

Cross-Modal Environment Self-Adaptation During Object Recognition in Artificial Cognitive Systems

David Miralles (✉ david.miralles@salle.url.edu)

La Salle - Universitat Ramon Llull

Guillem Garrofé

La Salle - Universitat Ramon Llull

Calota Parés

La Salle - Universitat Ramon Llull

Alejandro González

La Salle - Universitat Ramon Llull

Gerard Serra

La Salle - Universitat Ramon Llull

Alberto Soto

La Salle - Universitat Ramon Llull

Hans Op de Beeck

KU Leuven

Haemy Lee Masson

Johns Hopkins University

Xavier Sevillano

La Salle - Universitat Ramon Llull

Research Article

Keywords: Cross-modal, environment, ACS, Molyneux

Posted Date: July 30th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-754574/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Cross-modal environment self-adaptation during object recognition in artificial** 2 **cognitive systems**

3
4 David Miralles^{1*}, Guillem Garrofé¹, Carlota Parés¹, Alejandro González¹, Gerard Serra¹, Alberto Soto¹,
5 Xavier Sevillano¹, Hand Op de Beeck², Haemy Lee Masson³

- 6
7 (1) Grup de recerca en Tecnologies Media, La Salle-Universitat Ramon Llull, Barcelona, Catalonia,
8 Spain.
9 (2) Department of Brain and Cognition, Leuven Brain Institute, KU Leuven, Leuven, Belgium.
10 (3) Department of Cognitive Science, Johns Hopkins University, USA
11

12 **Abstract**

13
14 The cognitive connection between the senses of touch and vision is probably the best-known case of cross-
15 modality. Recent discoveries suggest that the mapping between both senses is learned rather than innate. These
16 evidences open the door to a dynamic cross-modality that allows individuals to adaptively develop within their
17 environment. Mimicking this aspect of human learning, we propose a new cross-modal mechanism that allows
18 artificial cognitive systems (ACS) to adapt quickly to unforeseen perceptual anomalies generated by the environment
19 or by the system itself. In this context, visual recognition systems have advanced remarkably in recent years thanks to
20 the creation of large-scale datasets together with the advent of deep learning algorithms. However, such advances
21 have not occurred on the haptic mode, mainly due to the lack of two-handed dexterous datasets that allow learning
22 systems to process the tactile information of human object exploration. This data imbalance limits the creation of
23 synchronized multimodal datasets that would enable the development of cross-modality in ACS during object
24 exploration. In this work, we use a multimodal dataset recently generated from tactile sensors placed on a collection
25 of objects that capture haptic data from human manipulation, together with the corresponding
26 visual counterpart. Using this data, we create a cross-modal learning transfer mechanism capable of detecting both
27 sudden and permanent anomalies in the visual channel and still maintain visual object recognition performance by
28 retraining the visual mode for a few minutes using haptic information. Here we show the importance of cross-
29 modality in perceptual awareness and its ecological capabilities to self-adapt to different environments.

30 **Introduction**

31
32 Humans perceive the environment through multiple senses. A set of sensory information, acquired
33 through each modality, is integrated and transformed into a supra-modal representation. This process
34 requires cross-modality (also referred to as cross-modal transfer or cross-modal matching) – the cognitive
35 ability to associate the sensory features acquired independently through multiple senses. Human
36 manipulation of objects, a natural example of cross-modality, connects the senses of sight and touch from
37 an early age, and this sensory connection is strengthened over the course of child development (1), and
38 stays throughout the lifespan (2). In particular, vision and haptics are complementary to each other,
39 improving the credibility of mental representation of object properties and recognition performance (3, 4,
40 5, 6). In this article, we present an Artificial Cognitive System (ACS) that builds on a cross-modality
41 ability using human manipulation data, achieving perceptual awareness and a dynamic capacity to adapt
42 to changing environments.
43

44 The question of how humans achieve cross-modality was sparked off by 17th-century natural philosopher
45 William Molyneux. In his letter to John Locke, he questioned whether a congenitally blind person, who
46 recently gained vision, would be able to visually recognize an object, previously known only by touch, or
47 she/he would need to learn to make the intermodal transfer from touch to vision (7). To this day, the
48 debate over whether this transfer ability is innate or acquired has led scientists to investigate cross-
49 modality in newborns, animals, and congenitally blind individuals (8, 9, 10). Of these studies, a recent
50 cross-modal matching experiment, conducted on congenitally blind individuals who later gained sight as
51 adults, suggests that cross-modality is acquired, dynamic, and moldable rather than innate and
52 predetermined (9). This dynamic nature of cross-modality allows human observers to modulate the
53 strength of intermodal connection based on the reliability of the information derived from each modality

54 (5), an ability that has also been observed in various animal species, including capuchin monkeys (11),
55 rodents (12), and even bumblebees (13).

56 In the past few years, there has been a growing interest in the development of cross-modality in artificial
57 agents (14), especially robots, as it may facilitate the creation of systems that autonomously adapt to
58 different environments. In particular, with the advancement of visuo-tactile sensors for haptic capture
59 systems, researchers have started investigating the cross-modal connection between vision and touch in
60 robotics (15, 16, 17). The tactile patterns that result from these sensors can be related to the images
61 obtained using the visual mode, thus creating a framework that enables the establishment of cross-modal
62 relationships. By gathering haptic data in the form of images, the dimensional gap between touch and
63 vision features can be successfully overcome (18). Furthermore, the maturity of image recognition
64 methods (e.g. deep learning algorithms) has allowed some progress in the study of the relationship
65 between these two modes. However, existing haptic data is still a long way from what would be a haptic
66 dexterous robot exploration similar to that of humans. To solve this problem, researchers have recently
67 designed a glove that through a mechanoreceptor network can provide tactile patterns to the system, as
68 well as information related to the dexterity of the human grasp (19). However, that work does not relate
69 haptic data to visual mode. In this paper, we use a haptic capture system (20) that also leverages
70 information from the human manipulation of objects, but with the ultimate goal of delving deeper into the
71 design of cross-modality in ACS.

72 To achieve this objective, we designed and printed novel 3D objects that collect human exploration data
73 with multiple capacitive touch sensors on the object surface. With this dataset, our ACS achieves cross-
74 modality via transfer learning from touch to vision. Unlike other approaches (17), where corrupted inputs
75 are incorporated during training time, we present a new mechanism that allows our ACS to use cross-
76 modality to continuously monitor whether the information received from visual and haptic modes
77 matches using cross-modality, hence being able to detect anomalies (e.g. blurred vision). Given that the
78 two sensory modalities are independent but collaborative, like those of a human, we examine how our
79 system dynamically changes the strength of the intermodal connection to better solve object recognition
80 problems when we degrade the quality of visual information at test time. If mismatches between vision
81 and touch channels persist over time, the ACS can autonomously retrain the faulty modality through
82 transfer knowledge within few minutes. Our findings suggest that with the implantation of biologically
83 inspired cross-modality, the ACS becomes perceptually aware of a faulty sensory modality and
84 autonomously adapts to changing environments without losing performance.

85

86

87 **Results**

88

89 **Haptic data and object recognition**

90

91 We designed a system that captures haptic information generated by humans during the object
92 manipulation process. This collected data is enough to create a haptic recognition system that outperforms
93 humans in a classification task with similar 3D shapes (20), both in accuracy and response time. As
94 illustrated in Fig. 1, the objects are six similar shapes we have digitally created and 3D printed (20, 21).
95 The external surface of each object is totally covered with 24 copper pads equally distributed and
96 connected to an electronic board placed inside, which also includes a gyroscope. During the object
97 manipulation, this system samples data from all the sensors at 40Hz and sends it to a computer through
98 wireless communication. Every sample is stored as a 24-bit array, h_j , called haptic state of the j sample,
99 one bit for each copper pad (Fig. 1B), hence the system has no information about the relation between the
100 location of the 24 sensors and the positions of their statuses in the array. Besides, since each sensor has its
101 position inside the array, the resulting state would vary if sensors were placed differently. For this reason,
102 and to make sure the presented algorithm does not use the sensors' order to recognize the objects, sensors

103 are placed in a way that every position in the haptic state corresponds to a sensor located in the same
 104 spatial location for every object. In the same way that touch receptors on the human hand would also
 105 receive input from corresponding locations on different objects when the relative orientation of the
 106 grasping hand and the objects would be held constant.

107 Our haptic dataset is based on these haptic states and their time evolution. The geometry of the objects
 108 affects their handling, and this is reflected in our data. To gather our haptic dataset, one participant is
 109 invited to manipulate each object with both hands and perform a random exploration task. Four series,
 110 each lasting for five minutes, have been recorded per object.

111 To perform automatic haptic object recognition, the dataset is divided into two parts: three series for
 112 training (15 minutes) and one for testing (5 minutes). To determine the probability of a set of n
 113 consecutive haptic states (h_1, \dots, h_n) belonging to a specific object S_i ($i \in [1, \dots, 6]$), i.e., $P(S_i |$
 114 $h_1, \dots, h_n)$, we adopt a naïve Bayes approach as follows:

$$115 \quad \hat{S}_i = \underset{S_i, i \in [1, \dots, 6]}{\operatorname{argmax}} \left[P(S_i) \prod_{j=1}^n P(h_j | S_i) \right]$$

116 Here the n -product $P(h_j | S_i)$ is the naïve condition, $P(S_i) = \frac{1}{6}$ is the probability of each object and \hat{S}_i the
 117 resulting prediction.

118 Our haptic object recognition system achieves an average accuracy of 89.63% after just 8 seconds of
 119 manipulation (average time for best accuracy in humans (20)), as shown in Fig. 1.

120

121 **Multimodal dataset generation and visual object recognition**

122 As stated earlier, the 3D printed objects have been produced from a 3D digital render. Moreover, from the
 123 data collected by the gyroscope located inside the objects and using the 3D renders, we synthesize a video
 124 with the movements of the objects caused by human manipulation (Fig. 2). We create one video frame for
 125 each haptic state; see Methods for details. We select the 5-minutes test data series mentioned above for
 126 this purpose. This opens the visual channel to our system. Now, the ACS can receive data from haptic and
 127 visual senses simultaneously (Fig. 2).

128 Since the ACS has previously had a haptic experience, it can autonomously tag the visual dataset through
 129 the results of the already trained haptic object recognition system, creating a cross-modal relationship.
 130 Obviously, after this transfer of knowledge, we can train a visual object recognition system from the new
 131 labelled visual dataset. We divide the 5-minute visual dataset into three parts: 60% for training, 20% for
 132 testing, and reserve the remaining 20% for a later experiment. Using this data, we train a set of
 133 convolutional neural networks (CNN) for visual object recognition. The visual object recognition system
 134 is based on a one-vs-all strategy combining six CNN classifiers, thus yielding a six-dimensional output,
 135 $v \in \mathbb{R}^6$. Each component of the output, v_i , is associated to the probability that the visual input
 136 corresponds to one of the 6 stimuli S_i (objects). To decide which is the corresponding object, the
 137 following criterion is adopted: a sample belongs to a certain class if only one of the values v_i exceeds a
 138 threshold, τ , ($\exists! v_i$ s.t. $v_i > \tau$). The value of τ is obtained from the relationship between the accuracy of
 139 the visual classifier and the different threshold values as shown in Fig. 3. Two other results can be
 140 obtained: a) Confusion (CF): there is more than one v_i value above the threshold ($\exists v_i, v_j$ s.t. $v_i > \tau$ and
 141 $v_j > \tau$ with $i \neq j$), that is, the visual object recognition system assigns the sample to more than one class.
 142 b) Ignorance (IG): there is no value of v_i above the threshold ($v_i < \tau, \forall i$), i.e., there is a lack of
 143 knowledge to decide which class it belongs to. Although there are several efficient methods for dealing
 144 with confusion, especially based on heuristics, at this stage we have preferred to include it as a false
 145 negative, even at the expense of visual object recognition performance. We will resolve the confusion
 146 with the help of the haptic mode, thus allowing true cross-modality, as we will see in the next section.

147 Visual object recognition results shown in Fig. 3A demonstrate that the visual channel of our ACS
148 achieves 75% accuracy including CF and IG as false negatives.

149

150

151 **Stressing the visual channel and the Molyneux mechanism**

152 The Molyneux problem addresses the following question: would a person born blind that later regains
153 sight as an adult, be able to visually recognize the shapes of objects previously experienced by touch?
154 Recent empirical studies have pointed out that upon recovery of sight, subjects are initially unable to
155 recognize these objects visually. However, after they experience the world with both senses, in a few days
156 a cross-modal link is created allowing them to pass the Molyneux test (9). In the present study, this
157 connection between the two senses equips the ACS with a cognitive mechanism that allows it to
158 autonomously detect a faulty channel.

159 The aforementioned mechanism, which we have called the Molyneux mechanism, allows the ACS to
160 continuously check if what it is seeing agrees with what it is touching. Hence, the ACS can detect sudden
161 anomalies by comparing the classification labels of the two independent channels, haptic and visual (Fig.
162 3B).

163 More specifically, the ACS can encounter four different situations while comparing the labels given by
164 the haptic and visual object recognition systems: i) both recognition systems agree, i.e., there is a match,
165 ii) do not agree, i.e., there is a mismatch, iii) the visual object recognition system does not have enough
166 knowledge to classify that sample (IG), which results in a mismatch, and iv) the visual channel object
167 recognition system assigns the sample to more than one class (CF) and the ACS checks if one of this
168 classes agrees with the haptic recognition (match) or, on the contrary, it does not (mismatch). It is worth
169 noting that when both recognition systems are working properly, some short duration mismatches can
170 appear, but the common situation is a continuous agreement between the visual and the haptic recognition
171 systems. As these mismatches are short in time, they can be easily removed with a low pass filter; see
172 Methods for details.

173 In order to test this mechanism and study the effectiveness of artificial cross-modality for perceptual
174 awareness, we stressed our visual channel by applying a blur filter. Once applied, the accuracy of the
175 visual object recognition system went down to approximately 20% for the 6 classes (see Fig. 3A). The
176 ACS detected this anomaly using the Molyneux mechanism, and did it quickly, with an average delay of
177 only 2.13 seconds (85.33 samples) after applying the blur filter to the visual input. In Fig. 4 we show the
178 results for each object.

179

180 **Self-adaptation to a new environment**

181 The Molyneux mechanism enables ACS to realize, in a fully automatic way, that the haptic and visual
182 object recognition systems lose coherence when a blur filter is applied to the visual channel. This
183 situation does not affect the haptic classifier, which remains stable and consistently gives trustworthy
184 information to the visual channel. Thus, the blurred images are tagged with the output of the haptic
185 classifier, and the ACS can retrain the CNNs of the visual object recognition system in order to classify
186 them correctly. Using the 20% (60 seconds) of the visual dataset that we had previously reserved after
187 going to a blur filter, we retrained the visual object recognition system by grouping the retraining samples
188 into 8-second batches (total of 7 batches) and iteratively feeding the CNNs with one batch at a time. For
189 each retraining batch iteration, we calculate the current accuracy of the visual object recognition system
190 in the blurred vision scenario. As shown in Fig. 5, the accuracy in the blurred vision scenario increases as
191 the ACS receives more blurred visual information in its retraining process.

192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237

Discussion

Visual object recognition and cross-modality

As described in the Results section, the visual mode benefits from the previous experience of the haptic mode. The haptic classifier recognizes the object being manipulated and the ACS tags the information of the visual channel in real-time. This would be equivalent to the process where for the first time a human becomes aware of the connection that exists between a known manipulated object and its visualization. It is important to highlight the autonomy that this cross-modal knowledge transfer gives to the ACS since it does not need an external agent to tag the visual channel (Fig. 2B).

On the other hand, heuristic arguments are commonly used to solve the confusion in the outputs of a one-vs-all visual classifiers. However, in our approach, cross-modality allows us to solve confusion through the haptic mode, i.e. through the interaction with the environment (object). In this study, if we use a highest value heuristic, our visual model reaches an accuracy of 79%, whereas if we use cross-modality to solve confusion, the accuracy goes up to 81% (see Fig. 5). Although in this case accuracies are similar, if the haptic classifier is stable, the cross-modality will always equal or improve the results of heuristic methods. This result resembles the exchange of multi-modal information to discriminate a stimulus in humans, which is a very common process (3).

Perceptual awareness

The Molyneux mechanism defined in this article allows the ACS to check the coherence between two synchronized samples (haptic-visual) (Fig. 3B). In other words, this mechanism allows the ACS to answer the question: is what I touch and what I see the same object? As shown in Fig. 4, in the Results and also in Methods section, the study of these anomalies has been detailed to differentiate them from those that lengthen in time. Filtering visual-haptic classification pairs in real-time allows the ACS to realize that the visual object recognition system is not working properly when the blur filter is applied. High and stable accuracy over time of the haptic classifier is assumed for this study as shown in Fig. 1. The goal of the filtering process is to identify changes in the visual channel. It can be observed (see Fig. 4) that for objects lat00 and lon05 the change to blurred vision is not detected. This is because the visual object recognition system continues to classify these two objects correctly despite the blur, and therefore there is no incoherence of any kind other than a decrease in visual accuracies for these two objects, i.e., the filter does not detect any change in the visual channel.

Resilience from cross-modality

Realizing that the environment has changed is the first step in the process of self-adapting to it. The proposed design with two independent object recognition systems and the cross-modality approach allows the ACS to autonomously adapt to changes in the environment that affect one of its sensory modes using the information of the mode that remains stable. The results from Fig. 5 show how after just one minute of retraining the visual classifier (CNN), the ACS adapts to the blurred scenario with an accuracy of 69.1%

It is important to highlight the ecological character of the transfer of knowledge between modes in the sense proposed in (9), since ACSs could adapt to transformations or changes in their perception systems during their lifespan.

238

239

240 **Limitations**

241 With the approach proposed in this article, we are aware that we are simplifying the problem by using
242 synthetic images and avoiding the occlusions caused by human hands during object manipulation.
243 Nonetheless, even though this work is not focused in solving this issue, this apparent problem could be
244 part of the solution, since these occlusions are strongly correlated with haptic data.

245 Another limitation of this approach is that, in this very first experiment where we have shown the benefits
246 of cross-modality for self-adaptation in a changing environment where suddenly the vision channel gets
247 blurred, we have not accounted for other errors that may cause the vision channel to stop working
248 properly. Moreover, we have used single frames instead of a group of states in the visual model. Our goal
249 with the presented experiment was to test the Molyneux mechanism and show the benefits of cross-
250 modality for environment self-adaptation for object recognition tasks. We believe this is the very first step
251 towards the development of perceptual awareness in ACSs for adaptation to changes in their perception
252 systems, and further studies are required.

253 **Opportunities for future research**

254 Although the current trend is to place touch sensors on robotic hands' end-effectors, the use of sensors on
255 objects is an equally important field of research, especially in obtaining data for ACS. In fact, it seems
256 reasonable to assume there would be a correlation between what was obtained by the introduced haptic
257 capture system and the point cloud that a robotic hand could generate if it could interact with that system.
258 We hypothesize that placing sensors directly on objects is equivalent to obtaining data from a human-like
259 robotic dexterous hands. This would allow us to integrate our ACS to a robotic hand such as Shadow
260 Hand (22). Other researchers have showed that this hand (23) could exhibit high levels of dexterity in
261 object manipulation tasks.

262 It would also be interesting to study other ways to stress the input channels. Two situations of special
263 interest are: a) the inducement of errors in the haptic channel to study cross-modality in the opposite
264 direction, b) the desynchronization between the haptic and visual channels.

265 Finally, to extend this work to unfamiliar objects, a deeper study of the Molyneux mechanism is
266 necessary. To this end, we propose an analysis of the haptic and visual perceptual spaces used in
267 neuroscience (24), which would allow us to understand how we can relate unknown objects with a
268 previous experience of the ACS.

269

270

271

272 **Methods**

273

274 In this section, we provide the methods and procedures used in this research article.

275 **Visual dataset generation**

276 To generate the 3D renders, the six similar shapes meshes are placed independently on a Processing 3D
277 scene with a 299x299 window to ensure generating square images during the rendering process. To this
278 end, we use a dataset consisting of four files per object (20 minutes in total). Then, an independent
279 rendering process has been performed for each of these files matching with our haptic dataset. In each
280 scene, the object of study is initially positioned in the world origin (0,0,0), which is the central point of

281 view of a camera that remains static during the whole process. This camera position is the same for all the
282 experiments.

283 Once the environment is set up, each object is texturized uniformly generating a UV Map with the same
284 colors on each side of the shape. Using this approach, each side can be easily identified despite the
285 symmetries present in all the shapes. Then, the data from the gyroscope collected during objects' human
286 manipulation is used to perform rotations on the object matching the ones from the human manipulation
287 experiment. Each haptic state has associated four values corresponding to each component of a quaternion
288 that are used to perform each rotation, considering the center of the figure the origin of the rotation. This
289 approach recreates the objects' original movements since the gyroscope in the 3D printed shapes is placed
290 inside the object on its center.

291 This rotation is performed for each haptic state and rendered from the world camera generating images
292 that capture each current object rotation forming the visual dataset. This synchronization between the
293 visual and the haptic dataset is what makes the experiments presented in this article possible.

294 **Blurred images generation**

295 In order to simulate a sudden loss of visual channel, we generate blurred images for all samples. These
296 blurred images are generated by applying over the original images an average filter, convoluting the
297 image with a 20x20 normalized box filter from OpenCV libraries.

298 **Visual object recognition system and training process**

299 In order to generate the visual recognition system, we follow a one-vs-all strategy. This strategy generates
300 an independent visual classifier for each one of the six classes ($Model_{[c]}$ where c defines the class).
301 $Model_{[c]}$ classifies a single image as belonging to class c or not. We adapt a pre-trained CNN model
302 based on InceptionV3 architecture (25), and we change the last layer for a fully connected dense layer
303 with 2 outputs, using a softmax activation function. Each of the $Model_{[c]}$ classifiers is trained using
304 categorical cross-entropy as loss function and a dropout of 0.4. The training is performed for 5 epochs
305 using 60/20 split of our own dataset for training/test the models (the remaining 20 percent is reserved for
306 the blur filter test). Samples of our dataset are split in two classes for each of the $Model_{[c]}$: i) positive
307 class: samples belonging to class c , and ii) negative class: samples that not belong to class c . Pre-trained
308 initial models are initialized using the weights of using ImageNet dataset (26).

309 **Visual retraining process**

310 By using the remaining 20% of samples that were not used in the model training/testing process (which
311 amount to around 60 seconds), each of the $Model_{[c]}$ classifiers is retrained using exactly the same
312 parameters used for training the original $Model_{[c]}$. The difference now is that the initial CNN weights are
313 the ones obtained after the previous training process. This retraining process is performed using one batch
314 (8 seconds of samples) at a time, and is repeated sequentially up to 7 times, as the samples are grouped in
315 7 batches. This retraining process simulates the gradual adaptation of the visual model to new conditions,
316 that in this work we simulate by a sudden loss of vision resulting in an input of blurred images. By using
317 the presented gradual retraining process, the model adapts to these new visual conditions for the blurred
318 vision scenario.

319 **Molyneux mechanism and filtering**

320 The haptic and visual classifiers constantly classify the haptic states and visual frames of the figures that
321 are acquired at a 40 Hz frequency. This means that every 25 ms, by applying the Molyneux mechanism,
322 the haptic classification is compared with the visual classification, checking if i) both channels are in
323 agreement, and ii) there is a failure in one of the channels or not.

324 It is normal that some short duration mismatches between both channels appear although both channels
325 are working properly. In order to provide a stable decision regardless of whether a channel is failing or
326 not, a low pass filter is applied to the output of the Molyneux mechanism.

327 The low pass filter consists in a 6th order Butterworth filter offering a flat output for the passband
328 frequencies and avoiding ripples. The first 1000 samples of each figure test file are used to study the
329 duration of the mismatches when there is no failure in neither of the two channels, obtaining a mean
330 duration of $\mu = 3.3$ samples and a standard deviation of $\sigma = 9.1$. Since it is desired to achieve a huge
331 attenuation for the mismatches frequencies, the cutoff frequency is set a decade before the frequency
332 corresponding to the $\mu + \sigma$ duration. Taking into account the sample rate of 40 Hz, the cutoff frequency
333 can be calculated as:

$$334 \quad f_{cutoff} = \frac{f_s}{10 \cdot (\mu + \sigma)} = 0.32 \text{ Hz}$$

335 After applying the filter, the output of the Molyneux mechanism fluctuates between 0 and 1 (see the
336 orange plot in Fig. 4), where 0 corresponds to channels not matching (failure in one channel) and 1 to
337 channels matching. In order to offer a binary response, as the one show in Fig. 4, the filter output goes
338 through an hysteresis cycle, where the output goes from 0 to 1 if the input is higher than 0.8 and from 1 to
339 0 if the input is lower than 0.2.

340

341 **References and Notes**

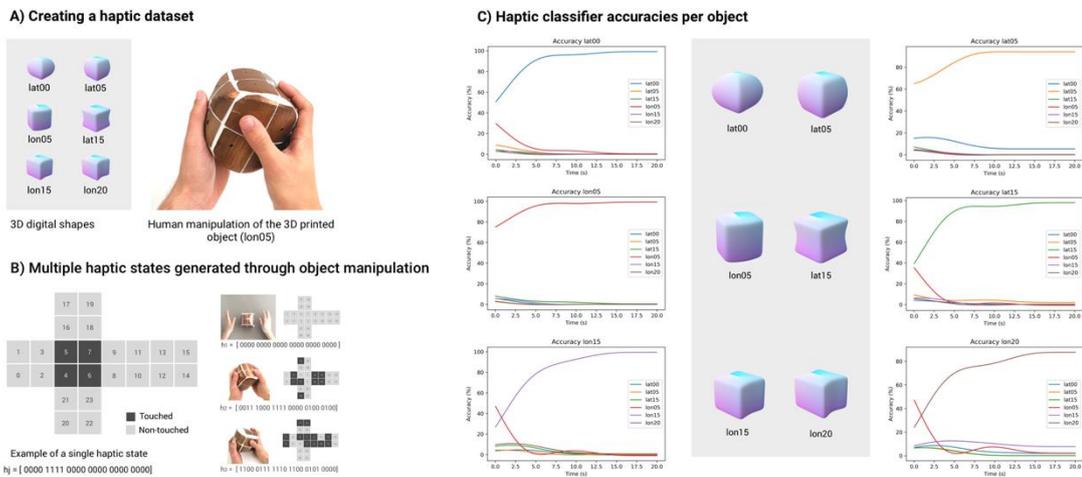
342

343

- 344 1. Purpura, G., Cioni, G., & Tinelli, F. Development of visuo-haptic transfer for object recognition
345 in typical preschool and school-aged children, *Child Neuropsychology* **24**, 657 (2018).
- 346 2. Norman, J. F., et al. Aging and the visual, haptic, and cross-modal perception of natural object
347 shape, *Perception*, **35**(10), 1383-1395 (2006)
- 348 3. Gibson, J. J. *The ecological approach to visual perception* (Taylor and Francis, 1979).
- 349 4. Gepshtein, S., & Banks, M. S. Viewing geometry determines how vision and haptics combine in
350 size perception. *Current Biology*, **13**(6), 483-488 (2003)
- 351 5. Ernst, M. O., & Bühlhoff, H. H. Merging the senses into a robust percept. *Trends in cognitive*
352 *sciences*, **8**(4), 162-169 (2004)
- 353 6. Lederman, S. J., & Klatzky, R. L. Haptic perception: A tutorial, *Attention, Perception, &*
354 *Psychophysics* **71**, 1439 (2009).
- 355 7. Locke, J. *An essay concerning human understanding*, (Hackett Publishing Company, 1996,
356 original 1689).
- 357 8. Cheselden, W. VII. An account of some observations made by a young gentleman, who was
358 born blind, or lost his sight so early, that he had no remembrance of ever having seen, and was
359 couch'd between 13 and 14 Years of age. *Philosophical Transactions of the Royal Society of*
360 *London*, **35**(402), 447-450 (1728)
- 361 9. Held, R. et al. The newly sighted fail to match seen with felt, *Nature neuroscience* **14**, 551
362 (2011).
- 363 10. Lewkowicz, D. J., & Lickliter, R. (Eds.). *The development of intersensory perception:*
364 *Comparative perspectives*. Psychology Press. (2013)
- 365 11. Carducci, P., Squillace, V., Manzi, G., & Truppa, V. Touch improves visual discrimination of
366 object features in capuchin monkeys (*Sapajus* spp.). *Behavioural processes*, **172**, 104044 (2020)
- 367 12. Hu, X., Urhie, O., Chang, K., Hostetler, R., & Agmon, A. A novel method for training mice in
368 visuo-tactile 3-D object discrimination and recognition. *Frontiers in behavioral*
369 *neuroscience*, **12**, 274 (2018)
- 370 13. Solvi, C., Al-Khudhairy, S. G., & Chittka, L. Bumble bees display cross-modal object
371 recognition between visual and tactile senses. *Science*, **367**(6480), 910-912 (2020)
- 372 14. Billard, A. & Kragic, D. Trends and challenges in robot manipulation, *Science* **364** (2019).
- 373 15. Lin, J., Calandra, R., & Levine, S. Learning to identify object instances by touch: Tactile
374 recognition via multimodal matching. *IEEE International Conference on Robotics and*
375 *Automation*, 3644-3650 (2019)
- 376 16. Falco, P., Lu, S., Natale, C., Pirozzi, S., & Lee, D. Connecting touch and vision via cross-modal
377 prediction, *IEEE Transactions on Robotics* **35**, 987 (2019).

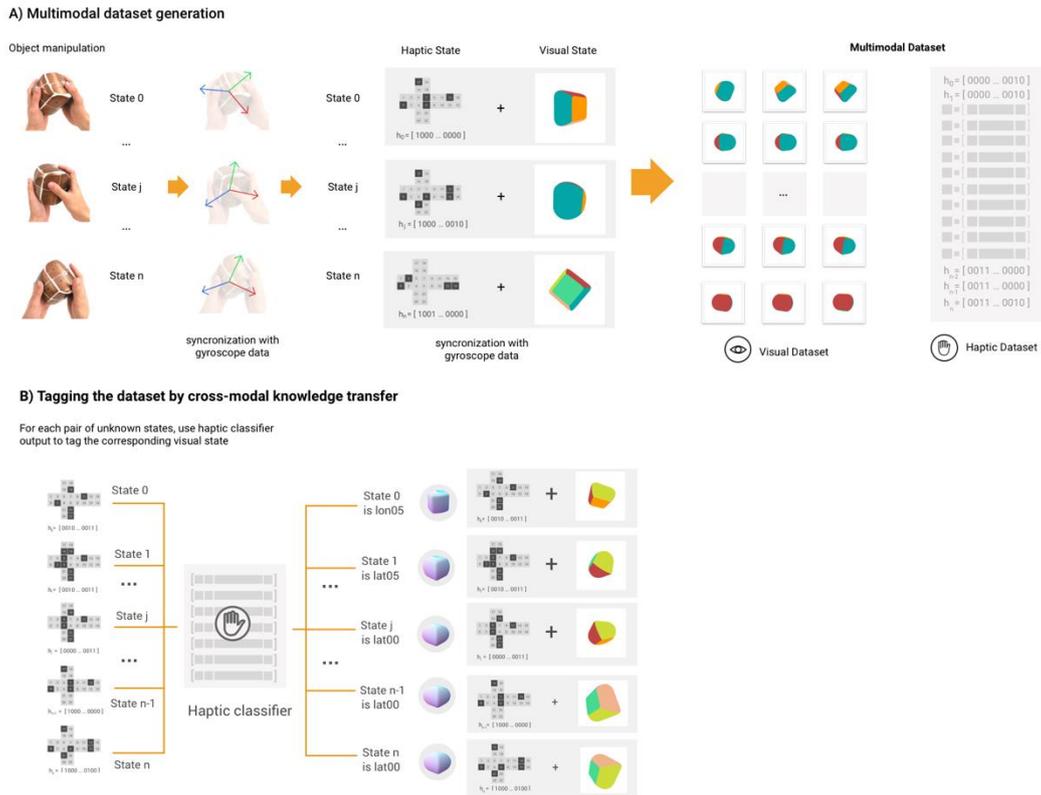
378 17. Lee, M. A., Tan, M., Zhu, Y., & Bohg, J. Detect, Reject, Correct: Crossmodal Compensation of
379 Corrupted Sensors. Preprint at <https://arxiv.org/abs/2012.00201> (2020)
380 18. Li, Y., Zhu, J. Y., Tedrake, R. & Torralba, A. Connecting touch and vision via cross-modal
381 prediction, *IEEE Conference on Computer Vision and Pattern Recognition* 10609-10618 (2019).
382 19. Sundaram, S. et al. Learning the signatures of the human grasp using a scalable tactile
383 glove. *Nature*, **569**(7758), 698-702 (2019)
384 20. Miralles, D. et al. Artificial haptic recognition through human manipulation of objects,
385 *Conference on Cognitive Computational Neuroscience*, (2019).
386 21. Gielis, J. A generic geometric transformation that unifies a wide range of natural and abstract
387 shapes, *American journal of botany*, **90**, 333 (2003).
388 22. Walker R., Shadow dexterous hand technical specification, *Shadow Robot Company* (2005).
389 23. Andrychowicz, O. M., et al., Learning dexterous in-hand manipulation, *The International*
390 *Journal of Robotics Research* **39**, 3 (2020).
391 24. Masson, H. L., Bulthé, J., De Beeck, H. P. O., & Wallraven C., Visual and haptic shape
392 processing in the human brain: unisensory processing, multisensory convergence, and top-down
393 influences, *Cerebral Cortex* **26**, 3402 (2016).
394 25. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J. & Wojna, Z. Rethinking the Inception
395 Architecture for Computer Vision, *IEEE Conference on Computer Vision and Pattern*
396 *Recognition*, 2818–2826 (2016)
397 26. Deng, W., et al. ImageNet: A Large-Scale Hierarchical Image Database, *IEEE Conference on*
398 *Computer Vision and Pattern Recognition*, 248–255 (2009)

403 **Figures and Tables**



405 **Fig. 1. From object human manipulation to automatic haptic classification.** (A) The objects of this
406 study are six 3D printed shapes labelled as lat00, lat05, lat15, lon05, lon15, and lon20. The participant sits
407 in front of a computer and follows instructions on how to manipulate the objects randomly. (B) The
408 human manipulation data of the objects are collected by the system and are stored as haptic states. A
409 haptic state is represented by a 24-bit array and indicates the status of each sensor (touch/not touched)
410 forty times per second. (C) ACS recognizes each object through a simple Bayes algorithm based on
411 haptic states. As show in this subfigure, accuracies higher than 80% are consistently reached after a few
412 seconds.

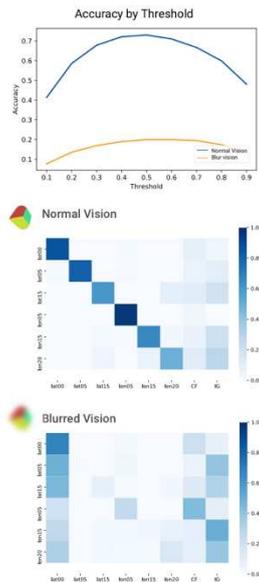
413
414
415



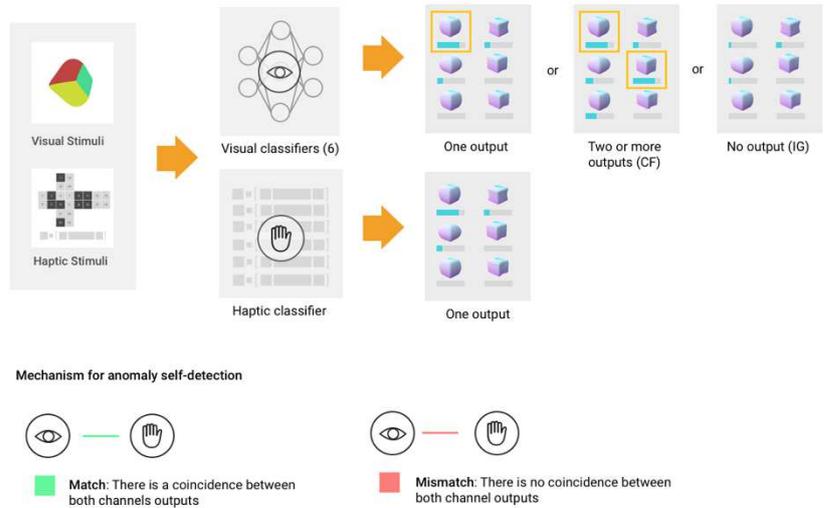
416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426

Fig. 2. Multimodal dataset generation through cross-modality. (A) Using the gyroscope, we can obtain the orientation of the object associated to each haptic state. From this information, we can draw (through 3D renders) a visual state of the object and associate it with its corresponding haptic state. From this, given any haptic dataset we can generate its corresponding visual dataset. With this method a new multi-modal dataset is created as a result of the combination of haptic and visual datasets. (B) Given that no manual annotation process has been carried out in the creation of the previous multi-modal dataset, the visual mode is not labelled. Here, we propose to tag the visual dataset from the results of the haptic object recognition system. This transfer of information generates a link between the two modes, visual and haptic, called cross-modality.

A) Visual classifier accuracies



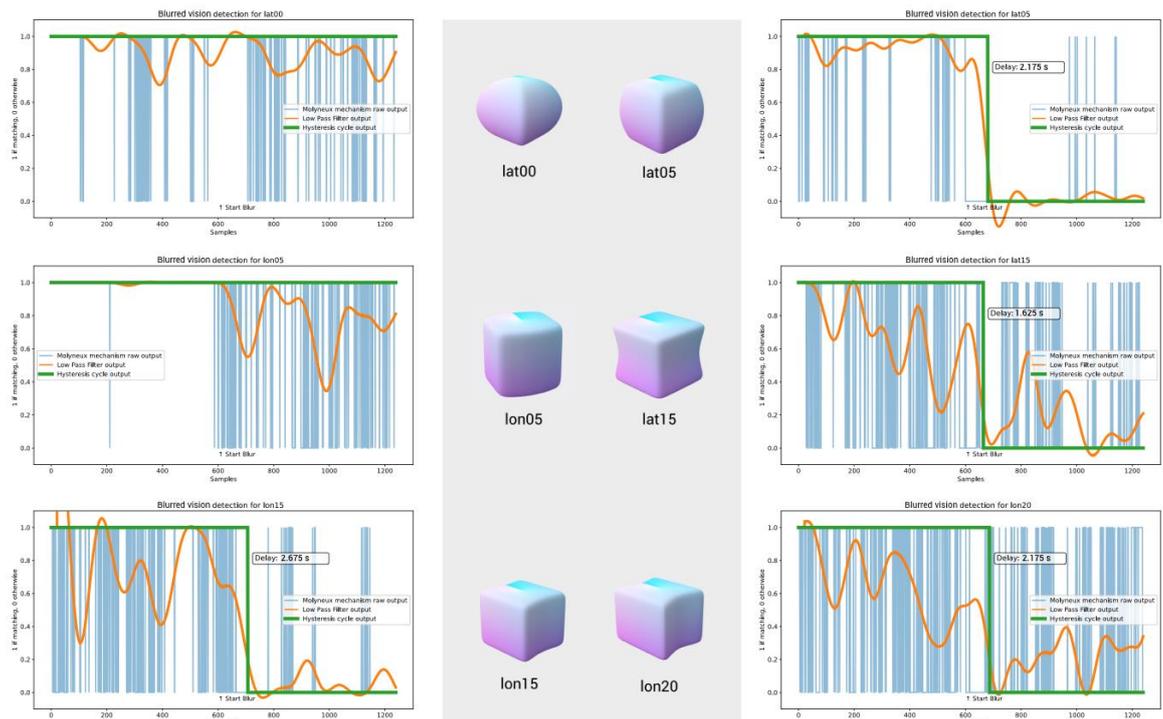
B) Molyneux mechanism in Robotic Cognitive System



427
428
429
430
431
432
433
434
435
436
437
438
439
440

Fig. 3. Visual classifier threshold and the Molyneux mechanism. (A) With the synthesized video generated from human object manipulation, we trained a one-vs-all based CNN model, generating an individual classifier for each class. In order to determine the output label of a classifier, we have studied which accuracy threshold results in the best performance of the visual model. By setting the threshold to 0.5, we obtain the confusion results in the best performance of the visual model. By setting the threshold to 0.5, we obtain the confusion matrices shown at the bottom after testing the model with normal vision images (like the ones used for training) and blurred vision images. (B) Every video frame and its corresponding haptic state are classified by the visual and haptic object recognition systems, respectively. The output of the visual recognition system can be i) a single class, ii) multiple classes (CF), or iii) none (IG). On the other hand, the label from the haptic recognition system is always univocal. By applying the Molyneux mechanism, we check if visual and haptic recognition agree, which we call a match. Otherwise, we call this disagreement a mismatch. In the case the visual recognition result is CF (ii), it will be a match if one of the possible classes agrees with the haptic classification.

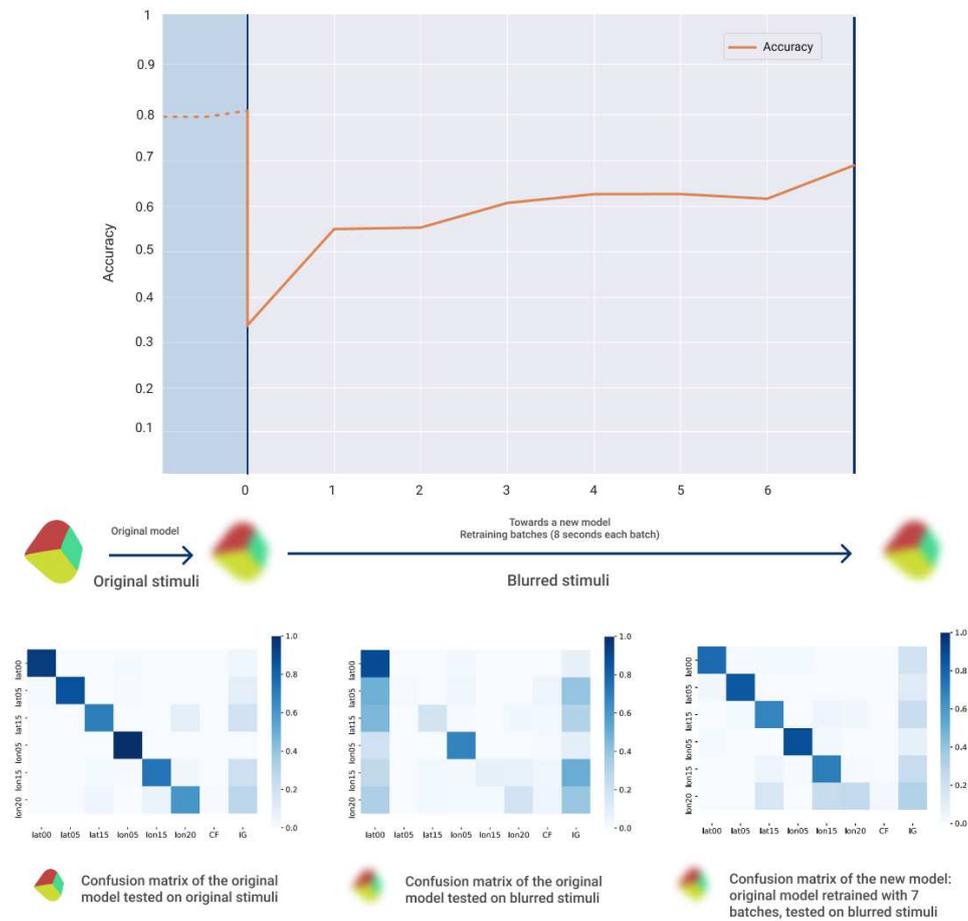
Blurred vision detection



441
442
443
444
445
446
447
448
449
450
451
452
453

Fig. 4. Blurred vision detection through Molyneux mechanism and filtering. By applying the Molyneux mechanism to every visual-haptic pair, the ACS can determine if the two channels are matching or not, hence it can detect failures in one channel. Although there is a clear trend towards the two channels matching when both are working properly, short mismatches can appear (blue plot). In order to obtain an accurate decision whether the channels are matching or not, a low pass filter is applied to attenuate these mismatches (orange plot). Finally, as a way to offer a binary response that decides if the vision is blurred or not, the filter output goes through a hysteresis cycle (green plot). With this entire process the detection of the blurred vision is not immediate and has an average delay of 2.13 seconds (85.33 samples). Note that there is no mismatch in lat00 and lon05, that is, for these two objects the blur filter does not affect the recognition system.

Accuracy evolution during self-adaptive process



454

455

456 **Fig. 5. Visual model adaptation to new stimuli acquisition conditions.** Top plot shows the decrease in

457 accuracy at the moment ($t = 0$) when blurred stimuli are first introduced. As blurred batches are

458 incorporated to retrain the system, previous accuracy with original stimuli is now achieved with blurred

459 stimuli ($t=7$). Bottom plots show the confusion matrices at the most relevant moments during this

460 adaptation process: (Left) Original model tested with original stimuli, (Center) Original model tested with

461 blurred stimuli ($t=0$), (Right) New model after the adaptation process with 7 batches of blurred images

462 ($t=7$) tested with blurred stimuli.