

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

**Experience Transforms Crossmodal Object Representations in the
Anterior Temporal Lobes**

Aedan Y. Li¹, Natalia Ladyka-Wojcik¹, Heba Qazilbash¹, Ali Golestani², Dirk B. Walther^{1,3},
Chris B. Martin⁴, Morgan D. Barense^{1,3}

Department of Psychology, University of Toronto¹
Department of Physics and Astronomy, University of Calgary²
Rotman Research Institute, Baycrest Health Sciences³
Department of Psychology, Florida State University⁴

Authors Note:

Earlier versions of this manuscript were presented as a talk at the Vision Sciences Society in 2022, as a poster at the Canadian Society for Brain, Behaviour and Cognitive Science in 2021, and as a poster at the Lake Ontario Visionary Establishment in 2020. The authors report no conflict of interest. Anonymized data are available on the Open Science Framework:

<https://osf.io/vq4wj/>. Univariate maps are available on NeuroVault:

<https://neurovault.org/collections/LFDCGMAY/>.

Correspondence concerning this article should be addressed to Aedan Li, Department of Psychology, University of Toronto, 100 St. George Street, Toronto, ON, Canada, M5S 3G3.

Contact: aedanyueli@gmail.com

Word Count: 5,658

34

Abstract

35 Combining information from multiple senses is essential to object recognition, core to the ability
36 to learn concepts, make new inferences, and generalize across distinct entities. Yet how the mind
37 combines sensory input into coherent crossmodal representations – the *crossmodal binding*
38 *problem* – remains poorly understood. Here, we applied multi-echo fMRI across a four-day
39 paradigm, in which participants learned 3-dimensional crossmodal representations created from
40 well-characterized unimodal visual shape and sound features. Our novel paradigm decoupled the
41 learned crossmodal object representations from their baseline unimodal shapes and sounds, thus
42 allowing us to track the emergence of crossmodal object representations as they were learned by
43 healthy adults. Critically, we found that two anterior temporal lobe structures – temporal pole
44 and perirhinal cortex – differentiated learned from non-learned crossmodal objects, even when
45 controlling for the unimodal features that composed those objects. These results provide
46 evidence for integrated crossmodal object representations in the anterior temporal lobes that were
47 different from the representations for the unimodal features. Furthermore, we found that
48 perirhinal cortex representations were by default biased towards visual shape, but this initial
49 visual bias was attenuated by crossmodal learning. Thus, crossmodal learning transformed
50 perirhinal representations such that they were no longer predominantly grounded in the visual
51 modality, which may be a mechanism by which object concepts gain their abstraction.

52 *Keywords:* Crossmodal binding problem, object representations, integrative coding,
53 distributed unimodal features, multi-echo fMRI

54

55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85

Acknowledgements

We are grateful to the Toronto Neuroimaging community for helpful feedback. In particular, the first author thanks Dr. Katherine Duncan and Dr. Massieh Moayedi for suggestions related to the experimental design, Dr. Michael Mack for initial guidance with 3D-printing, Dr. Rosanna Olsen for her tutorial on medial temporal lobe segmentation, as well as Dr. Andy Lee and Dr. Adrian Nestor for their neuroimaging and multivariate pattern analysis courses.

We thank Annie Kim and Katarina Savel for their assistance with participant recruitment, as well as Priya Abraham for her assistance with MRI scanning. Finally, we thank Dr. Rüdiger Stirnberg for sharing with us the multi-echo fMRI sequence used in this manuscript.

AYL is supported by an Alexander Graham Bell Canada Graduate Scholarship-Doctoral from the Natural Sciences and Engineering Research Council of Canada (NSERC CGS-D). This work is supported by a Scholar Award from the James S McDonnell Foundation, an Early Researcher Award from the Ontario Government, an NSERC Discovery grant, and a Canada Research Chair to MDB.

86 Experience Transforms Crossmodal Object Representations in the Anterior Temporal Lobes

87 The world is a great blooming, buzzing confusion¹ of the senses. Our ability to
88 understand “what is out there” depends on combining sensory features to form *crossmodal object*
89 *concepts*. A child, for example, might form the concept “frog” by learning that the visual
90 appearance of a four-legged creature goes with the sound of its croaking. Consequently, this
91 child has also learned that frogs do not produce barking sounds, as the child has created a unique
92 object association for a frog from specific unimodal shape and sound features. Forming coherent
93 crossmodal object representations is thus essential for human experience, allowing adaptive
94 behavior under changing environments. Yet how is it possible for the child to know that the
95 sound of croaking is associated with the visual shape of a frog, even when she might be looking
96 at a dog? How does the human mind form meaningful concepts from the vast amount of
97 unimodal feature information that bombards the senses, allowing us to interpret our external
98 world?

99 Known as the *crossmodal binding problem*, this unresolved question in the cognitive
100 sciences concerns how the mind combines unimodal sensory features into coherent crossmodal
101 object representations. Better characterization of how this computational challenge is solved will
102 not only improve our understanding of the human mind but will also have important
103 consequences for the design of future artificial neural networks. Current artificial machines do
104 not yet reach human performance on tasks involving crossmodal integration^{2,3} or generalization
105 beyond previous experience,^{4,5,6} which are limitations thought to be in part driven by the inability
106 of existing machines to resolve the binding problem.⁷

107 One theoretical view from the cognitive sciences suggests that crossmodal objects are
108 built from component unimodal features represented across distributed sensory regions.⁸ Under
109 this view, when a child thinks about “frog”, the visual cortex represents the appearance of the
110 shape of the frog, whereas the auditory cortex represents the croaking sound. Alternatively, other
111 theoretical views predict that multisensory objects are not only built from their component
112 unimodal sensory features, but that there is also a crossmodal integrative code that is different
113 from the sum of these parts.^{9,10,11,12,13} These latter views propose that anterior temporal lobe
114 structures can act as a polymodal “hub” that combines separate features into integrated
115 wholes.^{9,11,14,15}

116 Thus, a key theoretical challenge central to resolving the crossmodal binding problem is
117 understanding how anterior temporal lobe structures form object representations. Are crossmodal
118 objects entirely built from features distributed across sensory regions, or is there also integrative
119 coding in the anterior temporal lobes? Furthermore, the existing literature has predominantly
120 studied the neural representation of well-established object concepts from the visual domain
121 alone,⁸⁻²⁵ even though human experience is fundamentally crossmodal.

122 Here, we leveraged multi-echo fMRI²⁶ across a novel four-day task in which participants
123 learned to associate unimodal visual shape and sound features into 3D crossmodal object

124 representations. First, we characterized shape²⁷ and sound features in a separate validation
125 experiment, ensuring that the unimodal features were well-matched in terms of their subjective
126 similarity (*Figure 1*). On the learning task, participants independently explored the 3D-printed
127 shapes and heard novel experimenter-constructed sounds. The participants then learned specific
128 shape-sound associations (congruent objects), while other shape-sound associations were not
129 learned (incongruent objects).

130 Critically, our four-day learning task allowed us to isolate neural activity associated with
131 integrative coding in anterior temporal lobe structures that emerges with experience and differs
132 from the neural patterns recorded at baseline. The learned and non-learned crossmodal objects
133 were constructed from the same set of three validated shape and sound features, ensuring that
134 factors such as familiarity with the unimodal features, subjective similarity, and feature identity
135 were tightly controlled (*Figure 2*). If the mind represented crossmodal objects entirely as the
136 reactivation of unimodal shapes and sounds (i.e., objects are constructed from their parts), then
137 there should be no difference between the learned and non-learned objects (because they were
138 created from the same three shapes and sounds). By contrast, if the mind represented crossmodal
139 objects as something over and above their component features (i.e., representations for
140 crossmodal objects rely on integrative coding that is different from the sum of their parts), then
141 there should be behavioral and neural differences between learned and non-learned crossmodal
142 objects (because the only difference across the objects is the learned relationship between the
143 parts). Furthermore, this design allowed us to determine the relationship between the object
144 representation acquired *after* crossmodal learning and the unimodal feature representations
145 acquired *before* crossmodal learning. That is, we could examine whether learning led to
146 abstraction of the object representations such that it no longer resembled the unimodal feature
147 representations.

148 In brief, we found that crossmodal object concepts were represented as distributed
149 sensory-specific unimodal features along the visual and auditory processing pathways, as well as
150 integrative crossmodal combinations of those unimodal features in the anterior temporal lobes.
151 Intriguingly, the perirhinal cortex – an anterior temporal lobe structure – was biased towards the
152 visual modality before crossmodal learning at baseline, with greater activity towards shape over
153 sound features. Pattern similarity analyses revealed that the shape representations in perirhinal
154 cortex were initially unaffected by sound, providing evidence of a default visual shape bias.
155 However, crossmodal learning transformed the object representation in perirhinal cortex such
156 that it was no longer predominantly visual. These results are consistent with the idea that the
157 object representation had become abstracted away from the component unimodal features with
158 learning, such that perirhinal representations was no longer grounded in the visual modality.

159 **Results**

160 **Four-Day Crossmodal Object Learning Task**

161 *Measuring Within-Subject Changes After Crossmodal Learning*

162 We designed a 4-day learning task where each participant learned a set of shape-sound
163 associations that created crossmodal objects (*Figure 2*). There were two days involving only
164 behavioral measures (*Day 1* and *Day 3*). Before crossmodal learning on Day 1, participants
165 explored the 3D-printed shapes (*Visual*) and heard the sounds (*Sound*) separately. In blocks of
166 trials interleaved with these exploration phases, participants rated the similarity of the shapes and
167 sounds (see *Figure 2–figure supplement 1*). During crossmodal learning on Day 3, participants
168 explored specific shape-sound associations (*Congruent* objects) by pressing the button on each
169 3D-printed shape to play the associated sound, with pairings counterbalanced across observers.
170 Again, the participants rated the similarity of the shapes and sounds. Notably, all participants
171 could recognize their specific shape-sound associations at the end of Day 3, confirming that the
172 congruent shape-sound objects were successfully learned (performance = 100% for all
173 participants).

174 There were two neuroimaging days (*Day 2* and *Day 4*), during which we recorded brain
175 responses to unimodal features presented separately and to unimodal features presented
176 simultaneously using multi-echo fMRI (*Figure 2*). During Unimodal Feature runs, participants
177 either viewed images of the 3D-printed shapes or heard sounds. During Crossmodal Object runs,
178 participants experienced either the shape-sound associations learned on Day 3 (*Congruent*) or
179 shape-sound associations that had not been learned on Day 3 (*Incongruent*). We were especially
180 interested in neural differences between congruent and incongruent objects as evidence of
181 crossmodal integration; experience with the unimodal features composing congruent and
182 incongruent objects was equated and the only way to distinguish them was in terms of how the
183 features were integrated.

184

185 **Behavioral Pattern Similarity**

186 *Subjective Similarity Changes After Crossmodal Learning*

187 To understand how crossmodal learning impacts behaviour, we analyzed the within-
188 subject change in subjective similarity of the unimodal features *before* (Day 1) and *after* (Day 3)
189 participants learned their crossmodal pairings (*Figure 2*). In other words, we determined whether
190 the perceived similarity of the unimodal feature representations changed after participants had
191 experienced those unimodal features combined into crossmodal objects.

192 We conducted a linear mixed model which included learning day (before vs. after
193 crossmodal learning) and congruency (congruent vs. incongruent) as fixed effects. We observed
194 a robust learning-related behavioral change in terms of how participants experienced the
195 similarity of shape and sound features (*Figure 2–figure supplement 1*): there was a main effect of
196 learning day (before or after crossmodal learning: $F_{1,51} = 24.45$, $p < 0.001$, $\eta^2 = 0.32$), a main
197 effect of congruency (congruent or incongruent: $F_{1,51} = 6.93$, $p = 0.011$, $\eta^2 = 0.12$), and an

198 interaction between learning day and congruency ($F_{1,51} = 15.33, p < 0.001, \eta^2 = 0.23$). Before
199 crossmodal learning, there was no difference in similarity between congruent and incongruent
200 shape-sound features ($t_{17} = 0.78, p = 0.44$), whereas after crossmodal learning, participants rated
201 shapes and sounds associated with congruent objects to be more similar than shapes and sounds
202 associated with incongruent objects ($t_{17} = 5.10, p < 0.001, \text{Cohen's } d = 1.28$) (*Figure 2–figure*
203 *supplement 1*). Notably, this learning-related change in similarity was observed in 17 out of 18
204 participants. We confirmed this experience-dependent change in similarity structure in a separate
205 behavioral experiment with a larger sample size (observed in 38 out of 44 participants; learning
206 day x congruency interaction: $F_{1,129} = 13.74, p < 0.001; \eta^2 = 0.096$; *Figure 2–figure supplement*
207 *1*).

208

209 **Whole-brain Univariate Analysis**

210 *Unimodal Shape and Sound Representations are Distributed*

211 In the first set of neuroimaging analyses, we examined whether distributed brain regions
212 were involved in representing unimodal shapes and sounds. During unimodal runs (shapes and
213 sounds presented separately), we observed robust bilateral modality-specific activity across the
214 neocortex (*Figure 3a-c*). The ventral visual stream extending into the perirhinal cortex activated
215 more strongly to unimodal visual compared to sound information, indicating that perirhinal
216 cortex activity was by default biased towards visual information in the unimodal runs (i.e.,
217 towards complex visual shape configurations; *Figure 3a*). The auditory processing stream, from
218 the primary auditory cortex extending into the temporal pole along the superior temporal sulcus,
219 activated more strongly to unimodal sound compared to visual information (*Figure 3b*). These
220 results replicate the known representational divisions across the neocortex and show that regions
221 processing unimodal shapes and sounds are distributed across visual and auditory processing
222 pathways.^{29,30,31} Furthermore, the robust signal quality we observe in anterior temporal regions
223 demonstrates the improved quality of the multi-echo ICA pipeline employed in the current study,
224 as these anterior temporal regions are often susceptible to signal dropout with standard single
225 echo designs due to magnetic susceptibility issues near the sinus air/tissue boundaries (*Figure 3*
226 *–figure supplement 1*).

227

228 **Region-of-Interest Univariate Analysis**

229 *Anterior Temporal Lobes Differentiate Between Congruent and Incongruent Conditions*

230 We next examined univariate activity focusing on five *a priori* regions thought to be
231 important for representing unimodal features and their integration:^{9,11} temporal pole, perirhinal
232 cortex, lateral occipital complex (LOC), primary visual cortex (V1), and primary auditory cortex
233 (A1). For each ROI, we conducted a linear mixed model which included learning day (before vs.
234 after crossmodal learning) and modality (visual vs. sound feature) as fixed factors. Collapsing
235 across learning days, perirhinal cortex ($t_{67} = 5.53, p < 0.001, \text{Cohen's } d = 0.67$) and LOC ($t_{63} =$
236 $16.02, p < 0.001, \text{Cohen's } d = 2.00$) were biased towards visual information, whereas temporal

237 pole ($t_{67} = 6.73, p < 0.001, \text{Cohen's } d = 0.82$) and A1 ($t_{67} = 17.09, p < 0.001, \text{Cohen's } d = 2.07$)
238 were biased towards sound information (*Figure 3d*). Interestingly, we found a small overall bias
239 towards sound in V1, consistent with past work³² ($t_{67} = 2.26, p = 0.027, \text{Cohen's } d = 0.20$). Next,
240 we determined how neural responses in these regions changed following crossmodal learning.
241 We observed an interaction between learning day and modality in perirhinal cortex ($F_{1,48} = 5.24,$
242 $p = 0.027, \eta^2 = 0.098$) and LOC ($F_{1,45} = 25.89, p < 0.001, \eta^2 = 0.37$) (*Figure 3d*). These regions
243 activated more strongly to visual information at baseline before crossmodal learning compared to
244 after crossmodal learning, indicative of a visual bias that was attenuated with experience.

245 As a central goal of our study was to identify brain regions that were influenced by the
246 learned crossmodal associations, we next examined univariate differences between *Congruent*
247 *vs. Incongruent* for crossmodal object runs as a function of whether the crossmodal association
248 had been learned. We conducted a linear mixed model for each ROI which included learning day
249 (before vs. after crossmodal learning) and congruency (congruent vs. incongruent objects) as
250 fixed factors. We observed a significant interaction between learning day and congruency in the
251 temporal pole ($F_{1,48} = 7.63, p = 0.0081, \eta^2 = 0.14$). Critically, there was no difference in activity
252 between congruent and incongruent objects at baseline before crossmodal learning ($t_{33} = 0.37, p$
253 $= 0.72$), but there was more activation to incongruent compared to congruent objects after
254 crossmodal learning ($t_{33} = 2.42, p = 0.021, \text{Cohen's } d = 0.42$). As the unimodal shape-sound
255 *features* experienced by participants were the same before and after crossmodal learning (*Figure*
256 *2*), this finding reveals that the univariate signal in the temporal pole was differentiated between
257 congruent and incongruent objects that had been constructed from the same unimodal features.

258 By contrast, we did not observe a univariate difference between the congruent and
259 incongruent conditions in the perirhinal cortex, LOC, V1, or A1 ($F_{1,45-48}$ between 0.088 and 2.34,
260 p between 0.13 and 0.77). Similarly, the exploratory ROIs hippocampus (HPC: $F_{1,48} = 0.32, p =$
261 0.58) and inferior parietal lobe (IPL: $F_{1,48} = 0.094, p = 0.76$) did not distinguish between the
262 congruent and incongruent conditions.

263

264 **Neural Pattern Similarity**

265 *Congruent Associations Differ from Incongruent Associations in Anterior Temporal Lobes*

266 We next conducted a series of representational similarity analyses across Unimodal
267 Feature and Crossmodal Object runs before and after crossmodal learning. Here, we investigated
268 whether representations for unimodal features were changed after learning the crossmodal
269 associations between those features (i.e., learning the crossmodal pairings that comprised the
270 shape-sound objects). Such a finding could be taken as evidence that learning crossmodal *object*
271 concepts transforms the original representation of the component unimodal *features*. More
272 specifically, we compared the correlation between congruent and incongruent shape-sound
273 features within Unimodal Feature runs before and after crossmodal learning (*Figure 4a*).

274 We conducted a linear mixed model which included learning day (before vs. after
275 crossmodal learning) and congruency (congruent vs. incongruent) as fixed effects for each ROI.

276 Complementing the previous behavioral pattern similarity results (*Figure 2–figure supplement*
277 *1*), in the temporal pole we observed a main effect of learning day (before or after crossmodal
278 learning: $F_{1,32} = 4.63$, $p = 0.039$, $\eta^2 = 0.13$), a main effect of congruency (congruent or
279 incongruent object: $F_{1,64} = 7.60$, $p = 0.0076$, $\eta^2 = 0.11$), and an interaction between learning day
280 and congruency ($F_{1,64} = 6.09$, $p = 0.016$, $\eta^2 = 0.087$). At baseline before crossmodal learning,
281 there was no difference in pattern similarity between congruent features compared to incongruent
282 features in the temporal pole ($t_{33} = 0.22$, $p = 0.82$). After crossmodal learning, however, there
283 was lower pattern similarity for shape and sound features associated with congruent compared to
284 incongruent objects ($t_{33} = 3.47$, $p = 0.0015$, *Cohen's d* = 0.22; *Figure 4*). Thus, although in
285 behavior we observed that learning the crossmodal associations led to greater pattern similarity
286 between congruent compared to incongruent features (*Figure 2–figure supplement 1*), this
287 *greater behavioral similarity* was related to *reduced neural similarity* following crossmodal
288 learning in the temporal pole.

289 By contrast, the other four a priori determined ROIs (perirhinal cortex, LOC, V1, or A1)
290 did not show an interaction between learning day and congruency ($F_{1,60-64}$ between 0.039 and
291 1.30, p between 0.26 and 0.84; *Figure 4 – figure supplement 1*). Likewise, our 2 exploratory
292 ROIs (hippocampus, inferior parietal lobe) did not show an interaction between learning day and
293 congruency ($F_{1,64}$ between 0.68 and 0.91, p between 0.34 and 0.41; *Figure 5 – figure supplement*
294 *1*).

295

296 *The Visually-biased Code in Perirhinal Cortex was Attenuated with Learning*

297 The previous analyses found that the temporal pole differentiated between congruent and
298 incongruent shape-sound pairs after participants learned the crossmodal pairings. Next, we
299 characterized how the representations of these unimodal features changed after they had been
300 paired with features from another stimulus modality to form the crossmodal objects. Our key
301 question was whether learning crossmodal associations transformed the unimodal feature
302 representations.

303 First, the voxel-wise activity for unimodal feature runs was correlated to the voxel-wise
304 activity for crossmodal object runs at baseline before crossmodal learning (*Figure 5a*).
305 Specifically, we quantified the similarity in the patterns for the visual *shape features* with the
306 *crossmodal objects* that had that same shape, as well as between the *sound features* and the
307 *crossmodal objects* that had that same sound. We then conducted a linear mixed model which
308 included modality (visual vs. sound) as a fixed factor within each ROI. Consistent with the
309 univariate results (*Figure 3*), we observed greater pattern similarity when there was a match
310 between sound features in the temporal pole ($F_{1,32} = 15.80$, $p < 0.001$, $\eta^2 = 0.33$) and A1 ($F_{1,32} =$
311 145.73 , $p < 0.001$, $\eta^2 = 0.82$), and greater pattern similarity when there was a match in the visual
312 shape features in the perirhinal cortex ($F_{1,32} = 10.99$, $p = 0.0023$, $\eta^2 = 0.26$), LOC ($F_{1,30} = 20.09$,
313 $p < 0.001$, $\eta^2 = 0.40$), and V1 ($F_{1,32} = 22.02$, $p < 0.001$, $\eta^2 = 0.41$). Pattern similarity for each ROI
314 was higher for one of the two modalities, indicative of a baseline modality-specific bias towards
315 either visual or sound content.

316 We then examined whether the original representations would change after participants
317 learned how the features were paired together to make specific crossmodal objects, conducting
318 the same analysis described above after crossmodal learning had taken place (*Figure 5b*). With
319 this analysis, we sought to measure the relationship between the representation for the learned
320 crossmodal object and the original baseline representation for the unimodal features. More
321 specifically, the voxel-wise activity for unimodal feature runs *before* crossmodal learning was
322 correlated to the voxel-wise activity for crossmodal object runs *after* crossmodal learning
323 (*Figure 5b*). Another linear mixed model which included modality as a fixed factor within each
324 ROI revealed that the perirhinal cortex was no longer biased towards visual shape after
325 crossmodal learning ($F_{1,32} = 0.12, p = 0.73$), whereas the temporal pole, LOC, V1, and A1
326 remained biased towards either visual shape or sound ($F_{1,30-32}$ between 16.20 and 73.42, all $p <$
327 $0.001, \eta^2$ between 0.35 and 0.70).

328 To investigate this effect in perirhinal cortex more specifically, we conducted a linear
329 mixed model to directly compare the change in the visual bias of perirhinal representations from
330 before crossmodal learning to after crossmodal learning (green regions in *Figure 5a* vs. *5b*).
331 Specifically, the linear mixed model included learning day (before vs. after crossmodal learning)
332 and modality (visual feature match to crossmodal object vs. sound feature match to crossmodal
333 object). Results revealed a significant interaction between learning day and modality in the
334 perirhinal cortex ($F_{1,775} = 5.56, p = 0.019, \eta^2 = 0.071$), meaning that the baseline visual shape
335 bias observed in perirhinal cortex (green region of *Figure 5a*) was significantly attenuated with
336 experience (green region of *Figure 5b*). After crossmodal learning, a given shape no longer
337 invoked significant pattern similarity between objects that had the same shape but differed in
338 terms of what they sounded like. Taken together, these results suggest that prior to learning the
339 crossmodal objects, the perirhinal cortex had a default bias toward representing the visual shape
340 information and was not representing sound information of the crossmodal objects. After
341 crossmodal learning, however, the visual shape bias in perirhinal cortex was no longer present.
342 That is, with crossmodal learning, the representations within perirhinal cortex started to look less
343 like the visual features that comprised the crossmodal objects, providing evidence that the
344 perirhinal representations were no longer predominantly grounded in the visual modality.

345 To examine whether these results differed by congruency (i.e., whether any modality-
346 specific biases differed as a function of whether the object was congruent or incongruent), we
347 conducted exploratory linear mixed models for each of the five *a priori* ROIs across learning
348 days. More specifically, we correlated: 1) the voxel-wise activity for Unimodal Feature Runs
349 *before* crossmodal learning to the voxel-wise activity for Crossmodal Object Runs *before*
350 crossmodal learning (Day 2 vs. Day 2), 2) the voxel-wise activity for Unimodal Feature Runs
351 *before* crossmodal learning to the voxel-wise activity for Crossmodal Object Runs *after*
352 crossmodal learning (Day 2 vs Day 4), and 3) the voxel-wise activity for Unimodal Feature Runs
353 *after* crossmodal learning to the voxel-wise activity for Crossmodal Object Runs *after*
354 crossmodal learning (Day 4 vs Day 4). For each of the three analyses described, we then
355 conducted separate linear mixed models which included modality (visual feature match to

356 crossmodal object vs. sound feature match to crossmodal object) and congruency (congruent vs.
357 incongruent).

358 There was no significant relationship between modality and congruency in any ROI
359 between Day 2 and Day 2 ($F_{1,346-368}$ between 0.00 and 1.06, p between 0.30 and 0.99), between
360 Day 2 and Day 4 ($F_{1,346-368}$ between 0.021 and 0.91, p between 0.34 and 0.89), or between Day 4
361 and Day 4 ($F_{1,346-368}$ between 0.01 and 3.05, p between 0.082 and 0.93). However, exploratory
362 analyses revealed that perirhinal cortex was the only region without a modality-specific bias and
363 where the unimodal feature runs were not significantly correlated to the crossmodal object runs
364 *after crossmodal learning (Figure 5 – figure supplement 2)*.

365 Taken together, the overall pattern of results suggests that representations of the
366 crossmodal objects in perirhinal cortex were heavily influenced by their consistent visual
367 features *before* crossmodal learning. However, the crossmodal object representations were no
368 longer influenced by the component visual features *after* crossmodal learning (*Figure 5, Figure 5*
369 *– figure supplement 2*). Additional exploratory analyses did not find evidence of experience-
370 dependent changes in the hippocampus or inferior parietal lobes (*Figure 5 – figure supplement*
371 *1*).

372 Importantly, the change in pattern similarity in the perirhinal cortex across learning days
373 (*Figure 5*) is unlikely to be driven by noise, poor alignment of patterns across sessions, or
374 generally reduced responses. Other regions with numerically similar pattern similarity to
375 perirhinal cortex did not change across learning days (e.g., visual features x crossmodal objects
376 in A1 in *Figure 5*; the exploratory ROI hippocampus with numerically similar pattern similarity
377 to perirhinal cortex also did not change in *Figure 5 – figure supplement 1*).

378

379 *Representations in Perirhinal Cortex Change with Experience*

380 So far, we have shown that the perirhinal cortex was by default biased towards visual
381 shape features (*Figure 5a*), and that this visual shape bias was attenuated with experience
382 (*Figure 5b; Figure 5 – figure supplement 2*). In the final analysis, we tracked how the *individual*
383 *crossmodal object representations* themselves change after crossmodal learning.

384 We assessed the cross-day pattern similarity between Crossmodal Object Runs by
385 correlating the congruent and incongruent runs across learning days (*Figure 6*). We then
386 conducted a linear mixed model which included congruency (congruent vs. incongruent) as a
387 fixed factor for each *a priori* ROI. Perirhinal cortex was the only region that differentiated
388 between congruent and incongruent objects in this analysis (PRC: $F_{1,34} = 4.67$, $p = 0.038$, $\eta^2 =$
389 0.12 ; TP, LOC, V1, A1: $F_{1,32-34}$ between 0.67 and 2.83, p between 0.10 and 0.42). Pattern
390 similarity in perirhinal cortex did not differ from 0 for congruent objects across learning days (t_{35}
391 $= 0.39$, $p = 0.70$) but was significantly lower than 0 for incongruent objects ($t_{35} = 2.63$, $p = 0.013$,
392 *Cohen's d* = 0.44). By contrast, pattern similarity in temporal pole, LOC, V1, and A1 was
393 significantly correlated across learning days (pattern similarity > 0; t_{33-35} between 4.31 and 6.92

394 all $p < 0.001$) and did not differ between congruent and incongruent objects (temporal pole,
395 LOC, V1, and A1; $F_{1,32-34}$ between 0.67 and 2.83, p between 0.10 and 0.42). Thus, perirhinal
396 cortex was unique in that it not only differentiated between congruent and incongruent objects
397 that were built from the same unimodal features (i.e., representations of the whole crossmodal
398 object that was different than the unimodal features that composed it), but it also showed no
399 significant pattern similarity above 0 for the same representations across learning days (i.e.,
400 suggesting that the object representations were transformed after crossmodal learning).

401 No significant difference between the congruent and incongruent conditions were
402 observed for the hippocampus ($F_{1,34} = 0.34$, $p = 0.56$) or inferior parietal lobe ($F_{1,34} = 0.00$, $p =$
403 0.96) in a follow-up exploratory analysis (*Figure 5 – figure supplement 1*).

404

405 Discussion

406 Known as the *crossmodal binding problem*, a long-standing question in the cognitive
407 sciences has asked how the mind forms coherent concepts from multiple sensory modalities. To
408 study this problem, we designed a 4-day task to decouple the learned crossmodal object
409 representations (Day 3 and 4) from the baseline unimodal shape and sound features (Day 1 and
410 2). We equated the familiarity, subjective similarity, and identity of the unimodal feature
411 representations composing the learned (congruent) and unlearned (incongruent) objects, ensuring
412 that any differences between the two would not be driven by single features but rather by the
413 integration of those features (*Figure 2*). Paired with multi-echo fMRI to improve signal quality
414 in the anterior temporal lobes (*Figure 3 – figure supplement 1*), this novel paradigm tracked the
415 emergence of crossmodal object concepts from component baseline unimodal features in healthy
416 adults.

417 We found that the temporal pole and perirhinal cortex – two anterior temporal lobe
418 structures – came to represent new crossmodal object concepts with learning, such that the
419 acquired crossmodal object representations were different from the representation of the
420 constituent unimodal features (*Figure 5, 6*). Intriguingly, the perirhinal cortex was by default
421 biased towards visual shape, but that this initial visual bias was attenuated with experience
422 (*Figure 3c, 5, Figure 5 – figure supplement 2*). Within the perirhinal cortex, the acquired
423 crossmodal object concepts (measured after crossmodal learning) became less similar to their
424 original component unimodal features (measured at baseline before crossmodal learning); *Figure*
425 *5, 6, Figure 5 – figure supplement 2*. This is consistent with the idea that object representations
426 in perirhinal cortex integrate the component sensory features into a whole that is different from
427 the sum of the component parts, which might be a mechanism by which object concepts obtain
428 their abstraction.

429 As one solution to the crossmodal binding problem, we suggest that the temporal pole
430 and perirhinal cortex form unique crossmodal object representations that are different from the
431 distributed features in sensory cortex (*Figure 4, 5, 6, Figure 5 – figure supplement 2*). However,
432 the nature by which the integrative code is structured and formed in the temporal pole and
433 perirhinal cortex following crossmodal experience – such as through transformations, warping,

434 or other factors – is an open question and an important area for future investigation. Furthermore,
435 these distinct anterior temporal lobe structures may be involved with integrative coding in
436 different ways. For example, the crossmodal object representations measured after learning were
437 found to be related to the component unimodal feature representations measured before learning
438 in the temporal pole but not the perirhinal cortex (*Figure 5, 6, Figure 5 – figure supplement 2*).
439 Moreover, pattern similarity for congruent shape-sound pairs were lower than the pattern
440 similarity for incongruent shape-sound pairs after crossmodal learning in the temporal pole but
441 not the perirhinal cortex (*Figure 4b, Figure 4 – figure supplement 1*). As one interpretation of
442 this pattern of results, the temporal pole may represent new crossmodal objects by combining
443 previously learned knowledge.^{8,9,10,11,13,14,15,33} Specifically, research into *conceptual combination*
444 has linked the anterior temporal lobes to compound object concepts such as
445 “hummingbird”.^{34,35,36} For example, participants during our task may have represented the
446 sound-based “humming” concept and visually-based “bird” concept on Day 1, forming the
447 crossmodal “hummingbird” concept on Day 3; *Figure 1, 2*, which may recruit less activity in
448 temporal pole than an incongruent pairing such as “barking-frog”. For these reasons, the
449 temporal pole may form a crossmodal object code based on pre-existing knowledge, resulting in
450 reduced neural activity (*Figure 3d*) and pattern similarity towards features associated with
451 learned objects (*Figure 4b*).

452 By contrast, perirhinal cortex may be involved in pattern separation following crossmodal
453 experience. In our task, participants had to differentiate congruent and incongruent objects
454 constructed from the same three shape and sound features (*Figure 2*). An efficient way to solve
455 this task would be to form distinct object-level outputs from the overlapping unimodal feature-
456 level inputs such that congruent objects are made to be orthogonal from the representations
457 before learning (i.e., measured as pattern similarity equal to 0 in the perirhinal cortex; *Figure 5b,*
458 *6, Figure 5 – figure supplement 2*), whereas non-learned incongruent objects could be made to be
459 dissimilar from the representations before learning (i.e., anticorrelation, measured as pattern
460 similarity less than 0 in the perirhinal cortex; *Figure 6*). Because our paradigm could decouple
461 neural responses to the learned object representations (on Day 4) from the original component
462 unimodal features at baseline (on Day 2), these results could be taken as evidence of pattern
463 separation in the human perirhinal cortex.^{11,12} However, our pattern of results could also be
464 explained by other types of crossmodal integrative coding. For example, incongruent object
465 representations may be less stable than congruent object representations, such that incongruent
466 objects representation are warped to a greater extent than congruent objects (*Figure 6*).

467 Our results suggest that the temporal pole and perirhinal cortex are involved in
468 representing crossmodal objects after a period of crossmodal learning. Although this observation
469 is consistent with previous animal research³⁷ finding that a period of experience is necessary for
470 the perirhinal cortex to represent crossmodal objects, future work will need to determine whether
471 our findings are driven by *only* experience or by experience *combined with* sleep-dependent
472 consolidation.³⁸ Perhaps a future study could explore how separate unimodal features and the
473 integrative object representations change over the course of the same learning day compared to
474 multiple learning days after sleep. Nevertheless, perirhinal cortex was critically influenced by

475 experience, potentially explaining why findings in this literature have been at times mixed, as
476 stimulus history was not always controlled across different experiments.^{39,40} In our study, we
477 explicitly controlled for stimulus history (*Figure 2*), ensuring that participants extensively
478 explored individual features by the end of the first day and formed crossmodal objects by the end
479 of the third day.

480 Complementing seminal patient work causally linking anterior temporal lobe damage to
481 the loss of object concepts,⁴¹ we show that the formation of new crossmodal concepts also
482 recruits anterior temporal lobe structures like the temporal pole and perirhinal cortex. An
483 important direction of future work will be to investigate the fine-grained functional divisions
484 within the heterogeneous anterior temporal lobe region. One recent study has found that the
485 anterior temporal lobe can be separated into 34 distinct functional regions,⁴² suggesting that a
486 simple temporal pole versus perirhinal cortex division may not fully capture the complexity of
487 this region. Imaging the anterior temporal lobe has long been known to be challenging with
488 functional neuroimaging due to signal dropout.⁴³ We show that a multi-echo fMRI sequence²⁶
489 may be especially useful in future work, as multi-echo fMRI mitigates signal dropout better than
490 the standard single-echo fMRI (see *Figure 3 – figure supplement 1* for a visual comparison).

491 Importantly, the initial visual shape bias observed in the perirhinal cortex was attenuated
492 by experience (*Figure 5, Figure 5 – figure supplement 2*), suggesting that the perirhinal
493 representations had become abstracted and were no longer predominantly grounded in a single
494 modality after crossmodal learning. One possibility may be that the perirhinal cortex is by
495 default visually driven as an extension to the ventral visual stream,^{10,11,12} but can act as a
496 polymodal “hub” region for additional crossmodal input following learning. A complementary
497 possibility may be that our visual features contained tactile information (*Figure 1c*) that the
498 perirhinal cortex may be sensitive to following the initial exploration phase on our task (*Figure*
499 *2*).⁴⁰ Critically, other brain regions like the LOC also reduced in visual bias (*Figure 3c*), which
500 may reflect visual imagery or feedback connectivity between the anterior temporal lobes.
501 However, the perirhinal cortex was the only region where the visual bias was entirely attenuated
502 following crossmodal learning (*Figure 5b*).

503 An interesting future line of investigation may be to explore whether there exist similar
504 changes to the visual bias in artificial neural networks that aim to learn crossmodal object
505 concepts.^{2,3,7} Previous human neuroimaging has shown that the anterior temporal lobes are
506 important for intra-object configural representations,^{45,46} such that damage to the perirhinal
507 cortex^{20,47} leads to object discrimination impairment. For example, human participants with
508 perirhinal cortex damage are unable to resolve feature-level interference created by viewing
509 multiple objects with overlapping features. Certain types of errors made by deep learning
510 models⁴⁸ also seem to resemble the kinds of errors made by human patients,^{20,39,41,47} whereby
511 accurate object recognition can be disrupted by feature-level interference. Writing the word
512 “iPod” on an apple image, for instance, can lead to deep learning models falsely recognizing the
513 apple as an actual iPod.⁴⁹ As certain limitations of existing neural networks may be driven by an
514 inability to resolve the binding problem,⁷ future work to mimic the coding properties of anterior

515 temporal lobe structures may allow artificial machines to better mimic the remarkable human
516 ability to learn concepts, make new inferences, and generalize across distinct entities.

517 Notably, our perirhinal cortex mask overlaps with a key region of the ventral anterior
518 temporal lobe thought to be the central locus of crossmodal integration in the “hub and spokes”
519 model of semantic representations.^{9,50} However, additional work has also linked other brain
520 regions to the convergence of unimodal representations, such as the hippocampus^{51,52,53} and
521 inferior parietal lobes.^{54,55} This past work on the hippocampus and inferior parietal lobe does not
522 necessarily address the crossmodal binding problem that was the main focus of our present
523 study, as previous findings often do not differentiate *between* crossmodal integrative coding and
524 the convergence of unimodal feature representations *per se*. Furthermore, previous studies in the
525 literature typically do not control for stimulus-based factors such as experience with unimodal
526 features, subjective similarity, or feature identity that may complicate the interpretation of results
527 when determining regions important for crossmodal integration. Indeed, we found evidence
528 consistent with the convergence of unimodal feature-based representations in both the
529 hippocampus and inferior parietal lobes (*Figure 5 – figure supplement 1*), but no evidence of
530 crossmodal integrative coding different from the unimodal features. The hippocampus and
531 inferior parietal lobes were both sensitive to visual and sound features before and after
532 crossmodal learning (see *Figure 5 – figure supplement 1*). Yet the hippocampus and inferior
533 parietal lobes did not differentiate between the congruent and incongruent conditions or change
534 with experience (see *Figure 5 – figure supplement 1*).

535 In summary, forming crossmodal object concepts relies on the representations for the
536 whole crossmodal object in anterior temporal lobe structures different from the distributed
537 unimodal feature representations in sensory regions. It is this hierarchical architecture that
538 supports our ability to understand the external world, providing one solution to the age-old
539 question of how crossmodal concepts can be constructed from their component features.

540

541

542

543

References

- 544 1. James, W. (1890). *The principles of psychology, Vol. 1*. Henry Holt and Co.
545 <https://doi.org/10.1037/10538-000>
- 546 2. Guo, W., Wang, J., & Wang, S. (2019). Deep multimodal representation learning: A
547 survey. *IEEE Access*, 7, 63373-63394. doi: 10.1109/ACCESS.2019.2916887
- 548 3. Fei, N., Lu, Z., Gao, Y. et al. (2022). Towards artificial general intelligence via a
549 multimodal foundation model. *Nat Commun.*, 13, 3094. [https://doi.org/10.1038/s41467-](https://doi.org/10.1038/s41467-022-30761-2)
550 [022-30761-2](https://doi.org/10.1038/s41467-022-30761-2)
- 551 4. Keysers, D., Schärli, N., Scales, N., Buisman, H., Furrer, D., Kashubin, S., Momchev, N.,
552 Sinopalnikov, D., Stafiniak, L., Tihon, T., Tsarkov, D., Wang, X., Zee, M.V., &

- 553 Bousquet, O. (2020). Measuring Compositional Generalization: A Comprehensive
554 Method on Realistic Data. *ArXiv*, abs/1912.09713.
- 555 5. Hupkes, D., Dankers, V., Mul, M., & Bruni, E. (2020). Compositionality Decomposed:
556 How do Neural Networks Generalise? *J. Artif. Intell. Res.*, *67*, 757-795.
- 557 6. Santoro, A., Raposo, D., Barrett, D.G., Malinowski, M., Pascanu, R., Battaglia, P.W., &
558 Lillicrap, T.P. (2017). A simple neural network module for relational reasoning. *NIPS*.
- 559 7. Greff, K., Steenkiste, S.V., & Schmidhuber, J. (2020). On the Binding Problem in
560 Artificial Neural Networks. *ArXiv*, abs/2012.05208.
- 561 8. Barsalou, L. W. (2008). Grounded cognition. *Annual Review of Psychology*, *59*, 617-645.
562 <https://doi.org/10.1146/annurev.psych.59.103006.093639>
- 563 9. Patterson, K., Nestor, P. J., & Rogers, T. T. (2007). Where do you know what you know?
564 The representation of semantic knowledge in the human brain. *Nature Reviews*
565 *Neuroscience*, *8*, 976-987. <https://doi.org/10.1038/nrn2277>
- 566 10. Saksida, L. M., & Bussey, T. J. (2010). The representational–hierarchical view of
567 amnesia: Translation from animal to human. *Neuropsychologia*, *48*(8), 2370-2384.
568 <https://doi.org/10.1016/j.neuropsychologia.2010.02.026>
- 569 11. Cowell, R. A., Barense, M. D., & Sadil, P. S. (2019). A roadmap for understanding
570 memory: Decomposing cognitive processes into operations and representations. *eNeuro*,
571 *6*(4), ENEURO.0122-19.2019. <https://doi.org/10.1523/ENEURO.0122-19.2019>
- 572 12. Kent, B. A., Hvoslef-Eide, M., Saksida, L. M., & Bussey, T. J. (2016). The
573 representational–hierarchical view of pattern separation: Not just hippocampus, not just
574 space, not just memory? *Neurobiology of Learning and Memory*, *129*, 99-106.
575 <https://doi.org/10.1016/j.nlm.2016.01.006>
- 576 13. Damasio, A. R. (1989). Time-locked multiregional retroactivation: A systems-level
577 proposal for the neural substrates of recall and recognition. *Cognition*, *33*(1-2), 25-62.
578 [https://doi.org/10.1016/0010-0277\(89\)90005-X](https://doi.org/10.1016/0010-0277(89)90005-X)
- 579 14. Suzuki, W. A., & Naya, Y. (2014). The perirhinal cortex. *Annual Review of*
580 *Neuroscience*, *37*, 39–53. <https://doi.org/10.1146/annurev-neuro-071013-014207>
- 581 15. Ralph, M., Jefferies, E., Patterson, K. et al. (2017). The neural and computational bases
582 of semantic cognition. *Nat Rev Neurosci*, *18*, 42–55. <https://doi.org/10.1038/nrn.2016.150>
- 583 16. Ferko, K. M., Blumenthal, A., Martin, C. B., Proklova, D., Minos, A. N., Saksida, L. M.,
584 Bussey, T. J., Khan, A. R., & Kohler, S. (2022). Activity in perirhinal and entorhinal
585 cortex predicts perceived visual similarities among category exemplars with highest
586 precision. *eLife*, *11*, e66884. <https://doi.org/10.7554/eLife.66884>
- 587 17. Bausch, M., Niediek, J., Reber, T.P. et al. (2021). Concept neurons in the human medial
588 temporal lobe flexibly represent abstract relations between concepts. *Nat Commun*, *12*,
589 6164, <https://doi.org/10.1038/s41467-021-26327-3>
- 590 18. Pagan, M., Urban, L. S., Wohl, M. P., & Rust, N. C. (2013). Signals in inferotemporal
591 and perirhinal cortex suggest an “untangling” of visual target information. *Nat Neurosci*,
592 *16*(8), 1132-1139. doi:10.1038/nn.3433
- 593 19. Yee, E., & Thompson-Schill, S. L. (2016). Putting concepts into context. *Psychon. Bull.*
594 *Rev.*, *23*, 1015–1027. <https://doi.org/10.3758/s13423-015-0948-7>

- 595 20. Barense, M. D., Groen, I. I. I., Lee, A. C. H., et al. (2012). Intact memory for irrelevant
596 information impairs perception in amnesia. *Neuron*, 75(1), 157-167.
597 <https://doi.org/10.1016/j.neuron.2012.05.014>
- 598 21. Erez, J., Cusack, R., Kendall, W., & Barense, M. D. (2016). Conjunctive coding of
599 complex object features. *Cerebral Cortex*, 26(5), 2271-2282.
600 <https://doi.org/10.1093/cercor/bhv081>
- 601 22. Liang, J. C., Erez, J., Zhang, F., et al. (2020). Experience transforms conjunctive object
602 representations: Neural evidence for unitization after visual expertise. *Cerebral Cortex*,
603 30(5), 2721-2739. <https://doi.org/10.1093/cercor/bhz250>
- 604 23. Martin, C. B., Douglas, D., Newsome, R. N., et al. (2018). Integrative and distinctive
605 coding of visual and conceptual object features in the ventral visual stream. *eLife*, 7:
606 e31873. doi: 10.7554/eLife.31873
- 607 24. Hebart, M. N., Zheng, C. Y., Pereira, F., & Baker, C. I. (2020). Revealing the
608 multidimensional mental representations of natural objects underlying human similarity
609 judgements. *Nature Human Behaviour*, 4, 1173-1185. [https://doi.org/10.1038/s41562-](https://doi.org/10.1038/s41562-020-00951-3)
610 [020-00951-3](https://doi.org/10.1038/s41562-020-00951-3)
- 611 25. Li, A. Y., Fukuda, K., & Barense, M. D. (2022). Independent features form integrated
612 objects: Using a novel shape-color “conjunction task” to reconstruct memory resolution
613 for multiple object features simultaneously. *Cognition*, 223, 105024.
614 <https://doi.org/10.1016/j.cognition.2022.105024>
- 615 26. Kundu, P., Voon, V., Balchandani, P., et al. (2017). Multi-echo fMRI: A review of
616 applications in fMRI denoising and analysis of BOLD signals. *NeuroImage*, 154, 59-80.
617 <https://doi.org/10.1016/j.neuroimage.2017.03.033>
- 618 27. Li, A. Y., Liang, J. C., Lee, A. C. H., & Barense, M. D. (2020). The validated circular
619 shape space: Quantifying the visual similarity of shape. *J Exp Psychol Gen.*, 149(5): 949-
620 966. doi: 10.1037/xge0000693.
- 621 28. Shepard, R. N. (1980). Multidimensional Scaling, Tree-Fitting, and Clustering. *Science*,
622 210, 4468. doi: 10.1126/science.210.4468.390
- 623 29. Barense, M. D., Warren, J. D., Bussey, T. J., & Saksida, L. M. “The temporal lobes”.
624 *Oxford Textbook of Cognitive Neurology and Dementia*, edited by Masud Husain and
625 Jonathan M. Schott, Oxford University Press, 2016. doi:
626 10.1093/med/9780199655946.003.0004
- 627 30. Mishkin, M., Ungerleider, L. G., & Macko, K. A. (1983). Object vision and spatial
628 vision: Two cortical pathways. *Trends in Neurosciences*, 6(1983), 414-417.
629 [https://doi.org/10.1016/0166-2236\(83\)90190-X](https://doi.org/10.1016/0166-2236(83)90190-X)
- 630 31. Poremba, A., & Mishkin, M. (2007). Exploring the extent and function of higher-order
631 auditory cortex in rhesus monkeys. *Hearing Research*, 229(1-2), 14–23.
632 <https://doi.org/10.1016/j.heares.2007.01.003>
- 633 32. Vetter, P., Smith, F. W., & Muckli, L. (2014). Decoding sound and imagery content in
634 early visual cortex. *Current Biology*, 24(11), 1256-1262. doi: 10.1016/j.cub.2014.04.020
- 635 33. Binder, J. R., Desai, R. H. (2011). The neurobiology of semantic memory. *Trends Cogn.*
636 *Sci.*, 15(11), 527-36. doi: 10.1016/j.tics.2011.10.001.

- 637 34. Lynott, D., & Connell, L. (2010). Embodied conceptual combination. *Frontiers in*
638 *Psychology*. <https://doi.org/10.3389/fpsyg.2010.00212>
- 639 35. Coutanche, M. N., Solomon, S. H., & Thompson-Schill, S. L. (2020). Conceptual
640 combination. In D. Poeppel, G. R. Mangun and M. S. Gazzaniga (Eds.), *The Cognitive*
641 *Neurosciences, 6th edition*. Boston, MA: MIT Press.
- 642 36. Baron, S. G., & Osherson, D. (2011). Evidence for conceptual combination in the left
643 anterior temporal lobe. *NeuroImage*, 55(4), 1847-1852.
644 <https://doi.org/10.1016/j.neuroimage.2011.01.066>
- 645 37. Jacklin, D. L., Cloke, J. M., Potvin, A., et al. (2016). The dynamic multisensory engram:
646 Neural circuitry underlying crossmodal object recognition in rats changes with the nature
647 of object experience. *Journal of Neuroscience*, 36(4), 1273-1289.
648 <https://doi.org/10.1523/JNEUROSCI.3043-15.2016>
- 649 38. Schapiro, A.C., McDevitt, E.A., Chen, L. et al. (2017). Sleep Benefits Memory for
650 Semantic Category Structure While Preserving Exemplar-Specific Information. *Sci Rep*,
651 7, 14869. <https://doi.org/10.1038/s41598-017-12884-5>
- 652 39. Taylor, K. I., Moss, H. E., Stamatakis, E. A., & Tyler, L. K. (2006). Binding crossmodal
653 object features in perirhinal cortex. *PNAS*, 103(21).
654 <https://doi.org/10.1073/pnas.0509704103>
- 655 40. Holdstock, J. S., Hocking, J., Notley, et al. (2009). Integrating visual and tactile
656 information in the perirhinal cortex. *Cerebral Cortex*, 19(12), 2993–3000.
657 <https://doi.org/10.1093/cercor/bhp073>
- 658 41. Hodges, J. R., & Patterson, K. (1997). Semantic memory disorders. *Trends in Cognitive*
659 *Sciences*, 1(2), 68-72. [https://doi.org/10.1016/S1364-6613\(97\)01022-X](https://doi.org/10.1016/S1364-6613(97)01022-X)
- 660 42. Persichetti, A. S., Denning, J. M., Gotts, S. J., & Martin, A. (2021). A data-driven
661 functional mapping of the anterior temporal lobes. *Journal of Neuroscience*, 41(28),
662 6038-6049. <https://doi.org/10.1523/JNEUROSCI.0456-21.2021>
- 663 43. Visser, M., Jefferies, E., & Lambon Ralph, M. A. (2010). Semantic processing in the
664 anterior temporal lobes: A meta-analysis of the functional neuroimaging literature. *J*
665 *Cogn. Neurosci*, 22(6), 1083-94. doi: 10.1162/jocn.2009.21309.
- 666 44. Malach, R., Reppas, J. B., Benson, R. R., et al. (1995). Object-related activity revealed by
667 functional magnetic resonance imaging in human occipital cortex. *PNAS*, 92(18).
668 <https://doi.org/10.1073/pnas.92.18.8135>
- 669 45. Yeung, L. K., Olsen, R. K., Bild-Enkin, H., D'Angelo, M. C., Kacollja, A., McQuiggan,
670 D. A., Keshabyan, A., Ryan, J. D., & Barense, M. D. (2017). Anterolateral Entorhinal
671 Cortex Volume Predicted by Altered Intra-Item Configural Processing. *Journal of*
672 *Neuroscience*, 37(22), 5527–5538. <https://doi.org/10.1523/JNEUROSCI.3664-16.2017>
- 673 46. Watson, H. C., Lee, A. C. (2013). The perirhinal cortex and recognition memory
674 interference. *J Neurosci*, 33(9), 4192-4200. doi:10.1523/JNEUROSCI.2075-12.2013
- 675 47. Bonnen, T., Yamins, D. L. K., & Wagner, A. D. (2021). When the ventral visual stream
676 is not enough: A deep learning account of medial temporal lobe involvement in
677 perception. *Neuron*, 109(17), 2755-2766. <https://doi.org/10.1016/j.neuron.2021.06.018>

- 678 48. Guo, C., Lee, M. J., Leclerc, G., Dapello, J., Rao, Y., Madry, A., & DiCarlo, J. J. (2022).
679 Adversarially trained neural representations may already be as robust as corresponding
680 biological neural representations. *arXiv*, 2206.11228.
- 681 49. Goh, G., Cammarata, N., Voss, C., Carter, S., Petrov, M., Schubert, L., Radford, A., &
682 Olah, C. (2021). Multimodal neurons in artificial neural networks. *OpenAI*. doi:
683 10.23915/distill.00030
- 684 50. Ralph, M., Jefferies, E., Patterson, K. *et al.* (2017). The neural and computational bases
685 of semantic cognition. *Nat Rev Neurosci*, 18, 42-55.
686 <https://doi.org/10.1038/nrn.2016.150>
- 687 51. Butler, A. J., & James, K. H. (2011). Cross-modal versus within-modal recall:
688 Differences in behavioral and brain responses. *Behavioural Brain Research*, 224(2), 387-
689 396. <https://doi.org/10.1016/j.bbr.2011.06.017>
- 690 52. Clouter, A., Shapiro, K. L., & Hanslmayr, S. (2017). Theta phase synchronization is the
691 glue that binds human associative memory. *Current Biology*, 27(20), 3143-3148.
692 <https://doi.org/10.1016/j.cub.2017.09.001>
- 693 53. Vigano, S., & Piazza, M. (2020). Distance and direction codes underlie navigation of a
694 novel semantic space in the human brain. *Journal of Neuroscience*, 40(13), 2727-2736.
695 <https://doi.org/10.1523/JNEUROSCI.1849-19.2020>
- 696 54. Vigano, S., Rubino, V., Buiatti, M. *et al.* (2021). The neural representation of absolute
697 direction during mental navigation in conceptual spaces. *Commun Biol*, 4, 1294.
698 <https://doi.org/10.1038/s42003-021-02806-7>
- 699 55. Binder, J. R., & Desai, R. H. (2011). The neurobiology of semantic memory. *Trends in*
700 *Cognitive Sciences*, 15(11), 527-536. <https://doi.org/10.1016/j.tics.2011.10.001>
- 701 56. Stirnberg, R., & Stöcker, T. (2021). Segmented K-space blipped-controlled aliasing in
702 parallel imaging for high spatiotemporal resolution EPI. *Magnetic Resonance in*
703 *Medicine*, 85(3), 1540–1551. <https://doi.org/10.1002/mrm.28486>
- 704 57. Coutanche, M. N., & Thompson-Schill, S. L. (2012). The advantage of brief fMRI
705 acquisition runs for multi-voxel pattern detection across runs. *Neuroimage*, 61(4), 1113-
706 9. doi: 10.1016/j.neuroimage.2012.03.076
- 707 58. Ritchey, M., Montchal, M. E., Yonelinas, A. P., & Ranganath, C. (2015). Delay-
708 dependent contributions of medial temporal lobe regions to episodic memory retrieval.
709 *eLife*, 4, e05025. doi: 10.7554/eLife.05025
- 710 59. tedana Community, et al. (2021). ME-ICA/tedana:0.0.11. Zenodo. Available from
711 <https://doi.org/10.5281/zenodo.5541689>
- 712 60. Mumford, J. A. (2014). The impact of study design on pattern estimation for single-trial
713 multivariate pattern analysis. *NeuroImage*, 103, 130-138.
714 <https://doi.org/10.1016/j.neuroimage.2014.09.026>
- 715 61. Corey, D. M., Dunlap, W. P., & Burke, M. J. (1998). Averaging correlations: Expected
716 values and bias in combined Pearson *rs* and Fisher's *z* transformations. *Journal of*
717 *General Psychology*, 125(3), 245–262. <https://doi.org/10.1080/00221309809595548>
718
719

720 **Methods**

721 The experiments described in this study were approved by the University of Toronto
722 Ethics Review Board (protocols 37590 and 38856). Informed consent was obtained for all
723 participants in the study prior to their participation.

724 **Initial Stimulus Validation Experiment**

725 *Participants*

726 16 participants (*Females* = 11, *Age* = 18.63 years) were recruited from the University of
727 Toronto undergraduate participant pool and from the community. Course credit or \$10/hr CAD
728 was provided as compensation.

729 *Stimuli*

730 Three shape stimuli were sampled from the Validated Shape Space²⁷ at equidistant
731 positions, ensuring that the shapes were equated in their subjective similarity. The sound stimuli
732 were manually generated in a similar procedure to how the shape stimuli from the Validated
733 Shape Space²⁷ were originally created. More specifically, distinct sounds were morphed together
734 to create 5 complex, unrecognizable sounds that lasted for a duration of 2 seconds.

735 *Validation Procedure*

736 The stimulus validation procedure was based on previous work²⁷ (see *Figure 2–figure*
737 *supplement 1* for an example of the task). Across 9 trials, participants rated the similarity of each
738 of the 3 shapes in the context of every other shape, as well as 4 control trials in which each shape
739 was rated relative to itself. For this initial stimulus validation experiment we used line drawings
740 of the three shapes (for the 4-day crossmodal learning task we used images of the printed
741 objects). Afterwards, participants completed 40 trials in which they rated the similarity of each of
742 the 5 sounds in the context of every other sound, as well as 4 trials in which every sound was
743 rated relative to itself. In a self-timed manner, participants viewed pictures of shapes or clicked
744 icons to play the to-be-rated sounds from a headset.

745 For the shapes, we replicated the triangular geometry from participant similarity ratings
746 obtained in our past work²⁷ indicating that each shape was about as similar as every other shape
747 (*Figure 1a*). We then selected the three sounds that were best equated in terms of their perceived
748 similarity (*Figure 1a*). Thus, like the shapes, this procedure ensured that subjective similarity for
749 the sounds was explicitly controlled but the underlying auditory dimensions could vary (e.g.,
750 timbre, pitch, frequency). This initial validation experiment ensured that the subjective similarity
751 of the three features of each stimulus modality was equated within each modality prior to the
752 primary 4-day learning task.

753 *3D-Printed Shape-Sound Objects*

754 The three validated shapes were 3D-printed using a DREMEL Digilab 3D Printer 3D45-
755 01 with 1.75 mm gold-colored polymerized lactic acid filament. To create the 3D object models,
756 the original 2D images were imported into Blender and elongated to add depth. The face of the

757 shape image created a detachable lid, with a small circular opening to allow wiring to extend to a
758 playable button positioned on the exterior of the shape. An empty space was formed inside the
759 3D shape for the battery-powered embedded speaker. To ensure that the objects were graspable,
760 each shape was 3D-printed to be approximately the size of an adult hand (*Figure 1c*). The lid of
761 the shape was detached before each learning day (*Figure 2*), with the embedded speaker
762 programmed to play either no sound (Day 1) or to play the paired sound that formed the
763 congruent object (Day 3) (*Figure 1a*). After the speaker was programmed, the lid of the shape
764 was reattached using thermoplastic adhesive.

765 The sounds were played at an audible volume by the 3D-printed shapes during the
766 learning task (see next section). During the scanning sessions, we individually tailored the
767 volume until the participant could hear the sounds clearly when inside the MRI scanner.

768

769 **4-Day Crossmodal Object Learning Task**

770 *Participants*

771 Twenty new participants (*Females* = 13, *Age* = 23.15 years) were recruited and scanned
772 at the Toronto Neuroimaging Facility. All participants were right-handed, with normal or
773 corrected-to-normal vision, normal hearing, and no history of psychiatric illness. Of the 20
774 scanned participants, 1 participant dropped out after the first neuroimaging session. Severe
775 distortion was observed in a second participant from a metal retainer and data from this
776 participant was excluded from subsequent analyses. Due to technical difficulties, the functional
777 localizer scans were not saved for one participant and most feature runs could not be completed
778 for a second participant. Overall, the within-subject analyses described in the main text included
779 data from a minimum of 16 participants, with most analyses containing data from 17
780 participants. Critically, this within-subject learning design increases power to detect an effect.

781 Compensation was \$250 CAD for the 2 neuroimaging sessions and 2 behavioral sessions
782 (~approx. 6 hours total, which included set-up, consent, and debriefing), with a \$50 CAD
783 completion bonus.

784 *Behavioral Tasks*

785 On each behavioral day (Day 1 and Day 3; *Figure 2*), participants completed the
786 following tasks, in this order: Exploration Phase, one Unimodal Feature 1-back run (26 trials),
787 Exploration Phase, one Crossmodal 1-back run (26 trials), Exploration Phase, Pairwise Similarity
788 Task (24 trials), Exploration Phase, Pairwise Similarity Task (24 trials), Exploration Phase,
789 Pairwise Similarity Task (24 trials), and finally, Exploration Phase. To verify learning on Day 3,
790 participants also additionally completed a Learning Verification Task at the end of the session.
791 Details on each task are provided below.

792 The overall procedure ensured that participants extensively explored the unimodal
793 features on Day 1 and the crossmodal objects on Day 3. The Unimodal Feature and the

794 Crossmodal Object 1-back runs administered on Day 1 and Day 3 served as practice for the
795 neuroimaging sessions on Day 2 and Day 4, during which these 1-back tasks were completed.
796 Each behavioral session required less than 1 hour of total time to complete.

797 **Day 1 Exploration Phase.** On Day 1 (*Figure 2a*), participants separately learned the
798 shape and sound features in a random order. The 3D shapes were explored and physically
799 palpated by the participants. We also encouraged participants to press the button on each shape,
800 although the button was not operational on this day. Each 3D-printed shape was physically
801 explored for 1 minute and each sound was heard through a headset 7 times. There were 6
802 exploration phases in total, interleaved between the 1-back and pairwise similarity tasks (order
803 provided above). This procedure ensured that each individual stimulus was experienced
804 extensively by the end of the first day.

805 **Day 3 Exploration Phase.** On Day 3 (*Figure 2c*), participants experienced the 3D-printed
806 shape-sound objects in a random order. The sound was played over the embedded speakers by
807 pressing the now-operational button on each object. Participants were allotted 1 minute to
808 physically explore and palpate each shape-sound object, as well as to listen to the associated
809 sound by pressing the button. Like Day 1, there were 6 exploration phases in total, interleaved
810 between the 1-back and pairwise similarity tasks.

811 **Pairwise Similarity Task.** Using the same task as the stimulus validation procedure
812 (*Figure 2–figure supplement 1*), participants provided similarity ratings for all combinations of
813 the 3 validated shapes and 3 validated sounds (each of the six features were rated in the context
814 of every other feature in the set, with 4 repeats of the same feature, for a total of 72 trials). More
815 specifically, three stimuli were displayed on each trial, with one at the top and two at the bottom
816 of the screen in the same procedure as we have used previously²⁷. The 3D shapes were visually
817 displayed as a photo, whereas sounds were displayed on screen in a box that could be played
818 over headphones when clicked with the mouse. The participant made an initial judgment by
819 selecting the more similar stimulus on the bottom relative to the stimulus on the top. Afterwards,
820 the participant made a similarity rating between each bottom stimulus with the top stimulus from
821 0 being no similarity to 5 being identical. This procedure ensured that ratings were made relative
822 to all other stimuli in the set.

823 **Unimodal Feature and Crossmodal Object 1-back Tasks.** During fMRI scanning on
824 Days 2 and 4, participants completed 1-back tasks in which the target was an exact sequential
825 repeat of a feature (Unimodal Feature Task) or an exact sequential repeat of the shape-sound
826 object (Crossmodal Object Task). In total, there were 10 Unimodal Feature runs and 5
827 Crossmodal Object runs for each scanning session. Two Unimodal Feature runs were followed
828 by one Crossmodal Object run in an interleaved manner to participants until all 10 Unimodal
829 Feature runs and 5 Crossmodal Object runs were completed. Each run lasted 3 minutes and had
830 26 trials.

831 Each Unimodal Feature and Crossmodal Object run began with a blank screen appearing
832 for 6 seconds. For Unimodal Feature runs, either a shape or sound feature would then be
833 presented for two seconds, followed by a fixation cross appearing for 2 – 8 seconds (sampled

834 from the following probability distribution: 2 seconds = 30%, 4 seconds = 30%, 6 seconds =
835 30%, and 8 seconds = 10%). For Crossmodal Object runs, each shape appeared on the monitor at
836 the same time as a sound was played through the headset for two seconds, followed by a fixation
837 cross appearing for 2 – 8 seconds (sampled from the following probability distribution: 2 seconds
838 = 30%, 4 seconds = 30%, 6 seconds = 30%, and 8 seconds = 10%). Ensuring equal trial numbers,
839 three shape-sound pairings were congruent (learned by participants) and three shape-sound
840 pairings were incongruent (not learned by participants). Congruent and incongruent pairings
841 were built from different combinations of the same shape and sound features, with pairings
842 counterbalanced across participants.

843 Overall, each stimulus was presented four times in a random order per run, with two
844 repeats occurring at a random position for the corresponding 1-back task. The stimulus identity
845 and temporal position of any given 1-back repeat was random.

846 ***Learning Verification Task (Day 3 only).*** As the final task on Day 3, participants
847 completed a task to ensure that participants successfully formed their crossmodal pairing. All
848 three shapes and sounds were randomly displayed in 6 boxes on a display. Photos of the 3D
849 shapes were shown, and sounds were played by clicking the box with the mouse cursor. The
850 participant was cued with either a shape or sound, and then selected the corresponding paired
851 feature. At the end of Day 3, we found that all participants reached 100% accuracy on this task
852 (10 trials).

853 *Behavioral Pattern Similarity Analysis*

854 The pairwise similarity ratings for each stimulus were averaged into a single feature-level
855 RDM. We examined the magnitude of pattern similarity for congruent features compared to
856 incongruent features across learning days (see *Figure 2-figure supplement 1*).

857

858 **Neuroimaging Procedures**

859 Scanning was conducted using a 32-channel receiver head coil with the Siemens
860 Magnetom Prisma 3T MRI scanner at the Toronto Neuroimaging Facility. To record responses,
861 participants used a 4-button keypad (Current Designs, HHSC-1X4-CR). Stimulus materials were
862 displayed using an MR compatible screen at high resolution (1920 x 1080) with zero-delay
863 timing (32" BOLD screen) controlled by PsychToolbox-3 in MATLAB. At the start of each
864 neuroimaging session, we performed a sound check with a set of modified in-ear MR-compatible
865 headphones (Sensimetrics, model S14), followed by a functional localizer and then by the task-
866 related runs.

867 While in the scanner, participants completed the following: After an initial functional
868 localizer, we collected a resting state scan. After five 1-back runs, we acquired a whole-brain
869 high-resolution T1-weighted structural image. After an additional five 1-back runs, we acquired
870 a second resting-state scan, followed by the last five 1-back runs. The 15 total 1-back runs were

871 interleaved such that 2 Unimodal Feature runs would be presented, followed by 1 Crossmodal
872 Feature run until all 15 runs had been completed (see *Figure 2*).

873 *Multi-echo fMRI*

874 A 3D multi-echo echo-planer imaging (EPI) sequence with blipped-controlled aliasing in
875 parallel imaging (CAIPI) sampling⁵⁶ was used to acquire fMRI data on Day 2 and Day 4. For
876 task-related scans, the 3 echoes (TR = 2000 ms, TE 1 = 11 ms, TE 2 = 31.6 ms, and TE 3 = 52.2
877 ms) were each acquired with 90 images (210 x 210 field of view with a 100 x 100 matrix resize;
878 anterior to posterior phase encoding, 78 slices, slice thickness: 2.10 mm, flip angle: 17°,
879 interleaved multi-slice acquisition), resulting in an in-plane resolution of 2.10 x 2.10 mm. 3D
880 distortion correction and pre-scan normalization was enabled, with acceleration factor PE = 2
881 and acceleration factor 3D = 3. These parameters yielded coverage over the entire cortex, and a
882 B0 field map was collected at the completion of the experiment.

883 ***1-back Tasks (Unimodal Feature Runs and Crossmodal Object Runs)***. Rather than
884 collecting data from many different instances of a category as is common in a fMRI study using
885 multivariate pattern analysis, we collected data from many repetitions of the *same* stimulus using
886 a psychophysics-inspired approach. This paradigm ensured that the neural representations
887 specific to each unimodal feature and each crossmodal object was well-powered for subsequent
888 pattern similarity analyses.⁵⁷ Excluding 1-back repeats, each unimodal feature was displayed 4
889 times per run for a total of 40 instances per scanning session (80 instances of each unimodal
890 feature in total). Excluding 1-back repeats, each shape-sound pairing was displayed 4 times per
891 run for a total of 20 instances per scanning session (40 instances of each shape-sound object in
892 total). We designed our task-related runs to be 3 minutes in length, as “mini-runs” have been
893 shown to improve data quality in multivariate pattern analysis.⁵⁷ Details of the task can be found
894 in the section above.

895 ***Standard Functional Localizer***. Participants viewed intact visual features and phase
896 scrambled versions of the same features in separate 24 second blocks (8 functional volumes).⁴⁴
897 Each of the 32 images within a block were presented for 400 ms each with a 350 ms ISI. There
898 were 2 groups of 4 blocks, with each group separated by a 12 s fixation cross. Block order was
899 counterbalanced across participants. All stimuli were presented in the context of an 1-back task,
900 and the order of images within blocks was randomized with the 1-back repeat occurring once per
901 block. The identity and temporal position of the 1-back repeat was random.

902 *Structural and Resting State Scans*

903 A standard whole-brain high-resolution T1-weighted structural image was collected (TR
904 = 2000 ms, TE = 2.40 ms, flip angle = 9°, field of view = 256 mm, 160 slices, slice thickness =
905 1.00 mm, acceleration factor PE = 2), resulting in an in-place resolution of 1.00 mm x 1.00.

906 Two 6 minute 42 second resting state scans were also collected (TR = 2000 ms, TE = 30
907 ms; field of view: 220 mm, slice thickness: 2.00 mm; interleaved multi-slice acquisition, with
908 acceleration factor PE = 2).

909

910 **Neuroimaging Analysis**911 *ROI Definitions*

912 We conducted region-of-interest univariate (*Figure 3c, d*) and multivariate pattern
913 analysis (*Figure 4, 5, 6*) in five *a priori* masks: temporal pole, perirhinal cortex, lateral occipital
914 complex (LOC), primary visual cortex (V1), and primary auditory cortex (A1). These regions
915 were selected *a priori* given their hypothesized role in representing individual unimodal features
916 as well as their integrated whole.^{9,11} More specifically, we expected that the anterior temporal
917 lobe structures – temporal pole and perirhinal cortex – would differentiate between the congruent
918 and incongruent conditions. By contrast, we expected LOC, V1, and A1 to possess modality-
919 specific biases for either the visual or sound features. Temporal pole, V1, and A1 masks were
920 extracted from the Harvard-Oxford atlas. The perirhinal cortex mask was created from the
921 average of 55 manually-segmented T1 images from a previous publication.⁵⁸ The LOC mask was
922 extracted from the top 500 voxels in the lateral occipital region of each hemisphere that activated
923 more strongly to intact than phase scrambled objects in the functional localizer (uncorrected
924 voxel-wise $p < 0.001$).⁴⁴

925 Additionally, we conducted region-of-interest univariate and multivariate pattern analysis
926 in two *exploratory* masks: hippocampus and inferior parietal lobes (*Figure 5 – figure supplement*
927 *1*). These regions were selected given their hypothesized role in the convergence of unimodal
928 feature representations.⁵¹⁻⁵⁵

929 Probabilistic masks were thresholded at .5 (i.e., voxels labelled in 50% of participants),
930 with the masks transformed to subject space through the inverse warp matrix generated from
931 FNIRT nonlinear registration (see *Preprocessing*) then resampled from 1mm^3 to 2.1mm^3 . All
932 subsequent analyses were conducted in subject space.

933 *Multi-echo ICA-based Denoising*

934 For a detailed description of the overall ME-ICA pipeline, see the *tedana* Community.⁵⁹
935 The multi-echo ICA-based denoising approach was implemented using the function *meica.py* in
936 AFNI. We optimally averaged the three echoes, which weights the combination of echoes based
937 on the estimated T_2^* at each voxel for each echo. PCA then reduced the dimensionality of the
938 optimally-combined dataset and ICA decomposition was applied to remove non-BOLD noise.
939 TE-dependent components reflecting BOLD-like signal for each run were used as the dataset for
940 subsequent preprocessing in FSL (e.g., see *Figure 3 – figure supplement 1*).

941 *Preprocessing*

942 First, the anatomical image was skull-stripped. Data were high-pass temporally filtered
943 (50 s) and spatially smoothed (6 mm). Functional runs were registered to each participant's high-
944 resolution MPRAGE image using FLIRT boundary-based registration, with registration further

945 refined using FNIRT nonlinear registration. The resulting data were analyzed using first-level
946 FEAT Version 6.00 in each participant's native anatomical space.

947 *Univariate Analysis*

948 To obtain participant-level contrasts, we averaged the run-level Unimodal Feature (*Visual*
949 vs. *Sound*) and Crossmodal Object (*Congruent* vs. *Incongruent*) runs to produce the whole-brain
950 group-level contrasts in FSL FLAME. Whole-brain analyses were thresholded at voxel-level $p =$
951 0.001 with random field theory cluster correction at $p = 0.05$.

952 For ROI-based analyses (*Figure 3*), we estimated percent signal change using *featquery*.
953 The parameter estimates (beta weight) were scaled by the peak height of the regressor, divided
954 by the baseline intensity in the *Visual* vs. *Sound* and *Congruent* vs. *Incongruent* contrasts to
955 obtain a difference score. Inferential statistical analyses were performed with these difference
956 scores using a linear mixed model which included learning day (before vs. after crossmodal
957 learning) and hemisphere (left or right) as fixed effects for each ROI, with participants modelled
958 as random effects. All linear mixed model analyses were conducted using the *nlme* package in R
959 version 3.6.1.

960 *Single Trial Estimates*

961 We used the least squares single approach⁶⁰ with 2 mm smoothing on the raw data in a
962 separate set of analyses distinct from the univariate contrasts. Each individual stimulus, all other
963 repetitions of the stimulus, and all other individual stimuli were modelled as covariates, allowing
964 us to estimate whole-brain single-trial betas for each trial by run by mask by hemisphere by
965 subject. All pattern similarity analyses described in the main text were conducted using the
966 *CoSMoMVPA* package in MATLAB. After the single-trial betas were estimated, the voxel-wise
967 activity across runs were averaged into a single overall matrix.

968 *Neuroimaging Pattern Similarity Analysis*

969 Four comparisons were conducted for each a priori ROI: 1) the autocorrelation of the
970 average voxel-wise matrix during Unimodal Feature runs (*Figure 4a*, *Figure 5 – figure*
971 *supplement 1*), 2) the correlation between the RDM created from the Unimodal Feature runs
972 before crossmodal learning to the RDM created from the Crossmodal Object runs before
973 crossmodal learning (*Figure 5a*), 3) the correlation between the RDM created from the Unimodal
974 Feature runs before crossmodal learning to the RDM created from the Crossmodal Object runs
975 after crossmodal learning (*Figure 5b*), and 4) the correlation between the RDM created from the
976 Crossmodal Object runs before crossmodal learning to the RDM created from the Crossmodal
977 Object runs after crossmodal learning (*Figure 6*).

978 The z-transformed Pearson's correlation coefficient was used as the distance metric for
979 all pattern similarity analyses. More specifically, each individual Pearson correlation was Fisher
980 z-transformed and then averaged (see⁶¹). Inferential statistical analyses were performed for each
981 individual ROI using linear mixed models which could include congruency (congruent or

982 incongruent), learning day (before or after crossmodal learning), modality (visual or sound), and
983 hemisphere (left or right) as fixed factors, with participant modelled as random effects allowing
984 intercepts to vary by learning day when appropriate. One-sample t-tests also compared the z-
985 transformed pattern similarity scores relative to 0. All linear mixed model analyses were
986 conducted using the *nlme* package in R version 3.6.1.

987

988 **Crossmodal Object Learning Task: Behavioral Replication**

989 *Participants*

990 44 new participants (*Females* = 34, *M_{age}* = 23.95 years) were recruited from the
991 University of Toronto undergraduate participant pool and from the community. Course credit or
992 \$10/hr CAD was provided as compensation.

993 *Procedure*

994 We conducted a same-day behavioural-only variant of the four-day task described in the
995 main text (*Figure 2*), excluding neuroimaging sessions. Participants first explored the 3D-printed
996 shapes and heard the sounds separately (the button-activated speaker was not operational on this
997 day). Each 3D-printed shape was physically explored for 1 minute and each sound was heard
998 through a headset 7 times. On a separate pairwise similarity rating task, participants then
999 provided similarity ratings for all combinations of the 3 shapes and 3 sounds (rated in the context
1000 of each other stimulus in the set, with 4 repeats of the same item; 72 total trials). Every 24 trials,
1001 participants again explored the same shapes and sounds (separately before crossmodal learning,
1002 in a counterbalanced order across participants).

1003 Next, participants learned that certain shapes are associated with certain sounds, such that
1004 the 3D-printed shapes now played a sound when the button was pressed. Participants were
1005 allotted 1 minute to physically explore and palpate each shape-sound object, as well as to listen
1006 to the associated sound by pressing the button. Participants repeated the pairwise similarity rating
1007 task, and every 24 trials, participants explored the 3D-printed shape-sound objects.

1008 The behavioral similarity judgments before and after crossmodal learning were analyzed
1009 in the same pattern similarity approach described in the main text (*Figure 2-figure supplement*
1010 *I*).

1011

1012

1013

1014

1015

1016

1017 **Figure Captions**

1018 **Figure 1. 3D-printed objects.** An independent validation experiment ensured that the similarity of the selected
1019 shapes and sounds were well-matched. (a) Three shapes were sampled from the *Validated Circular Shape (VCS)*
1020 *Space* (shown as black points on VCS space),²⁷ a stimulus space whereby angular distance corresponds to subjective
1021 shape similarity. Three sounds were sampled from a set of five experimenter-created sounds. This independent
1022 validation experiment ensured that we could characterize the change in similarity structure following crossmodal
1023 learning, because we knew the baseline similarity structure (i.e., two triangular representational geometries
1024 visualized using multidimensional scaling²⁸; also see *Figure 2–figure supplement 1*). Furthermore, this procedure
1025 ensured that the subjective similarity of the three features was equated within each modality. (b) The shapes were
1026 then 3D-printed with a hollow space and embedded with a button-activated speaker. (c) Participants could
1027 physically explore and palpate the 3D shape-sound objects. Critically, we manipulated whether the button-activated
1028 speaker was operational across learning days (see Methods/Figure 2).

1029
1030 **Figure 2. Four-day crossmodal object learning task.** On **Day 1** (behavior), participants heard sounds through a
1031 headset and explored 3D-printed shapes while the button-activated speakers were not operational. During a separate
1032 task (*Figure 2–figure supplement 1*), participants rated the similarity of the visual shapes and sound features. On
1033 **Day 2** (neuroimaging), participants completed (i) 10 Unimodal Feature runs in which they performed a 1-back task
1034 involving the shape and sound features experienced separately and (ii) 5 Crossmodal Object runs in which they
1035 performed a 1-back task for the shapes and sounds experienced simultaneously. As participants at this point have not
1036 yet learned the congruent shape-sound pairings, the Day 2 neuroimaging session serves as a within-subject neural
1037 baseline for how the unimodal features were represented before crossmodal learning. On **Day 3** (behavior),
1038 participants again explored the shape and sound features. Participants now learned to make crossmodal associations
1039 between the specific visual and sound features that composed the shape-sound object by pressing the button to play
1040 an embedded speaker, thus forming congruent object representations (i.e., crossmodal learning). Shape-sound
1041 associations were counterbalanced across participants, and we again collected similarity ratings between the shapes
1042 and sounds on a separate task. On **Day 4** (neuroimaging), participants completed the same task as on Day 2. In
1043 summary, across four days, we characterized the neural and behavioral changes that occurred before and after
1044 shapes and sounds were paired together to form crossmodal object representations. As the baseline similarity
1045 structure of the shape and sound features were *a priori* defined (see *Figure 1*) and measured on the first day of
1046 learning (see *Figure 2–figure supplement 1*), changes to the within-subject similarity structure provide insight into
1047 whether the crossmodal object representations (acquired after crossmodal learning) differed from component
1048 unimodal representations (acquired before crossmodal learning).

1049
1050 **Figure 2 – figure supplement 1.** Pairwise similarity task and results. In the initial stimulus validation experiment,
1051 participants provided pairwise ratings for 5 sounds and 3 shapes. The shapes were equated in their subjective similarity that
1052 had been selected from a well-characterized perceptually uniform stimulus space²⁷ and the pairwise ratings followed the
1053 same procedure as described in ref²⁷. Based on this initial experiment, we then selected the 3 sounds from the that were
1054 most closely equated in their subjective similarity. (a) 3D-printed shapes were displayed as images, whereas sounds were
1055 displayed in a box that could be played when clicked by the participant. Ratings were averaged to produce a similarity
1056 matrix for each participant, and then averaged to produce a group-level similarity matrix. Shown as triangular
1057 representational geometries recovered from multidimensional scaling in the above, shapes (blue) and sounds (orange) were
1058 approximately equated in their subjective similarity. These features were then used in the four-day crossmodal learning
1059 task. (b) Behavioral results from the four-day crossmodal learning task paired with multi-echo fMRI described in the main
1060 text. Before crossmodal learning, there was no difference in similarity between shape and sound features associated with
1061 congruent objects compared to incongruent objects – indicating that similarity was controlled at the unimodal feature-level.
1062 After crossmodal learning, we observed a robust shift in the magnitude of similarity. The shape and sound features
1063 associated with congruent objects were now significantly more similar than the same shape and sound features associated
1064 with incongruent objects ($p < 0.001$), evidence that crossmodal learning changed how participants experienced the

1065 unimodal features (observed in 17/18 participants). (c) We replicated this learning-related shift in pattern similarity with a
 1066 larger sample size ($n = 44$; observed in 38/44 participants). *** denotes $p < 0.001$. Horizontal lines denote the comparison
 1067 of congruent vs. incongruent conditions.

1068

1069 **Figure 3.** (a-b) Univariate analyses superimposed on MNI-152 standard space. All contrasts were thresholded at voxel-
 1070 wise $p = 0.001$ and cluster-corrected at $p = 0.05$ (random-effects, FSL FLAME; 6-mm spatial smoothing). Collapsing
 1071 across learning days, robust modality-specific activity was observed across the neocortex. (c-d) Five ROIs were *a priori*
 1072 selected based on existing theory:^{9,11} temporal pole – TP, perirhinal cortex – PRC, lateral occipital complex – LOC,
 1073 primary visual cortex – V1, and primary auditory cortex – A1. (c) Consistent with the whole-brain results, LOC was biased
 1074 towards visual features whereas A1 and TP were biased towards sound features. Activation in PRC and LOC showed
 1075 learning-related shifts, with the magnitude of visual bias decreasing after crossmodal learning. (d) TP was the only brain
 1076 region to show an experience-dependent change in univariate activity to the learned shape-sound associations during
 1077 crossmodal object runs. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Asterisks above or below bars indicate a significant
 1078 difference from zero. Horizontal lines within brain regions reflect an interaction between modality or congruency with
 1079 learning day (e.g., reduction in visual bias after crossmodal learning in PRC).

1080

1081 **Figure 3 – figure supplement 1. Signal quality comparison from a representative participant.** (a) The multi-
 1082 echo sequence we used acquired 3 measurements after every radiofrequency pulse, compared to the standard single-
 1083 echo EPI which acquires a single measurement (usually at a TE around 30 ms). A multi-echo sequence with 3
 1084 echoes acquires 3 times as much data as the current standard single-echo approach, and accounts for differences in
 1085 measured T_2^* across brain regions. For example, better signal is obtained at high TE values for the anterior
 1086 temporal lobes, which would otherwise reveal substantial signal dropout due to susceptibility artifacts at TE = 30
 1087 ms. (b) We optimally averaged the three echoes, using a method that weighs the combination of echoes based on the
 1088 estimated T_2^* at each voxel for each echo, and then applied ICA decomposition to remove non-BOLD noise. We
 1089 found that the multi-echo approach better recovered signal from the anterior temporal lobe structures compared to
 1090 the standard single-echo approach (shown in the Echo 2 column).

1091

1092 **Figure 4.** (a) Contrast matrix comparing the effect of congruency on feature representations. The voxel-wise matrix
 1093 averaged across unimodal runs were autocorrelated using the z-transformed Pearson's correlation, creating a
 1094 unimodal feature-level contrast matrix. We examined the average pattern similarity between unimodal features
 1095 associated with congruent objects (green) compared to the same unimodal features associated with incongruent
 1096 objects (yellow). (b) Pattern similarity analysis revealed an interaction between learning day and congruency in the
 1097 temporal pole (TP). At baseline before crossmodal learning, there was no difference in neural similarity between
 1098 unimodal features that paired to create congruent objects compared to the same unimodal features that paired to
 1099 create incongruent objects. After crossmodal learning, however, there was *less* neural similarity between the
 1100 unimodal features of pairs comprising congruent objects compared to the unimodal features of pairs comprising
 1101 incongruent objects. Because congruent and incongruent objects were built from the same shapes and sounds, this
 1102 result provides evidence that learning about crossmodal object associations influenced the representations of the
 1103 component features in the temporal pole. There was no difference between the congruent and incongruent pairings
 1104 in any other ROI (*Figure 4 – figure supplement 1*). ** $p < 0.01$.

1105

1106 **Figure 4 – figure supplement 1.** Pattern similarity analyses between unimodal features associated with congruent
 1107 objects and incongruent objects, before and after crossmodal learning (analysis visualized in *Figure 4* in the main
 1108 text). (a-c) Interestingly, the perirhinal cortex, LOC, and V1 – primarily visually-biased regions (see main text) –
 1109 reduced in pattern similarity after crossmodal learning. (d) By contrast, there was no change across learning days in
 1110 A1. No region displayed a difference between congruent and incongruent feature pairings other than the temporal

1111 pole (see *Figure 4*). * denotes $p < 0.05$, ** denotes $p < 0.01$, *** denotes $p < 0.001$. Horizontal lines denote the
 1112 main effect of learning day.

1113

1114 **Figure 5.** Contrast matrices and pattern similarity analyses investigating the effect of crossmodal learning on
 1115 modality-specific biases. The voxel-wise matrix for unimodal feature runs on Day 2 were correlated to the voxel-
 1116 wise matrix for crossmodal object runs on (a) Day 2 and (b) Day 4, creating a contrast matrix between visual and
 1117 auditory unimodal features to crossmodal objects that contained those features. We compared the average pattern
 1118 similarity (z-transformed Pearson correlation) between shape (blue) and sound (orange) features across learning
 1119 days. (a) Robust modality-specific feature biases were observed in all examined regions before crossmodal learning.
 1120 That is, pattern similarity for each brain region was higher for one of the two modalities, indicative of a modality-
 1121 specific bias. For example, pattern similarity in perirhinal cortex (PRC) preferentially tracked the visual features of
 1122 the crossmodal objects, evidence of a default visual shape bias *before crossmodal learning*. (b) Critically, we found
 1123 that perirhinal representations were transformed with experience, such that the initial visual bias was attenuated *after*
 1124 *crossmodal learning* (i.e., denoted by a significant interaction, shown by shaded green regions), evidence that
 1125 representations were no longer predominantly grounded in the visual modality. * $p < 0.05$, ** $p < 0.01$, *** $p <$
 1126 0.001 . Horizontal lines within brain regions indicate a significant main effect of modality. Vertical asterisks denote
 1127 pattern similarity comparisons relative to 0.

1128

1129 **Figure 5 – figure supplement 1.** Analyses for the hippocampus (HPC) and inferior parietal lobe (IPL). (a) In the
 1130 visual vs. auditory univariate analysis, there was no visual or sound bias in HPC, but there was a bias towards
 1131 sounds that increased numerically after crossmodal learning in the IPL. (b) Pattern similarity analyses between
 1132 unimodal features associated with congruent objects and incongruent objects. Similar to *Figure 4 – figure*
 1133 *supplement 1*, there was no main effect of congruency in either region. (c) When we looked at the pattern similarity
 1134 between Unimodal Feature runs on Day 2 to Crossmodal Object runs on Day 2, we found that there was significant
 1135 pattern similarity when there was a match between the unimodal feature and the crossmodal object (e.g., pattern
 1136 similarity > 0). This pattern of results held when (d) correlating the Unimodal Feature runs on Day 2 to Crossmodal
 1137 Object runs on Day 4, and (e) correlating the Unimodal Feature runs on Day 4 to Crossmodal Object runs on Day 4.
 1138 Finally, (f) there was no significant pattern similarity between Crossmodal Object runs before learning correlated to
 1139 Crossmodal Object after learning in HPC, but there was significant pattern similarity in IPL ($p < 0.001$). Taken
 1140 together, these results suggest that both HPC and IPL are sensitive to visual and sound content, as the (c, d, e)
 1141 unimodal feature-level representations were correlated to the crossmodal object representations irrespective of
 1142 learning day. However, there was no difference between congruent and incongruent pairings in any analysis,
 1143 suggesting that HPC and IPL did not represent crossmodal objects differently from the component unimodal
 1144 features. For these reasons, HPC and IPL may represent the convergence of unimodal feature representations (i.e.,
 1145 because HPC and IPL were sensitive to both visual and sound features), but our results do not seem to support these
 1146 regions in forming crossmodal integrative coding distinct from the unimodal features (i.e., because representations
 1147 in HPC and IPL did not differentiate the congruent and incongruent conditions and did not change with experience).
 1148 * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Asterisks above or below bars indicate a significant difference from zero.
 1149 Horizontal lines within brain regions in (a) reflect an interaction between modality and learning day, whereas
 1150 horizontal lines within brain regions in reflect main effects of (b) learning day, (c-e) modality, or (f) congruency.

1151

1152 **Figure 5 – figure supplement 2.** The voxel-wise matrix for Unimodal Feature runs on Day 4 were correlated to the
 1153 voxel-wise matrix for Crossmodal Object runs on Day 4 (see *Figure 5* in the main text for an example). We
 1154 compared the average pattern similarity (z-transformed Pearson correlation) between shape (blue) and sound
 1155 (orange) features specifically after crossmodal learning. Consistent with *Figure 5b*, perirhinal cortex was the only
 1156 region without a modality-specific bias. Furthermore, perirhinal cortex was the only region where the
 1157 representations of both the visual and sound features were not significantly correlated to the crossmodal objects. By

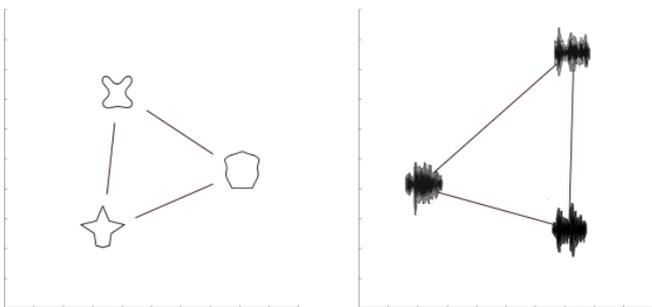
1158 contrast, every other region maintained a modality-specific bias for either the visual or sound features. These results
1159 suggest that perirhinal cortex representations were transformed with experience, such that the initial visual shape
1160 representations (*Figure 5a*) were no longer grounded in a single modality after crossmodal learning. Furthermore,
1161 these results suggest that crossmodal learning formed an integrative code different from the unimodal features in
1162 perirhinal cortex, as the visual and sound features were not significantly correlated with the crossmodal objects. * p
1163 < 0.05 , ** $p < 0.01$, *** $p < 0.001$. Horizontal lines within brain regions indicate a significant main effect of
1164 modality. Vertical asterisks denote pattern similarity comparisons relative to 0.

1165

1166 **Figure 6.** Contrast matrix shown on the left panel, with actual results shown on the right panel. We compared the
1167 average pattern similarity across learning days between crossmodal object runs on Day 2 with crossmodal object
1168 runs on Day 4 (z-transformed Pearson correlation). We observed lower average pattern similarity for incongruent
1169 objects (yellow) compared to congruent (green) objects in perirhinal cortex (PRC). These results suggest that
1170 perirhinal cortex differentiated congruent and incongruent objects constructed from the same features. Furthermore,
1171 pattern similarity was never above 0 for the perirhinal cortex. By contrast, there was no significant difference
1172 between congruent and incongruent objects in any other examined region, and pattern similarity was always above
1173 0. * denotes $p < 0.05$, ** denotes $p < 0.01$, *** denotes $p < 0.001$. Horizontal lines within brain regions denote a
1174 main effect of congruency. Vertical asterisks denote pattern similarity comparisons relative to 0.

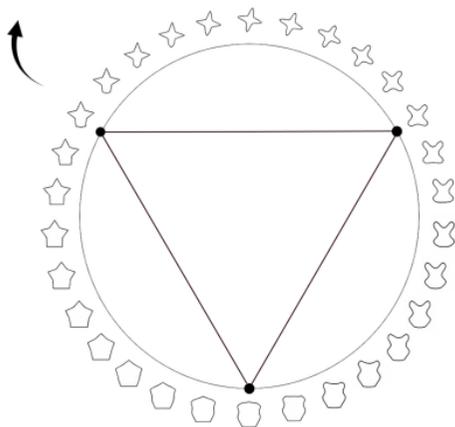
a. Characterizing Subjective Similarity

Validation Experiment ($n = 16$)



Shape Similarity

Sound Similarity



VCS space

b. 3D-Printed Shapes with Embedded Speakers



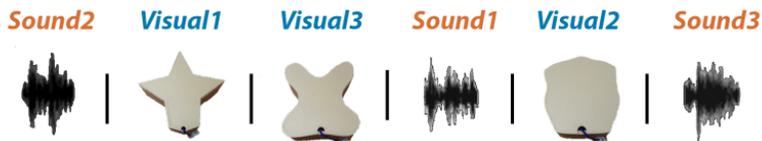
c. Example 3D Shape-Sound Object



Before Crossmodal Learning

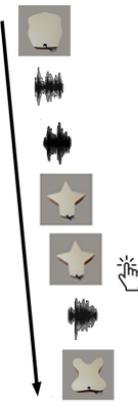
Day 1. Unimodal Learning (Behavior)

Unimodal Exploration:  



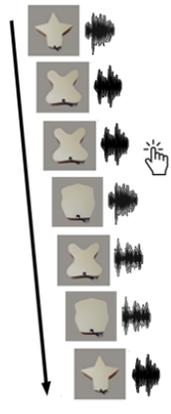
Day 2. Multi-echo fMRI (Neuroimaging)

Unimodal Feature Runs



 1-back target

Crossmodal Object Runs



CONGRUENT

-   *Visual1+Sound2*
-   *Visual3+Sound1*
-   *Visual2+Sound3*

INCONGRUENT

-   *Visual1+Sound3*
-   *Visual3+Sound2*
-   *Visual2+Sound1*

After Day 1: congruent objects not yet learned

After Crossmodal Learning

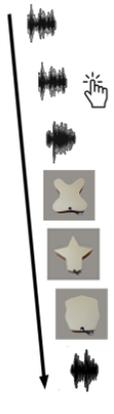
Day 3. Crossmodal Learning (Behavior)

Crossmodal Exploration:  



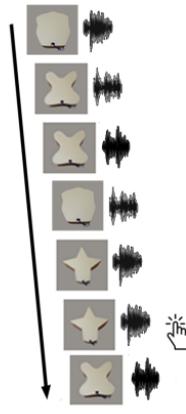
Day 4. Multi-echo fMRI (Neuroimaging)

Unimodal Feature Runs



 1-back target

Crossmodal Object Runs



CONGRUENT

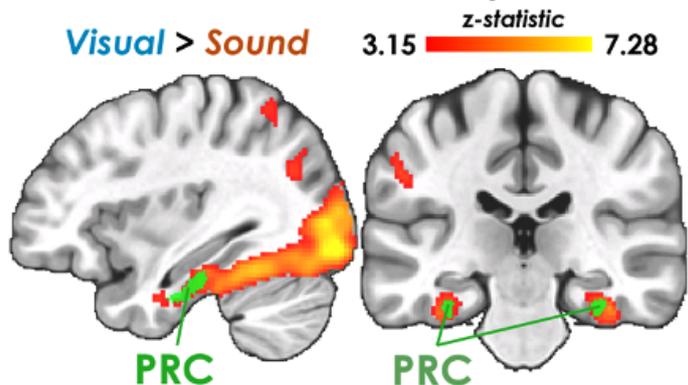
-   *Visual1+Sound2*
-   *Visual3+Sound1*
-   *Visual2+Sound3*

INCONGRUENT

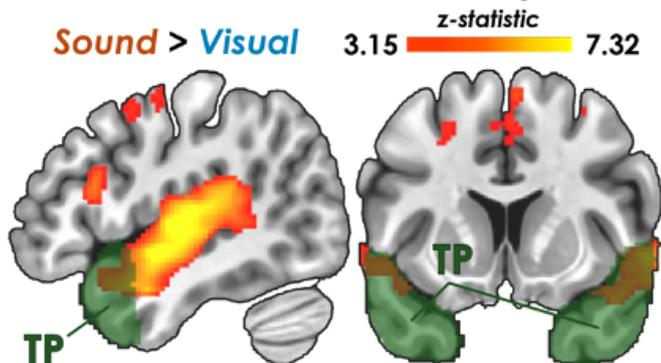
-   *Visual1+Sound3*
-   *Visual3+Sound2*
-   *Visual2+Sound1*

After Day 3: congruent objects learned

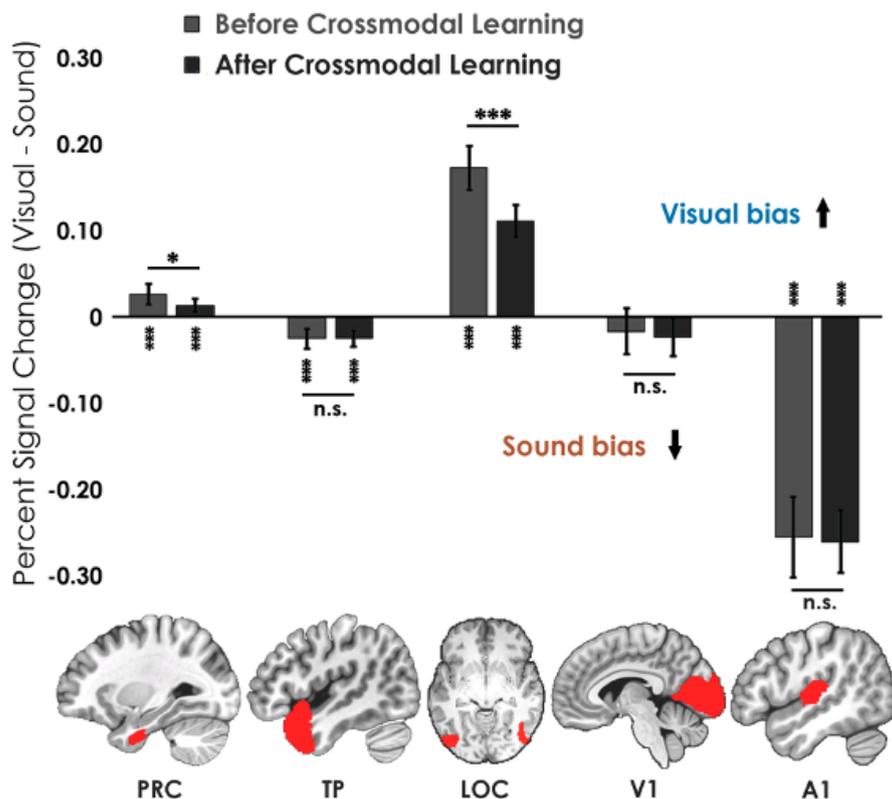
a. Whole-brain: Visual Activity



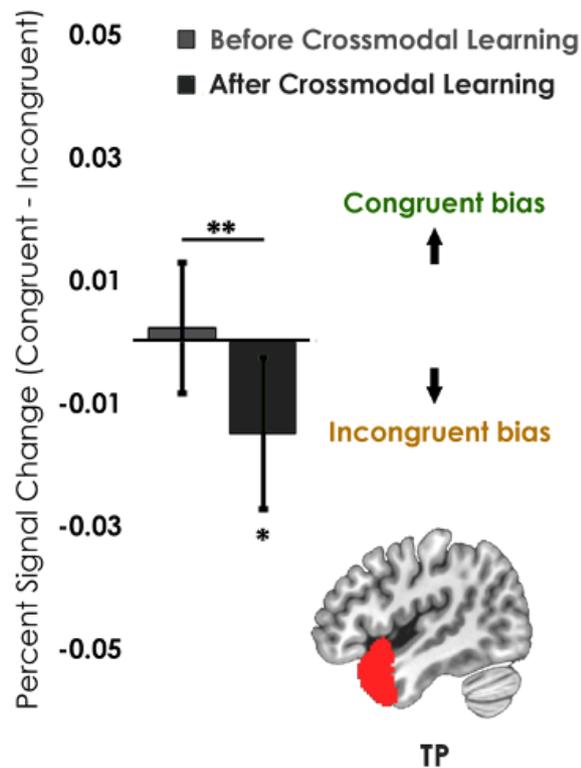
b. Whole-brain: Sound Activity



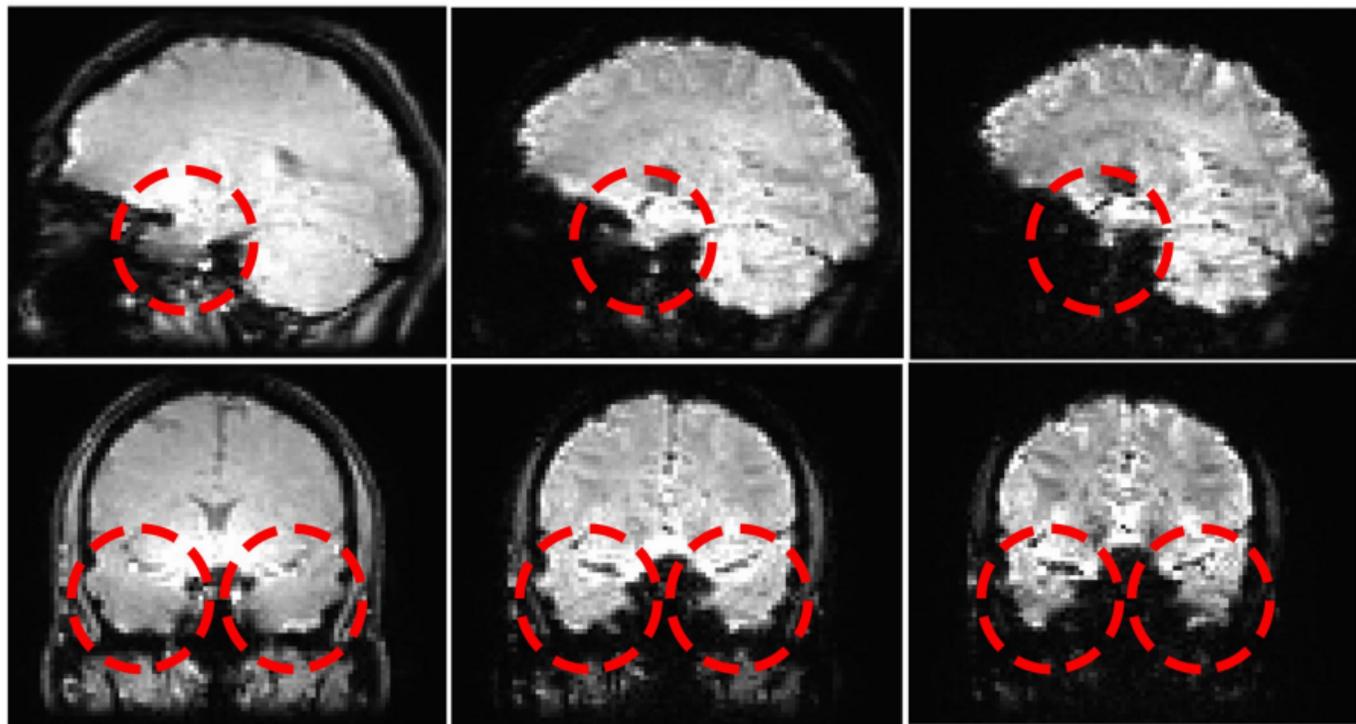
c. Unimodal Runs: Visual vs Auditory



d. Crossmodal Object Runs: Congruent vs Incongruent



a. Multi-echo fMRI Sequence

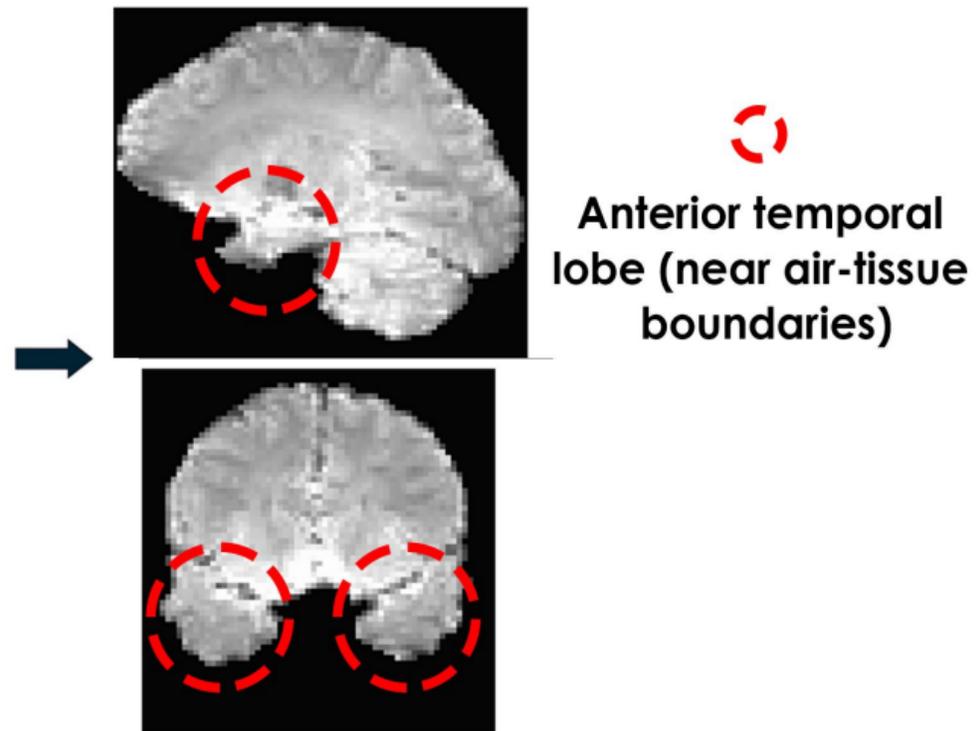


Echo 1
TE = 11 ms

Echo 2 (Standard EPI)
TE = 31.6 ms

Echo 3
TE = 52.2 ms

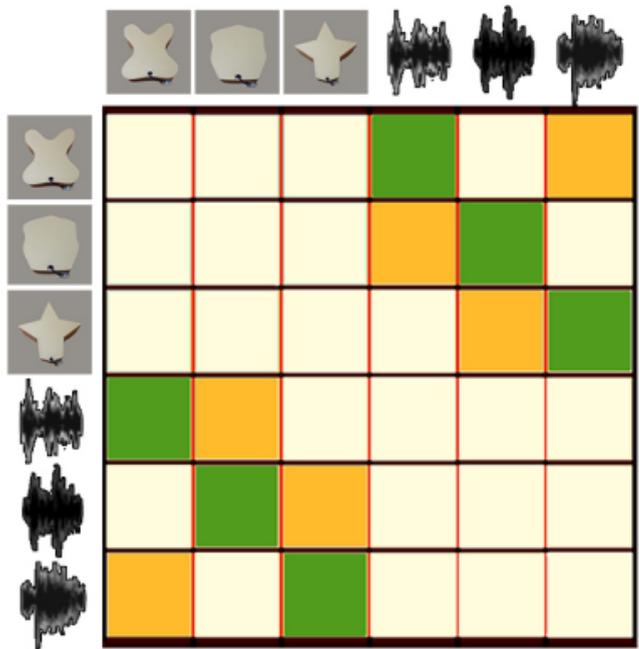
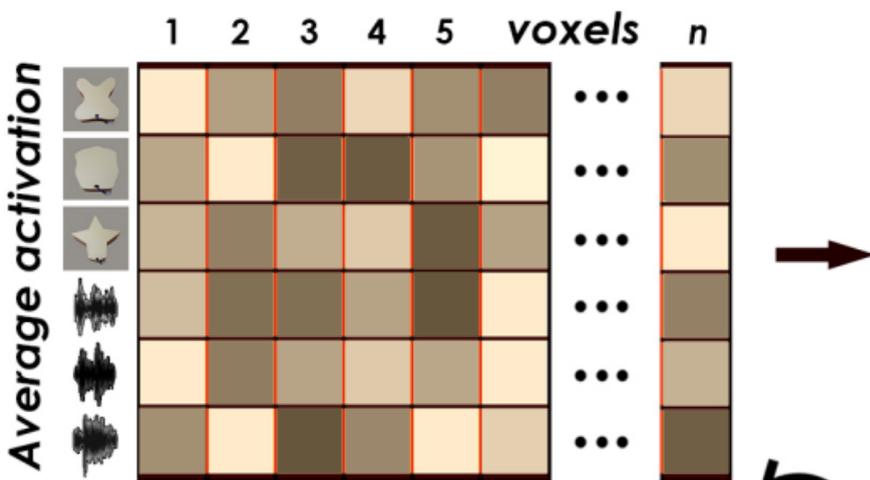
b. Optimal Combination with ICA



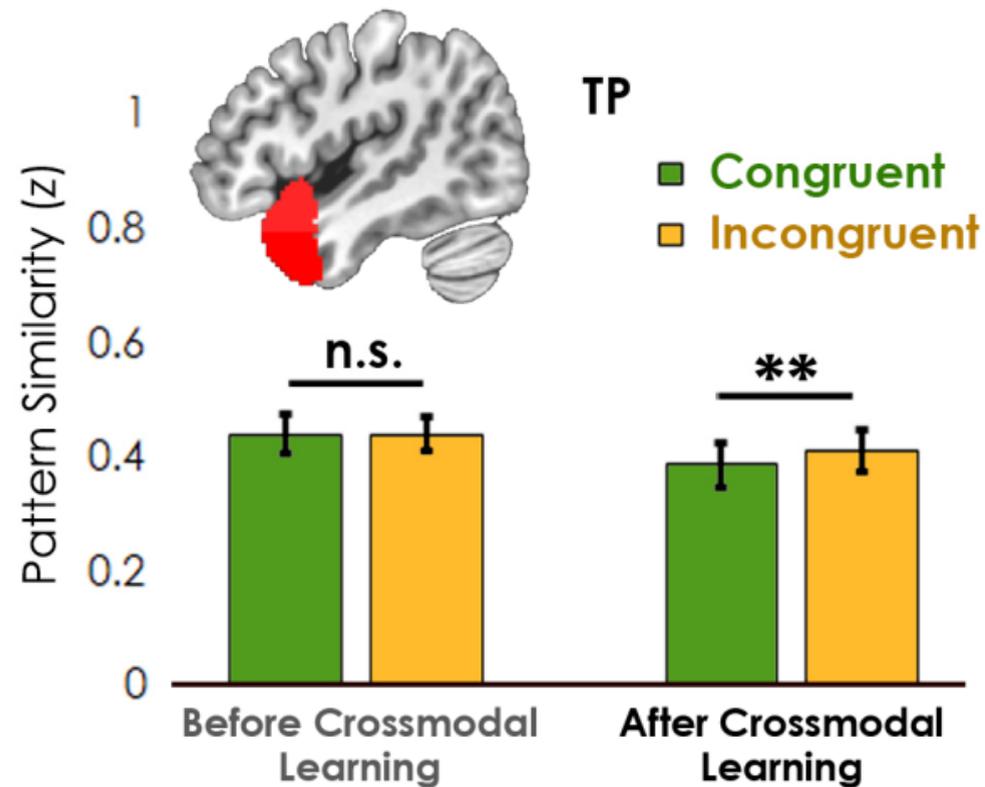
Combined Echoes

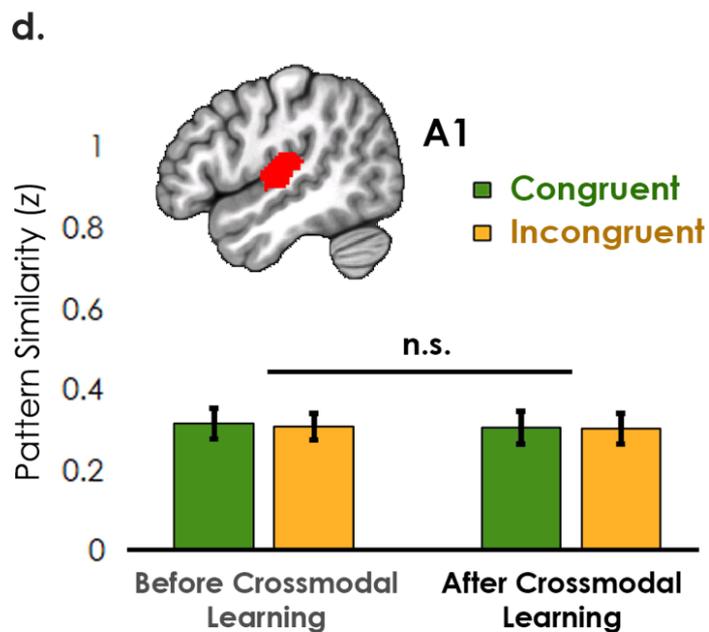
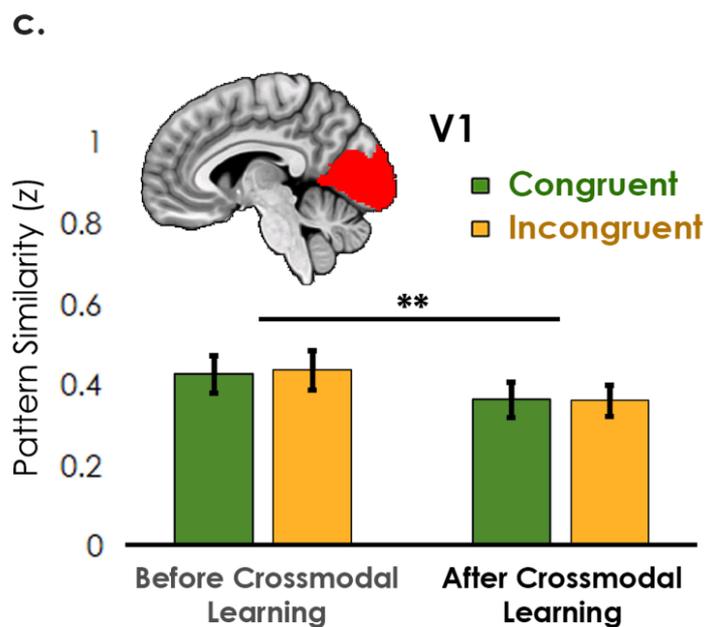
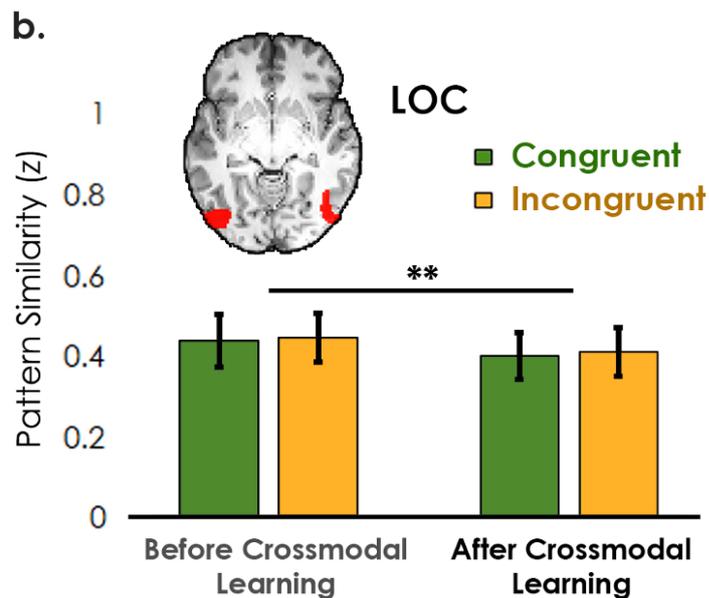
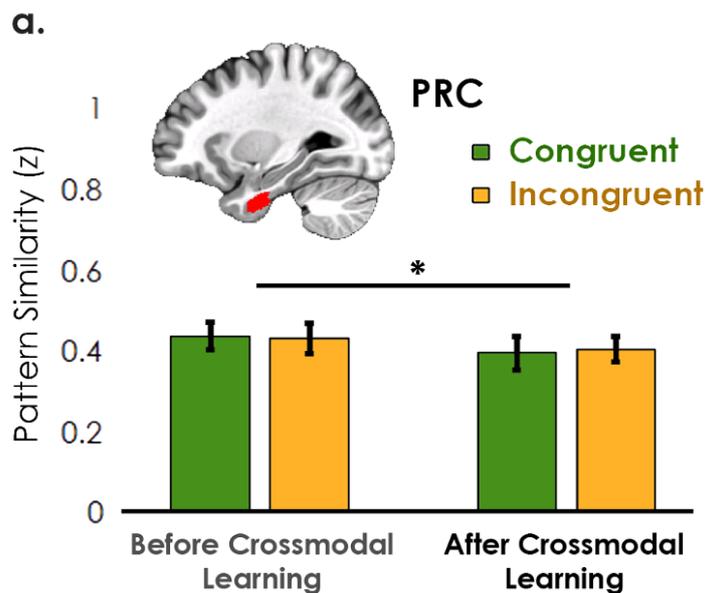
a. Pattern similarity analysis (Unimodal Feature Runs)

ROI: Average voxel-wise matrix

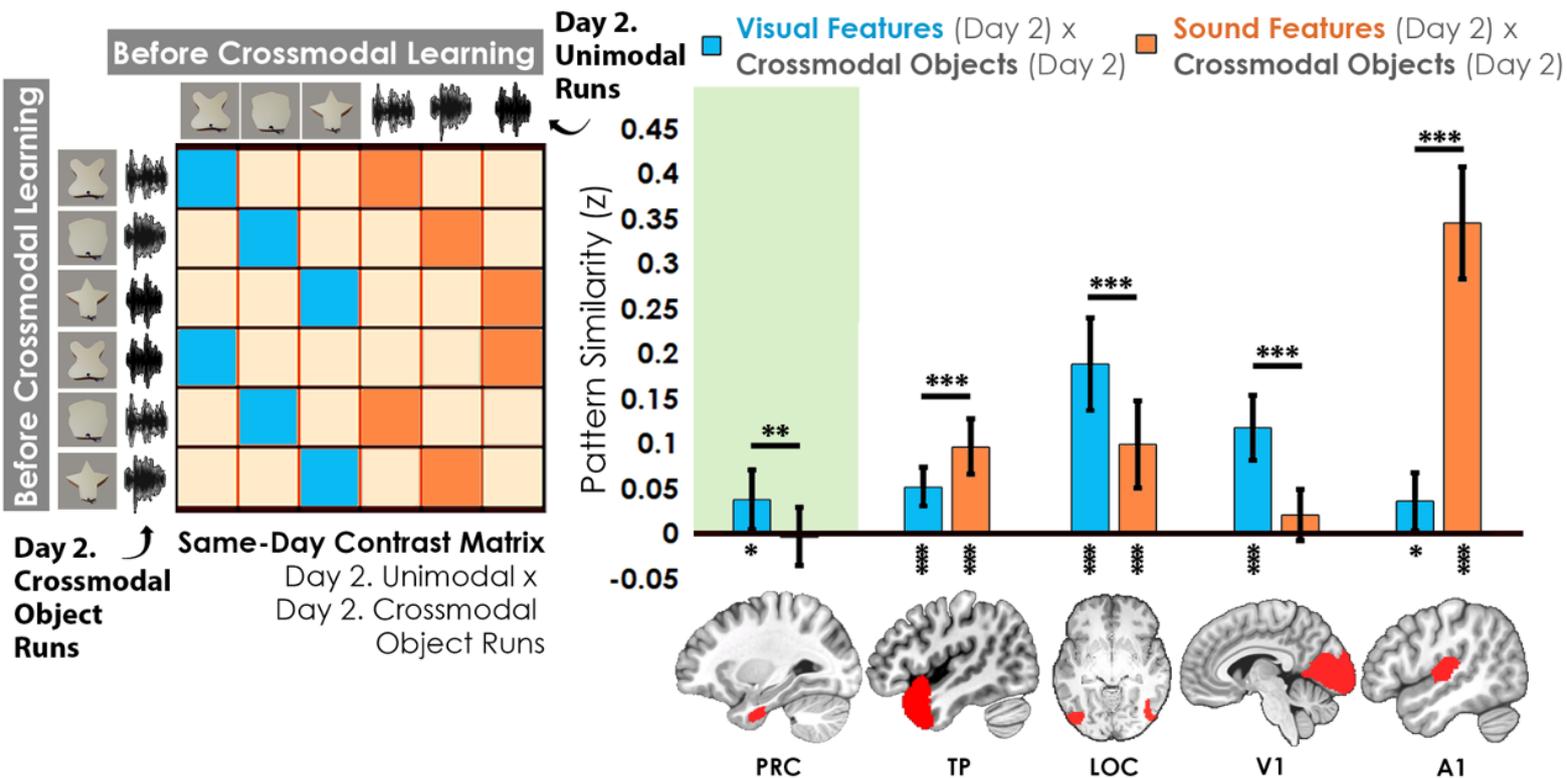


b. Temporal pole differentiates congruent and incongruent pairings

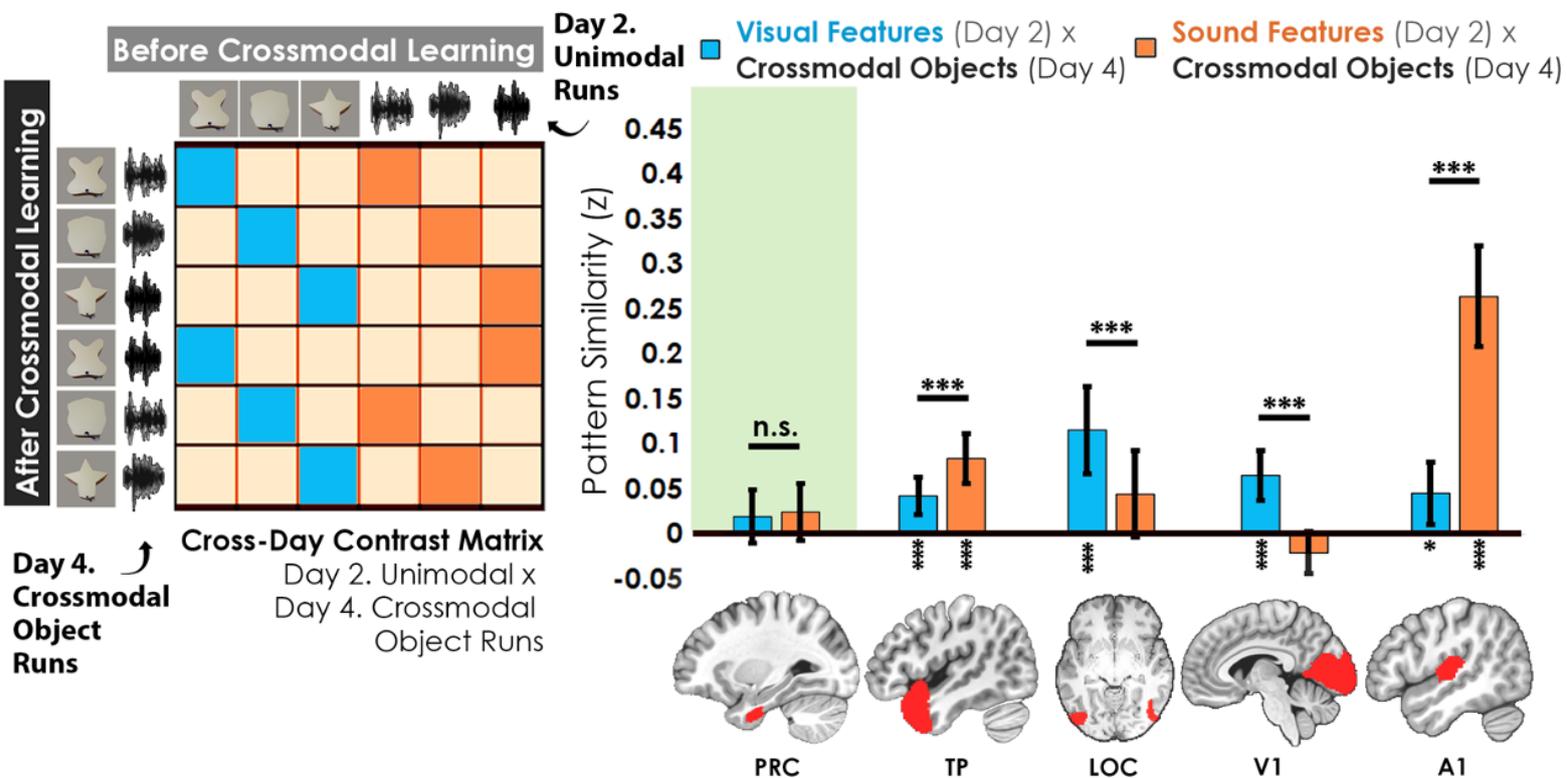




a. All regions show a modality-specific bias prior to crossmodal learning

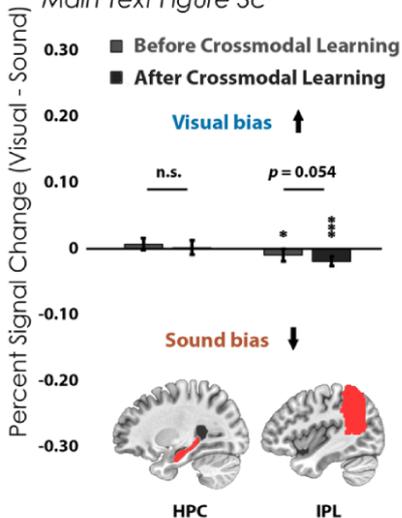


b. Perirhinal cortex is the only region to lose its modality-specific bias after crossmodal learning



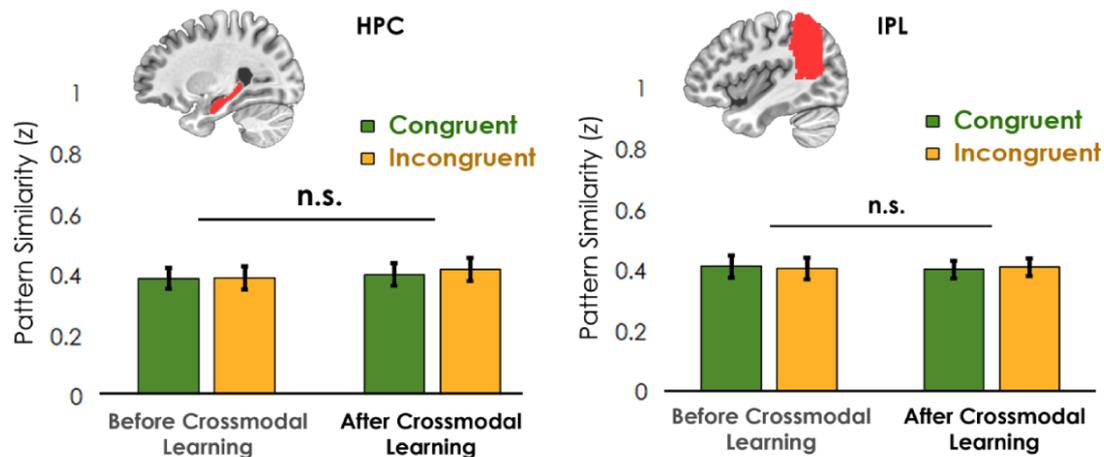
a. Unimodal Runs: Visual vs Auditory

Main Text Figure 3c



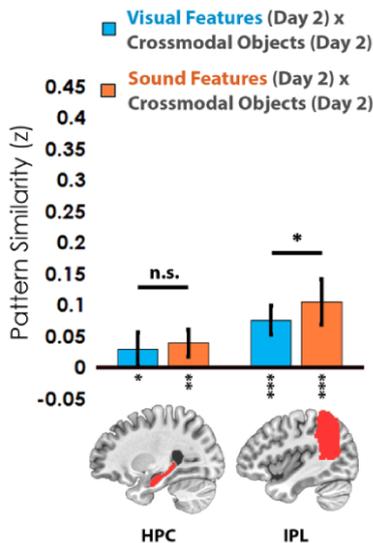
b. Pattern Similarity Analysis (Unimodal Runs)

Main Text Figure 4



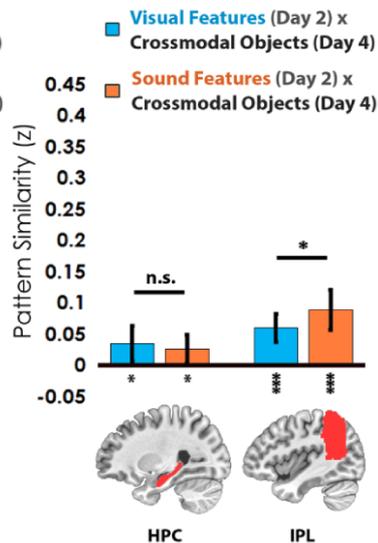
c. Day 2 x Day 2

Main Text Figure 5a



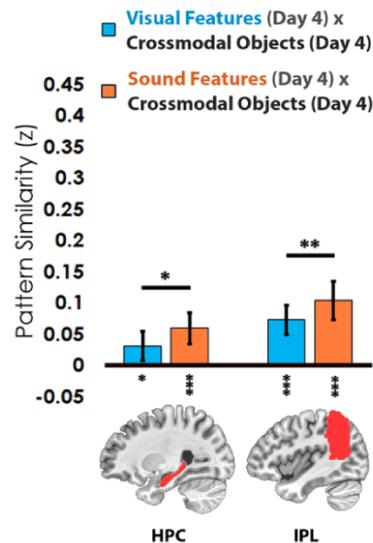
d. Day 2 x Day 4

Main Text Figure 5b



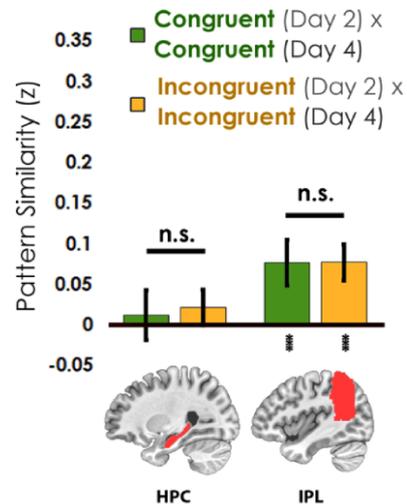
e. Day 4 x Day 4

Supplemental Figure S5



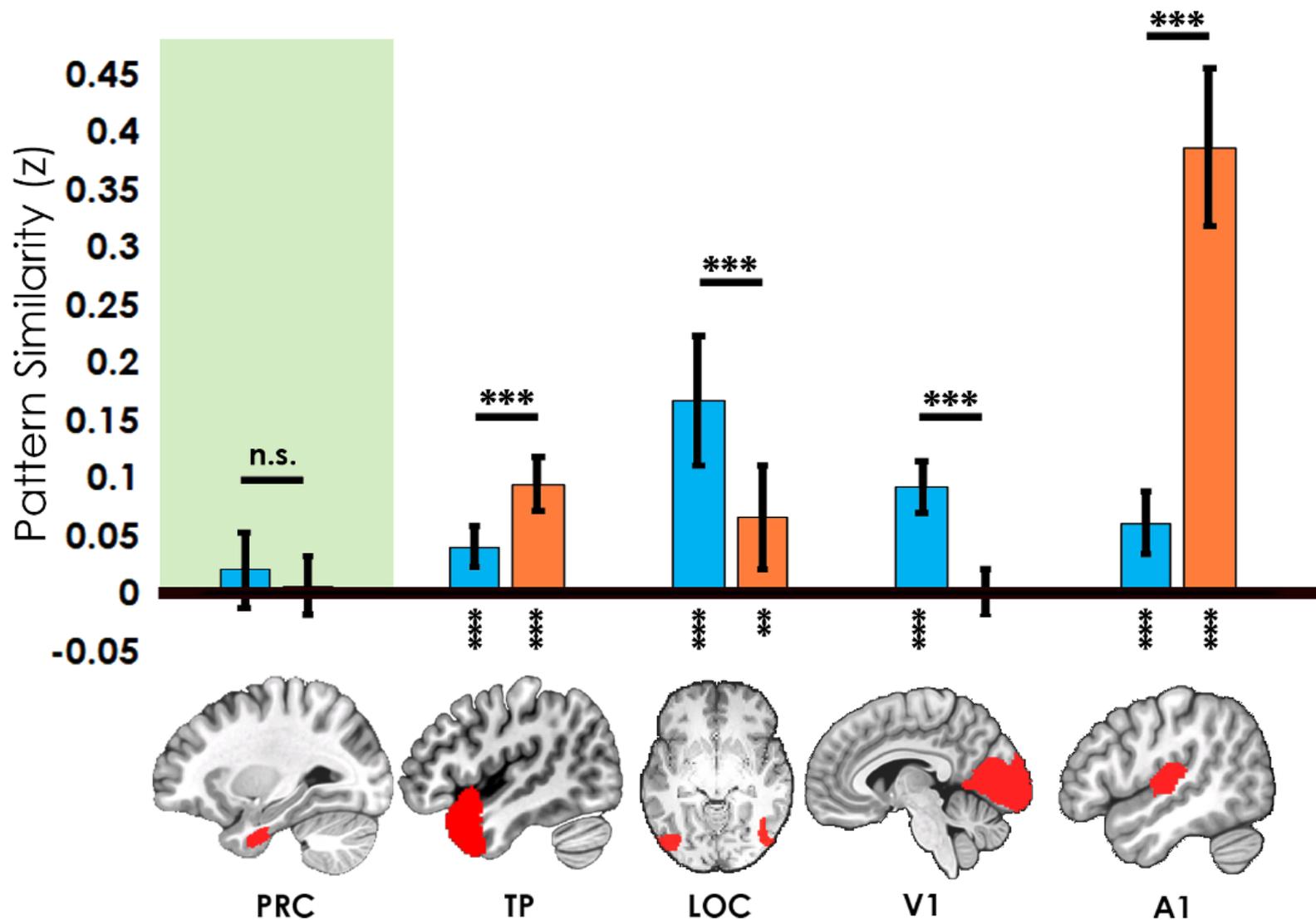
f. Crossmodal Object Runs Before and After Crossmodal Learning

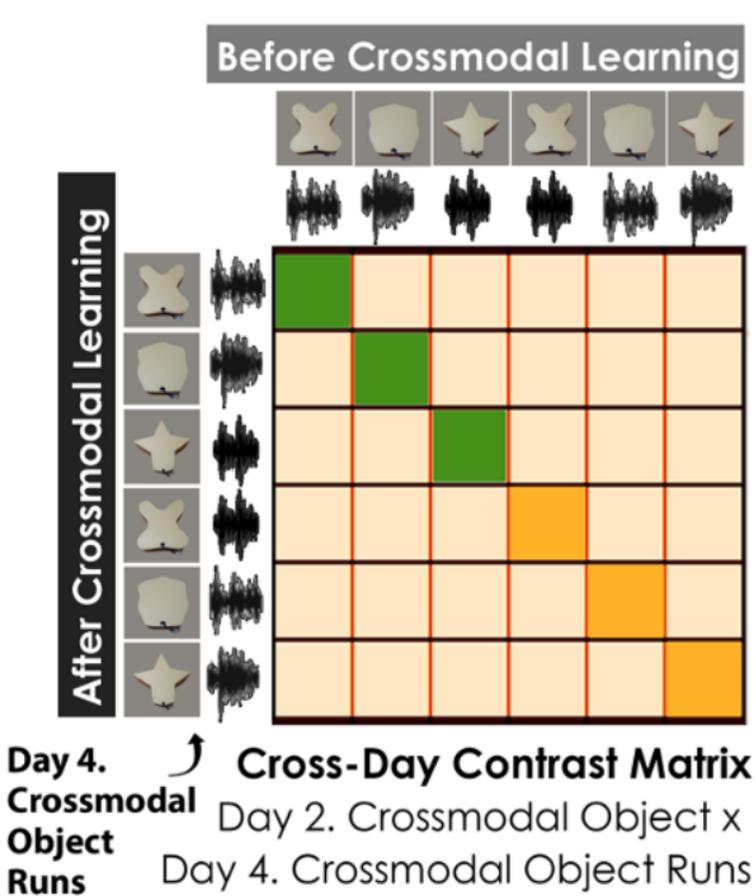
Main Text Figure 6



■ **Visual Features** (Day 4) x
Crossmodal Objects (Day 4)

■ **Sound Features** (Day 4) x
Crossmodal Objects (Day 4)





Day 2. Crossmodal Object Runs

■ **Congruent** (Day 2) x **Congruent** (Day 4)

■ **Incongruent** (Day 2) x **Incongruent** (Day 4)

Pattern Similarity (z)

