. We observed such
497 reduction in early lexical-semantic representations (i.e., from ~130 ms post adjective-onset),
498 supporting previous research that reported effects of negation as soon as lexical-semantic
499 representations of words are formed^{12,29–31,71}, and not exclusively at later processing stages (e.g.,
500 P600^{72,73}).

501 Our behavioral and neural data jointly point to a *mitigation* rather than an *inversion* effect
502 of negation: initial interpretations and neural representations of negated adjectives are similar to
503 that of affirmative adjectives, but weakened; final interpretations do not overlap with neither
504 affirmative extreme of the semantic scale. While previous fMRI studies on sentential negation have
505 shown that negation reduces hemodynamic brain activations related to verb processing^{40,41}, the
506 current study offers novel time-resolved behavioral and neural data on how negation selectively
507 operates on abstract concepts. Previous research has highlighted that negation might behave
508 differently depending on the pragmatics of discourse interpretation, e.g., when presented in
509 isolation as compared to when presented in context (“not wrong” vs. “your theory is not wrong”
510^{9,10}), or when used ironically (“they are not really good” said ironically to mean that they are
511 “mediocre”, e.g.,^{11,71}). Within this pragmatic framework, it has been suggested that the opposite
512 meaning of a scalar adjective would be more simply conveyed by the affirmative counterpart than
513 by negation^{11,44,74}; thus, to convey the opposite meaning of “bad”, it would be more appropriate to
514 use “good” as opposed to “not bad”. Following this logic, negation would be purposefully used
515 (and understood) to convey a different, mitigated meaning of the adjective (e.g., “not bad” = “less
516 than bad”). Although we did not directly manipulate sentential or pragmatic contexts, our findings

517 provide behavioral and neural evidence that negation acts as a mitigator. Here we only tested
518 adjective pairs that form *contraries* (which lie on a continuum, e.g., “bad” and “good”); thus
519 inherently different patterns of results could emerge in the case of *contradictories* (which form a
520 dichotomy, e.g., “dead” and “alive”,⁴⁴), where there is no continuum for mitigation to have an
521 effect.

522 Overall, evidence that negation weakens adjective representations invites the hypothesis
523 that negation operates as a suppression mechanism, possibly through general-purpose inhibitory
524 systems^{36,37}. To address this, we compared beta power modulations in affirmative and negated
525 phrases (**Fig.6**). In addition to subserving motor processing, beta-power modulation (12-30 Hz) is
526 associated with multiple aspects of language processing^{35,75–78}; for a review, see^{57,79}). We
527 evaluated differences between negated and affirmative phrases separately in the low- and high-beta
528 bands. We found greater power for negated than affirmative phrases in both bands, during the
529 processing of the modifier and throughout the processing of the adjective up to ~700 ms, localized
530 in left-lateralized sensorimotor areas. The timing and spatial correlates of beta-power in relation to
531 negation align with studies that examined the effect of negation on (mental and motor) action
532 representation³⁶. Strikingly, we demonstrated that negation recruits brain areas and
533 neurophysiological mechanisms similar to that recruited by response inhibition - however in the
534 absence of action-related language material. Within a framework that recognizes two interactive
535 neural systems, i.e., a semantic representation and a semantic control system⁵³, negation would
536 operate through the latter, modulating how activation propagates through the (ventral) language
537 semantic network wherein meaning is represented. The precise connectivity that underpins
538 mitigation of lexical-semantic representations remains to be investigated.

539 Collectively, we demonstrated that, by characterizing subtle changes of linguistic meaning
540 through negation, using time-resolved behavioral and neuroimaging methods and multivariate
541 decoding, we can tease apart different possible representation outcomes of combinatorial
542 operations, above and beyond the sum of the processing of individual word meanings.

543

544

545 **Materials and Methods**

546 *Participants*

547 *Experiment 1: continuous behavioral tracking.* 101 participants (46 females; mean age = 29.6 years;
548 range 18-67 years) completed an online mouse tracking experiment. Participants were recruited via
549 Amazon Mechanical Turk and via the platform SONA (a platform for students’ recruitment). All
550 participants were native English speakers with self-reported normal hearing, normal or corrected to

551 normal vision, and no neurological deficits. 97 participants were right-handed. Participants were
552 paid or granted university credits for taking part in the study, which was performed online. All
553 participants provided written informed consent, as approved by the local institutional review board
554 (New York University's Committee on Activities Involving Human Subjects). The data of 23
555 participants were excluded from the data analysis due to (i) number of "incorrect" feedback (based
556 on the warnings) > 30%, (ii) mean RTs > 2SD from the group mean, or (iii) response trajectory
557 always ending within 1/4 from the center of the scale, regardless of condition (i.e., participants who
558 did not pay attention to the instructions of the task). Thus, 78 participants were included in the
559 analysis. The sample size was determined based on previous studies using a similar behavioral
560 approach (~30 participants^{15,45,80}) and was increased to account for the exclusion rate reported for
561 online crowdsourcing experiments^{81,82}. For participants in *Experiment 1 (replication)* see **Fig.S2**.

562
563 *Experiment 2: MEG.* A new group of 28 participants (17 females; mean age = 28.7 years; range 19-
564 53 years) took part in the in-lab MEG experiment. All participants were native English speakers
565 with self-reported normal hearing, normal or corrected to normal vision, and no neurological
566 deficits. 24 participants were right-handed. They were paid or granted university credits for taking
567 part in the study. All participants provided written informed consent, as approved by the local
568 institutional review board (New York University's Committee on Activities Involving Human
569 Subjects). The data of 2 participants were excluded from the data analysis because their accuracy
570 scores in the behavioral task was < 60%. Thus, 26 participants were included in the analysis. The
571 sample size was determined based on previous studies investigating negation using EEG (17 to 33
572 participants^{26,35,37}), investigating semantic representation using MEG (25 to 27 participants^{7,8}), or
573 employing decoding methods with MEG data (17 to 20 participants^{83,84}).

574

575 ***Stimuli, Design, and Procedure***

576 *Experiment 1 (and replication): continuous mouse tracking.*

577 *Stimuli and Design.* The linguistic stimulus set comprises 108 unique adjective phrases (for the
578 complete list, see **Table S1**). Adjectives were selected to be antonyms (i.e., *low* and *high* poles of
579 the scale) in the following six cognitive or sensory dimensions: *quality* ("bad", "good"), *beauty*
580 ("ugly", "beautiful"), *mood* ("sad", "happy"), *temperature* ("cold", "hot"), *speed* ("slow", "fast"),
581 and *size* ("small", "big"). These antonyms are all *contraries* (i.e., adjectives that lie on a continuum
582⁴⁴). Lexical characteristics of the antonyms were balanced according to the English Lexicon Project
583⁸⁵; mean (SD) HAL log frequency of *low* adjectives: 10.69 (1.09), *high* adjectives: 11.51 (1.07),
584 mean (SD) bigram frequency of *low* adjectives: 1087.10 (374), *high* adjectives: 1032 (477.2); mean

585 (SD) lexical decision RTs of *low* adjectives: 566 (37), *high* adjectives: 586 ms (70)). Adjectives
586 were combined with zero (e.g., “### ##”), one (e.g., “really ##”), or two modifiers (e.g., “really
587 not”). Modifiers were either the intensifier “really” or the negation “not” (see³³ for a similar choice
588 of modifiers). A sequence of dashes was used to indicate the absence of a modifier, e.g., “really
589 ## good”. Each of the 12 adjectives was preceded by each of the nine possible combinations of
590 modifiers: “### ##”, “### really”, “really ##”, “### not”, “not ##”, “really not”, “not really”,
591 “really really” and “not not” (“not not” was included to achieve a full experimental design, even if
592 it is not a frequent combination in natural language and its cognitive and linguistic representations
593 are still under investigation⁸⁶). Each dimension (e.g., quality) was presented in two blocks (one
594 block for each scale orientation, e.g., low to high and high to low) for a total of 12 blocks. Each
595 phrase was repeated three times within each block (note that “### really”/“really ##” were
596 repeated an overall of three times, and so were “### not”/“not ##”). Thus, the overall experiment
597 comprised 504 trials. The order of phrases was randomized within each block for each participant.
598 The order of pairs of blocks was randomized across participants.

599
600 *Procedure.* Behavioral trajectories provide time-resolved dynamic data that reflect changes in
601 representation^{15,45,46}. The online experiment was developed using oTree, a Python-based
602 framework for the development of controlled experiments on online platforms⁸⁷. Participants
603 performed this study remotely, using their own monitor and mouse (touchpads were not allowed).
604 They were instructed to read affirmative or negated adjective phrases (e.g., “really really good”,
605 “really not bad”) and rate the overall meaning of each phrase on a scale, e.g., from “really really
606 bad” to “really really good”. Participants were initially familiarized with the experiment through
607 short videos and a short practice block (18 trials with feedback). They were instructed that the poles
608 of the scale (e.g., “bad” and “good”) would be reversed in half of the trials and warned that (i) they
609 could not cross the vertical borders of the response space, (ii) they had to maintain a constant
610 velocity, by following an horizontal line moving vertically, and (iii) they could not rate the meaning
611 of the phrase before the third word was presented. At the beginning of each trial, a response area of
612 600 (horizontal) x 450 (vertical) pixels and a solid line at the top of the rectangle were presented
613 (**Fig.1A**). Participants were informed about the scale (e.g., quality) and the direction of the scale
614 (e.g., “bad” to “good” or “good” to “bad”, i.e., 1 to 10 or 10 to 1). Participants were instructed to
615 click on the “start” button and move the cursor of the mouse to the portion of the scale that best
616 represented the overall meaning of the phrase. The “start” button was placed in the center portion
617 of the bottom of the response space (i.e., in a neutral position). Once “start” was clicked on,
618 information about the scale and scale direction disappeared, leaving only the solid line on screen.

619 Phrases were presented at the top of the response space, from the time when participants clicked on
620 “start”, one word at a time, each word for 250 ms (inter-word-interval: 50 ms). After each trial,
621 participants were provided the “incorrect” feedback if the cursor’s movement violated the warnings
622 provided during the familiarization phase, and an explanation was provided (e.g., “you crossed the
623 vertical borders”). To keep participants engaged, we provided feedback also based on the final
624 interpretation: “incorrect” if the response was in the half of the scale opposite to the adjective (for
625 the conditions: “### ##”, “#### really”, “really ###” and “really really”), or in the same half of
626 the scale of the adjective (for the conditions: “### not” or “not ##”), or in the outer 20% left and
627 right portions of the scale (for the conditions: “really not”, “not really” and “not not”); feedback
628 was “correct” otherwise. In case of an “incorrect” trial, the following trial was delayed for 4
629 seconds. For each trial, we collected continuous mouse trajectories and RTs. The overall duration
630 of the behavioral experiment was approximately 90 minutes. To verify that the feedback did not
631 affect our results, we ran a replication study with 55 online participants where no feedback was
632 provided based on the final interpretation (**Fig.S2**).

633
634 *Experiment 2: MEG.*
635 *Stimuli and Design.* The linguistic stimulus set comprised 72 unique adjective phrases (for the
636 complete list, see **Table S2**). Similar to the Experiment 1, adjectives were selected for being
637 antonyms (and *contraries*) in the following cognitive or sensory dimensions: *quality* (“bad”,
638 “good”), *temperature* (“cool”, “warm”), *loudness* (“quiet”, “loud”), and *brightness* (“dark”,
639 “bright”). Lexical characteristics of the antonyms were balanced according to the English Lexicon
640 Project⁽⁸⁵⁾; mean (SD) HAL log frequency of “low” adjectives: 10.85 (1.03), “high” adjectives:
641 10.55 (1.88); mean (SD) bigram frequency of “low” adjectives: 1196.5 (824.6), “high” adjectives:
642 1077.5 (376.3); mean (SD) lexical decision RTs of “low” adjectives: 594 ms (39), “high” adjectives:
643 594 (33)). Adjectives were combined with zero (e.g., “### ##”), one (e.g., “really ##”) or two
644 modifiers (e.g., “really not”). Modifiers were either the intensifier “really” or the negation “not”. A
645 sequence of dashes was used to indicate the absence of a modifier, e.g., “really ### good”. Each of
646 the eight adjectives was preceded by each of the nine possible combinations of modifiers: “###
647 ##”, “#### really”, “really ###”, “### not”, “not ##”, “really not”, “not really”, “really really”
648 and “not not” (“not not” was included to achieve a full experimental design, even if it is not a
649 frequent combination in natural language). To avoid possible differences in neural representation
650 of phrases with and without syntactic/semantic composition, the condition with no modifiers (“###
651 ##”) was exclusively employed as a baseline comparison in the time-frequency analysis and was
652 excluded from all other analyses. Each dimension (e.g., quality) was presented in two blocks, one

653 block for each yes/no key orientation (8 blocks in total, see Procedure). Each phrase (e.g., “really
654 really bad”) was repeated four times within one block. Thus, the overall experiment comprised 576
655 trials. The order of phrases was randomized within each block for each participant. The order of
656 blocks was randomized across participants within the first and second half of the experiment. The
657 yes/no order was randomized across participants.

658
659 *Procedure.* Participants were familiarized with the linguistic stimuli through a short practice block
660 that mimicked the structure of the experimental blocks. They were instructed to read affirmative or
661 negated adjective phrases (e.g., “really really good”, “really not bad”) and derive the overall
662 meaning of each adjective phrase, on a scale from 0 to 8, e.g., from “really really bad” to “really
663 really good”. Each trial started with a fixation cross (duration: 750 ms), followed by each phrase
664 presented one word at a time, each word for 100 ms (inter-word-interval: 250 ms, **Fig.1B**). After
665 each phrase, a fixation cross was presented for 1500 ms. A number (i.e., probe) was then presented,
666 which did or did not correspond to the overall meaning of the adjective phrase on the scale.
667 Participants were required to indicate whether the probe number correctly represented the meaning
668 of the phrase on the scale (*yes/no* answer). The yes/no order was swapped halfway through the
669 experiment. Responses had no time limit. If correct (+/- one step on the scale), a green fixation
670 cross was presented; if incorrect, a red fixation cross was presented, and feedback was provided.
671 While performing the experiment, participants lay supine in a magnetically shielded room while
672 continuous MEG data were recorded through a 157-channel whole-head axial gradiometer system
673 (Kanazawa Institute of Technology, Kanazawa, Japan). Sampling rate was 1000 Hz, and online
674 high-pass filter of 1 Hz and low-pass filter of 200 Hz were applied. Five electromagnetic coils were
675 attached to the forehead of the participants and their position was measured twice, before the first
676 and after the last block. Instructions, visual stimuli and visual feedback were back-projected onto a
677 Plexiglas screen using a Hitachi projector. Stimuli were presented using Psychtoolbox v3⁽⁸⁸⁾;
678 www.psychtoolbox.org), running under MATLAB R2019a (MathWorks) on an Apple iMac model
679 10.12.6. Participants responded to the yes/no question with their index finger of their left and right
680 hand, using a keypad. For each trial, we also collected accuracy and RTs. The overall duration of
681 the MEG experiment was approximately 60 minutes.

682

683

684 *Data analysis*

685 *Experiment 1 (and replication): RTs and mouse trajectories data.*

686 The RTs and mouse trajectory analyses were limited to correct trials (group mean accuracy: 82%,
687 SD: 13%), and RTs were limited within the range of participant median RTs \pm 2 SD.

688 To evaluate differences in RTs between antonyms (“small”, “cold”, “ugly”, “bad”, “sad” vs. “big”,
689 “hot”, “beautiful”, “good”, “happy”, “fast”, i.e., *low* vs. *high* poles in each scalar dimension), and
690 between negated and affirmative phrases (e.g., “really really good” vs. “really not good”), and their
691 interactions, median RTs of each participant were entered into 2 (*antonym*: low vs. high) x 2
692 (*negation*: negated vs. affirmative) repeated-measures ANOVA.

693 To evaluate differences in the final interpretations between antonyms in each scale, between
694 negated and affirmative phrases, and their interactions, mean and standard deviation of the final
695 responses of each participant were entered into a 2 (*antonym*: low vs. high) x 2 (*negation*: negated
696 vs. affirmative) repeated-measures ANOVA. Post-hoc tests were conducted for significant
697 interactions (correction = Holm). Effect sizes were calculated using partial eta squared (η_p^2).

698 To compare mouse trajectories over time across participants, we resampled participants’
699 mouse trajectories at 100 Hz using linear interpolation, up to 2 seconds, to obtain 200 time points
700 for each trial. Furthermore, trajectories were normalized between -1 and 1. For visualization
701 purposes, we computed the median of trajectories across trials for each participant, dimension (e.g.,
702 quality), antonym (e.g., “bad”) and modifier (e.g., “really not”), and at each timepoint.

703 Finally, to quantitatively evaluate how the interpretation of each phrase changed over time,
704 for every participant we carried out regression analyses per each time point, for affirmative and
705 negated phrases separately (for a similar approach, see ⁴⁵). The dependent variable was the mouse
706 coordinate along the scale (note that the scale which was swapped in half of the trials was swapped
707 back for data analysis purposes), and the predictor was whether the adjective was a low or high
708 antonym (e.g., “bad” vs. “good”). To identify the time windows where predictors were significantly
709 different from 0 at the group level, we performed permutation cluster tests on beta values (10,000
710 permutations) in the time window from the onset of the adjective up to 1.4 s from adjective onset
711 (i.e., 2 s from the onset of word 1).

712

713 *Experiment 2: Accuracy and RTs data.*

714 To evaluate differences in accuracy between *low* and *high* antonyms (“bad”, “cool”, “quiet”, “dark”
715 vs. “good”, “warm”, “loud”, “bright”), and between negated and affirmative phrases (e.g., “really
716 really good” vs. “really not good”), and their interactions, mean accuracies in the yes/no task of
717 each participant were entered into 2 (*antonym*: low vs. high) x 2 (*negation*: negated vs. affirmative)
718 repeated-measures ANOVA.

719 The response time analysis was limited to correct trials. RTs outside the range of participant
720 median RTs ± 2 SD were removed. To evaluate differences in RTs between *low* and *high* antonyms
721 in each scale and between negated and affirmative phrases, and their interactions, median RTs of
722 each participant in the yes/no task were entered into a 2 (*antonym*: low vs. high) x 2 (*negation*:
723 negated vs. affirmative) repeated-measures ANOVA.

724

725 *Experiment 2: MEG data.*

726 *Preprocessing.*

727 MEG data preprocessing was performed using MNE-python⁸⁹ and Eelbrain
728 (10.5281/zenodo.438193). First, bad channels (i.e., below the 3rd or above the 97th percentile
729 across all channels, for more than 20% of the entire recording) were interpolated. The MEG
730 responses were denoised by applying least square projections of the reference channels and
731 removing the corresponding components from the data⁹⁰. Denoised data were lowpass-filtered at
732 20 Hz for the decoding analyses and at 40 Hz for the time-frequency analyses. FastICA was used
733 to decompose the signal into independent components, to visually inspect and remove artifacts
734 related to eye-blinks, heartbeat and external noise sources. MEG recordings were then epoched into
735 epochs of -300 ms and 2550 ms around the onset of the first, second, or third word (or probe) for
736 the decoding analyses, and into epochs of -800 and 3000 ms around the onset of the first word for
737 the time-frequency analyses (and then cut between -300 and 2550 ms for group analyses). Note
738 that, for visualization purposes, only 1700 ms from the onset of the first word (i.e., 1000 ms from
739 adjective onset) were included in most figures (as no significant results were observed for control
740 analyses run for later time windows). Finally, epochs with amplitudes greater than an absolute
741 threshold of 3000 fT were removed and a baseline between -300 to 0 ms was applied to all epochs.

742

743 *Source reconstruction.*

744 Structural magnetic resonance images (MRIs) were collected for 10 out of 26 participants. For the
745 remaining 16 participants, we manually scaled and co-registered the “fsaverage” brain to the
746 participant’s head-digitalized shape and fiducials^{89,91}.

747 For every participant, an ico-4 source space was computed, containing 2562 vertices per hemisphere
748 and the forward solution was calculated using the Boundary Element Model (BEM). A noise
749 covariance matrix was estimated from the 300 ms before the onset of the first word. The inverse
750 operator was created and applied to the neuromagnetic data to estimate the source time courses at
751 each vertex using dynamic statistical parametric mapping (dSPM:⁹²). The results were then
752 morphed to the ico-5 “fsaverage” brain, yielding to time courses for 10242 vertices per hemisphere.

753 We then estimated the magnitude of the activity at each vertex (signal to noise ratio: 3, lambda2:
754 0.11, with orientation perpendicular to the cortical surface), which was used in the decoding
755 analyses (*Spatial decoders*).

756

757 *Decoding analyses.*

758 Decoding analyses were limited to correct trials and were performed with the MNE⁸⁹ and Scikit-
759 Learn packages⁴⁷. First, X (or the selected principal components) were set to have zero mean and
760 unit variance (i.e., using a standard scaler). Second, we fitted a l2 linear estimator to a subset of the
761 epochs (training set, X_{train}) and estimated y on a separate group of epochs (test set, \hat{y}_{test}). We then
762 computed the accuracy (AUC, see below) of the decoder, by comparing \hat{y}_{test} with the ground truth
763 y. For this analysis, we used the default values provided by the Scikit-Learn package and set the
764 class-weight parameter to “balanced”.

765

766 *Temporal decoders.* Temporal decoding analyses were performed in sensor-space. Before fitting
767 the estimators, linear dimensionality reduction (principal component analysis, PCA) was performed
768 on the channel amplitudes to project them to a lower dimensional space (i.e., to new virtual channels
769 that explained more than 99% of the feature variance). We then fitted the linear estimator on each
770 participant separately, across all selected components, at each time-point separately. Time was
771 subsampled to 100 Hz. We then employed a 5-fold stratified cross-validation (or 10-fold, depending
772 on the number of trials per class), that fitted the linear estimator to 80% (or 90%) of the epochs and
773 generated predictions on 20% (or 10%) of the epochs, while keeping the distributions of the training
774 and test set maximally homogeneous. This decoding approach was used for analyses of **Fig.3B**,
775 **Fig.3C**, **Fig.4A**, **Fig.4C** and decoding approach (ii) in **Fig.5C**. To investigate whether the
776 representation of antonyms was comparable between affirmative and negated phrases, in a different
777 set of analyses (i.e., decoding approach (i), **Fig.5B**) we fitted the linear estimator to all epochs
778 corresponding to affirmative phrases and generated predictions on all epochs corresponding to
779 negated phrases. In both decoding approaches, accuracy and probability estimates for each class
780 were then computed. Decoding accuracy is summarized with an empirical area under the curve
781 (rocAUC, 0 to 1, chance at 0.5).

782 At the group level, we extracted the clusters of time where AUC across participants was
783 significantly higher than chance using a one-sample permutation cluster test, as implemented in
784 MNE-python (10000 permutations⁹³). We performed separate permutation cluster tests for the
785 following time windows: -700 to -350 ms from adjective onset (i.e., word 1), -350 to 0 ms from
786 adjective onset (i.e., word 2), 0 to 500 ms from adjective onset (i.e., time window for lexical-

787 semantic processes^{65,66}) and 500 to 1000 ms from adjective onset (i.e., to account for potential later
788 processes).

789

790 *Expected outcome for the effect of negation on the representation of antonyms.* Temporal decoding
791 approach (i) and (ii) described above allow us to make specific predictions about the effect of
792 negation on the representation of antonyms (**Fig.5A**).

793 *Approach (i)* train set: affirmative phrases; test set: negated phrases. For our results to
794 support predictions (1) *No effect* or (2) *Mitigation*, this decoding approach should show probability
795 estimates of high and low adjectives significantly above the computed chance level and in the
796 direction of the respective classes, indicating that the initial representation of adjectives in negated
797 phrases is similar to that in affirmative phrases (left column, first and second row under *decoding*
798 *approach* in **Fig.5A**). Conversely, for our results to support prediction (3) *Inversion*, this decoding
799 approach should show probability estimates of high and low adjectives significantly above the
800 computed chance level but in the direction of the opposite classes (i.e., swapped), as adjective
801 representations would be systematically inverted in negated phrases (left column, third row under
802 *decoding approach* in **Fig.5A**). Finally, we should observe at chance probability estimates in the
803 case of (4) *Change*, where adjective representations in negated phrases are not predictable from the
804 corresponding representations in affirmative phrases (left column, fourth row under *decoding*
805 *approach* in **Fig.5A**).

806 *Approach (ii)* train set: affirmative and negated phrases together; test set: affirmative and
807 negated phrases separately. This decoding analysis allows us to disentangle predictions (1) *No effect*
808 from (2) *Mitigation*. For the results of this analysis to support prediction (1) *No effect*, we should
809 observe quantitatively comparable probability estimates in affirmative and negated phrases,
810 suggesting that negation does not change the representation of adjectives (right column, first row
811 under *decoding approach* in **Fig.5A**). Conversely, in support of prediction (2) *Mitigation*, we
812 should observe significantly reduced probability estimates for negated relative to affirmative
813 phrases, suggesting less robust differences between low and high antonyms in negated phrases
814 (right column, second row under *decoding approach* in **Fig.5A**). The outcome of predictions (3)
815 *Inversion* and (4) *Change* would be at chance probability estimates (as the model is trained on
816 different representations within the same class; right column, third and fourth row under *decoding*
817 *approach* in **Fig.5A**).

818 *Spatial decoders.* Spatial decoding analyses were performed in source-space. We fitted each
819 estimator on each participant separately, across 50 to 650 ms time samples relative to the onset of
820 the adjective (to include the three significant time windows that emerge from the temporal decoding

821 analysis in **Fig.3B**), at each brain source separately, after morphing individual participant's source
822 estimates to the ico-5 "fsaverage" common reference space. We employed a 5-fold stratified cross-
823 validation, which fitted the linear estimator to 80% of the epochs and generated predictions on 20%
824 of the epochs, while keeping the distributions of the training and test set maximally homogeneous.
825 Decoding accuracy is summarized with an empirical area under the curve (AUC, 0 to 1, chance at
826 0.5). At the group level, we extracted the brain areas where the AUC across participants was
827 significantly higher than chance, using a one-sample permutation cluster test as implemented in
828 MNE-python (10000 permutations; adjacency computed from the "fsaverage" brain⁹³).

829

830 *Time-frequency analysis.*

831 We extracted time-frequency power of the epochs (-800 to 3000 ms from the onset of word 1) using
832 Morlet wavelets of 3 cycles per frequency, in frequencies between 3.9 and 37.2 Hz, logarithmically
833 spaced (19 frequencies overall). Power estimates were then cut between -300 and 2550 ms from
834 onset of word 1 and baseline corrected using a window of -300 to -100 ms from the onset of word
835 1, by subtracting the mean of baseline values and dividing by the mean of baseline values (mode =
836 'percent'). Power in the low-beta frequency range (12 to 20 Hz) and in the high-beta frequency
837 range (21 to 30 Hz^{57,79}) was averaged to obtain a time course of power in low and high-beta
838 rhythms. We then subtracted the beta power of affirmative phrases from that of negated phrases. At
839 the group level, we extracted the clusters of time where this difference in power across participants
840 was significantly greater than 0, using a one-sample permutation cluster test as implemented in
841 MNE-python (10000 permutations⁹³). We performed separate permutation cluster tests in the same
842 time windows used for the decoding analysis: -700 to -350 ms, -350 to 0 ms, 0 to 500 ms, and 500
843 to 1000 ms from the onset of the adjective (note that no significant differences were observed in
844 analyses ran for time windows after 1000 ms). We then computed the induced power in source
845 space (method: dSPM and morphing individual participant's source estimates to the ico-5
846 "fsaverage" reference space) for the significant clusters of time in the low- and high-beta range
847 separately and averaged over time. At the group level, we extracted the brain areas where the power
848 difference across participants was significantly greater than 0, using a one-sample permutation
849 cluster test as implemented in MNE-python (10000 permutations; adjacency computed from the
850 "fsaverage" brain⁹³).

851

852

853

854

855 **References**

- 856 1. Ding, N., Melloni, L., Zhang, H., Tian, X. & Poeppel, D. Cortical tracking of hierarchical linguistic structures in
857 connected speech. *Nature Neuroscience* **19**, 158–164 (2015).
- 858 2. Fedorenko, E. *et al.* Neural correlate of the construction of sentence meaning. *Proceedings of the National*
859 *Academy of Sciences of the United States of America* **113**, E6256–E6262 (2016).
- 860 3. Martin, A. E. & Baggio, G. Modelling meaning composition from formalism to mechanism. 1–7 (2019).
- 861 4. Matchin, W. & Hickok, G. The Cortical Organization of Syntax. *Cerebral Cortex* **30**, 1481–1498 (2020).
- 862 5. Oseki, Y. & Marantz, A. Modeling morphological processing in human magnetoencephalography. *Proceedings of*
863 *the Society for Computation in Linguistics* **3**, (2020).
- 864 6. Pallier, C., Devauchelle, A.-D. & Dehaene, S. Cortical representation of the constituent structure of sentences.
865 *Proceedings of the National Academy of Sciences* **108**, 2522–2527 (2011).
- 866 7. Pylkkänen, L. The neural basis of combinatory syntax and semantics. *Science* **366**, 62–66 (2019).
- 867 8. Ziegler, J. & Pylkkänen, L. Scalar adjectives and the temporal unfolding of semantic composition: An MEG
868 investigation. *Neuropsychologia* **89**, 161–171 (2016).
- 869 9. Tian, Y., Ferguson, H. & Breheny, R. Processing negation without context – why and when we represent the
870 positive argument. *Language, Cognition and Neuroscience* **31**, 683–698 (2016).
- 871 10. Tian, Y., Breheny, R. & Ferguson, H. J. Why we simulate negated information: A dynamic pragmatic account.
872 *Quarterly Journal of Experimental Psychology* **63**, 2305–2312 (2010).
- 873 11. Giora, R. Anything negatives can do affirmatives can do just as well, except for some metaphors. *Journal of*
874 *Pragmatics* **38**, 981–1014 (2006).
- 875 12. Horn, L. R. *A natural history of negation*. (University of Chicago Press, 1989).
- 876 13. Ettinger, A. What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models.
877 *Transactions of the Association for Computational Linguistics* **8**, 34–48 (2020).
- 878 14. Dale, R. & Duran, N. D. The cognitive dynamics of negated sentence verification. *Cognitive Science* **35**, 983–996
879 (2011).
- 880 15. Darley, E. J., Kent, C. & Kazanina, N. A ‘no’ with a trace of ‘yes’: A mouse-tracking study of negative sentence
881 processing. *Cognition* **198**, 104084 (2020).
- 882 16. Dudschig, C. & Kaup, B. How does ‘not left’ become ‘right’? Electrophysiological evidence for a dynamic
883 conflict-bound negation processing account. *Journal of Experimental Psychology: Human Perception and*
884 *Performance* **44**, 716–728 (2018).

- 885 17. Just, M. A. & Carpenter, P. A. Comprehension of negation with quantification. *Journal of Verbal Learning and*
886 *Verbal Behavior* **10**, 244–253 (1971).
- 887 18. Kaup, B., Yaxley, R. H., Madden, C. J., Zwaan, R. A. & Ldtke, J. Experiential simulations of negated text
888 information. *Quarterly Journal of Experimental Psychology* **60**, 976–990 (2007).
- 889 19. Dudschig, C., Kaup, B., Liu, M. & Schwab, J. The processing of negation and polarity: An overview. *Journal of*
890 *Psycholinguistic Research* **50**, 1199–1213 (2021).
- 891 20. Sherman, M. A. Adjectival negation and the comprehension of multiply negated sentences. *Journal of Verbal*
892 *Learning and Verbal Behavior* **15**, 143–157 (1976).
- 893 21. Kaup, B. Negation and its impact on the accessibility of text information. *Memory and Cognition* **29**, 960–967
894 (2001).
- 895 22. Kaup, B. & Zwaan, R. A. Effects of negation and situational presence on the accessibility of text information.
896 *Journal of Experimental Psychology: Learning Memory and Cognition* **29**, 439–446 (2003).
- 897 23. MacDonald, M. C. & Just, M. A. Changes in activation levels with negation. *Journal of Experimental Psychology:*
898 *Learning, Memory, and Cognition* **15**, 633–642 (1989).
- 899 24. Carpenter, P. A. & Just, M. A. Sentence comprehension: A psycholinguistic processing model of verification.
900 *Psychological Review* **82**, 45–73 (1975).
- 901 25. Clark, H. H. & Chase, W. G. On the process of comparing sentences against pictures. *Cognitive Psychology* **3**,
902 472–517 (1972).
- 903 26. Lüdtke, J., Friedrich, C. K., De Filippis, M. & Kaup, B. Event-related potential correlates of negation in a
904 sentence-picture verification paradigm. *Journal of Cognitive Neuroscience* **20**, 1355–1370 (2008).
- 905 27. Kaup, B. & Dudschig, C. Understanding negation: Issues in the processing of negation. in *The Oxford Handbook*
906 *of Negation* (eds. Déprez, V. & Espinal, M. T.) 634–655 (Oxford University Press, 2020).
- 907 28. Papeo, L. & de Vega, M. The neurobiology of lexical and sentential negation. *The Oxford Handbook of Negation*
908 739–756 (2020).
- 909 29. Papeo, L., Hochmann, J.-R. & Battelli, L. The default computation of negated meanings. *Journal of Cognitive*
910 *Neuroscience* **28**, 1980–1986 (2016).
- 911 30. Lyons, J. *Linguistic semantics: An introduction*. (Cambridge University Press, 1995).
- 912 31. Mayo, R., Schul, Y. & Burnstein, E. ‘I am not guilty’ vs ‘I am innocent’: Successful negation may depend on the
913 schema used for its encoding. *Journal of Experimental Social Psychology* **40**, 433–449 (2004).
- 914 32. Orenes, I., Beltrán, D. & Santamaría, C. How negation is understood: Evidence from the visual world paradigm.
915 *Journal of Memory and Language* **74**, 36–45 (2014).

- 916 33. van Gaal, S. *et al.* Can the meaning of multiple words be integrated unconsciously? *Phil. Trans. R. Soc. B* **369**,
917 20130212 (2014).
- 918 34. Bartoli, E. *et al.* The disembodiment effect of negation: Negating action-related sentences attenuates their
919 interference on congruent upper limb movements. *Journal of Neurophysiology* **109**, 1782–1792 (2013).
- 920 35. Beltrán, D., Morera, Y., García-Marco, E. & De Vega, M. Brain inhibitory mechanisms are involved in the
921 processing of sentential negation, regardless of its content. Evidence from EEG theta and beta rhythms. *Frontiers*
922 *in Psychology* **10**, 1–14 (2019).
- 923 36. Beltrán, D., Liu, B. & de Vega, M. Inhibitory mechanisms in the processing of negations: A neural reuse
924 hypothesis. *Journal of Psycholinguistic Research* **50**, 1243–1260 (2021).
- 925 37. De Vega, M. *et al.* Sentential negation might share neurophysiological mechanisms with action inhibition.
926 Evidence from frontal theta rhythm. *Journal of Neuroscience* **36**, 6002–6010 (2016).
- 927 38. Djokic, V., Maillard, J., Bulat, L. & Shutova, E. Modeling affirmative and negated action processing in the brain
928 with lexical and compositional semantic models. 5155–5165 (2019).
- 929 39. Gallese, V. & Lakoff, G. The brain's concepts: The role of the sensory-motor system in conceptual knowledge.
930 *Cognitive Neuropsychology* **22**, 455–479 (2005).
- 931 40. Tettamanti, M. *et al.* Negation in the brain: Modulating action representations. *NeuroImage* **43**, 358–367 (2008).
- 932 41. Tomasino, B., Weiss, P. H. & Fink, G. R. To move or not to move: Imperatives modulate action-related verb
933 processing in the motor system. *Neuroscience* **169**, 246–258 (2010).
- 934 42. Bianchi, I., Savardi, U., Burro, R. & Torquati, S. Negation and psychological dimensions. *Journal of Cognitive*
935 *Psychology* **23**, 275–301 (2011).
- 936 43. Colston, H. L. 'Not good' is 'bad,' but 'not bad' is not 'good': an analysis of three accounts of negation
937 asymmetry. *Discourse Processes* **28**, 237–256 (1999).
- 938 44. Fraenkel, T. & Schul, Y. The meaning of negated adjectives. *Intercultural Pragmatics* **5**, 517–540 (2008).
- 939 45. Dotan, D. & Dehaene, S. How do we convert a number into a finger trajectory? *Cognition* **129**, 512–529 (2013).
- 940 46. Maldonado, M., Dunbar, E. & Chemla, E. Mouse tracking as a window into decision making. *Behav Res* **51**,
941 1085–1101 (2019).
- 942 47. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–
943 2830 (2011).
- 944 48. Caucheteux, C. & King, J.-R. Brains and algorithms partially converge in natural language processing. *Commun*
945 *Biol* **5**, 134 (2022).

- 946 49. Binder, J. R., Desai, R. H., Graves, W. W. & Conant, L. L. Where is the semantic system? A critical review and
947 meta-analysis of 120 functional neuroimaging studies. *Cerebral Cortex* **19**, 2767–2796 (2009).
- 948 50. Caucheteux, C., Gramfort, A. & King, J.-R. Disentangling syntax and semantics in the brain with deep networks.
949 (2021).
- 950 51. Huth, A. G., de Heer, W. A., Griffiths, T. L., Theunissen, F. E. & Gallant, J. L. Natural speech reveals the
951 semantic maps that tile human cerebral cortex. *Nature* **532**, 453–458 (2016).
- 952 52. Lau, E. F., Gramfort, A., Hämäläinen, M. S. & Kuperberg, G. R. Automatic semantic facilitation in anterior
953 temporal cortex revealed through multimodal neuroimaging. *Journal of Neuroscience* **33**, 17174–17181 (2013).
- 954 53. Lambon Ralph, M. A., Jefferies, E., Patterson, K. & Rogers, T. T. The neural and computational bases of semantic
955 cognition. *Nature Reviews Neuroscience* **18**, 42–55 (2016).
- 956 54. Hagoort, P., Hald, L., Bastiaansen, M. & Petersson, K. M. Integration of word meaning and world knowledge in
957 language comprehension. *Science* **304**, 438–441 (2004).
- 958 55. Popham, S. F. *et al.* Visual and linguistic semantic representations are aligned at the border of human visual
959 cortex. *Nat Neurosci* **24**, 1628–1636 (2021).
- 960 56. Parrish, A. & Pyllkkänen, L. Conceptual combination in the LATL with and without syntactic composition.
961 *Neurobiology of Language* **3**, 46–66 (2022).
- 962 57. Weiss, S. & Mueller, H. M. ‘Too many betas do not spoil the broth’: The role of beta brain oscillations in
963 language processing. *Frontiers in Psychology* **3**, 1–15 (2012).
- 964 58. Wagner, J., Wessel, J. R., Ghahremani, A. & Aron, A. R. Establishing a right frontal beta signature for stopping
965 action in scalp EEG: Implications for testing inhibitory control in other task contexts. *Journal of Cognitive*
966 *Neuroscience* **30**, 107–118 (2018).
- 967 59. Alemanno, F. *et al.* Action-related semantic content and negation polarity modulate motor areas during sentence
968 reading: An event-related desynchronization study. *Brain Research* **1484**, 39–49 (2012).
- 969 60. Aravena, P. *et al.* Grip force reveals the context sensitivity of language-induced motor activity during “action
970 words” processing: evidence from sentential negation. *PLoS ONE* **7**, e50287 (2012).
- 971 61. Foroni, F. & Semin, G. R. Comprehension of action negation involves inhibitory simulation. *Frontiers in Human*
972 *Neuroscience* **7**, 1–7 (2013).
- 973 62. Ghio, M., Haegert, K., Vaghi, M. M. & Tettamanti, M. Sentential negation of abstract and concrete conceptual
974 categories: A brain decoding multivariate pattern analysis study. *Philosophical Transactions of the Royal Society*
975 *B: Biological Sciences* **373**, 7–10 (2018).

- 976 63. Liuzza, M. T., Candidi, M. & Aglioti, S. M. Do not resonate with actions: Sentence polarity modulates cortico-
977 spinal excitability during action-related sentence reading. *PLoS ONE* **6**, 38–41 (2011).
- 978 64. Hauk, O., Davis, M. H., Ford, M., Pulvermüller, F. & Marslen-Wilson, W. D. The time course of visual word
979 recognition as revealed by linear regression analysis of ERP data. *NeuroImage* **30**, 1383–1400 (2006).
- 980 65. Kutas, M. & Federmeier, K. D. Thirty years and counting: Finding meaning in the N400 component of the event-
981 related brain potential (ERP). *Annual Review of Psychology* **62**, 621–647 (2011).
- 982 66. Pulvermüller, F., Shtyrov, Y. & Hauk, O. Understanding in an instant: Neurophysiological evidence for
983 mechanistic language circuits in the brain. *Brain and Language* **110**, 81–94 (2009).
- 984 67. Pulvermüller, F., Assadollahi, R. & Elbert, T. Neuromagnetic evidence for early semantic access in word
985 recognition. *European Journal of Neuroscience* **13**, 201–205 (2001).
- 986 68. Teige, C. *et al.* Dynamic semantic cognition: Characterising coherent and controlled conceptual retrieval through
987 time using magnetoencephalography and chronometric transcranial magnetic stimulation. *Cortex* **103**, 329–349
988 (2018).
- 989 69. Zhang (张琳敏), L. & Pykkänen, L. Semantic composition of sentences word by word: MEG evidence for shared
990 processing of conceptual and logical elements. *Neuropsychologia* **119**, 392–404 (2018).
- 991 70. Fyshe, A., Sudre, G., Wehbe, L., Rafidi, N. & Mitchell, T. M. The lexical semantics of adjective–noun phrases in
992 the human brain. *Human Brain Mapping* **40**, 4457–4469 (2019).
- 993 71. Nieuwland, M. S. & Kuperberg, G. R. When the truth is not too hard to handle: An event-related potential study
994 on the pragmatics of negation. *Psychological Science* **19**, 1213–1218 (2008).
- 995 72. Palaz, B., Rhodes, R. & Hestvik, A. Informative use of “not” is N400-blind. *Psychophysiology* **57**, (2020).
- 996 73. Xiang, M., Grove, J. & Giannakidou, A. Semantic and pragmatic processes in the comprehension of negation: An
997 event related potential study of negative polarity sensitivity. *Journal of Neurolinguistics* **38**, 71–88 (2016).
- 998 74. Grice, H. P. Logic and Conversation. in *Syntax and Semantics* vol. 3 41–58 (New York: Academic Press, 1975).
- 999 75. Bastiaansen, M. C. M., van der Linden, M., ter Keurs, M., Dijkstra, T. & Hagoort, P. Theta responses are involved
1000 in lexical—semantic retrieval during language processing. *Journal of Cognitive Neuroscience* **17**, 530–541
1001 (2005).
- 1002 76. Luo, Y., Zhang, Y., Feng, X. & Zhou, X. Electroencephalogram oscillations differentiate semantic and prosodic
1003 processes during sentence reading. *Neuroscience* **169**, 654–664 (2010).
- 1004 77. Supp, G. G. *et al.* Lexical memory search during N400: cortical couplings in auditory comprehension:
1005 *NeuroReport* **15**, 1209–1213 (2004).

- 1006 78. Weiss, S. & Rappelsberger, P. EEG coherence within the 13–18 Hz band as a correlate of a distinct lexical
1007 organisation of concrete and abstract nouns in humans. *Neuroscience Letters* **209**, 17–20 (1996).
- 1008 79. Schaller, F., Weiss, S. & Müller, H. M. EEG beta-power changes reflect motor involvement in abstract action
1009 language processing. *Brain and Language* **168**, 95–105 (2017).
- 1010 80. Pinheiro-Chagas, P., Dotan, D., Piazza, M. & Dehaene, S. Finger tracking reveals the covert stages of mental
1011 arithmetic. *Open Mind* **1**, 30–41 (2017).
- 1012 81. Peer, E., Brandimarte, L., Samat, S. & Acquisti, A. Beyond the Turk: Alternative platforms for crowdsourcing
1013 behavioral research. *Journal of Experimental Social Psychology* **70**, 153–163 (2017).
- 1014 82. Simcox, T. & Fiez, J. A. Collecting response times using Amazon Mechanical Turk and Adobe Flash. *Behavior*
1015 *Research Methods* **46**, 95–111 (2014).
- 1016 83. Gwilliams, L. & King, J. R. Recurrent processes support a cascade of hierarchical decisions. *eLife* **9**, 1–20 (2020).
- 1017 84. King, J. R., Pescetelli, N. & Dehaene, S. Brain mechanisms underlying the brief maintenance of seen and unseen
1018 sensory information. *Neuron* **92**, 1122–1134 (2016).
- 1019 85. Balota, D. A. *et al.* The English Lexicon Project. *Behavior Research Methods* **39**, 445–459 (2007).
- 1020 86. Schiller, N. O. *et al.* Solving the problem of double negation is not impossible: electrophysiological evidence for
1021 the cohesive function of sentential negation. *Language, Cognition and Neuroscience* **32**, 147–157 (2017).
- 1022 87. Chen, D. L., Schonger, M. & Wickens, C. oTree—An open-source platform for laboratory, online, and field
1023 experiments. *Journal of Behavioral and Experimental Finance* **9**, 88–97 (2016).
- 1024 88. Brainard, D. H. The Psychophysics Toolbox. *Spatial Vision* **10**, 433–436 (1997).
- 1025 89. Gramfort, A. *et al.* MEG and EEG data analysis with MNE-Python. *Frontiers in Neuroscience* **7**, 1–13 (2013).
- 1026 90. Adachi, Y., Shimogawara, M., Higuchi, M., Haruta, Y. & Ochiai, M. Reduction of non-periodic environmental
1027 magnetic noise in MEG measurement by Continuously Adjusted Least squares Method. *IEEE Transactions on*
1028 *Applied Superconductivity* **11**, 669–672 (2001).
- 1029 91. Andersen, L. M. Group analysis in MNE-python of evoked responses from a tactile stimulation paradigm: A
1030 pipeline for reproducibility at every step of processing, going from individual sensor space representations to an
1031 across-group source space representation. *Frontiers in Neuroscience* **12**, (2018).
- 1032 92. Dale, A. M. *et al.* Dynamic statistical parametric mapping: Combining fMRI and MEG for high-resolution
1033 imaging of cortical activity. *Neuron* **26**, 55–67 (2000).
- 1034 93. Maris, E. & Oostenveld, R. Nonparametric statistical testing of EEG- and MEG-data. *Journal of Neuroscience*
1035 *Methods* **164**, 177–190 (2007).
- 1036

1037 **Acknowledgements**

1038 This work was supported by the Leon Levy Foundation (A.Z.), the European Union's Horizon 2020
1039 research and innovation program under grant agreement No 660086 (J.R.K.), the Bettencourt-
1040 Schueller Foundation (J.R.K.), the Fondation Roger de Spoelberch (J.K.R.), the Philippe
1041 Foundation (J.R.K.), the FrontCog grant ANR-17-EURE-0017 (J.R.K.), and the Ernst Struengmann
1042 Foundation (D.P.).

1043

1044 **Author contributions**

1045 AZ, PR, JRK, and DP conceptualized the experiment; AZ, PR, and WML collected the data; AZ
1046 analyzed the data; PR, LG, and JRK contributed to analysis; AZ wrote the paper; AZ, PR, LG, JRK,
1047 and DP discussed the results and edited the paper.

1048

1049 **Competing interests**

1050 The authors declare no competing interests.

1051

1052 **Data and materials availability**

1053 All data needed to evaluate the conclusions in the paper are present in the paper and/or the
1054 Supplementary Materials. Additional data related to this paper may be requested from the
1055 corresponding author.

1056

1057 **Supplementary Materials**

1058 **Tables**

1059

List of linguistic stimuli employed in Experiment 1 (behavior)					
### ###	small	really really	small	not not	small
### ###	big	really really	big	not not	big
### ###	cold	really really	cold	not not	cold
### ###	hot	really really	hot	not not	hot
### ###	ugly	really really	ugly	not not	ugly
### ###	beautiful	really really	beautiful	not not	beautiful
### ###	bad	really really	bad	not not	bad
### ###	good	really really	good	not not	good
### ###	sad	really really	sad	not not	sad
### ###	happy	really really	happy	not not	happy
### ###	slow	really really	slow	not not	slow
### ###	fast	really really	fast	not not	fast
### really	small	### not	small	really not	small
### really	big	### not	big	really not	big
### really	cold	### not	cold	really not	cold
### really	hot	### not	hot	really not	hot
### really	ugly	### not	ugly	really not	ugly
### really	beautiful	### not	beautiful	really not	beautiful
### really	bad	### not	bad	really not	bad
### really	good	### not	good	really not	good
### really	sad	### not	sad	really not	sad
### really	happy	### not	happy	really not	happy
### really	slow	### not	slow	really not	slow
### really	fast	### not	fast	really not	fast
really ###	small	not ###	small	not really	small
really ###	big	not ###	big	not really	big
really ###	cold	not ###	cold	not really	cold
really ###	hot	not ###	hot	not really	hot
really ###	ugly	not ###	ugly	not really	ugly
really ###	beautiful	not ###	beautiful	not really	beautiful
really ###	bad	not ###	bad	not really	bad
really ###	good	not ###	good	not really	good
really ###	sad	not ###	sad	not really	sad
really ###	happy	not ###	happy	not really	happy
really ###	slow	not ###	slow	not really	slow
really ###	fast	not ###	fast	not really	fast

1060

1061 **Table S1.** Comprehensive list of the 108 stimuli used in the behavioral experiment, color coded for
 1062 each experimental condition; purple: low adjectives, orange: high adjectives; green: affirmative
 1063 phrases, red: negated phrases.

List of linguistic stimuli employed in Experiment 2 (MEG)					
### ###	quiet	really really	quiet	not not	quiet
### ###	loud	really really	loud	not not	loud
### ###	cool	really really	cool	not not	cool
### ###	warm	really really	warm	not not	warm
### ###	dark	really really	dark	not not	dark
### ###	bright	really really	bright	not not	bright
### ###	bad	really really	bad	not not	bad
### ###	good	really really	good	not not	good
### really	quiet	### not	quiet	really not	quiet
### really	loud	### not	loud	really not	loud
### really	cool	### not	cool	really not	cool
### really	warm	### not	warm	really not	warm
### really	dark	### not	dark	really not	dark
### really	bright	### not	bright	really not	bright
### really	bad	### not	bad	really not	bad
### really	good	### not	good	really not	good
really ###	quiet	not ###	quiet	not really	quiet
really ###	loud	not ###	loud	not really	loud
really ###	cool	not ###	cool	not really	cool
really ###	warm	not ###	warm	not really	warm
really ###	dark	not ###	dark	not really	dark
really ###	bright	not ###	bright	not really	bright
really ###	bad	not ###	bad	not really	bad
really ###	good	not ###	good	not really	good

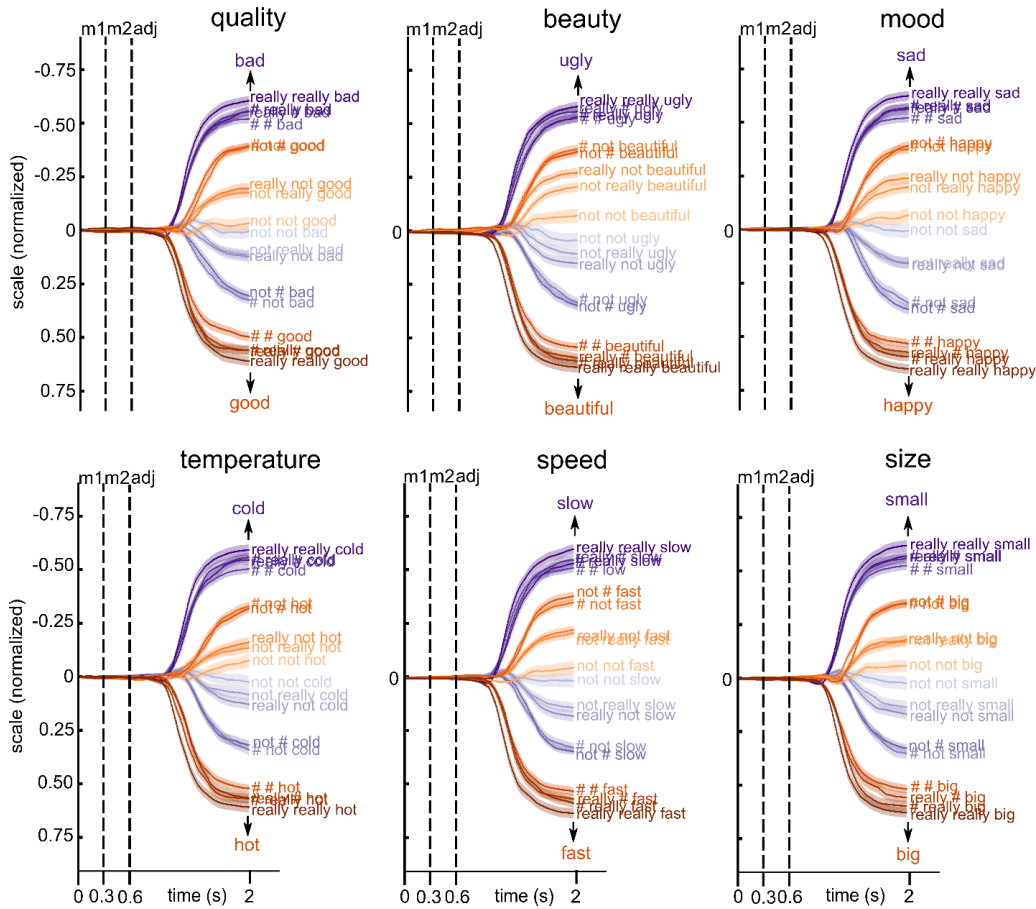
1064

1065 **Table S2.** Comprehensive list of the 72 stimuli used in the MEG experiment, color coded for each
 1066 experimental condition; purple: low adjectives, orange: high adjectives; green: affirmative phrases,
 1067 red: negated phrases. Note that the condition with no modifiers (“### ###”) was only employed as
 1068 a baseline condition in the time-frequency analysis.

1069 **Figures**

1070

Trajectories for affirmative and negated phrases, for each scalar dimension

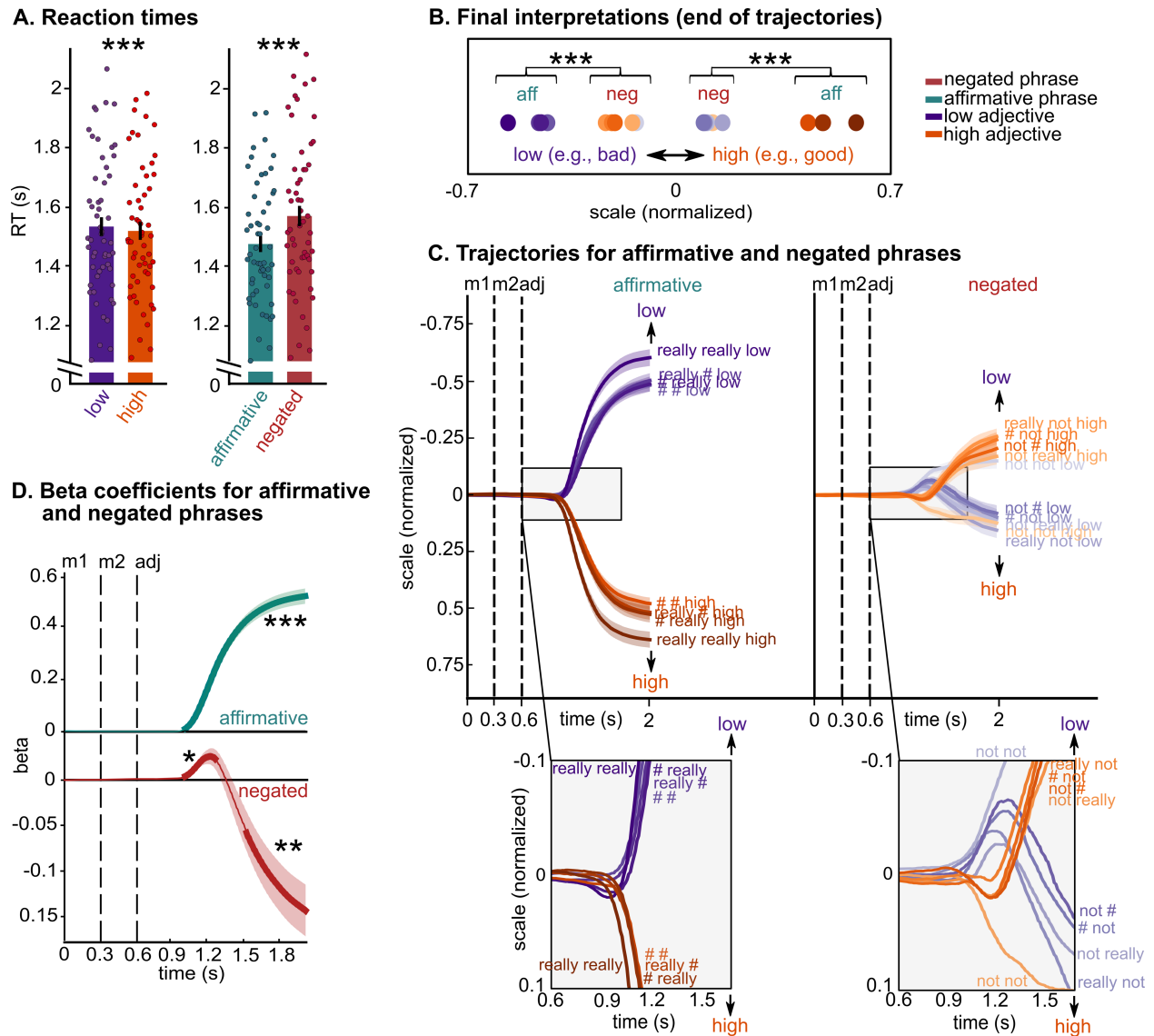


1071

1072

1073 **Fig. S1. Trajectories for each scalar dimension.**

1074 Behavioral trajectories for low (purples) and high (oranges) antonyms over time, for each scalar
1075 dimension (i.e., quality, beauty, mood, temperature, speed and size), for each modifier (shades of
1076 orange and purple), and for affirmative and negated phrases. Black vertical dashed lines indicate
1077 the presentation onset of each word: modifier1, modifier2 and adjective.



1078

1079

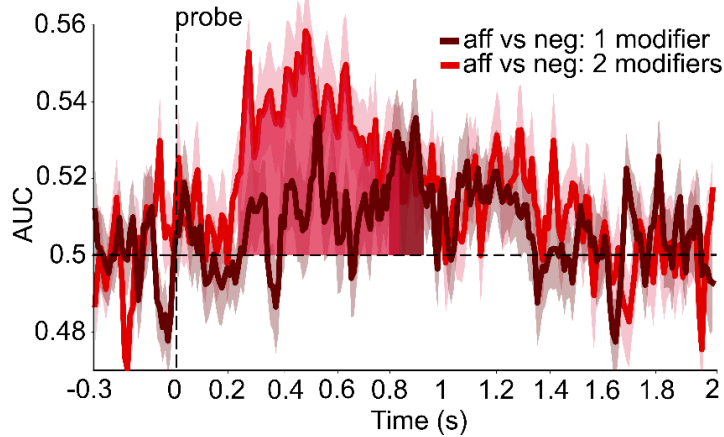
1080 **Fig. S2. Replication of Experiment 1, without feedback on interpretation.**

1081 A new group of 60 participants (37 females; mean age = 19.26 years; range 18-23 years) completed
 1082 the online mouse tracking experiment. Participants were recruited via the platform SONA (a
 1083 platform for students' recruitment). All participants were native English speakers with self-reported
 1084 normal hearing and no neurological deficits. 59 participants were right-handed. Participants were
 1085 granted university credits for taking part in the study, which was performed online. All participants
 1086 provided written informed consent, as approved by the local institutional review board (New York
 1087 University's Committee on Activities Involving Human Subjects). The data of 5 participants were
 1088 excluded from the data analysis due to (i) number of "incorrect" feedback based on the warnings >
 1089 30%, (ii) mean RTs > 2SD from the group mean, or (iii) response trajectory always ending within
 1090 1/4 from the center of the scale, regardless of condition (i.e., participants who did not pay attention
 1091 to the instructions of the task). Thus, 55 participants were included in the analysis. The experimental

1092 procedure was the same as that of Experiment 1, except that no feedback was provided to
1093 participants based on the final interpretation, but only if the cursor's movement violated the
1094 warnings provided during the familiarization phase (e.g., "you crossed the vertical borders", see
1095 Procedure of Experiment 1). We performed the same data analyses performed for Experiment 1:
1096 **(A) Reaction times:** To evaluate the specific effect of antonyms and of negation, we performed a 2
1097 (*antonym*: low vs high) x 2 (*negation*: negated vs affirmative) repeated-measures ANOVA. The
1098 results reveal a significant main effect of antonyms ($F(1,54) = 36.90, p < 0.001, \eta_p^2 = 0.40$) and a
1099 significant main effect of negation ($F(1,54) = 73.04, p < 0.001, \eta_p^2 = 0.57$). Moreover, a significant
1100 crossover interaction between antonyms and negation was found ($F(1,54) = 16.40, p < 0.001, \eta_p^2 =$
1101 0.23). These results replicate Experiment 1, showing that participants were faster for high adjectives
1102 (e.g., "good") than for low adjectives (e.g., "bad") and for affirmative phrases (e.g., "really really
1103 good") than for negated phrases (e.g., "really not good"). A further analysis including the number
1104 of modifiers as factor (i.e., *complexity*) indicates that participants were faster for phrases with two
1105 modifiers, e.g., "not really", than phrases with one modifier, e.g., "not ####" ($F(1,54) = 28.87, p <$
1106 $0.001, \eta_p^2 = 0.35$, especially in affirmative phrases: complexity by negation interaction $F(1,54) =$
1107 $6.26, p = 0.015, \eta_p^2 = 0.10$), again replicating results of Experiment 1. Bars represent the
1108 participants' mean \pm SEM and dots represent individual participants.
1109 **(B) Continuous mouse trajectories:** To investigate how negation changes the interpretation of scalar
1110 adjectives, we performed a 2 (*antonym*: low vs high) x 2 (*negation*: negated vs affirmative)
1111 repeated-measures ANOVA for participants' final interpretations (filled circles, purple = low,
1112 orange = high, averaged across dimensions and participants), which revealed a significant main
1113 effect of antonyms ($F(1,54) = 166.40, p < 0.001, \eta_p^2 = 0.47$), a significant main effect of negation
1114 ($F(1,54) = 48.62, p < 0.001, \eta_p^2 = 0.47$), and a significant interaction between antonyms and
1115 negation ($F(1,54) = 210.13, p < 0.001, \eta_p^2 = 0.80$). Post-hoc tests show that the final interpretation
1116 of negated phrases was located at a more central portion of the semantic scale than that of
1117 affirmative phrases (affirmative low $<$ negated high, and affirmative high $>$ negated low, $p_{\text{holm}} <$
1118 0.001), indicating that negation never inverts the interpretation of adjectives to that of their
1119 antonyms. Results also show that the final interpretations of negated phrases was significantly more
1120 variable (measured as standard deviations) than that of affirmative phrases ($F(1,54) = 15.43, p <$
1121 $0.001, \eta_p^2 = 0.22$). These results replicate Experiment 1. **(C) and (D).** To quantify the degree of
1122 deviation towards each side of the scale, we performed regression analyses with antonyms as the
1123 predictor and mouse trajectories as the dependent variable. Trials with "not not" were not included
1124 in this analysis as, in this experiment, the trajectories pattern was different compared to the other
1125 conditions with negation. Our results indicate that, while mouse trajectories of affirmative phrases

1126 branched towards either side of the scale and remained on that side until the final interpretation
1127 (lines in the left, gray, zoomed-in panel in **C**), the trajectories of negated phrases first deviated
1128 towards the side of the adjective and then towards the side of the antonym (lines in the right, gray,
1129 zoomed-in panel in **C**). The results of the regression analyses show that (1) in affirmative phrases,
1130 betas are positive (i.e., mouse trajectories moving towards the adjective) starting from 400 ms from
1131 the adjective onset ($p < 0.001$, green line in **D**); and that (2) in negated phrases, betas are positive
1132 (i.e., mouse trajectories moving towards the adjective) between 400 and 650 ms from the adjective
1133 onset ($p = 0.02$), and only became negative (i.e., mouse trajectories moving towards the antonym)
1134 from 910 ms from the adjective onset ($p = 0.003$, i.e., red line in **D**). Thicker lines indicate
1135 significant time windows. These results again replicate Experiment 1. For panels C and D: black
1136 vertical dashed lines indicate the presentation onset of each word: modifier 1, modifier 2 and
1137 adjective; each line and shading represent participants' mean \pm SEM; *** $p < 0.001$; ** $p < 0.01$; *
1138 $p < 0.05$.

Temporal decoding of negation as a function of complexity



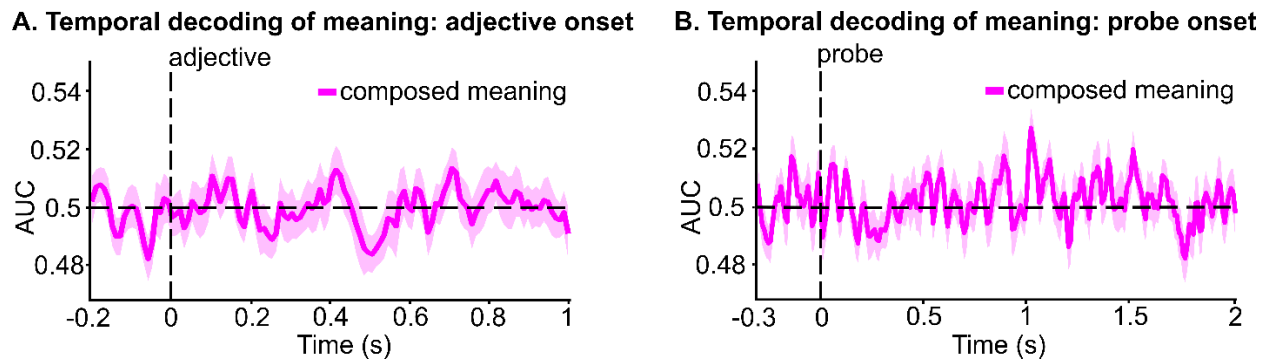
1139

1140

1141 **Fig. S3. Temporal decoding of negation as a function of number of modifiers (i.e., complexity),**
1142 **time-locked to the onset of the probe.**

1143 Decoding accuracy of negation over time, as a function of the number of modifiers (1 modifier:
1144 dark red line and shading; 2 modifiers: light red line and shading). Significant time windows are
1145 indicated by dark red (1 modifier) and light red (2 modifiers) areas. These results show that we
1146 could significantly decode the difference between affirmative and negated phrases between 230 and
1147 930 ms after the onset of the probe, especially when the phrase included two modifiers (1 modifier:
1148 between 790 and 930 ms: $p < 0.001$; 2 modifiers: between 230 and 840 ms: $p < 0.001$). This suggests
1149 that the representation of modifiers is reactivated at the stage when participants have to perform the
1150 yes/no task. AUC = area under the receiver operating characteristic curve, chance = 0.5 (black
1151 dashed horizontal line); the black vertical dashed line indicates the presentation onset of the probe;
1152 aff = affirmative, neg = negated; each line and shading represent participants mean \pm SEM.

1153

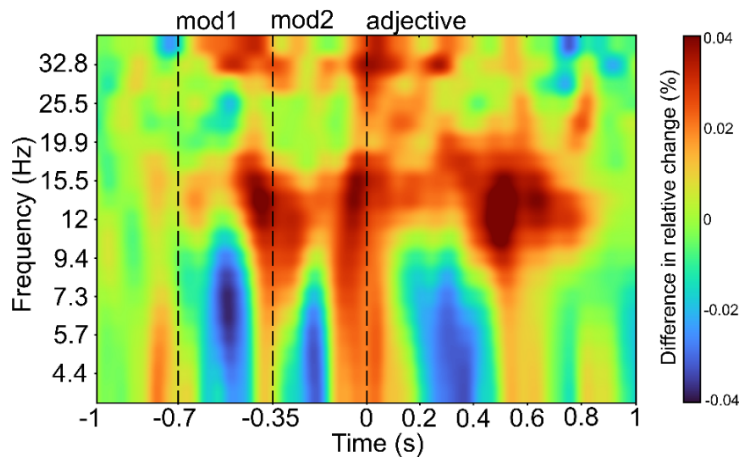


1154
1155

1156 **Fig. S4. Temporal decoding of composed meaning.**

1157 We trained estimators on phrases where the predicted composed meaning was “low” vs. “high” in
1158 90% of the trials and computed the accuracy of the model in predicting the representation of the
1159 meaning “low” vs. “high” in the remaining 10% of the trials. For instance, for the *quality* dimension,
1160 classes are: [0: *bad*] “### really bad”, “really ### bad”, “really really bad”, “### not good”, “not
1161 ### good”, “not not good”, “really not good”, “not really good”; and [1: *good*] “### really good”,
1162 “really ### good”, “really really good”, “### not bad”, “not ### bad”, “not not bad”, “really not
1163 bad”, “not really bad”. The composed meaning was derived from the behavioral results of
1164 Experiment 1. **(A)** Temporal decoding analyses time-locked to the onset of the adjective do not
1165 reveal any significant temporal cluster, suggesting that negation does not invert the representation
1166 of the adjective to that of its antonym (e.g., “bad” to “good”), as would be predicted by prediction
1167 **(3) Inversion**. **(B)** Temporal decoding analyses time-locked to the onset of the probe do not reveal
1168 any significant temporal cluster, suggesting that negation does not invert the representation of the
1169 adjective to that of its antonym (e.g., “bad” to “good”) after the presentation of the probe number.
1170 For all panels: AUC = area under the receiver operating characteristic curve, chance = 0.5 (black
1171 horizontal dashed line); black vertical dashed lines indicate the presentation onset of the adjective
1172 in **A** and the probe in **B**; each line and shading represent participants' mean \pm SEM.

(Negated - affirmative) power across times and frequencies

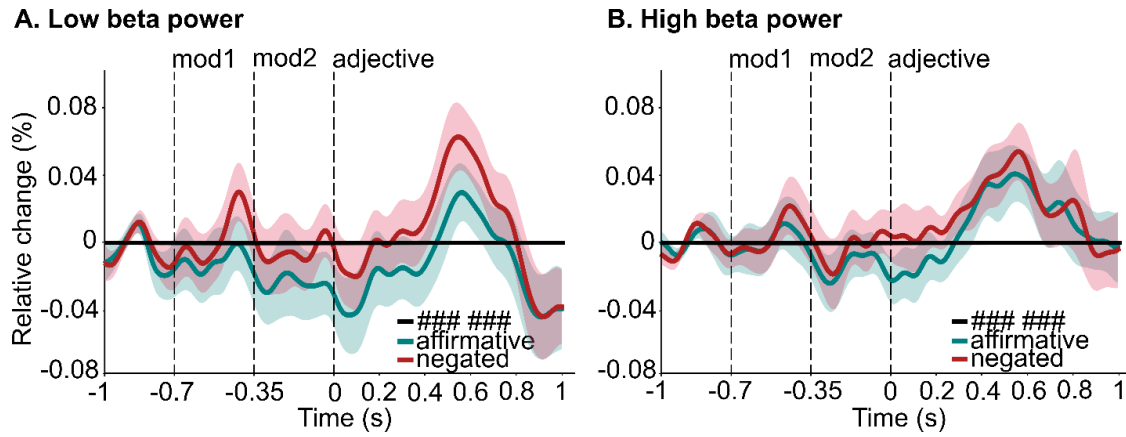


1173

1174

1175 **Fig. S5. Differences between negated and affirmative phrases across time and frequencies.**

1176 Time-frequency spectrum of the differences between negated and affirmative phrases averaged
1177 across all sensors and all participants. Frequencies are between 3.9 and 37.2 Hz, logarithmically
1178 spaced. Black vertical dashed lines indicate the presentation onset of each word: modifier1,
1179 modifier2 and adjective; colors indicate % differences in change relative to a baseline of -300 to -
1180 100 ms from the onset of word 1 (modifier1).



1181

1182

1183 **Fig. S6. Low- and high-beta power for negated and affirmative phrases across time.**

1184 The mean beta power for the no modifier condition was subtracted from the mean beta power of
1185 affirmative and negated phrases, separately for low-beta (12-20 Hz, **(A)**) and high-beta (21-30 Hz,
1186 **(B)**). The horizontal solid black line represents the no modifier condition (i.e., ### ##) after
1187 subtraction (thus = 0), and the green and red lines represent beta power over time for affirmative
1188 and negated phrases, respectively. Relative change (%) was obtained by subtracting the mean of
1189 baseline values (-300 to -100 ms from the onset of word1) and dividing by the mean of baseline
1190 values. Black vertical dashed lines indicate the presentation onset of each word: modifier1,
1191 modifier2 and adjective; each line and shading represent participants' mean \pm SEM.