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

2    Short title: **Semantic representations through negation**

3

4    Arianna Zuanazzi[1*], Pablo Ripollés[1,2,3], Wy Ming Lin[4], Laura Gwilliams[5], Jean-Rémi King[1,6†],

5    David Poeppel[1,3,7†]

6

7    1 Department of Psychology, New York University, New York, NY, USA.

8    2 Music and Audio Research Lab (MARL), New York University, New York, NY, USA.

9    3 Center for Language, Music and Emotion (CLaME), New York University, New York, NY, USA.

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

11   Tübingen, Germany.

12   5 Department of Psychology, Stanford University, Stanford, CA, USA.

13   6 Ecole Normale Supérieure, PSL University, Paris, France.

14   7 Ernst Strüngmann Institute for Neuroscience, Frankfurt, Germany.

15

16   †These authors contributed equally to this work.

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

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

**Abstract**

Combinatoric linguistic operations underpin human language processes, but how meaning is composed and refined in the mind of the reader is not well understood. We address this puzzle by exploiting the ubiquitous function of negation. We track the online effects of negation ("not") and intensifiers ("really") on the representation of scalar adjectives (e.g., "good") in parametrically designed behavioral and neurophysiological (MEG) experiments. The behavioral data show that participants first interpret negated adjectives as affirmative and later modify their interpretation towards, but never exactly as, the opposite meaning. Decoding analyses of neural activity further reveal significant above chance decoding accuracy for negated adjectives within 600 ms from adjective onset, suggesting that negation does not invert the representation of adjectives (i.e., "not bad" represented as "good"); furthermore, decoding accuracy for negated adjectives is found to be significantly lower than that for affirmative adjectives. Overall, these results suggest that negation mitigates rather than inverts the neural representations of adjectives. This putative suppression mechanism of negation is supported by increased synchronization of beta-band neural activity in sensorimotor areas. The analysis of negation provides a steppingstone to understand how the human brain represents changes of meaning over time.

**Introduction**

A hallmark of language processing is that we combine elements of the stored inventory - informally speaking, words - and thereby flexibly generate new meanings or change current meanings. The final representations derive in systematic ways from the combination of individual pieces. The composed meanings can be extracted in relatively straightforward ways, such as by sequentially combining individual meanings of words and phrases (e.g., "this theory is correct") or stem from more subtle inferential processes, where further operations are required to achieve understanding (e.g., "this theory is not even wrong", meaning "this theory is incoherent"). A mechanistic understanding of the underlying processes requires characterization of how meaning representations are constructed in real time. There has been steady progress and productive debate on syntactic structure building [1–6]. In contrast, how novel semantic configurations are represented over time is less widely investigated. In the experimental approach pursued here, we build on the existing literature on precisely controlled *minimal* linguistic environments [7,8]. We deploy a new, simple parametric experimental paradigm that capitalizes on the powerful role that *negation* plays in shaping semantic representations of words. While negation is undoubtfully a complex linguistic operation that can affect comprehension as a function of other linguistic factors (such as discourse and pragmatics [9–11]), our investigation specifically focuses on how negation operates in phrasal structures. Combining behavioral and neurophysiological data, we show how word meaning is (and is not) modulated in controlled contexts that contrast affirmative (e.g., "really good") and negated (e.g., "not good") phrases. The results identify models and mechanisms of how negation, a compelling window into semantic representation, operates in real time.
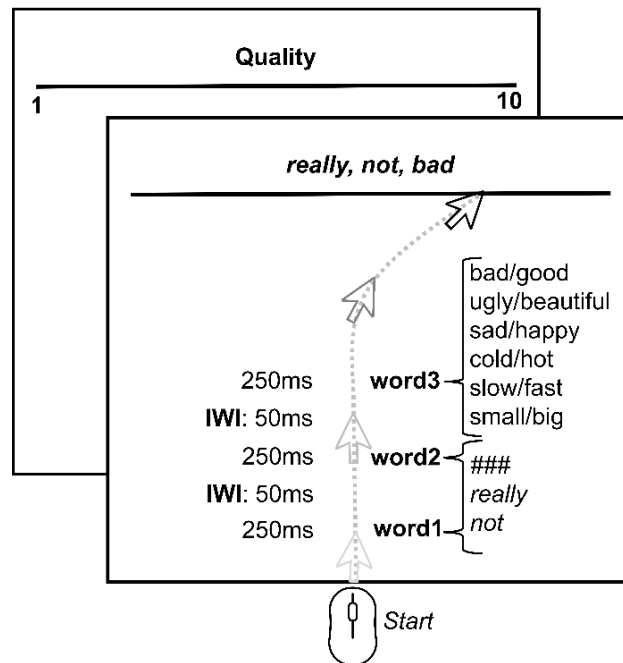
Negation is ubiquitous – and therefore interesting in its own right. Furthermore, it offers a compelling linguistic framework to understand how the human brain builds meaning through combinatoric processes. Intuitively, negated concepts (e.g., "not good") entertain some relation with the affirmative concept (e.g., "good") as well as their counterpart (e.g., "bad"). The function of negation in natural language has been a matter of longstanding debate among philosophers, psychologists, logicians, and linguists [12]. In spite of its intellectual history and relevance (interpreting negation was, famously, a point of debate between Bertrand Russell and Ludwig Wittgenstein), comparatively little research investigates the cognitive and neural mechanisms underpinning negation. Previous work shows that negated phrases/sentences are processed with more difficulty (slower, with more errors) than the affirmative counterparts, suggesting an asymmetry between negated and affirmative representations; furthermore, state-of-the-art artificial neural networks appear to be largely insensitive to the contextual impacts of negation [13–20]. This asymmetry motivates one fundamental question: *how* does negation operate?

103  Studies addressing this question suggest that negation operates as a suppression mechanism

104  by reducing the extent of available information [21–23], either in two steps [18,24–28] or in one

105  incremental step [12,29–31]; other studies demonstrate that negation is rapidly and dynamically

106  integrated into meaning representations [10,32], even unconsciously [33]. Within the context of

107  action representation (e.g., "cut", "wish"), previous research suggests that negation recruits general-

108  purpose inhibitory and cognitive control systems [34–41].

109  While the majority of neuroimaging studies focused on how negation affects action

110  representation, psycholinguistic research shows that scalar adjectives (e.g., "bad-good", "close-

111  open", "empty-full") offer insight into how negation operates on semantic representations of single

112  words. These studies provide behavioral evidence that negation can either *eliminate* the negated

113  concept and convey the opposite meaning ("not good" = "bad") or *mitigate* the meaning of its

114  antonym along a semantic continuum ("not good" = "less good", "average", or "somehow bad";

115  [11,12,42–44]). Thus, the system of polar opposites generated by scalar adjectives provides an

116  especially useful testbed to investigate changes in representation of abstract concepts along a

117  semantic scale (e.g., "bad" to "good"), as a function of negation (e.g., "bad" vs. "not good").

118  Here, we capitalize on the semantic continuum offered by scalar adjectives to investigate

119  *how* negation operates on the representation of abstract concepts (e.g., "bad" vs. "good"). First, we

120  track how negation affects semantic representations over time in a behavioral mouse tracking study

121  (and a replication study; **Fig.1A**). Next, we use magnetoencephalography (MEG) and a decoding

122  approach to track the evolution of neural representations of target adjectives in affirmative and

123  negated phrases (**Fig.1B**). Mouse tracking and decoding approaches allow us to quantify and

124  compare dynamic changes in participants' interpretations and neural representations of adjectives

125  over time (e.g., [45,46]). We test four hypotheses: (1) negation does not change the representation

126  of adjectives (e.g., "not good" = "good"), (2) negation weakens the representation of adjectives

127  (e.g., "not good" < "good"), (3) negation inverts the representation of adjectives (e.g., "not good"

128  = "bad"), and (4) negation changes the representation of adjectives to another representation (e.g.,

129  "not good" = e.g., "unacceptable"). The combined behavioral and neurophysiological data

130  adjudicate among these hypotheses and identify potential mechanisms that underlie how negation

131  functions in online meaning construction. Emerging temporal dynamics clarify how the effect of

132  negation on adjective meaning unfolds over time, whether incrementally (i.e., in parallel to

133  adjective processing) or serially (i.e., in a second step after adjective processing).

134

135

136

4

**A. Behavioral experiment: mouse trajectories**

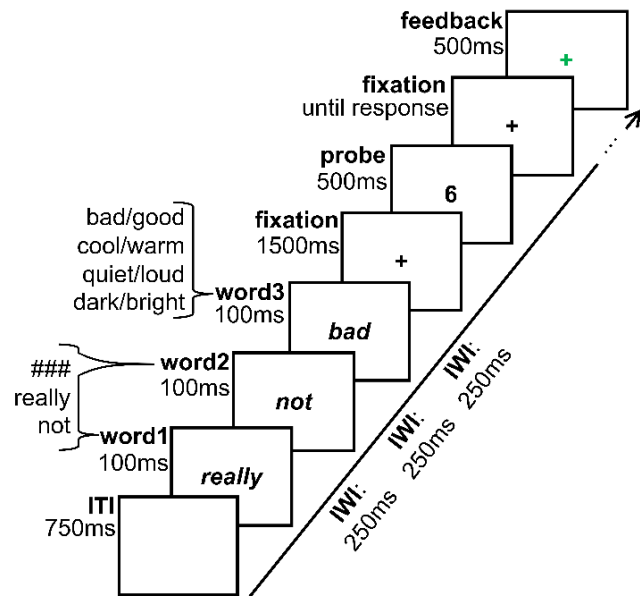**B. MEG experiment: behavioral task**

137 **Figure 1. Experimental procedures.**

138 (**A**) Behavioral procedure. Participants read affirmative or negated adjective phrases (e.g., "really really good", "###

139 not bad") word by word and rated the overall meaning of each phrase on a scale. Each trial consisted of combinations

140 of "###", "really", and "not" in word positions 1 and 2, followed by an adjective representing the low or high pole

141 across six possible scalar dimensions. Before each trial, participants were informed about the scale direction, e.g., "bad"

142 to "good", i.e., 1 to 10. Scale direction was pseudorandomized across blocks. Feedback was provided at the end of each

143 trial (to which 1 and 0 was assigned to compute the average feedback score). For each trial, we collected continuous

144 mouse trajectories throughout the entire trial as well as reaction times. (**B**) MEG procedure. Participants read

145 affirmative or negated adjective phrases and were instructed to derive the overall meaning of each adjective phrase on

146 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

147 presented, and participants were required to indicate whether the probe number represented the overall meaning of the

148 phrase on the scale (*yes/no* answer, using a keypad). Feedback was provided at the end of each trial (green or red cross,

149 to which 1 and 0 was assigned to compute the average feedback score). While performing the task, participants lay

150 supine in a magnetically shielded room while continuous MEG data were recorded through a 157-channel whole-head

151 axial gradiometer system. Panels A and B: "###" = no modifier; IWI = inter-word-interval.

152

153

154 **Results**

155 *Experiment 1: Continuous mouse tracking reveals a two-stage representation of negated*

156 *adjectives*

157 Experiment 1 (online behavioral experiment; N = 78) aimed to track changes in representation over

158 time of scalar adjectives in affirmative and negated phrases. Participants read two-to-three-word

phrases comprising one or two modifiers ("not" and "really") and a scalar adjective (e.g., "really really good", "really not quiet", "not ### fast"). The number and position of modifiers were manipulated to allow for a characterization of negation in simple and complex phrasal contexts, above and beyond single word processing. Adjectives were selected to represent opposite poles (i.e., antonyms) of the respective semantic scales: *low* pole of the scale (e.g., "bad", "ugly", "sad", "cold", "slow", and "small") and *high* pole of the scale (e.g., "good", "beautiful", "happy", "hot", "fast", and "big"). A sequence of dashes was used to indicate the absence of a modifier. **Fig. 1A** and **Table S1** provide a comprehensive list of the linguistic stimuli. On every trial, participants rated the overall meaning of each phrase on a scale defined by each antonym pair (**Fig. 1A**). Feedback was provided at the end of each trial (to which 1 and 0 were assigned to compute the average feedback score). We analyzed reaction times and continuous mouse trajectories, which consist of the positions of the participant's mouse cursor while rating the phrase meaning. Continuous mouse trajectories offer the opportunity to measure the unfolding of word and phrase comprehension over time, thus providing time-resolved dynamic data that reflect changes in meaning representation [15,45,47].

*Reaction times.* To evaluate the effect of antonyms and of negation on reaction times in behavioral Experiment 1, we performed a 2 (*antonym*: low vs. high) x 2 (*negation*: negated vs. affirmative) repeated-measures ANOVA. The results reveal a significant main effect of antonyms ($F(1,77) = 60.83$, $p < 0.001$, $\eta_p^2 = 0.44$) and a significant main effect of negation ($F(1,77) = 104.21$, $p < 0.001$, $\eta_p^2 = 0.57$, **Fig.2A**). No significant crossover interaction between antonyms and negation was observed ($p > 0.05$). Participants were faster for high adjectives (e.g., "good") than for low adjectives (e.g., "bad") and for affirmative phrases (e.g., "really really good") than for negated phrases (e.g., "really not good"). These results support previous behavioral data showing that negation is associated with increased processing difficulty [15,16]. A further analysis including the number of modifiers as factor (i.e., *complexity*) indicates that participants were faster for phrases with two modifiers, e.g., "not really", than phrases with one modifier, e.g., "not ###" ($F(1,77) = 16.02$, $p < 0.001$, $\eta_p^2 = 0.17$; see **Table S3A** for pairwise comparisons between each pair of modifiers), suggesting that the placeholder "###" may induce some processing slow-down. To confirm this hypothesis, further research should investigate the specific effect of placeholders (e.g., "###" or "xkq") on word and phrase representation and semantic composition.
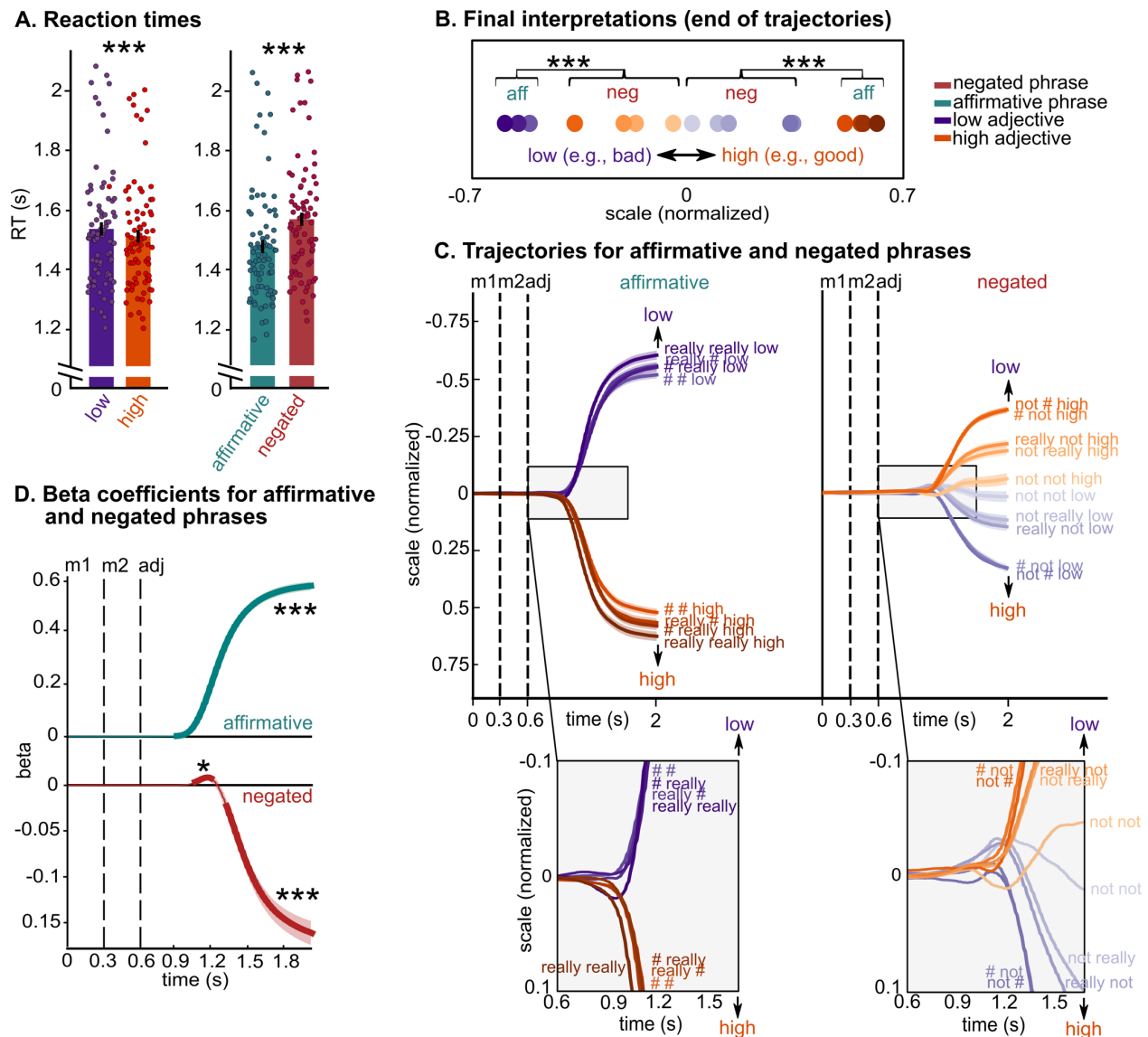
**Figure 2. Behavioral results.**

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

205     lines indicate the presentation onset of each word: modifier 1, modifier 2 and adjective; each line and shading represent

206     participants' mean ± SEM; Panels A,B,D: *** $p < 0.001$; * $p < 0.05$.

207

208     *Continuous mouse trajectories.* Continuous mouse trajectories across all adjective pairs and across

209     all participants are depicted in **Fig.2B** and **Fig.2C** (*low* and *high* summarize the two antonyms

210     across all scalar dimensions, see **Fig.S1** for each adjective dimension separately).

211             To quantify how the final interpretation of scalar adjectives changes as a function of

212     negation, we first performed a 2 (*antonym*: low vs. high) x 2 (*negation*: negated vs. affirmative)

213     repeated-measures ANOVA for participants' ends of trajectories (filled circles in **Fig.2B**), which

214     reveal a significant main effect of antonyms ($F(1,77) = 338.57$, $p < 0.001$, $\eta_p^2 = 0.83$), a significant

215     main effect of negation ($F(1,77) = 65.50$, $p < 0.001$, $\eta_p^2 = 0.46$), and a significant antonyms by

216     negation interaction ($F(1,77) = 1346.07$, $p < 0.001$, $\eta_p^2 = 0.95$). Post-hoc tests show that the final

217     interpretation of negated phrases is located at a more central portion on the semantic scale than that

218     of affirmative phrases (affirmative low < negated high, and affirmative high > negated low, $p_{holm} <$

219     0.001). Furthermore, the final interpretation of negated phrases is significantly more variable

220     (measured as standard deviations) than that of affirmative phrases ($F(1,77) = 78.14$, $p < 0.001$, $\eta_p^2$

221     = 0.50). Taken together, these results suggest that negation shifts the final interpretation of

222     adjectives towards the antonyms, but never to a degree that overlaps with the interpretation of the

223     affirmative antonym.

224             Second, we explored the temporal dynamics of adjective representation as a function of

225     negation (i.e., from the presentation of word 1 to the final interpretation; lines in **Fig.2C**). While

226     mouse trajectories of affirmative phrases branch towards either side of the scale and remain on that

227     side until the final interpretation (lines in the left, gray, zoomed-in panel in **Fig.2C**), trajectories of

228     negated phrases first deviate towards the side of the adjective and then towards the side of the

229     antonym, to reach the final interpretation (i.e., "not low" first towards "low" and then towards

230     "high"; right, gray, zoomed-in panel in **Fig.2C**; see **Fig.S1** for each adjective dimension separately).

231     To characterize the degree of deviation towards each side of the scale, we performed regression

232     analyses with antonyms as the predictor and mouse trajectories as the dependent variable (see

233     **Methods**). The results confirm this observation, showing that (1) in affirmative phrases, betas are

234     positive (i.e., mouse trajectories moving towards the adjective) starting at 300 ms from adjective

235     onset ($p < 0.001$, green line in **Fig.2D**); and that (2) in negated phrases, betas are positive between

236     450 and 580 ms from adjective onset (i.e., mouse trajectories moving towards the adjective, $p =$

237     0.04), and only become negative (i.e., mouse trajectories moving towards the antonym, $p < 0.001$)

238     from 700 ms from adjective onset (red line in **Fig.2D**). Note that beta values of negated phrases are

239  smaller than that for affirmative phrases, again suggesting that negation does not invert the
240  interpretation of the adjective to that of the antonym.

241

242  ***Replication of Experiment 1: Continuous mouse tracking reveals a two-stage representation of***
243  ***negated adjectives, in the absence of feedback***

244  We replicated Experiment 1 in a new group of online participants (N=55; **Fig.3**). The experimental
245  procedure was the same as that of Experiment 1, except that no feedback was provided to
246  participants based on the final interpretation, but only if the cursor's movement violated the
247  warnings provided during the familiarization phase (e.g., "you crossed the vertical borders", see
248  **Methods**). We performed the same data analyses performed for Experiment 1.

249

250  *Reaction times.* The 2 (*antonym*: low vs high) x 2 (*negation*: negated vs affirmative) repeated-
251  measures ANOVA reveal a significant main effect of antonyms ($F(1,54) = 36.90$, $p < 0.001$, $\eta_p^2 =$
252  0.40) and a significant main effect of negation ($F(1,54) = 73.04$, $p < 0.001$, $\eta_p^2 = 0.57$). Moreover,
253  a significant crossover interaction between antonyms and negation was found ($F(1,54) = 16.40$, $p$
254  $< 0.001$, $\eta_p^2 = 0.23$, **Fig.3A**). These results replicate Experiment 1, showing that participants were
255  faster for high adjectives (e.g., "good") than for low adjectives (e.g., "bad") and for affirmative
256  phrases (e.g., "really really good") than for negated phrases (e.g., "really not good"). Results on
257  *complexity* reveal that participants were faster for phrases with two modifiers, e.g., "not really",
258  than phrases with one modifier, e.g., "not ###" ($F(1,54) = 28.87$, $p < 0.001$, $\eta_p^2 = 0.35$, especially
259  in affirmative phrases: complexity by negation interaction $F(1,54) = 6.26$, $p = 0.015$, $\eta_p^2 = 0.10$),
260  again replicating results of Experiment 1 (see **Table S3B** for pairwise comparisons between each
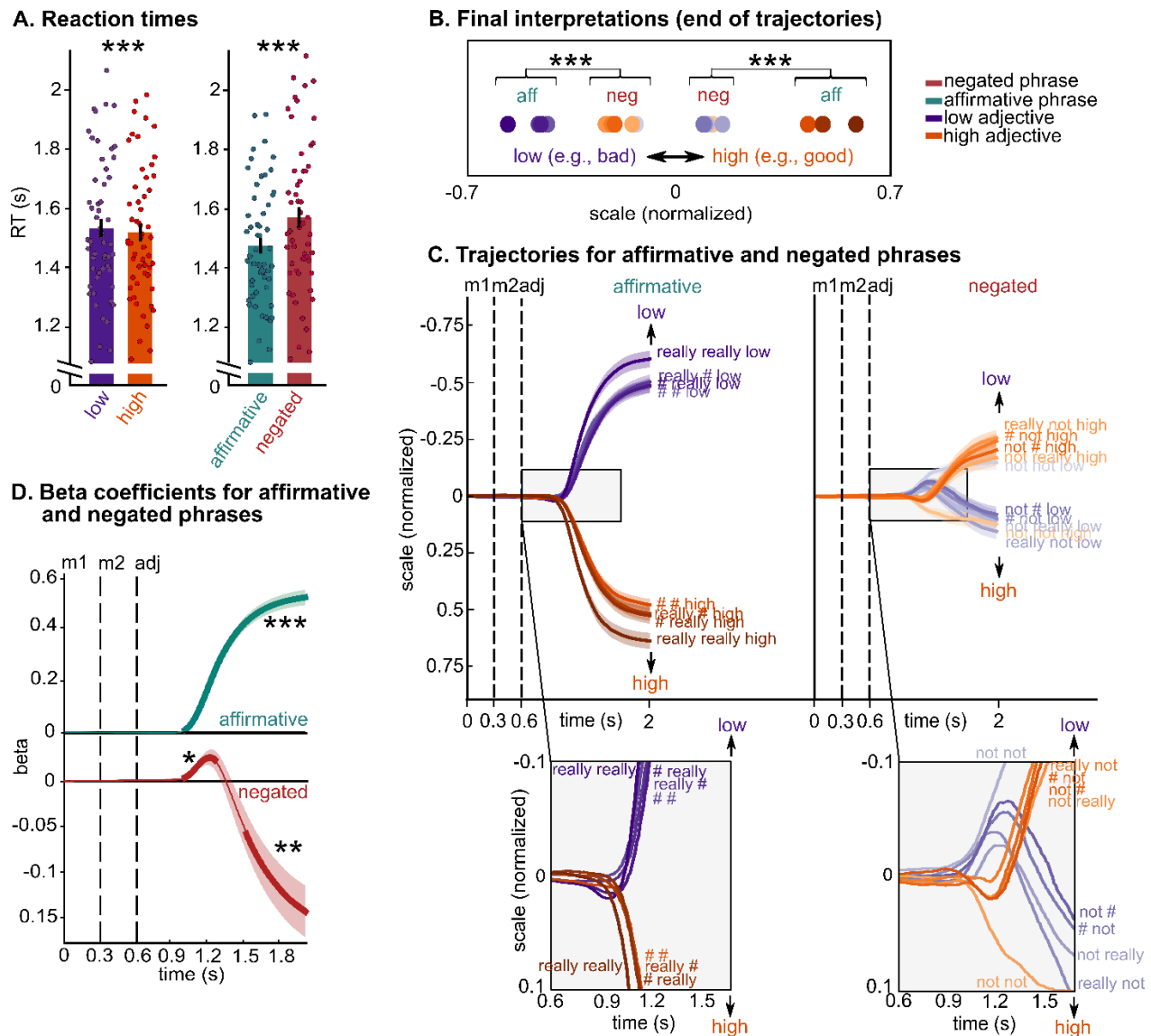261  pair of modifiers).

262

**Figure 3. Replication of Experiment 1, without feedback on interpretation.**

**(A)** Reaction times results for the online behavioral study (N=55). Bars represent the participants' mean ± SEM and dots represent individual participants. Participants were faster for high adjectives (e.g., "good") than for low adjectives (e.g., "bad") and for affirmative phrases (e.g., "really really good") than for negated phrases (e.g., "really not good"). These results replicate Experiment 1. **(B)** Final interpretations (i.e., end of trajectories) of each phrase, represented by filled circles (purple = low, orange = high), averaged across adjective dimensions and participants, showing that negation never inverts the interpretation of adjectives to that of their antonyms. **(C)** Mouse trajectories for low (purple) and high (orange) antonyms, for each modifier (shades of orange and purple) and for affirmative (left panel) and negated (right panel) phrases. Zoomed-in panels at the bottom demonstrate that mouse trajectories of affirmative phrases branch towards the adjective's side of the scale and remain on that side until the final interpretation; in contrast, the trajectories of negated phrases first deviate towards the side of the adjective and subsequently towards the side of the antonym (except for "not not"). This result is confirmed by linear models fitted to the data at each timepoint in **D**. These results also replicate Experiment 1. **(D)** Beta values (average over 55 participants) over time, separately for affirmative and negated phrases. Thicker lines indicate significant time windows. Trials with "not not" were not included in this analysis

10

279    as the trajectories pattern was different compared to the other conditions with negation. Panels **C**, **D**: black vertical
280    dashed lines indicate the presentation onset of each word: modifier 1, modifier 2 and adjective; each line and shading
281    represent participants' mean ± SEM; Panels A,B,D: *** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$.

282

283    *Continuous mouse trajectories.* The 2 (*antonym*: low vs high) x 2 (*negation*: negated vs affirmative)
284    repeated-measures ANOVA for participants' final interpretations reveal a significant main effect
285    of antonyms (F(1,54) = 166.40, $p < 0.001$, $\eta_p^2 = 0.75$), a significant main effect of negation (F(1,54)
286    = 48.62, $p < 0.001$, $\eta_p^2 = 0.47$), and a significant interaction between antonyms and negation
287    (F(1,54) = 210.13, $p < 0.001$, $\eta_p^2 = 0.80$). Post-hoc tests show that the final interpretation of negated
288    phrases was located at a more central portion of the semantic scale than that of affirmative phrases
289    (affirmative low < negated high, and affirmative high > negated low, $p_{holm} < 0.001$, **Fig.3B**),
290    indicating that negation never inverts the interpretation of adjectives to that of their antonyms.
291    Results also show that the final interpretations of negated phrases was significantly more variable
292    (measured as standard deviations) than that of affirmative phrases (F(1,54) = 15.43, $p < 0.001$, $\eta_p^2$
293    = 0.22). These results again replicate Experiment 1. As for Experiment 1, we then performed
294    regression analyses with antonyms as the predictor and mouse trajectories as the dependent
295    variable. For this analysis, trials with "not not" were not included as, in this experiment, the
296    trajectories pattern was different compared to the other conditions with negation (**Fig.3C**). The
297    results of the regression analyses show that (1) in affirmative phrases, betas are positive (i.e., mouse
298    trajectories moving towards the adjective) starting from 400 ms from the adjective onset ($p < 0.001$,
299    green line in **Fig.3D**); and that (2) in negated phrases, betas are positive (i.e., mouse trajectories
300    moving towards the adjective) between 400 and 650 ms from the adjective onset ($p = 0.02$), and
301    only became negative (i.e., mouse trajectories moving towards the antonym) from 910 ms from the
302    adjective onset ($p = 0.003$, i.e., red line in **Fig.3D**). This pattern replicates that of Experiment 1.

303        The replication of Experiment 1 illustrates the robustness of the behavioral mouse tracking
304    findings, even in the absence of feedback. Taken together, these results suggest that participants
305    initially interpreted negated phrases as affirmative (e.g., "not good" interpreted along the "good"
306    side of the scale) and later as a mitigated interpretation of the opposite meaning (e.g., the antonym
307    "bad").

308

309    ***Experiment 2: MEG shows that negation weakens the representation of adjectives and recruits***
310    ***response inhibition networks***
311    In this study (MEG experiment, N = 26), participants read adjective phrases comprising one or two
312    modifiers ("not" and "really") and scalar adjectives across different dimensions (e.g., "really really

11

good", "really not quiet", "not ### dark"). Adjectives were selected to represent opposite poles (i.e., the antonyms) of the respective semantic scales: *low* pole of the scale (e.g., "bad", "cool", "quiet", "dark") and *high* pole of the scale (e.g., "good", "warm", "loud", "bright"). A sequence of dashes was used to indicate the absence of a modifier. **Fig.1B** and **Table S2** provide the comprehensive list of the linguistic stimuli. Participants were asked to indicate whether a probe (e.g., 6) represented the meaning of the phrase on a scale from "really really low" (0) to "really really high" (8) (*yes/no* answer, **Fig.1B**). Feedback consisted of a green or red cross, to which 1 and 0 was assigned to compute the average feedback score. Behavioral data of Experiment 2 replicate that of Experiment 1: negated phrases are processed slower and with lower feedback score than affirmative phrases (main effect of negation for RTs: $F(1,25) = 26.44$, $p < 0.001$, $\eta_p^2 = 0.51$; main effect of negation for feedback score: $F(1,25) = 8.03$, $p = 0.009$, $\eta_p^2 = 0.24$).

The MEG analyses, using largely temporal and spatial decoding approaches [48], comprise four incremental steps: (1) we first identify the temporal correlates of simple word representation (i.e., the words "really" and "not" in the modifier position, and each pair of scalar adjectives in the second word position, i.e., the head position; see **Table S2**); (2) we test lexical-semantic representations of adjectives over time beyond the single word level, by entering *low* ("bad", "cool", "quiet" and "dark") and *high* ("good", "warm", "loud" and "bright") antonyms in the same model (adjectives in purple vs. orange in **Table S2**). We then test the representation of the negation operator over time (modifiers in green vs. red in **Table S2**); (3) we then ask how negation operates on the representation of adjectives, by teasing apart four possible mechanisms (i.e., *No effect*, *Mitigation*, *Inversion*, *Change*; adjectives in purple vs. orange for modifiers in green and red separately in **Table S2**); (4) we explore changes in beta power as a function of negation (motivated by the literature implicating beta-band neural activity in linguistic processing).

*(1)    Temporal decoding of single word processing*

The butterfly (bottom) and topography plots (top) in **Fig.4A** illustrate the grand average of the event-related fields elicited by the presentation of all words, as well as the probe, regardless of condition. Results of decoding analyses performed on these preprocessed MEG data (after performing linear dimensionality reduction; see **Methods**) show that the temporal decoding of "really" vs. "not" is significant between 120 and 430 ms and between 520 and 740 ms from the onset of the first modifier (dark gray shading, $p < 0.001$ and $p = 0.001$) and between 90 and 640 ms from the onset of the second modifier (light gray shading, $p < 0.001$, **Fig.4B**). Pairs of antonyms from different scales (regardless of specific modifier) were similarly decodable between 90 and 410 ms from adjective onset (quality: 110 to 200 ms, $p = 0.002$ and 290 to 370 ms, $p = 0.018$;

12

347    temperature: 140 to 280 ms, *p* < 0.001; loudness: 110 to 410 ms, *p* < 0.001; brightness: 90 to 350

348    ms, *p* < 0.001, **Fig.4C**), reflecting time windows during which the brain represents visual, lexical,

349    and semantic information (e.g., [7,49]). These results further show that single words can be decoded
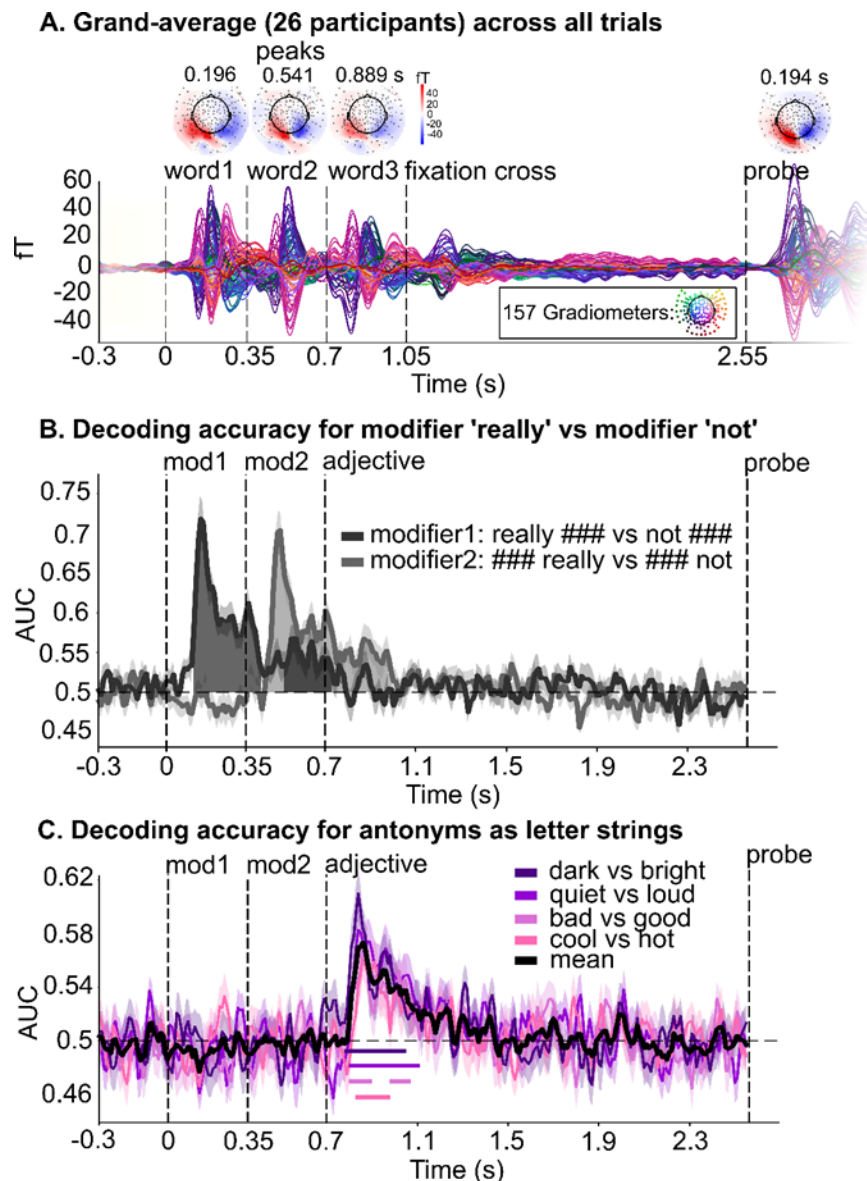
350    with relatively high accuracy (~70%).



351    **Figure 4. Evoked activity and temporal decoding of modifiers and adjectives as letter strings.**
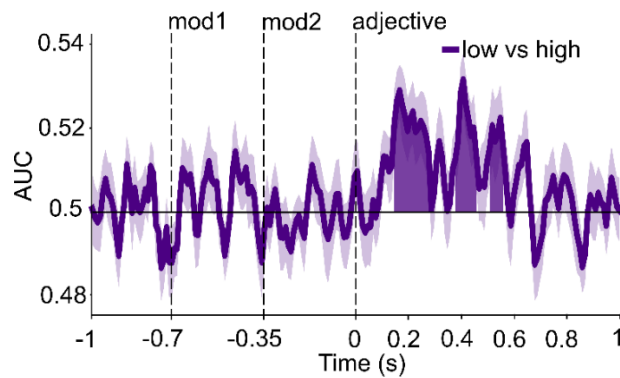
352    **(A)** The butterfly (bottom) and topo plots (top) illustrate the event-related fields elicited by the presentation of each

353    word as well as the probe, with a primarily visual distribution of neural activity right after visual onset (i.e., letter string

354    processing). We performed multivariate decoding analyses on these preprocessed MEG data, after performing linear

355    dimensionality reduction (see **Methods**). Detector distribution of MEG system in inset box. fT: femtoTesla magnetic

356    field strength. **(B)** We estimated the ability of the decoder to discriminate "really" vs. "not" separately in the first and

357    second modifier's position, from all MEG sensors. We contrasted phrases with modifiers "really ###" and "not ###",

358    and phrases with modifiers "### not" and "### really". **(C)** We evaluated whether the brain encodes representational

13

359 differences between each pair of antonyms (e.g., "bad" vs. "good"), in each of the four dimensions (quality,

360 temperature, loudness, and brightness). The mean across adjective pairs is represented as a solid black line; significant

361 windows are indicated by horizontal solid lines below. For panels B and C: AUC = area under the receiver operating

362 characteristic curve, chance = 0.5 (black horizontal dashed line); For all panels: black vertical dashed lines indicate the

363 presentation onset of each word: modifier 1, modifier 2, and adjective; each line and shading represent participants'

364 mean ± SEM.

365

366 *(2)    Temporal and spatial decoding of adjectives and negation*
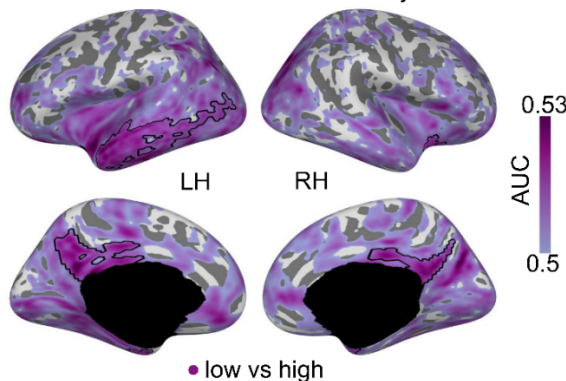
367 After establishing that single words' features can be successfully decoded in sensible time windows

368 (see **Fig.4**), we moved beyond single word representation and clarified the temporal patterns of

369 adjective and negation representation independently from their interaction and identified temporal

370 windows where to expect changes in adjective representation as a function of negation. First, we

371 selectively evaluated lexical-semantic differences between *low* ("bad", "cool", "quiet" and "dark")

372 and *high* ("good", "warm", "loud" and "bright") adjectives, regardless of the specific scale (i.e.,

373 pooling over *quality*, *temperature*, *loudness*, and *brightness*) and by pooling over all modifiers.

374 Temporal decoding analyses (see **Methods**) reveal significant decodability of *low* vs. *high*

375 antonyms in three time windows between 140 and 560 ms from adjective onset (140 to 280 ms, $p$

376 $< 0.001$; 370 to 460 ms: $p = 0.009$; 500 to 560 ms: $p = 0.044$, purple shading in **Fig.5A**). No

377 significant differences in lexical-semantic representation between *low* and *high* antonyms were

378 observed in later time windows (i.e., after 560 ms from adjective onset). The spatial decoding

379 analysis illustrated in **Fig.5B** (limited to 50-650 ms from adjective onset, see **Methods**) show that

380 decoding accuracy for *low* vs. *high* antonyms is significantly above chance in a widespread left-

381 lateralized brain network, encompassing the anterior portion of the superior temporal lobe, the

382 middle, and the inferior temporal lobe (purple shading in **Fig.5B**, significant clusters are indicated

383 by a black contour: left temporal lobe cluster, $p = 0.002$). A significant cluster was also found in

384 the right temporal pole, into the insula ($p = 0.007$). Moreover, we found significant clusters in the

385 bilateral cingulate gyri (posterior and isthmus) and precunei (left precuneus/cingulate cluster, $p =$

386 $0.009$; right precuneus/cingulate cluster, $p = 0.037$). Overall, these regions are part of the

387 (predominantly left-lateralized) frontotemporal brain network that underpins lexical-semantic

388 representation and composition [7,8,46,49–55].

14

**A. Temporal decoding of antonyms: word meaning**

**B. Spatial decoding of antonyms**

50-650 ms from the onset of the adjective

LH    RH

low vs high

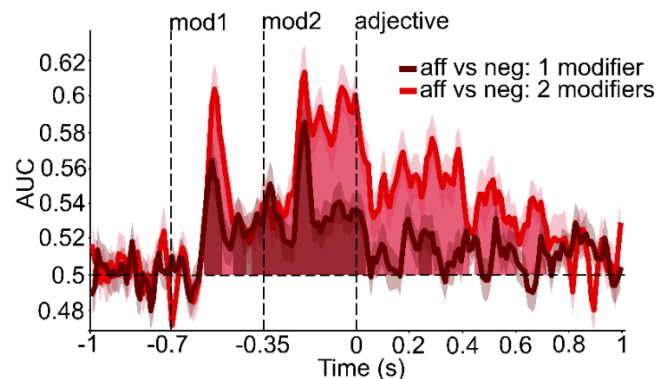**C. Temporal decoding of negation as a function of complexity**

**Figure 5. Temporal and spatial decoding of antonyms across all scales and temporal decoding of negation.**
(**A**) Decoding accuracy (purple line) of lexical-semantic differences between antonyms across all scales (i.e., pooling over "bad", "cool", "quiet" and "dark"; and "good", "warm", "loud" and "bright" before fitting the estimators) over time, regardless of modifier; significant time windows are indicated by purple shading; (**B**) Decoding accuracy (shades of purple) for antonyms across all scales over brain sources (after pooling over the four dimensions), between 50 and 650 ms from adjective onset. Significant spatial clusters are indicated by a black contour. (**C**) Decoding accuracy of negation over time, as a function of the number of modifiers (1 modifier: dark red line and shading; 2 modifiers: light red line and shading). 1 modifier: "really ###", "### really", "not ###", "### not"; 2 modifiers: "really really", "really not", "not really", "not not". Significant time windows are indicated by dark red (1 modifier) and light red (2 modifiers)

15

398    shading. For all panels: AUC: area under the receiver operating characteristic curve, chance = 0.5 (black horizontal

399    dashed line); black vertical dashed lines indicate the presentation onset of each word: modifier1, modifier2 and

400    adjective; each line and shading represent participants' mean ± SEM; aff = affirmative, neg = negated; LH = left

401    hemisphere; RH = right hemisphere.

402

403    Next, we turn to representations of negation over time. We performed a temporal decoding analysis

404    for phrases containing "not" vs. phrases not containing "not", separately for phrases with one and

405    two modifiers (to account for phrase complexity; see **Table S2** for a list of all trials). For phrases

406    with one modifier, the decoding of negation is significantly higher than chance throughout word 1

407    (-580 to -500 ms from adjective onset, $p = 0.005$), then again throughout word 2 (-470 to 0 ms from

408    adjective onset, $p < 0.001$). After the presentation of the adjective, negation decodability is again

409    significantly above chance between 0 and 40 ms ($p = 0.034$) and between 230 and 290 ms from

410    adjective onset ($p = 0.018$; dark red line and shading in **Fig.5C**). Similarly, for phrases with two

411    modifiers, the decoding of negation is significantly higher than chance throughout word 1 (-580 to

412    -410 ms from adjective onset, $p = 0.002$), throughout word 2 (-400 to 0 ms from adjective onset, $p$

413    $< 0.001$), and for a longer time window from adjective onset compared to phrases with one modifier,

414    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, $p <$

415    $0.001$; 620 to 720 ms, $p < 0.001$; light red line and shading in **Fig.5C**). The same analysis time-

416    locked to the onset of the probe shows that negation is once again significantly decodable between

417    230 and 930 ms after the probe, likely being reinstated when participants perform the task (**Fig.S2**).

418           Cumulatively, these results suggest that the brain encodes negation every time a "not" is

419    presented and maintains this information up to 720 ms after adjective onset. Further, they show that

420    the duration of negation maintenance is amplified by the presence of a second modifier, highlighting

421    combinatoric effects [2,6,56].
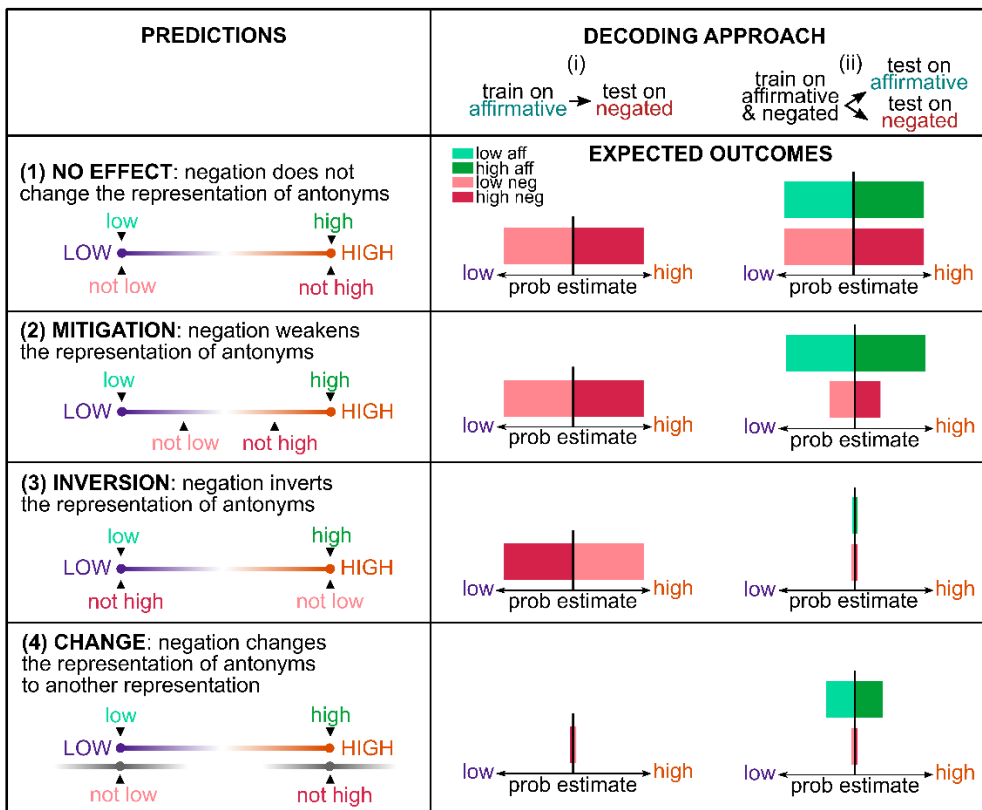
422

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

424    The temporal decoding analyses performed separately for adjectives and for negation demonstrate

425    that the brain maintains the representation of the modifiers available throughout the presentation of

426    the adjective. Here we ask how negation *operates on* the representation of the antonyms at the

427    neural level, leveraging theoretical accounts of negation [11,12,42–44], behavioral results of

428    Experiment 1, and two complementary decoding approaches. We test four hypotheses (see

429    *Predictions* in **Fig.6A**): (1) *No effect of negation*: negation does not change the representation of

430    adjectives (i.e., "not low" = "low"). We included this hypothesis based on the two-step theory of

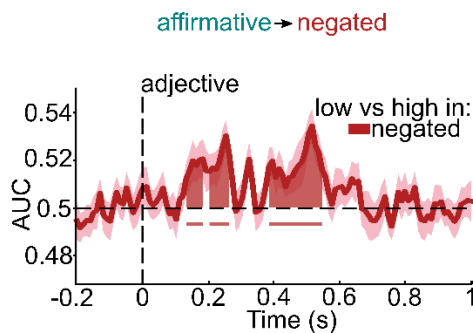431    negation, wherein the initial representation of negated adjectives would not be affected by negation

16

432    [27]. (2) *Mitigation*: negation weakens the representation of adjectives (i.e., "not low" < "low"). (3)

433    *Inversion*: negation inverts the representation of adjectives (i.e., "not low" = "high"). Hypotheses

434    (2) and (3) are derived from previous linguistics and psycholinguistics accounts on comprehension

435    of negated adjectives [42–44]. Finally, (4) *Change*: we evaluated the possibility that negation might

436    change the representation of adjectives to another representation outside the semantic scale defined

437    by the two antonyms (e.g., "not low" = e.g., "fair"). Importantly, these predictions focus on *how*

438    negation affects representations rather than on *when*. Thus, a combination of mechanisms may be

439    observed over time (e.g., first *no effect* and then *inversion*).

440    To adjudicate between these four hypotheses, we performed two complementary sets of decoding

441    analyses. Decoding approach (i): we computed the accuracy with which estimators trained on *low*

442    vs. *high* antonyms in affirmative phrases (e.g., "really really bad" vs. "really really good")

443    generalize to the representation of *low* vs. *high* antonyms in negated phrases (e.g., "really not bad"

444    vs. "really not good") at each time sample time-locked to adjective onset (see **Methods**); decoding

445    approach (ii): we trained estimators on *low* vs. *high* antonyms in affirmative and negated phrases

446    together (in 90% of the trials) and computed the accuracy of the model in predicting the

447    representation of *low* vs. *high* antonyms in affirmative and negated phrases separately (in the

448    remaining 10% of the trials; see **Methods**). Decoding approach (ii) allows for a direct comparison

449    between AUC and probability estimates in affirmative and negated phrases and to disentangle

450    predictions (1) *No effect* from (2) *Mitigation*. Expected probability estimates (i.e., the averaged

451    class probabilities for *low* and *high* classes) as a result of decoding approach (i) and (ii) are depicted

452    as light and dark, green and red bars under *Decoding approach* in **Fig.6A**.

17

**Figure 6. Predictions, decoding approaches, and results of the effect of negation on the representation of adjectives.**

(**A**) We tested four possible effects of negation on the representation of adjectives: (1) *No effect*, (2) *Mitigation*, (3) *Inversion*, (4) *Change* (left column). Note that we depicted predictions of (3) *Inversion* on the extremes of the scale, but a combination of inversion and mitigation would have the same expected outcomes. We performed two sets of decoding analyses (right column): (i) We trained estimators on low (purple) vs. high (orange) antonyms in affirmative

18

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

Temporal decoding approach (i) reveals that the estimators trained on the representation of *low* vs. *high* antonyms in affirmative phrases significantly generalize to the representation of *low* vs. *high* antonyms in negated phrases, in four time windows between 130 and 550 ms from adjective 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 ms: $p = 0.008$; red shading in **Fig.6B**). **Fig.6D** depicts the probability estimates averaged over the significant time windows for *low* and *high* antonyms in negated phrases. These results only support predictions (1) *No effect* and (2) *Mitigation*, thus invalidating predictions (3) *Inversion* and (4) *Change*. **Fig.S3** illustrates a different approach that similarly leads to the exclusion of prediction *(3) Inversion*.

Temporal decoding approach (ii) shows significant above chance decoding accuracy for affirmative phrases between 130 and 280 ms ($p < 0.001$) and between 370 and 420 ms ($p = 0.035$) from adjective onset. Conversely, decoding accuracy for negated phrases is significantly above chance only between 380 and 450 ms after the onset of the adjective ($p = 0.004$). Strikingly, negated phrases are associated with significantly lower decoding accuracy than affirmative phrases in the time window between 130 and 190 ms from adjective onset ($p = 0.040$; black horizontal line in **Fig.6C**). **Fig.6E** represents the probability estimates averaged over this 130-190 ms significant time window for *low* and *high* antonyms, separately in affirmative and negated phrases, illustrating

19

496    reduced probability estimates for negated compared to affirmative phrases. No significant

497    difference between decoding accuracy of affirmative and negative phrases was found for later time

498    windows (500-1000 ms from adjective onset, $p > 0.05$). A follow-up analysis where we trained and

499    tested on *low* vs. *high* antonyms in affirmative and negated phrases separately shows similar results

500    (**Fig.S4A**). Furthermore, the analysis including all trials, regardless of feedback score, also shows

501    similar results (**Fig.S4B**).

502        Overall, the generalization of representation from affirmative to negated phrases and the

503    higher decoding accuracy (and probability estimates) for affirmative than negated phrases within

504    the first 500 ms from adjective onset (i.e., within the time window of lexical-semantic processing

505    shown in **Fig.5A**) provide direct evidence in support of prediction (2) *Mitigation*, wherein negation

506    weakens the representation of adjectives. The alternative hypotheses did not survive the different
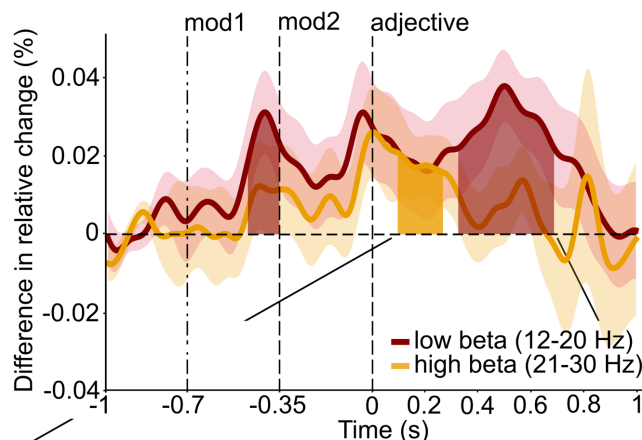
507    decoding approaches.

508

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

510    We distinguished among four possible mechanisms of how negation could operate on the

511    representation of adjectives and demonstrated that negation does not invert or change the

512    representation of adjectives but rather weakens the decodability of *low* vs. *high* antonyms within

513    the first ~300 ms from adjective onset (**Fig.6C**; with AUC for affirmative and negated adjectives

514    being significantly different for about 60 ms within this time window). The availability of negation

515    upon the processing of the adjective (**Fig.5A** and **Fig.5C**) and the reduced decoding accuracy for

516    antonyms in negated phrases (**Fig.6C**) raise the question of whether negation operates through

517    inhibitory mechanisms, as suggested by previous research employing action-related verbal material

518    [35–37]. We therefore performed time-frequency analyses, focusing on beta power (including low-

519    beta: 12 to 20 Hz, and high-beta: 20 to 30 Hz, [57], see **Methods**), which has been previously

520    associated with inhibitory control [58] (see **Fig.S5** for comprehensive time-frequency results). We

521    reasoned that, if negation operates through general-purpose inhibitory systems, we should observe

522    higher beta power for negated than affirmative phrases in sensorimotor brain regions.

523        Our results are consistent with this hypothesis, showing significantly higher low-beta power

524    (from 229 to 350 ms from the onset of modifier1: $p = 0.036$; from 326 to 690 ms from adjective

525    onset: $p = 0.012$; red line in **Fig.7A**) and high-beta power (from 98 to 271 ms from adjective onset:

526    $p = 0.044$; yellow line in **Fig.7A**) for negated than affirmative phrases. **Fig.S6** further shows low

527    and high-beta power separately for negated and affirmative phrases, compared to phrases with no

528    modifier (i.e., with "### ###").

529  Our whole-brain source localization analysis shows significantly higher low-beta power for
530  negated than affirmative phrases in the left precentral, postcentral, and paracentral gyri ($p = 0.012$;
531  between 326 and 690 ms from adjective onset, red cluster in **Fig.7C**). For high-beta power, similar
532  (albeit not significant) sensorimotor spatial patterns emerge (yellow cluster in **Fig.7B**).



533  **Figure 7. Differences in beta power over time between negated and affirmative phrases.**
534  (**A**) Differences in low (12-20 Hz, red) and high (21-30 Hz, yellow) beta power over time between negated (i.e., "###
535  not", "not ###", "really not", "not really", "not not") and affirmative phrases (i.e., "### really", "really ###", "really
536  really"). Negated phrases show higher beta power compared to affirmative phrases throughout the presentation of the
537  modifiers and for a sustained time window from adjective onset up to ~700 ms; significant time windows are indicated
538  by red (low-beta) and yellow (high-beta) shading; black vertical dashed lines indicate the presentation onset of each
539  word: modifier1, modifier2 and adjective; each line and shading represent participants' mean ± SEM. (**B**) Differences
540  (however not reaching statistical significance, $\alpha = 0.05$) in high-beta power between negated and affirmative phrases
541  (restricted between 97 and 271 ms from adjective onset, yellow cluster). (**C**) Significant differences in low-beta power
542  between negated and affirmative phrases (restricted between 326 and 690 ms from adjective onset) in the left precentral,
543  postcentral and paracentral gyrus (red cluster). Note that no significant spatial clusters were found in the right
544  hemisphere.

21

545

546

**Discussion**

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

The longer RTs and lower feedback score for negated phases shown in Experiment 1 (**Fig.2A**), in the replication experiment (**Fig.3A**), and in Experiment 2, are consistent with data demonstrating that negation incurs increased processing costs [13–18,27,32]. More significantly, mouse trajectories show that participants initially interpreted negated phrases as affirmative (e.g., "not good" is located on the "good" side of the scale, for ~130 ms, **Fig.2C** and **Fig.3C**), indicating that initial representations of negated scalar adjectives are closer to the representations of the adjectives rather than that of their antonyms. Similarly, participants' final interpretations of negated adjectives (e.g., "not good", "really not good") never overlapped with the final interpretations of the corresponding affirmative antonyms (e.g., "bad", "really bad", "really really bad"; **Fig.2B** and **Fig.3B**) highlighting how negation never inverts the meaning of an adjective to that of its antonym, even when participants are making decisions on a binary semantic scale (9,37-40).

Continuous mouse trajectories allowed us to quantify dynamic changes in participants' interpretations. MEG provided a means to directly track neural representations over time. We first identified the temporal correlates of lexical-semantic processing *separately* for scalar adjectives and for the negation operator. The time window of adjective representation (~140-560 ms from adjective onset, **Fig.5A**) is consistent with previous studies investigating lexical-semantic processing in language comprehension (130–200 ms up to ~550 ms from adjective onset [64–68]). Spatial decoding results corroborate temporal results, highlighting the involvement of the left-lateralized frontotemporal brain network in adjective processing (**Fig.5B,** [7,8,46,49–55]). Our data further show that negation is processed up to ~700 ms from adjective onset (**Fig.5C**). Overall, these data demonstrate that both scalar adjectives and negation are represented between 140 and 560 ms from adjective onset (compare **Fig.5A** and **Fig.5C**), suggesting that they are represented in parallel

579 and not serially (i.e., one after the other; see [69,70] for related patterns in the context of negation
580 + auxiliary verb and adjective + noun). Finally, they show that the decodability of negation
581 increases in phrases with two modifiers (e.g., "really not", "not really", **Fig.5C, Fig.S2**),
582 highlighting compositional effects [6].

583     We then evaluated the effects of the negation operator *on* adjective representation, to
584 address the question of *how* negation operates on lexical-semantic representations of antonyms. We
585 contrasted four hypotheses (**Fig.6A**): negation (1) does not change the representation of scalar
586 adjectives (e.g., "not good" = "good", *No effect*), (2) weakens the representation of scalar adjectives
587 (e.g., "not good" < "good", *Mitigation*), (3) inverts the representation of scalar adjectives (e.g., "not
588 good" = "bad", *Inversion*), or (4) changes the representation of scalar adjectives to another
589 representation (e.g., "not good" = e.g., "unacceptable", *Change*). These four hypotheses make
590 predictions about how negation operates on scalar adjectives at any given time. It is thus possible
591 that multiple mechanisms may unfold over time when looking at time-resolved data (e.g., first *no*
592 *effect* and then *inversion*). Using two complementary decoding approaches, we demonstrated that,
593 within the time window of adjective encoding, the representation of affirmative adjectives
594 generalizes to that of negated adjectives (**Fig.6B** and **Fig.6D**). This finding rules out predictions (3)
595 *Inversion* and (4) *Change*. Moreover, these findings complement our behavioral data that show that
596 negated adjectives are initially interpreted by participants as affirmative. Second, we showed that
597 the representation of adjectives in affirmative and negated phrases is not identical but is weakened
598 by negation (**Fig.6C** and **Fig.6E**). This result rules out prediction (1) *No effect* and supports
599 prediction (2) *Mitigation*, wherein negation weakens the representation of adjectives. We observed
600 such a reduction in early representations (i.e., within ~300 ms from adjective onset). This finding
601 is consistent with previous research that reported effects of negation as soon as lexical-semantic
602 representations of words are formed [12,29–31,71], and not exclusively at later processing stages
603 (e.g., P600 [72,73]). In addition, the fact that *low* vs. *high* adjectives are decodable ~400 ms after
604 the adjective onset in negated phrases (**Fig.6C**, **Fig.S4A**, **Fig.S4B) raises** two novel questions:
605 First, is the mitigation effect of negation stable over time? Second, at what exact stages does it
606 operate upon (see [64,66])? Using a masked priming paradigm, van Gaal et al. [33] analyzed
607 participants' EEG responses to sequences of words that were either consciously or unconsciously
608 perceived. Their findings indicate that the meaning of multiple words, including negation, can be
609 integrated even when subjects report not seeing them, but that conscious perception is required for
610 later grammatical integration. Future research remains necessary to more precisely tease apart the
611 lexical, semantic, and syntactic features that are selectively affected by the negation operator over
612 time.

613       Taken together, our behavioral and neural data jointly point to a *mitigation* rather than an

614    *inversion* effect of negation at early semantic processing stages, and exclude the hypothesis

615    according to which negation does not change the representation of antonyms. Specifically, these

616    results show that initial interpretations and early neural representations of negated adjectives are

617    similar to that of affirmative adjectives, but weakened. The comparison between MEG and

618    behavioral results also reveals interesting differences. Behavioral data reveal that, in negated

619    phrases, participants later modify their initial interpretation towards, but never exactly as, the

620    opposite meaning. Our MEG data do not show an inversion of adjective representation as a function

621    of negation, at early or later lexico-semantic processing stages. Differences between our behavioral

622    and neural results could be ascribed to the fact that the behavioral task had to be adapted to the

623    MEG environment. In the behavioral experiment (and its replication), participants were

624    continuously and explicitly indicating their interpretation, while in the MEG experiment they were

625    required to make a decision on their interpretation only after the probe was presented (1850 ms after

626    the adjective presentation, **Fig.2B** and **Fig.4A**), which could have hindered later effects of negation.

627       While previous fMRI studies on sentential negation have shown that negation reduces

628    hemodynamic brain activations related to verb processing [40,41], the current study offers novel

629    time-resolved behavioral and neural data on how negation selectively operates on abstract concepts.

630    Previous research has highlighted that negation might behave differently depending on the

631    pragmatics of discourse interpretation, e.g., when presented in isolation as compared to when

632    presented in context ("not wrong" vs. "this theory is not wrong" [9,10]), or when used ironically

633    ("they are not really good" said ironically to mean that they are "mediocre", e.g., [11,71]). Within

634    this pragmatic framework, it has been suggested that the opposite meaning of a scalar adjective

635    would be more simply conveyed by the affirmative counterpart than by negation [11,44,74]; thus,

636    to convey the opposite meaning of "bad", it would be more appropriate to use "good" as opposed

637    to "not bad". Following this logic, negation would be purposefully used (and understood) to convey

638    a different, mitigated meaning of the adjective (e.g., "not bad" = "less than bad"). Although we did

639    not directly manipulate sentential or pragmatic contexts, our findings provide behavioral and neural

640    evidence that negation acts as a mitigator. Here we only tested adjective pairs that form *contraries*

641    (which lie on a continuum, e.g., "bad" and "good"); thus inherently different patterns of results

642    could emerge in the case of *contradictories* (which form a dichotomy, e.g., "dead" and "alive",

643    [44]), where there is no continuum for mitigation to have an effect.

644       Overall, evidence that negation weakens adjective representations invites the hypothesis

645    that negation operates as a suppression mechanism, possibly through general-purpose inhibitory

646    systems [36,37]. To address this, we compared beta power modulations in affirmative and negated

phrases (**Fig.7**). In addition to subserving motor processing, beta-power modulation (12-30 Hz) has been associated with attention and expectancy violation and with multiple aspects of language processing, such as semantic memory and syntactic binding, as well as feedback processing ([35,75–78]; for a review, see [57,79]). We evaluated differences between negated and affirmative phrases separately in the low- and high-beta bands. We found greater power for negated than affirmative phrases in both bands, during the processing of the modifier and throughout the processing of the adjective up to ~700 ms, localized in left-lateralized sensorimotor areas. The timing and spatial correlates of beta-power in relation to negation align with studies that examined the effect of negation on (mental and motor) action representation [36]. Strikingly, we demonstrated that negation recruits brain areas and neurophysiological mechanisms similar to that recruited by response inhibition - however in the absence of action-related language material. Within a framework that recognizes two interactive neural systems, i.e., a semantic representation and a semantic control system [53], negation would operate through the latter, modulating how activation propagates through the (ventral) language semantic network wherein meaning is represented. The precise connectivity that underpins mitigation of lexical-semantic representations remains to be investigated.

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

**Materials and Methods**

***Participants***

*Experiment 1 (and replication): continuous behavioral tracking.* 101 participants (46 females; mean age = 29.6 years; range 18-67 years) completed an online mouse tracking experiment. Participants were recruited via Amazon Mechanical Turk and via the platform SONA (a platform for students' recruitment). All participants were native English speakers with self-reported normal hearing, normal or corrected to normal vision, and no neurological deficits. 97 participants were right-handed. Participants were paid or granted university credits for taking part in the study, which was performed online. All participants provided written informed consent, as approved by the local institutional review board (New York University's Committee on Activities Involving Human Subjects). The data of 23 participants were excluded from the data analysis due to (i) number of "incorrect" feedback (based on the warnings) > 30%, (ii) mean RTs > 2SD from the group mean,

681  or (iii) response trajectory always ending within 1/4 from the center of the scale, regardless of

682  condition (i.e., participants who did not pay attention to the instructions of the task). Thus, 78

683  participants were included in the analyses. The sample size was determined based on previous

684  studies using a similar behavioral approach (~30 participants [15,45,80]) and was increased to

685  account for the exclusion rate reported for online crowdsourcing experiments [81,82].

686  A new group of 60 participants (37 females; mean age = 19.26 years; range 18-23 years) completed

687  the online mouse tracking replication experiment. Participants were recruited via the platform

688  SONA. All participants were native English speakers with self-reported normal hearing and no

689  neurological deficits. 59 participants were right-handed. Participants were granted university credits

690  for taking part in the study, which was performed online. All participants provided written informed

691  consent, as approved by the local institutional review board (New York University's Committee on

692  Activities Involving Human Subjects). The data of 5 participants were excluded from the data

693  analysis due to (i) number of "incorrect" feedback based on the warnings > 30%, (ii) mean RTs >

694  2SD from the group mean, or (iii) response trajectory always ending within 1/4 from the center of

695  the scale, regardless of condition (i.e., participants who did not pay attention to the instructions of

696  the task). Thus, 55 participants were included in the analyses.

697

698  *Experiment 2: MEG.* A new group of 28 participants (17 females; mean age = 28.7 years; range 19-

699  53 years) took part in the in-lab MEG experiment. All participants were native English speakers

700  with self-reported normal hearing, normal or corrected to normal vision, and no neurological

701  deficits. 24 participants were right-handed. They were paid or granted university credits for taking

702  part in the study. All participants provided written informed consent, as approved by the local

703  institutional review board (New York University's Committee on Activities Involving Human

704  Subjects). The data of 2 participants were excluded from the data analysis because their feedback

705  scores in the behavioral task was < 60%. Thus, 26 participants were included in the analysis. The

706  sample size was determined based on previous studies investigating negation using EEG (17 to 33

707  participants [26,35,37]), investigating semantic representation using MEG (25 to 27 participants

708  [7,8]), or employing decoding methods with MEG data (17 to 20 participants [83,84]).

709

710  **Stimuli, Design, and Procedure**

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

712  *Stimuli and Design.* The linguistic stimulus set comprises 108 unique adjective phrases (for the

713  complete list, see **Table S1**). Adjectives were selected to be antonyms (i.e., *low* and *high* poles of

714  the scale) in the following six cognitive or sensory dimensions: *quality* ("bad", "good"), *beauty*

26

715    ("ugly", "beautiful"), *mood* ("sad", "happy"), *temperature* ("cold", "hot"), *speed* ("slow", "fast"),

716    and *size* ("small", "big"). These antonyms are all *contraries* (i.e., adjectives that lie on a continuum

717    [44]). Lexical characteristics of the antonyms were balanced according to the English Lexicon

718    Project [85]; mean (SD) HAL log frequency of *low* adjectives: 10.69 (1.09), *high* adjectives: 11.51

719    (1.07), mean (SD) bigram frequency of *low* adjectives: 1087.10 (374), *high* adjectives: 1032

720    (477.2); mean (SD) lexical decision RTs of *low* adjectives: 566 (37), *high* adjectives: 586 ms (70)).

721    Adjectives were combined with zero (e.g., "### ###"), one (e.g., "really ###"), or two modifiers

722    (e.g., "really not"). Modifiers were either the intensifier "really" or the negation "not" (see [33] for

723    a similar choice of modifiers; "really" was preferred to "very" as it more strongly intensifies the

724    meaning of the adjective, e.g., "really hot" > "very hot"). A sequence of dashes was used to indicate

725    the absence of a modifier, e.g., "really ### good". Each of the 12 adjectives was preceded by each

726    of the nine possible combinations of modifiers: "### ###", "### really", "really ###", "### not",

727    "not ###", "really not", "not really", "really really" and "not not", to diversify modifiers' sequences

728    and measure how negation affects adjective representation above and beyond the specific effects of

729    the words "really" and "not". Note that "not not" was included to achieve a full experimental design,

730    even if it is not a frequent combination in natural language and its cognitive and linguistic

731    representations are still under investigation (see [86]). Each dimension (e.g., quality) was presented

732    in two blocks (one block for each scale orientation, e.g., *low* to *high* and *high* to *low*) for a total of

733    12 blocks. Each phrase was repeated three times within each block (note that "### really"/"really

734    ###" were repeated an overall of three times, and so were "### not"/"not ###"). Thus, the overall

735    experiment comprised 504 trials. The order of phrases was randomized within each block for each

736    participant. The order of pairs of blocks was randomized across participants.

737

738    *Procedure.* Behavioral trajectories provide time-resolved dynamic data that reflect changes in

739    representation [15,45,47]. The online experiment was developed using oTree, a Python-based

740    framework for the development of controlled experiments on online platforms [87]. Participants

741    performed this study remotely, using their own monitor and mouse (touchpads were not allowed).

742    They were instructed to read affirmative or negated adjective phrases (e.g., "really really good",

743    "really not bad") and rate the overall meaning of each phrase on a scale, e.g., from "really really

744    bad" to "really really good". Participants were initially familiarized with the experiment through

745    short videos and a short practice block (18 trials with feedback). They were instructed that the poles

746    of the scale (e.g., "bad" and "good") would be reversed in half of the trials and warned that (i) they

747    could not cross the vertical borders of the response space, (ii) they had to maintain a constant

748    velocity, by following an horizontal line moving vertically, and (iii) they could not rate the meaning

27

749   of the phrase before the third word was presented. At the beginning of each trial, a response area of
750   600 (horizontal) x 450 (vertical) pixels and a solid line at the top of the rectangle were presented
751   (**Fig.1A**). Participants were informed about the scale (e.g., quality) and the direction of the scale
752   (e.g., "bad" to "good" or "good" to "bad", i.e., 1 to 10 or 10 to 1). Participants were instructed to
753   click on the "start" button and move the cursor of the mouse to the portion of the scale that best
754   represented the overall meaning of the phrase. The "start" button was placed in the center portion
755   of the bottom of the response space (i.e., in a neutral position). Once "start" was clicked on,
756   information about the scale and scale direction disappeared, leaving only the solid line on screen.
757   Phrases were presented at the top of the response space, from the time when participants clicked on
758   "start", one word at a time, each word for 250 ms (inter-word-interval: 50 ms). After each trial,
759   participants were provided the "incorrect" feedback if the cursor's movement violated the warnings
760   provided during the familiarization phase, and an explanation was provided (e.g., "you crossed the
761   vertical borders"). To keep participants engaged, we provided feedback also based on the final
762   interpretation: "negative" if the response was in the half of the scale opposite to the adjective (for
763   the conditions: "### ###", "#### really", "really ###" and "really really"), or in the same half of
764   the scale of the adjective (for the conditions: "### not" or "not ###"), or in the outer 20% left and
765   right portions of the scale (for the conditions: "really not", "not really" and "not not"); feedback
766   was "positive" otherwise. In case of a trial with negative feedback, the following trial was delayed
767   for 4 seconds. For each trial, we collected continuous mouse trajectories and RTs. The overall
768   duration of the behavioral experiment was approximately 90 minutes. To verify that the feedback
769   did not affect our results, we ran a replication study with a new group of 55 online participants
770   where no feedback was provided based on the final interpretation.

771

772   *Experiment 2: MEG.*
773   *Stimuli and Design.* The linguistic stimulus set comprised 72 unique adjective phrases (for the
774   complete list, see **Table S2**). Similar to Experiment 1, adjectives were selected for being antonyms
775   (and *contraries*) in the following cognitive or sensory dimensions (touch, audition, vision): *quality*
776   ("bad", "good"), *temperature* ("cool", "warm"), *loudness* ("quiet", "loud"), and *brightness* ("dark",
777   "bright"). The number of semantic scales (4) represents a tradeoff between stimulus variability,
778   number of stimuli within each condition - which is essential to achieve a reliable decoding accuracy
779   -, and experiment duration for attention maintenance. Lexical characteristics of the antonyms were
780   balanced according to the English Lexicon Project ([85]; mean (SD) HAL log frequency of "low"
781   adjectives: 10.85 (1.03), "high" adjectives: 10.55 (1.88); mean (SD) bigram frequency of "low"
782   adjectives: 1196.5 (824.6), "high" adjectives: 1077.5 (376.3); mean (SD) lexical decision RTs of

28

783    "low" adjectives: 594 ms (39), "high" adjectives: 594 (33)). Adjectives were combined with zero

784    (e.g., "### ###"), one (e.g., "really ###") or two modifiers (e.g., "really not"). Modifiers were either

785    the intensifier "really" or the negation "not". A sequence of dashes was used to indicate the absence

786    of a modifier, e.g., "really ### good". Each of the eight adjectives was preceded by each of the nine

787    possible combinations of modifiers: "### ###", "#### really", "really ###", "### not", "not ###",

788    "really not", "not really", "really really" and "not not" ("not not" was included to achieve a full

789    experimental design, even if it is not a frequent combination in natural language. See **Fig.S4C**,

790    **Fig.S4D** and **Fig.S4E** where we speculate that two "not", i.e., double negation, do not cancel each

791    other out but rather have mitigation effects similar to that of "really not"). To avoid possible

792    differences in neural representation of phrases with and without syntactic/semantic composition,

793    the condition with no modifiers ("### ###") was exclusively employed as a baseline comparison in

794    the time-frequency analysis and was excluded from all other analyses. Each dimension (e.g.,

795    quality) was presented in two blocks, one block for each yes/no key orientation (8 blocks in total,

796    see Procedure). Each phrase (e.g., "really really bad") was repeated four times within one block.

797    Thus, the overall experiment comprised 576 trials. The order of phrases was randomized within

798    each block for each participant. The order of blocks was randomized across participants within the

799    first and second half of the experiment. The yes/no order was randomized across participants.

800

801    *Procedure.* Participants were familiarized with the linguistic stimuli through a short practice block

802    that mimicked the structure of the experimental blocks. They were instructed to read affirmative or

803    negated adjective phrases (e.g., "really really good", "really not bad") and derive the overall

804    meaning of each adjective phrase, on a scale from 0 to 8, e.g., from "really really bad" to "really

805    really good". Each trial started with a fixation cross (duration: 750 ms), followed by each phrase

806    presented one word at a time, each word for 100 ms (inter-word-interval: 250 ms, **Fig.1B**). After

807    each phrase, a fixation cross was presented for 1500 ms. A number (i.e., probe) was then presented.

808    To keep the task engaging, participants were required to indicate whether the probe number

809    represented the meaning of the phrase on the scale (*yes/no* answer). The order of the yes/no response

810    keys was swapped halfway through the experiment. Responses had no time limit. If matching (+/-

811    one step on the scale from a likely predefined value), a green fixation cross was presented; if not, a

812    red fixation cross was presented, and feedback was provided.

813    While performing the experiment, participants lay supine in a magnetically shielded room while

814    continuous MEG data were recorded through a 157-channel whole-head axial gradiometer system

815    (Kanazawa Institute of Technology, Kanazawa, Japan). Sampling rate was 1000 Hz, and online

816    high-pass filter of 1 Hz and low-pass filter of 200 Hz were applied. Five electromagnetic coils were

817    attached to the forehead of the participants and their position was measured twice, before the first

818    and after the last block. Instructions, visual stimuli and visual feedback were back-projected onto a

819    Plexiglas screen using a Hitachi projector. Stimuli were presented using Psychtoolbox v3 ([88];

820    www.psychtoolbox.org), running under MATLAB R2019a (MathWorks) on an Apple iMac model

821    10.12.6. Participants responded to the yes/no question with their index finger of their left and right

822    hand, using a keypad. For each trial, we also collected feedback scores and RTs. The overall

823    duration of the MEG experiment was approximately 60 minutes.

824

825

826    ***Data analysis***

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

828    The RTs and mouse trajectory analyses were limited to trials with positive feedback (group mean

829    feedback scores: 82%, SD: 13%), and RTs were limited within the range of participant median RTs

830    ± 2 SD.

831    To evaluate differences in RTs between antonyms ("small", "cold", "ugly", "bad", "sad" vs. "big",

832    "hot", "beautiful", "good", "happy", "fast", i.e., *low* vs. *high* poles in each scalar dimension), and

833    between negated and affirmative phrases (e.g., "really really good" vs. "really not good"), and their

834    interactions, median RTs of each participant were entered into 2 (*antonym*: low vs. high) x 2

835    (*negation*: negated vs. affirmative) repeated-measures ANOVA.

836    To evaluate differences in the final interpretations between antonyms in each scale, between

837    negated and affirmative phrases, and their interactions, mean and standard deviation of the final

838    responses of each participant were entered into a 2 (*antonym*: low vs. high) x 2 (*negation*: negated

839    vs. affirmative) repeated-measures ANOVA. Post-hoc tests were conducted for significant

840    interactions (correction = Holm). Effect sizes were calculated using partial eta squared ($\eta_p^2$).

841    To compare mouse trajectories over time across participants, we resampled participants' mouse

842    trajectories at 100 Hz using linear interpolation, up to 2 seconds, to obtain 200 time points for each

843    trial. Furthermore, trajectories were normalized between -1 and 1. For visualization purposes, we

844    computed the median of trajectories across trials for each participant, dimension (e.g., quality),

845    antonym (e.g., "bad") and modifier (e.g., "really not"), and at each timepoint.

846    Finally, to quantitatively evaluate how the interpretation of each phrase changed over time, for

847    every participant we carried out regression analyses per each time point, for affirmative and negated

848    phrases separately (for a similar approach, see [45]). Note that, for the replication of Experiment 1,

849    trials with "not not" were not included in this analysis, as the trajectories pattern was different

850    compared to the other conditions with negation. The dependent variable was the mouse coordinate

851 along the scale (the scale which was swapped in half of the trials was swapped back for data analysis

852 purposes), and the predictor was whether the adjective was a low or high antonym (e.g., "bad" vs.

853 "good"). To identify the time windows where predictors were significantly different from 0 at the

854 group level, we performed permutation cluster tests on beta values (10,000 permutations) in the

855 time window from the onset of the adjective up to 1.4 s from adjective onset (i.e., 2 s from the onset

856 of word 1).

857

858 *Experiment 2: Feedback scores and RTs data.*

859 To evaluate differences in feedback scores between *low* and *high* antonyms ("bad", "cool", "quiet",

860 "dark" vs. "good", "warm", "loud", "bright"), and between negated and affirmative phrases (e.g.,

861 "really really good" vs. "really not good"), and their interactions, mean feedback score in the yes/no

862 task of each participant, computed as an average of 0 (red cross) and 1 (green cross), were entered

863 into 2 (*antonym*: low vs. high) x 2 (*negation*: negated vs. affirmative) repeated-measures ANOVA.

864 The response time analysis was limited to trials with positive feedback. RTs outside the

865 range of participant median RTs ± 2 SD were removed. To evaluate differences in RTs between

866 *low* and *high* antonyms in each scale and between negated and affirmative phrases, and their

867 interactions, median RTs of each participant in the yes/no task were entered into a 2 (*antonym*: low

868 vs. high) x 2 (*negation*: negated vs. affirmative) repeated-measures ANOVA.

869

870 *Experiment 2: MEG data.*

871 *Preprocessing.*

872 MEG data preprocessing was performed using MNE-python [89] and Eelbrain

873 (10.5281/zenodo.438193). First, bad channels (i.e., below the 3rd or above the 97th percentile

874 across all channels, for more than 20% of the entire recording) were interpolated. The MEG

875 responses were denoised by applying least square projections of the reference channels and

876 removing the corresponding components from the data [90]. Denoised data were lowpass-filtered

877 at 20 Hz for the decoding analyses and at 40 Hz for the time-frequency analyses. FastICA was used

878 to decompose the signal into 20 independent components, to visually inspect and remove artifacts

879 related to eye-blinks, heartbeat, and external noise sources (removed components across blocks and

880 participants: mean = 5.98, SD = 1.73). MEG recordings were then epoched into epochs of -300 ms

881 and 2550 ms around the onset of the first, second, or third word (or probe) for the decoding analyses,

882 and into epochs of -800 and 3000 ms around the onset of the first word for the time-frequency

883 analyses (and then cut between -300 and 2550 ms for group analyses). Note that, for visualization

884 purposes, only 1700 ms from the onset of the first word (i.e., 1000 ms from adjective onset) were

885 included in most figures (as no significant results were observed for control analyses run for later

886 time windows). Finally, epochs with amplitudes greater than an absolute threshold of 3000 fT were

887 removed and a baseline between -300 to 0 ms was applied to all epochs.

888

889 *Source reconstruction.*

890 Structural magnetic resonance images (MRIs) were collected for 10 out of 26 participants. For the

891 remaining 16 participants, we manually scaled and co-registered the "fsaverage" brain to the

892 participant's head-digitized shape and fiducials [89,91].

893 For every participant, an ico-4 source space was computed, containing 2562 vertices per hemisphere

894 and the forward solution was calculated using the Boundary Element Model (BEM). A noise

895 covariance matrix was estimated from the 300 ms before the onset of the first word up to the onset

896 of the first word presentation. The inverse operator was created and applied to the neuromagnetic

897 data to estimate the source time courses at each vertex using dynamic statistical parametric mapping

898 (dSPM: [92]). The results were then morphed to the ico-5 "fsaverage" brain, yielding to time

899 courses for 10242 vertices per hemisphere. We then estimated the magnitude of the activity at each

900 vertex (signal to noise ratio: 3, lambda2: 0.11, with orientation perpendicular to the cortical

901 surface), which was used in the decoding analyses (*Spatial decoders*).

902

903 *Decoding analyses.*

904 Decoding analyses were limited to trials with positive feedback and were performed with the MNE

905 [89] and Scikit-Learn packages [48]. First, X (or the selected principal components) were set to

906 have zero mean and unit variance (i.e., using a standard scaler). Second, we fitted an l2-regularized

907 logistic regression model as estimator to a subset of the epochs (training set, $X_{train}$) and estimated y

908 on a separate group of epochs (test set, $\hat{y}_{test}$). We then computed the accuracy (AUC, see below) of

909 the decoder, by comparing $\hat{y}_{test}$ with the ground truth y. For this analysis, we used the default values

910 provided by the Scikit-Learn package and set the class-weight parameter to "balanced".

911

912 *Temporal decoders.* Temporal decoding analyses were performed in sensor-space. Before fitting

913 the estimators, linear dimensionality reduction (principal component analysis, PCA) was performed

914 on the channel amplitudes to project them to a lower dimensional space (i.e., to new virtual channels

915 that explained more than 99% of the feature variance). We then fitted the estimator on each

916 participant separately, across all selected components, at each time-point separately. Time was

917 subsampled to 100 Hz. We then employed a 5-fold (for analyses in **Fig.4B** and **Fig.4C**) or 10-fold

918 stratified cross-validation (for analyses in **Fig.5A**, **Fig.5C**, and **Fig.6C**) that fitted the estimator to

919     80% or 90% of the epochs and generated predictions on 20% or 10% of the epochs, while keeping

920     the distributions of the training and test set maximally homogeneous. To investigate whether the

921     representation of antonyms was comparable between affirmative and negated phrases, in a different

922     set of analyses (i.e., decoding approach (i), **Fig.6B**) we fitted the estimator to all epochs

923     corresponding to affirmative phrases and generated predictions on all epochs corresponding to

924     negated phrases. In both decoding approaches, accuracy and probability estimates for each class

925     were then computed. Decoding accuracy is summarized with an empirical area under the curve

926     (rocAUC, 0 to 1, chance at 0.5).

927     At the group level, we extracted the clusters of time where AUC across participants was

928     significantly higher than chance using a one-sample permutation cluster test, as implemented in

929     MNE-python (10000 permutations [93]). We performed separate permutation cluster tests for the

930     following time windows: -700 to -350 ms from adjective onset (i.e., word 1), -350 to 0 ms from

931     adjective onset (i.e., word 2), 0 to 500 ms from adjective onset (i.e., time window for lexical-

932     semantic processes [65,66]) and 500 to 1000 ms from adjective onset (i.e., to account for potential

933     later processes).

934

935     *Expected outcome for the effect of negation on the representation of antonyms.* Temporal decoding

936     approach (i) and (ii) described above allow us to make specific predictions about the effect of

937     negation on the representation of antonyms (**Fig.6A**).

938     *Approach (i)* train set: affirmative phrases (in green in **Table S2**); test set: negated phrases

939     (in red in **Table S2**). For our results to support predictions (1) *No effect* or (2) *Mitigation*, this

940     decoding approach should show probability estimates of high and low adjectives significantly

941     above the computed chance level and in the direction of the respective classes, indicating that the

942     initial representation of adjectives in negated phrases is similar to that in affirmative phrases (left

943     column, first and second row under *decoding approach* in **Fig.6A**). Conversely, for our results to

944     support prediction (3) *Inversion*, this decoding approach should show probability estimates of high

945     and low adjectives significantly above the computed chance level but in the direction of the opposite

946     classes (i.e., swapped), as adjective representations would be systematically inverted in negated

947     phrases (left column, third row under *decoding approach* in **Fig.6A**). Finally, we should observe at

948     chance probability estimates in the case of (4) *Change*, where adjective representations in negated

949     phrases are not predictable from the corresponding representations in affirmative phrases (left

950     column, fourth row under *decoding approach* in **Fig.6A**).

951     *Approach (ii)* train set: affirmative and negated phrases together (in green/red in **Table S2**);

952     test set: affirmative and negated phrases separately (in green and red in **Table S2**). This decoding

953　analysis allows us to disentangle predictions (1) *No effect* from (2) *Mitigation*. For the results of

954　this analysis to support prediction (1) *No effect*, we should observe quantitatively comparable

955　probability estimates in affirmative and negated phrases, suggesting that negation does not change

956　the representation of adjectives (right column, first row under *decoding approach* in **Fig.6A**).

957　Conversely, in support of prediction (2) *Mitigation*, we should observe significantly reduced

958　probability estimates for negated relative to affirmative phrases, suggesting less robust differences

959　between low and high antonyms in negated phrases (right column, second row under *decoding*

960　*approach* in **Fig.6A**). The outcome of predictions (3) *Inversion* would be at chance probability

961　estimates for affirmative and negated phrases (as the model is trained on opposite representations

962　within the same class; right column, third row under *decoding approach* in **Fig.6A**) and the outcome

963　of (4) *Change* would be at chance probability estimates for negated phrases (as the model is trained

964　on different representations within the same class; right column, fourth row under *decoding*

965　*approach* in **Fig.6A**).

966　*Spatial decoders.* Spatial decoding analyses were performed in source-space. We fitted each

967　estimator on each participant separately, across 50 to 650 ms time samples relative to the onset of

968　the adjective (to include the three significant time windows that emerge from the temporal decoding

969　analysis in **Fig.4B**), at each brain source separately, after morphing individual participant's source

970　estimates to the ico-5 "fsaverage" common reference space. We employed a 5-fold stratified cross-

971　validation, which fitted the estimator to 80% of the epochs and generated predictions on 20% of the

972　epochs, while keeping the distributions of the training and test set maximally homogeneous.

973　Decoding accuracy is summarized with an empirical area under the curve (AUC, 0 to 1, chance at

974　0.5). At the group level, we extracted the brain areas where the AUC across participants was

975　significantly higher than chance, using a one-sample permutation cluster test as implemented in

976　MNE-python (10000 permutations; adjacency computed from the "fsaverage" brain [93]).

977

978　*Time-frequency analysis.*

979　We extracted time-frequency power of the epochs (-800 to 3000 ms from the onset of word 1) using

980　Morlet wavelets of 3 cycles per frequency, in frequencies between 3.9 and 37.2 Hz, logarithmically

981　spaced (19 frequencies overall). Power estimates where then cut between -300 and 2550 ms from

982　onset of word 1 and baseline corrected using a window of -300 to -100 ms from the onset of word

983　1, by subtracting the mean of baseline values and dividing by the mean of baseline values (mode =

984　'percent'). Power in the low-beta frequency range (12 to 20 Hz) and in the high-beta frequency

985　range (21 to 30 Hz [57,79]) was averaged to obtain a time course of power in low and high-beta

986　rhythms. We then subtracted the beta power of affirmative phrases from that of negated phrases. At

34

987  the group level, we extracted the clusters of time where this difference in power across participants

988  was significantly greater than 0, using a one-sample permutation cluster test as implemented in

989  MNE-python (10000 permutations [93]). We performed separate permutation cluster tests in the

990  same time windows used for the decoding analysis: -700 to -350 ms, -350 to 0 ms, 0 to 500 ms, and

991  500 to 1000 ms from the onset of the adjective (note that no significant differences were observed

992  in analyses ran for time windows after 1000 ms). We then computed the induced power in source

993  space (method: dSPM and morphing individual participant's source estimates to the ico-5

994  "fsaverage" reference space) for the significant clusters of time in the low- and high-beta range

995  separately and averaged over time. At the group level, we extracted the brain areas where the power

996  difference across participants was significantly greater than 0, using a one-sample permutation

997  cluster test as implemented in MNE-python (10000 permutations; adjacency computed from the

998  "fsaverage" brain [93]).

999

1000

1001  **References**

1002  1.   Ding N, Melloni L, Zhang H, Tian X, Poeppel D. Cortical tracking of hierarchical linguistic structures in
1003       connected speech. Nature Neuroscience. 2015;19: 158–164. doi:10.1038/nn.4186

1004  2.   Fedorenko E, Scott TL, Brunner P, Coon WG, Pritchett B, Schalk G, et al. Neural correlate of the construction
1005       of sentence meaning. Proceedings of the National Academy of Sciences of the United States of America.
1006       2016;113: E6256–E6262. doi:10.1073/pnas.1612132113

1007  3.   Martin AE, Baggio G. Modelling meaning composition from formalism to mechanism. 2019; 1–7.

1008  4.   Matchin W, Hickok G. The Cortical Organization of Syntax. Cerebral Cortex. 2020;30: 1481–1498.
1009       doi:10.1093/cercor/bhz180

1010  5.   Oseki Y, Marantz A. Modeling morphological processing in human magnetoencephalography. Proceedings of
1011       the Society for Computation in Linguistics. 2020;3.

1012  6.   Pallier C, Devauchelle A-D, Dehaene S. Cortical representation of the constituent structure of sentences.
1013       Proceedings of the National Academy of Sciences. 2011;108: 2522–2527. doi:10.1073/pnas.1018711108

1014  7.   Pylkkänen L. The neural basis of combinatory syntax and semantics. Science. 2019;366: 62–66.
1015       doi:10.1126/science.aax0050

1016  8.   Ziegler J, Pylkkänen L. Scalar adjectives and the temporal unfolding of semantic composition: An MEG
1017       investigation. Neuropsychologia. 2016;89: 161–171. doi:10.1016/j.neuropsychologia.2016.06.010

1018  9.   Tian Y, Ferguson H, Breheny R. Processing negation without context – why and when we represent the
1019       positive argument. Language, Cognition and Neuroscience. 2016;31: 683–698.
1020       doi:10.1080/23273798.2016.1140214

1021  10.  Tian Y, Breheny R, Ferguson HJ. Why we simulate negated information: A dynamic pragmatic account.
1022       Quarterly Journal of Experimental Psychology. 2010;63: 2305–2312. doi:10.1080/17470218.2010.525712

1023  11.  Giora R. Anything negatives can do affirmatives can do just as well, except for some metaphors. Journal of
1024       Pragmatics. 2006;38: 981–1014. doi:10.1016/j.pragma.2005.12.006

35

1025   12.   Horn LR. A natural history of negation. University of Chicago Press; 1989.

1026   13.   Ettinger A. What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models.
1027         Transactions of the Association for Computational Linguistics. 2020;8: 34–48. doi:10.1162/tacl_a_00298

1028   14.   Dale R, Duran ND. The cognitive dynamics of negated sentence verification. Cognitive Science. 2011;35: 983–
1029         996. doi:10.1111/j.1551-6709.2010.01164.x

1030   15.   Darley EJ, Kent C, Kazanina N. A 'no' with a trace of 'yes': A mouse-tracking study of negative sentence
1031         processing. Cognition. 2020;198: 104084. doi:10.1016/j.cognition.2019.104084

1032   16.   Dudschig C, Kaup B. How does "not left" become "right"? Electrophysiological evidence for a dynamic
1033         conflict-bound negation processing account. Journal of Experimental Psychology: Human Perception and
1034         Performance. 2018;44: 716–728. doi:10.1037/xhp0000481

1035   17.   Just MA, Carpenter PA. Comprehension of negation with quantification. Journal of Verbal Learning and
1036         Verbal Behavior. 1971;10: 244–253. doi:10.1016/S0022-5371(71)80051-8

1037   18.   Kaup B, Yaxley RH, Madden CJ, Zwaan RA, Ldtke J. Experiential simulations of negated text information.
1038         Quarterly Journal of Experimental Psychology. 2007;60: 976–990. doi:10.1080/17470210600823512

1039   19.   Dudschig C, Kaup B, Liu M, Schwab J. The processing of negation and polarity: An overview. Journal of
1040         Psycholinguistic Research. 2021;50: 1199–1213. doi:10.1007/s10936-021-09817-9

1041   20.   Sherman MA. Adjectival negation and the comprehension of multiply negated sentences. Journal of Verbal
1042         Learning and Verbal Behavior. 1976;15: 143–157.

1043   21.   Kaup B. Negation and its impact on the accessibility of text information. Memory and Cognition. 2001;29:
1044         960–967. doi:10.3758/BF03195758

1045   22.   Kaup B, Zwaan RA. Effects of negation and situational presence on the accessibility of text information.
1046         Journal of Experimental Psychology: Learning Memory and Cognition. 2003;29: 439–446. doi:10.1037/0278-
1047         7393.29.3.439

1048   23.   MacDonald MC, Just MA. Changes in activation levels with negation. Journal of Experimental Psychology:
1049         Learning, Memory, and Cognition. 1989;15: 633–642. doi:10.1037/0278-7393.15.4.633

1050   24.   Carpenter PA, Just MA. Sentence comprehension: A psycholinguistic processing model of verification.
1051         Psychological Review. 1975;82: 45–73. doi:10.1037/h0076248

1052   25.   Clark HH, Chase WG. On the process of comparing sentences against pictures. Cognitive Psychology. 1972;3:
1053         472–517. doi:10.1016/0010-0285(72)90019-9

1054   26.   Lüdtke J, Friedrich CK, De Filippis M, Kaup B. Event-related potential correlates of negation in a sentence-
1055         picture verification paradigm. Journal of Cognitive Neuroscience. 2008;20: 1355–1370.
1056         doi:10.1162/jocn.2008.20093

1057   27.   Kaup B, Dudschig C. Understanding negation: Issues in the processing of negation. In: Déprez V, Espinal MT,
1058         editors. The Oxford Handbook of Negation. Oxford University Press; 2020. pp. 634–655. Available:
1059         http://oxfordhandbooks.com/view/10.1093/oxfordhb/9780198830528.001.0001/oxfordhb-9780198830528-e-33

1060   28.   Papeo L, de Vega M. The neurobiology of lexical and sentential negation. The Oxford Handbook of Negation.
1061         2020; 739–756.

1062   29.   Papeo L, Hochmann J-R, Battelli L. The default computation of negated meanings. Journal of Cognitive
1063         Neuroscience. 2016;28: 1980–1986. doi:10.1162/jocn_a_01016

1064   30.   Lyons J. Linguistic semantics: An introduction. New York, NY: Cambridge University Press; 1995.

31. Mayo R, Schul Y, Burnstein E. "I am not guilty" vs "I am innocent": Successful negation may depend on the schema used for its encoding. Journal of Experimental Social Psychology. 2004;40: 433–449. doi:10.1016/j.jesp.2003.07.008

32. Orenes I, Beltrán D, Santamaría C. How negation is understood: Evidence from the visual world paradigm. Journal of Memory and Language. 2014;74: 36–45. doi:10.1016/j.jml.2014.04.001

33. van Gaal S, Naccache L, Meuwese JDI, van Loon AM, Leighton AH, Cohen L, et al. Can the meaning of multiple words be integrated unconsciously? Phil Trans R Soc B. 2014;369: 20130212. doi:10.1098/rstb.2013.0212

34. Bartoli E, Tettamanti A, Farronato P, Caporizzo A, Moro A, Gatti R, et al. The disembodiment effect of negation: Negating action-related sentences attenuates their interference on congruent upper limb movements. Journal of Neurophysiology. 2013;109: 1782–1792. doi:10.1152/jn.00894.2012

35. Beltrán D, Morera Y, García-Marco E, De Vega M. Brain inhibitory mechanisms are involved in the processing of sentential negation, regardless of its content. Evidence from EEG theta and beta rhythms. Frontiers in Psychology. 2019;10: 1–14. doi:10.3389/fpsyg.2019.01782

36. Beltrán D, Liu B, de Vega M. Inhibitory mechanisms in the processing of negations: A neural reuse hypothesis. Journal of Psycholinguistic Research. 2021;50: 1243–1260. doi:10.1007/s10936-021-09796-x

37. De Vega M, Morera Y, León I, Beltrán D, Casado P, Martín-Loeches M. Sentential negation might share neurophysiological mechanisms with action inhibition. Evidence from frontal theta rhythm. Journal of Neuroscience. 2016;36: 6002–6010. doi:10.1523/JNEUROSCI.3736-15.2016

38. Djokic V, Maillard J, Bulat L, Shutova E. Modeling affirmative and negated action processing in the brain with lexical and compositional semantic models. 2019; 5155–5165.

39. Gallese V, Lakoff G. The brain's concepts: The role of the sensory-motor system in conceptual knowledge. Cognitive Neuropsychology. 2005;22: 455–479. doi:10.1080/02643290442000310

40. Tettamanti M, Manenti R, Della Rosa PA, Falini A, Perani D, Cappa SF, et al. Negation in the brain: Modulating action representations. NeuroImage. 2008;43: 358–367. doi:10.1016/j.neuroimage.2008.08.004

41. Tomasino B, Weiss PH, Fink GR. To move or not to move: Imperatives modulate action-related verb processing in the motor system. Neuroscience. 2010;169: 246–258. doi:10.1016/j.neuroscience.2010.04.039

42. Bianchi I, Savardi U, Burro R, Torquati S. Negation and psychological dimensions. Journal of Cognitive Psychology. 2011;23: 275–301. doi:10.1080/20445911.2011.493154

43. Colston HL. "Not good" is "bad," but "not bad" is not "good": an analysis of three accounts of negation asymmetry. Discourse Processes. 1999;28: 237–256. doi:10.1080/01638539909545083

44. Fraenkel T, Schul Y. The meaning of negated adjectives. Intercultural Pragmatics. 2008;5: 517–540. doi:10.1515/IPRG.2008.025

45. Dotan D, Dehaene S. How do we convert a number into a finger trajectory? Cognition. 2013;129: 512–529. doi:10.1016/j.cognition.2013.07.007

46. Caucheteux C, Gramfort A, King J-R. Disentangling syntax and semantics in the brain with deep networks. 2021. Available: http://arxiv.org/abs/2103.01620

47. Maldonado M, Dunbar E, Chemla E. Mouse tracking as a window into decision making. Behav Res. 2019;51: 1085–1101. doi:10.3758/s13428-018-01194-x

48. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine learning in Python. Journal of Machine Learning Research. 2011;12: 2825–2830.

1106    49.    Caucheteux C, King J-R. Brains and algorithms partially converge in natural language processing. Commun
1107            Biol. 2022;5: 134. doi:10.1038/s42003-022-03036-1

1108    50.    Binder JR, Desai RH, Graves WW, Conant LL. Where is the semantic system? A critical review and meta-
1109            analysis of 120 functional neuroimaging studies. Cerebral Cortex. 2009;19: 2767–2796.
1110            doi:10.1093/cercor/bhp055

1111    51.    Huth AG, de Heer WA, Griffiths TL, Theunissen FE, Gallant JL. Natural speech reveals the semantic maps
1112            that tile human cerebral cortex. Nature. 2016;532: 453–458. doi:10.1038/nature17637

1113    52.    Lau EF, Gramfort A, Hämäläinen MS, Kuperberg GR. Automatic semantic facilitation in anterior temporal
1114            cortex revealed through multimodal neuroimaging. Journal of Neuroscience. 2013;33: 17174–17181.
1115            doi:10.1523/JNEUROSCI.1018-13.2013

1116    53.    Lambon Ralph MA, Jefferies E, Patterson K, Rogers TT. The neural and computational bases of semantic
1117            cognition. Nature Reviews Neuroscience. 2016;18: 42–55. doi:10.1038/nrn.2016.150

1118    54.    Hagoort P, Hald L, Bastiaansen M, Petersson KM. Integration of word meaning and world knowledge in
1119            language comprehension. Science. 2004;304: 438–441. doi:10.1126/science.1095455

1120    55.    Popham SF, Huth AG, Bilenko NY, Deniz F, Gao JS, Nunez-Elizalde AO, et al. Visual and linguistic semantic
1121            representations are aligned at the border of human visual cortex. Nat Neurosci. 2021;24: 1628–1636.
1122            doi:10.1038/s41593-021-00921-6

1123    56.    Parrish A, Pylkkänen L. Conceptual combination in the LATL with and without syntactic composition.
1124            Neurobiology of Language. 2022;3: 46–66. doi:10.1162/nol_a_00048

1125    57.    Weiss S, Mueller HM. "Too many betas do not spoil the broth": The role of beta brain oscillations in language
1126            processing. Frontiers in Psychology. 2012;3: 1–15. doi:10.3389/fpsyg.2012.00201

1127    58.    Wagner J, Wessel JR, Ghahremani A, Aron AR. Establishing a right frontal beta signature for stopping action
1128            in scalp EEG: Implications for testing inhibitory control in other task contexts. Journal of Cognitive
1129            Neuroscience. 2018;30: 107–118. doi:10.1162/jocn_a_01183

1130    59.    Alemanno F, Houdayer E, Cursi M, Velikova S, Tettamanti M, Comi G, et al. Action-related semantic content
1131            and negation polarity modulate motor areas during sentence reading: An event-related desynchronization study.
1132            Brain Research. 2012;1484: 39–49. doi:10.1016/j.brainres.2012.09.030

1133    60.    Aravena P, Delevoye-Turrell Y, Deprez V, Cheylus A, Paulignan Y, Frak V, et al. Grip force reveals the
1134            context sensitivity of language-induced motor activity during "action words" processing: evidence from
1135            sentential negation. Paterson K, editor. PLoS ONE. 2012;7: e50287. doi:10.1371/journal.pone.0050287

1136    61.    Foroni F, Semin GR. Comprehension of action negation involves inhibitory simulation. Frontiers in Human
1137            Neuroscience. 2013;7: 1–7. doi:10.3389/fnhum.2013.00209

1138    62.    Ghio M, Haegert K, Vaghi MM, Tettamanti M. Sentential negation of abstract and concrete conceptual
1139            categories: A brain decoding multivariate pattern analysis study. Philosophical Transactions of the Royal
1140            Society B: Biological Sciences. 2018;373: 7–10. doi:10.1098/rstb.2017.0124

1141    63.    Liuzza MT, Candidi M, Aglioti SM. Do not resonate with actions: Sentence polarity modulates cortico-spinal
1142            excitability during action-related sentence reading. PLoS ONE. 2011;6: 38–41.
1143            doi:10.1371/journal.pone.0016855

1144    64.    Hauk O, Davis MH, Ford M, Pulvermüller F, Marslen-Wilson WD. The time course of visual word recognition
1145            as revealed by linear regression analysis of ERP data. NeuroImage. 2006;30: 1383–1400.
1146            doi:10.1016/j.neuroimage.2005.11.048

1147    65.    Kutas M, Federmeier KD. Thirty years and counting: Finding meaning in the N400 component of the event-
1148            related brain potential (ERP). Annual Review of Psychology. 2011;62: 621–647.
1149            doi:10.1146/annurev.psych.093008.131123

66. Pulvermüller F, Shtyrov Y, Hauk O. Understanding in an instant: Neurophysiological evidence for mechanistic language circuits in the brain. Brain and Language. 2009;110: 81–94. doi:10.1016/j.bandl.2008.12.001

67. Pulvermüller F, Assadollahi R, Elbert T. Neuromagnetic evidence for early semantic access in word recognition. European Journal of Neuroscience. 2001;13: 201–205. doi:10.1046/j.0953-816X.2000.01380.x

68. Teige C, Mollo G, Millman R, Savill N, Smallwood J, Cornelissen PL, et al. Dynamic semantic cognition: Characterising coherent and controlled conceptual retrieval through time using magnetoencephalography and chronometric transcranial magnetic stimulation. Cortex. 2018;103: 329–349. doi:10.1016/j.cortex.2018.03.024

69. Zhang (张琳敏) L, Pylkkänen L. Semantic composition of sentences word by word: MEG evidence for shared processing of conceptual and logical elements. Neuropsychologia. 2018;119: 392–404. doi:10.1016/j.neuropsychologia.2018.08.016

70. Fyshe A, Sudre G, Wehbe L, Rafidi N, Mitchell TM. The lexical semantics of adjective–noun phrases in the human brain. Human Brain Mapping. 2019;40: 4457–4469. doi:10.1002/hbm.24714

71. Nieuwland MS, Kuperberg GR. When the truth is not too hard to handle: An event-related potential study on the pragmatics of negation. Psychological Science. 2008;19: 1213–1218. doi:10.1111/j.1467-9280.2008.02226.x

72. Palaz B, Rhodes R, Hestvik A. Informative use of "not" is N400-blind. Psychophysiology. 2020;57. doi:10.1111/psyp.13676

73. Xiang M, Grove J, Giannakidou A. Semantic and pragmatic processes in the comprehension of negation: An event related potential study of negative polarity sensitivity. Journal of Neurolinguistics. 2016;38: 71–88. doi:10.1016/j.jneuroling.2015.11.001

74. Grice HP. Logic and Conversation. P. Cole, J. L. Morgan. Syntax and Semantics. P. Cole, J. L. Morgan. New York: Academic Press; 1975. pp. 41–58.

75. Bastiaansen MCM, van der Linden M, ter Keurs M, Dijkstra T, Hagoort P. Theta responses are involved in lexical—semantic retrieval during language processing. Journal of Cognitive Neuroscience. 2005;17: 530–541. doi:10.1162/0898929053279469

76. Luo Y, Zhang Y, Feng X, Zhou X. Electroencephalogram oscillations differentiate semantic and prosodic processes during sentence reading. Neuroscience. 2010;169: 654–664. doi:10.1016/j.neuroscience.2010.05.032

77. Supp GG, Schlögl A, Gunter TC, Bernard M, Pfurtscheller G, Petsche H. Lexical memory search during N400: cortical couplings in auditory comprehension: NeuroReport. 2004;15: 1209–1213. doi:10.1097/00001756-200405190-00026

78. Weiss S, Rappelsberger P. EEG coherence within the 13–18 Hz band as a correlate of a distinct lexical organisation of concrete and abstract nouns in humans. Neuroscience Letters. 1996;209: 17–20. doi:10.1016/0304-3940(96)12581-7

79. Schaller F, Weiss S, Müller HM. EEG beta-power changes reflect motor involvement in abstract action language processing. Brain and Language. 2017;168: 95–105. doi:10.1016/j.bandl.2017.01.010

80. Pinheiro-Chagas P, Dotan D, Piazza M, Dehaene S. Finger tracking reveals the covert stages of mental arithmetic. Open Mind. 2017;1: 30–41. doi:10.1162/opmi_a_00003

81. Peer E, Brandimarte L, Samat S, Acquisti A. Beyond the Turk: Alternative platforms for crowdsourcing behavioral research. Journal of Experimental Social Psychology. 2017;70: 153–163. doi:10.1016/j.jesp.2017.01.006

82. Simcox T, Fiez JA. Collecting response times using Amazon Mechanical Turk and Adobe Flash. Behavior Research Methods. 2014;46: 95–111. doi:10.3758/s13428-013-0345-y

1192    83.    Gwilliams L, King JR. Recurrent processes support a cascade of hierarchical decisions. eLife. 2020;9: 1–20.
1193           doi:10.7554/ELIFE.56603

1194    84.    King JR, Pescetelli N, Dehaene S. Brain mechanisms underlying the brief maintenance of seen and unseen
1195           sensory information. Neuron. 2016;92: 1122–1134. doi:10.1016/j.neuron.2016.10.051

1196    85.    Balota DA, Yap MJ, Cortese MJ, Hutchison KA, Kessler B, Loftis B, et al. The English Lexicon Project.
1197           Behavior Research Methods. 2007;39: 445–459. doi:10.3758/BF03193014

1198    86.    Schiller NO, Van Lenteren L, Witteman J, Ouwehand K, Band GPH, Verhagen A. Solving the problem of
1199           double negation is not impossible: electrophysiological evidence for the cohesive function of sentential
1200           negation. Language, Cognition and Neuroscience. 2017;32: 147–157. doi:10.1080/23273798.2016.1236977

1201    87.    Chen DL, Schonger M, Wickens C. oTree—An open-source platform for laboratory, online, and field
1202           experiments. Journal of Behavioral and Experimental Finance. 2016;9: 88–97. doi:10.1016/j.jbef.2015.12.001

1203    88.    Brainard DH. The Psychophysics Toolbox. Spatial Vision. 1997;10: 433–436. doi:10.1163/156856897X00357

1204    89.    Gramfort A, Luessi M, Larson E, Engemann DA, Strohmeier D, Brodbeck C, et al. MEG and EEG data
1205           analysis with MNE-Python. Frontiers in Neuroscience. 2013;7: 1–13. doi:10.3389/fnins.2013.00267

1206    90.    Adachi Y, Shimogawara M, Higuchi M, Haruta Y, Ochiai M. Reduction of non-periodic environmental
1207           magnetic noise in MEG measurement by Continuously Adjusted Least squares Method. IEEE Transactions on
1208           Applied Superconductivity. 2001;11: 669–672. doi:10.1109/77.919433

1209    91.    Andersen LM. Group analysis in MNE-python of evoked responses from a tactile stimulation paradigm: A
1210           pipeline for reproducibility at every step of processing, going from individual sensor space representations to
1211           an across-group source space representation. Frontiers in Neuroscience. 2018;12.
1212           doi:10.3389/fnins.2018.00006

1213    92.    Dale AM, Liu AK, Fischl BR, Buckner RL, Belliveau JW, Lewine JD, et al. Dynamic statistical parametric
1214           mapping: Combining fMRI and MEG for high-resolution imaging of cortical activity. Neuron. 2000;26: 55–67.
1215           doi:10.1016/S0896-6273(00)81138-1

1216    93.    Maris E, Oostenveld R. Nonparametric statistical testing of EEG- and MEG-data. Journal of Neuroscience
1217           Methods. 2007;164: 177–190. doi:10.1016/j.jneumeth.2007.03.024

1218

## Acknowledgements

1225

## Author contributions

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

1230

1231 **Competing interests**

1232 The authors declare no competing interests.

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

**Supplementary Materials**

**Tables**

| 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 |

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

42

| 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 |

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

1294

**A. Experiment 1: Post Hoc Comparisons for RTs - Modifiers**

| | | Mean Difference | 95% CI for Mean Difference Lower | Upper | SE | t | $p_{tukey}$ |
|---|---|---|---|---|---|---|---|
| # # | # not | -0.120 | -0.214 | -0.026 | 0.030 | -3.956 | **0.003** |
| | # really | -0.025 | -0.119 | 0.069 | 0.030 | -0.817 | 0.996 |
| | not # | -0.110 | -0.204 | -0.016 | 0.030 | -3.642 | **0.009** |
| | not not | -0.079 | -0.173 | 0.015 | 0.030 | -2.617 | 0.181 |
| | not really | -0.087 | -0.181 | 0.007 | 0.030 | -2.876 | 0.096 |
| | really # | -0.006 | -0.100 | 0.089 | 0.030 | -0.182 | 1.000 |
| | really not | -0.090 | -0.185 | 0.004 | 0.030 | -2.988 | 0.071 |
| | really really | 0.020 | -0.074 | 0.114 | 0.030 | 0.654 | 0.999 |
| # not | # really | 0.095 | 8.348e-4 | 0.189 | 0.030 | 3.139 | **0.046** |
| | not # | 0.010 | -0.085 | 0.104 | 0.030 | 0.314 | 1.000 |
| | not not | 0.041 | -0.054 | 0.135 | 0.030 | 1.338 | 0.920 |
| | not really | 0.033 | -0.061 | 0.127 | 0.030 | 1.080 | 0.977 |
| | really # | 0.114 | 0.020 | 0.208 | 0.030 | 3.774 | **0.005** |
| | really not | 0.029 | -0.065 | 0.123 | 0.030 | 0.968 | 0.989 |
| | really really | 0.140 | 0.045 | 0.234 | 0.030 | 4.610 | **< .001** |
| # really | not # | -0.086 | -0.180 | 0.009 | 0.030 | -2.825 | 0.110 |
| | not not | -0.055 | -0.149 | 0.040 | 0.030 | -1.801 | 0.682 |
| | not really | -0.062 | -0.157 | 0.032 | 0.030 | -2.059 | 0.502 |
| | really # | 0.019 | -0.075 | 0.113 | 0.030 | 0.634 | 0.999 |
| | really not | -0.066 | -0.160 | 0.028 | 0.030 | -2.171 | 0.426 |
| | really really | 0.045 | -0.050 | 0.139 | 0.030 | 1.470 | 0.869 |
| not # | not not | 0.031 | -0.063 | 0.125 | 0.030 | 1.024 | 0.984 |
| | not really | 0.023 | -0.071 | 0.117 | 0.030 | 0.766 | 0.998 |
| | really # | 0.105 | 0.011 | 0.199 | 0.030 | 3.460 | **0.017** |
| | really not | 0.020 | -0.074 | 0.114 | 0.030 | 0.654 | 0.999 |
| | really really | 0.130 | 0.036 | 0.224 | 0.030 | 4.296 | **< .001** |
| not not | not really | -0.008 | -0.102 | 0.086 | 0.030 | -0.258 | 1.000 |
| | really # | 0.074 | -0.020 | 0.168 | 0.030 | 2.435 | 0.266 |
| | really not | -0.011 | -0.105 | 0.083 | 0.030 | -0.371 | 1.000 |
| | really really | 0.099 | 0.005 | 0.193 | 0.030 | 3.271 | **0.031** |
| not really | really # | 0.082 | -0.013 | 0.176 | 0.030 | 2.693 | 0.152 |
| | really not | -0.003 | -0.098 | 0.091 | 0.030 | -0.112 | 1.000 |
| | really really | 0.107 | 0.013 | 0.201 | 0.030 | 3.529 | **0.013** |
| really # | really not | -0.085 | -0.179 | 0.009 | 0.030 | -2.806 | 0.115 |
| | really really | 0.025 | -0.069 | 0.120 | 0.030 | 0.836 | 0.996 |
| really not | really really | 0.110 | 0.016 | 0.204 | 0.030 | 3.642 | **0.009** |

1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307

**B. Replication of Experiment 1: Post Hoc Comparisons for RTs - Modifiers**

| | | Mean Difference | 95% CI for Mean Difference | | SE | t | $p_{tukey}$ |
|---|---|---|---|---|---|---|---|
| | | | Lower | Upper | | | |
| # # | # not | -0.122 | -0.265 | 0.020 | 0.046 | -2.669 | 0.162 |
| | # really | -0.027 | -0.170 | 0.116 | 0.046 | -0.589 | 1.000 |
| | not # | -0.115 | -0.258 | 0.028 | 0.046 | -2.512 | 0.229 |
| | not not | -0.095 | -0.238 | 0.048 | 0.046 | -2.069 | 0.496 |
| | not really | -0.123 | -0.266 | 0.019 | 0.046 | -2.692 | 0.153 |
| | really # | -0.016 | -0.159 | 0.127 | 0.046 | -0.352 | 1.000 |
| | really not | -0.123 | -0.266 | 0.020 | 0.046 | -2.689 | 0.154 |
| | really really | 0.022 | -0.121 | 0.164 | 0.046 | 0.474 | 1.000 |
| # not | # really | 0.095 | -0.047 | 0.238 | 0.046 | 2.079 | 0.489 |
| | not # | 0.007 | -0.136 | 0.150 | 0.046 | 0.157 | 1.000 |
| | not not | 0.027 | -0.115 | 0.170 | 0.046 | 0.600 | 1.000 |
| | not really | -0.001 | -0.144 | 0.142 | 0.046 | -0.023 | 1.000 |
| | really # | 0.106 | -0.037 | 0.249 | 0.046 | 2.317 | 0.334 |
| | really not | -9.364e-4 | -0.144 | 0.142 | 0.046 | -0.020 | 1.000 |
| | really really | 0.144 | 0.001 | 0.287 | 0.046 | 3.143 | **0.046** |
| # really | not # | -0.088 | -0.231 | 0.055 | 0.046 | -1.922 | 0.599 |
| | not not | -0.068 | -0.211 | 0.075 | 0.046 | -1.479 | 0.865 |
| | not really | -0.096 | -0.239 | 0.046 | 0.046 | -2.102 | 0.473 |
| | really # | 0.011 | -0.132 | 0.154 | 0.046 | 0.238 | 1.000 |
| | really not | -0.096 | -0.239 | 0.047 | 0.046 | -2.100 | 0.475 |
| | really really | 0.049 | -0.094 | 0.191 | 0.046 | 1.063 | 0.979 |
| not # | not not | 0.020 | -0.122 | 0.163 | 0.046 | 0.443 | 1.000 |
| | not really | -0.008 | -0.151 | 0.135 | 0.046 | -0.180 | 1.000 |
| | really # | 0.099 | -0.044 | 0.242 | 0.046 | 2.160 | 0.434 |
| | really not | -0.008 | -0.151 | 0.135 | 0.046 | -0.178 | 1.000 |
| | really really | 0.137 | -0.006 | 0.280 | 0.046 | 2.986 | 0.072 |
| not not | not really | -0.029 | -0.171 | 0.114 | 0.046 | -0.623 | 0.999 |
| | really # | 0.079 | -0.064 | 0.221 | 0.046 | 1.717 | 0.736 |
| | really not | -0.028 | -0.171 | 0.114 | 0.046 | -0.620 | 0.999 |
| | really really | 0.117 | -0.026 | 0.259 | 0.046 | 2.543 | 0.214 |
| not really | really # | 0.107 | -0.036 | 0.250 | 0.046 | 2.340 | 0.320 |
| | really not | 1.182e-4 | -0.143 | 0.143 | 0.046 | 0.003 | 1.000 |
| | really really | 0.145 | 0.002 | 0.288 | 0.046 | 3.166 | **0.043** |
| really # | really not | -0.107 | -0.250 | 0.036 | 0.046 | -2.338 | 0.322 |
| | really really | 0.038 | -0.105 | 0.181 | 0.046 | 0.826 | 0.996 |
| really not | really really | 0.145 | 0.002 | 0.288 | 0.046 | 3.163 | **0.043** |

1308
1309

1310 **Table S3.** We performed a one-way ANOVA and Tukey post-hoc tests on the average RTs across

1311 trials per subject and per each modifier condition. Each line represents pairwise comparisons

1312 between each pair of modifiers, for Experiment 1 (i.e., behavioral experiment, **A**) and its replication

1313 (**B**). *p*-value and confidence intervals are adjusted for comparing a family of 9 estimates. Significant

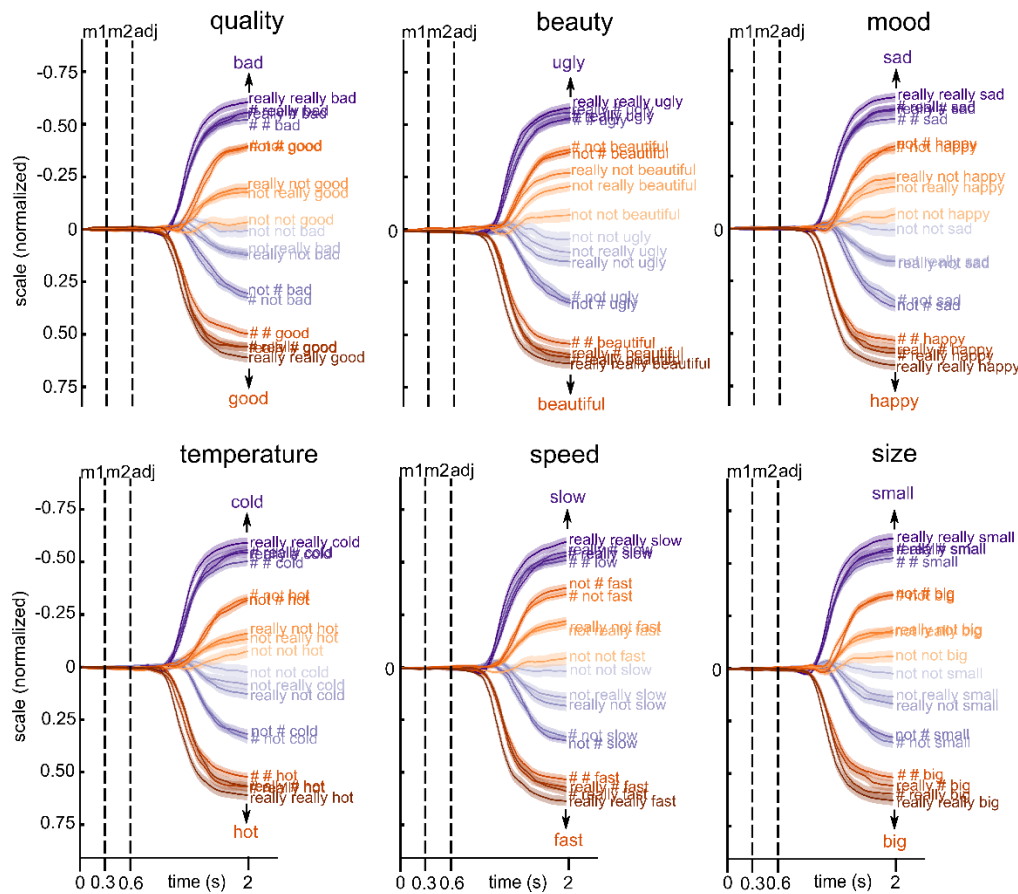1314 *p*-values are highlighted in bold.

**Figures**



**Figure S1. Trajectories for each scalar dimension.**

Behavioral trajectories for low (purples) and high (oranges) antonyms over time, for each scalar dimension (i.e., quality, beauty, mood, temperature, speed and size), for each modifier (shades of orange and purple), and for affirmative and negated phrases. Black vertical dashed lines indicate the presentation onset of each word: modifier1, modifier2 and adjective.
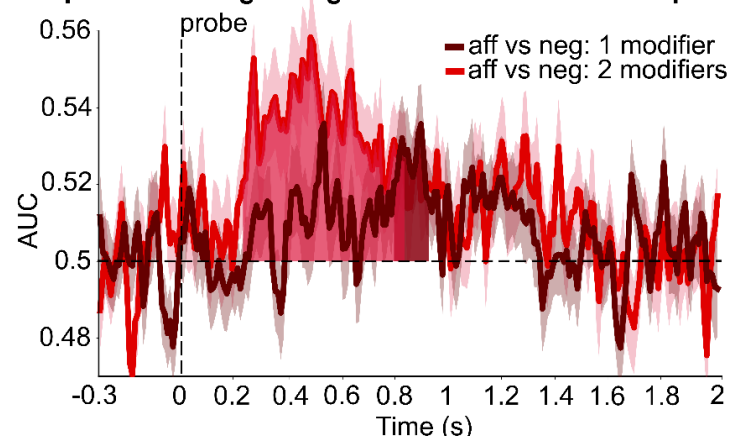
**Figure S2. Temporal decoding of negation as a function of number of modifiers (i.e., complexity), time-locked to the onset of the probe.**

Decoding accuracy of negation over time, as a function of the number of modifiers (1 modifier: dark red line and shading; 2 modifiers: light red line and shading). Significant time windows are indicated by dark red (1 modifier) and light red (2 modifiers) shading. These results show that we could significantly decode the difference between affirmative and negated phrases between 230 and 930 ms after the onset of the probe, especially when the phrase included two modifiers (1 modifier: between 790 and 930 ms: $p < 0.001$; 2 modifiers: between 230 and 840 ms: $p < 0.001$). This suggests that the representation of modifiers is reactivated at the stage when participants have to perform the yes/no task. 1 modifier: "really ###", "### really", "not ###", "### not"; 2 modifiers: "really really", "really not", "not really", "not not". AUC = area under the receiver operating characteristic curve, chance = 0.5 (black dashed horizontal line); the black vertical dashed line indicates the presentation onset of the probe; aff = affirmative, neg = negated; each line and shading represent participants mean ± SEM.
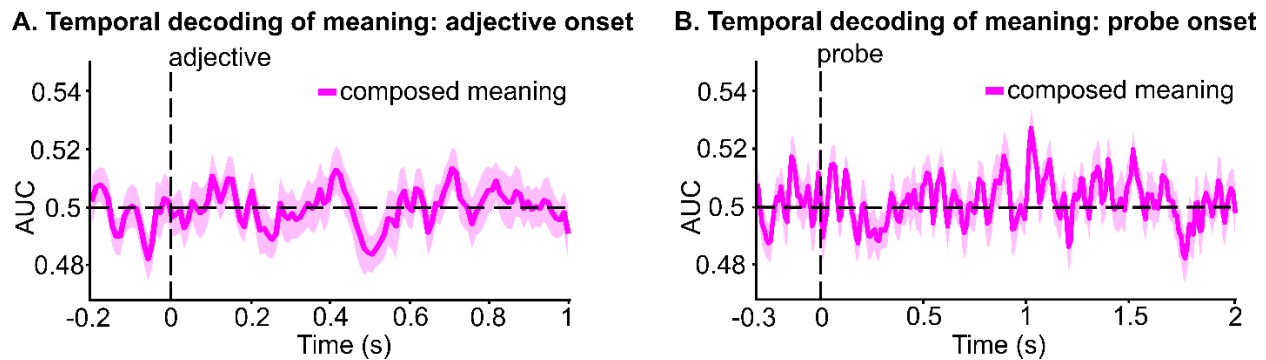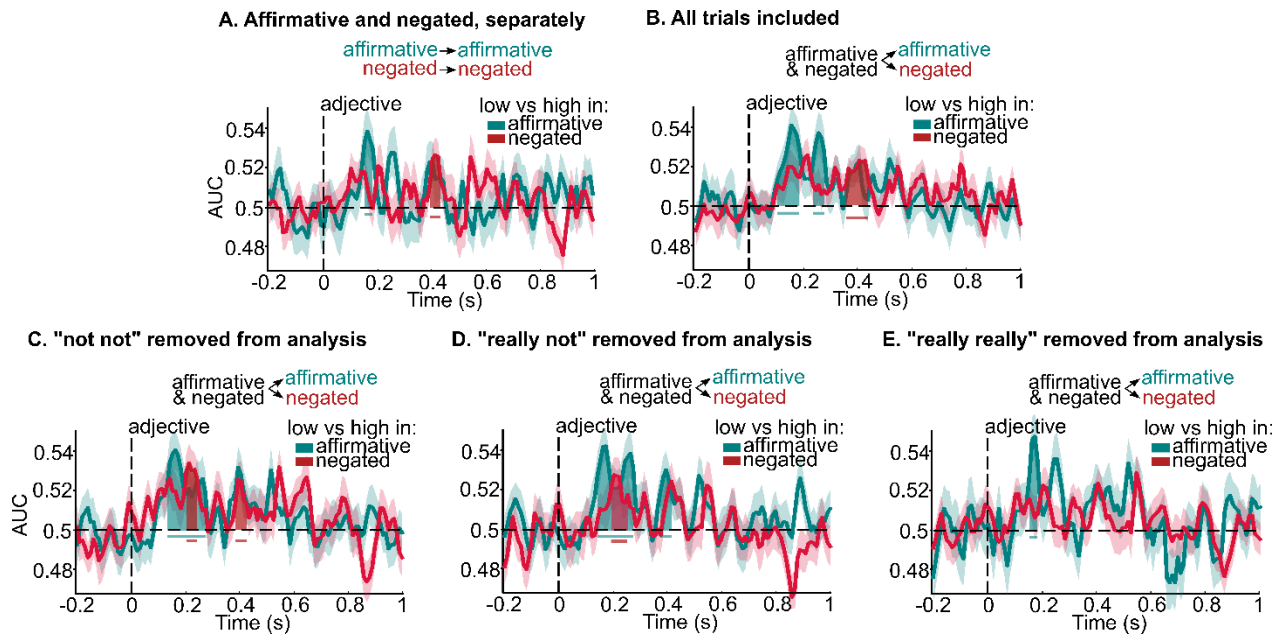
**Figure S3. Temporal decoding of composed meaning.**

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

48

1369



1370

1371

1372 **Figure S4. Follow-up analyses of Fig.6C.**

1373 **A.** We conducted a follow-up analysis where we trained and tested on "low" vs. "high" antonyms

1374 in affirmative and negated phrases separately, to further investigate lowering in decoding accuracy

1375 when representations are closer on the semantic scale, as predicted by the mitigation hypothesis for

1376 negated phrases. We found similar patterns to our main analysis. Results show that affirmative

1377 phrases (green line) are associated with significantly above-chance decoding accuracy between 150

1378 and 190 ms ($p = 0.026$; green shading and horizontal solid line) from adjective onset. No significant

1379 above-chance decoding accuracy was found for negated phrases before ~400 ms from adjective

1380 onset (390 to 440 ms, $p = 0.009$; red shading and horizontal solid line). **B.** We conducted a follow-

1381 up analysis where no trials were removed due to the feedback score. We found similar patterns to

1382 our main analysis. Results show that affirmative phrases (green line) are associated with

1383 significantly above-chance decoding accuracy between 100 and 190 ms and 230 and 280 ms from

1384 adjective onset ($p = 0.001$ and $p = 0.032$ respectively, green shading and horizontal solid lines).

1385 Negative phrases (red line) are associated with significantly above-chance decoding accuracy

1386 between 350 to 440 ms from adjective onset ($p < 0.001$, red shading and horizontal solid line).

1387 **C.D.E.** We conducted a series of follow-up analyses where we removed one condition (i.e., one

1388 modifiers combination) at a time to evaluate its specific effect on adjective representation. **C.** "not

1389 not" is removed from the analysis: affirmative phrases (green line) are associated with significantly

1390 above-chance decoding accuracy between 130 and 280 ms from adjective onset ($p < 0.001$, green

1391 shading and horizontal solid line), negative phrases (red line) are associated with significantly

above-chance decoding accuracy between 200 to 250 ms and between 380 to 430 ms from adjective onset ($p = 0.011$ and $p = 0.049$, red shading and horizontal solid lines). **D.** "really not" is removed from the analysis: affirmative phrases (green line) are associated with significantly above-chance decoding accuracy between 140 and 280 ms and between 370 and 420 ms from adjective onset ($p = 0.001$ and $p = 0.038$, green shading and horizontal solid lines), negative phrases (red line) are associated with significantly above-chance decoding accuracy between 190 to 260 ms from adjective onset ($p = 0.009$, red shading and horizontal solid lines). **E.** "really really" is removed from the analysis. affirmative phrases (green line) are associated with significantly above-chance decoding accuracy between 150 and 190 ms from adjective onset ($p = 0.025$, green shading and horizontal solid line), no significant above-chance decoding accuracy was found for negated phrases. Overall, these results suggest that "not not" and "really not" have similar mitigation effects. Conversely, and as expected, "really really" does not have mitigation effects on adjective representation.
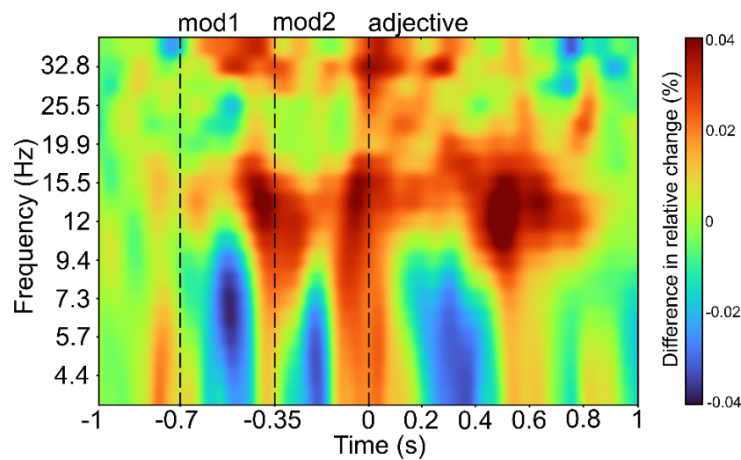
**Figure S5. Differences between negated and affirmative phrases across time and frequencies.** Time-frequency spectrum of the differences between negated and affirmative phrases averaged across all sensors and all participants. Frequencies are between 3.9 and 37.2 Hz, logarithmically spaced. Black vertical dashed lines indicate the presentation onset of each word: modifier1, modifier2 and adjective; colors indicate % differences in change relative to a baseline of -300 to -100 ms from the onset of word 1 (modifier1).
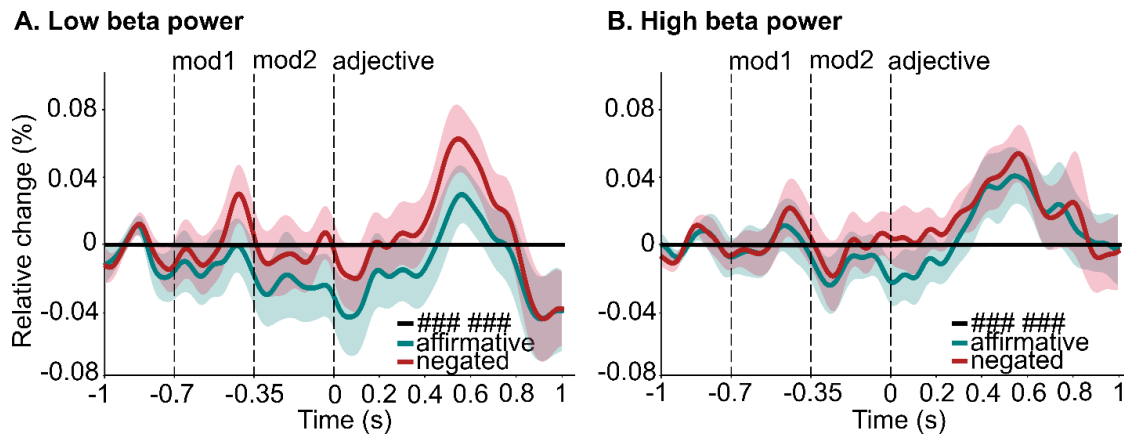
**Figure S6. Low- and high-beta power for negated and affirmative phrases across time.**

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