

# Multimodal Representation Learning for Place Recognition using Deep Hebbian Predictive Coding

Martin J. Pearson,<sup>1\*</sup> Shirin Dora,<sup>2 5</sup> Oliver Struckmeier,<sup>3</sup> Thomas C. Knowles,<sup>1</sup>  
Ben Mitchinson,<sup>4</sup> Kshitij Tiwari,<sup>3</sup> Ville Kyrki,<sup>3</sup> Sander Bohte,<sup>5 6</sup> and Cyriel  
M.A. Pennartz<sup>6</sup>

<sup>1</sup>*Bristol Robotics Laboratory, University of The West England Bristol, UK*

<sup>2</sup>*Department of Computer Science, Loughborough University, UK*

<sup>3</sup>*Intelligent Robotics Group, Aalto University, Finland*

<sup>4</sup>*Dept. of Computer Science, University of Sheffield, UK*

<sup>5</sup>*Center for Mathematics and Informatics, Netherlands*

<sup>6</sup>*Dept. of Cognitive and Systems Neuroscience, University of Amsterdam, Netherlands*

Correspondence\*:

Martin J. Pearson

[martin.pearson@brl.ac.uk](mailto:martin.pearson@brl.ac.uk)

## 2 ABSTRACT

3

4 Recognising familiar places is a competence required in many engineering applications that  
5 interact with the real world such as robot navigation. Combining information from different sensory  
6 sources promotes robustness and accuracy of place recognition. However, mismatch in data  
7 registration, dimensionality, and timing between modalities remain challenging problems in  
8 multisensory place recognition. Spurious data generated by sensor drop-out in multisensory  
9 environments is particularly problematic and often resolved through adhoc and brittle solutions.  
10 An effective approach to these problems is demonstrated by animals as they gracefully move  
11 through the world. Therefore, we take a neuro-ethological approach by adopting self-supervised  
12 representation learning based on a neuroscientific model of visual cortex known as predictive  
13 coding. We demonstrate how this parsimonious network algorithm which is trained using a local  
14 learning rule can be extended to combine visual and tactile sensory cues from a biomimetic robot  
15 as it naturally explores a visually aliased environment. The place recognition performance  
16 obtained using joint latent representations generated by the network is significantly better  
17 than contemporary representation learning techniques. Further, we see evidence of improved  
18 robustness at place recognition in face of unimodal sensor drop-out. The proposed multimodal  
19 deep predictive coding algorithm presented is also linearly extensible to accommodate more than  
20 two sensory modalities, thereby providing an intriguing example of the value of neuro-biologically  
21 plausible representation learning for multimodal navigation.

22 **Keywords:** Predictive coding, Multisensory integration, Place recognition, Sensory reconstruction, Whisker tactile

## 1 INTRODUCTION

23 The study of biology and the brain has inspired many innovative and robust solutions to hard problems  
24 in engineering. Biologically inspired machine learning has great potential for robotics and automation  
25 Sunderhauf et al. (2018) with significant progress being made in perception Giusti et al. (2016); Eitel  
26 et al. (2015) and scene understanding Eslami et al. (2018); Badrinarayanan et al. (2017); Gu et al. (2019);  
27 Sheppard and Rahnemoonfar (2017). Supervised deep learning takes biological inspiration from layered  
28 neural connectivity, synaptic plasticity, and distributed computation to learn non-linear mappings between  
29 inputs and desired outputs. These approaches usually rely on biologically implausible learning principles.  
30 Closer to neurobiology, deep reinforcement learning also leverages these bio-inspired architectural  
31 properties but instead learns against a task specific cost function Mnih et al. (2015). Both approaches  
32 require an error signal that is either back propagated or otherwise distributed through the layered network  
33 weight space during training. Unsupervised learning in neural networks does not typically require a globally  
34 distributed error signal for training, instead they find and exacerbate patterns in the input space by learning  
35 correlations or through local competition, typically to enable a useful reduction in dimensionality. These  
36 low dimensional latent representations of input are often used to perform clustering of complex data  
37 or serve as efficient pre-processing for a supervised or reinforcement learning back-end. All of these  
38 approaches to machine learning adopt the same assumed flow of information through the network, namely,  
39 from sensory input toward an appropriate output representation. The flow of information in a Predictive  
40 Coding Network (PCN) is both from sensory input to output and the opposite, i.e., each layer in the  
41 network predicts representations of the previous layer in parallel, ultimately predicting the actual input  
42 being passed into the lowest layer of the network Rao and Ballard (1999); Spratling (2017). Prior to  
43 learning layer-wise predictions are randomly initialised, during weight learning and inference they are  
44 compared to the predictions received from the previous layers. Local learning rules are then applied to  
45 update weights and infer neural activity in each layer to minimize the error in predictions (which is related  
46 to the free-energy or ‘surprise’ in the system Friston (2010)) on subsequent exposures to similar input.  
47 This approach to learning is more biologically plausible as a globally distributed error is not required to  
48 update the weights, instead local Hebbian-like rules are applied Dora et al. (2018). Moreover, predictive  
49 coding is also an unsupervised learning approach and, hence, does not require labelled datasets for training.  
50 It has been used in robotics for learning sensorimotor models Park et al. (2012); Nagai (2019); Lanillos  
51 and Cheng (2018) and for goal directed planning in visuomotor tasks Hwang et al. (2017); Choi and Tani  
52 (2018). To the best of our knowledge, it has not been used for place recognition. Place recognition is  
53 the ability to interpret and recall sensory views of the world to inform an estimate of location or pose,  
54 with visual place recognition being its most common and well-studied form Lowry et al. (2016). A place  
55 recognition system is typically decomposed into two sub-systems, an image/sensory processing module,  
56 and a mapping module which stores either a metric or relative association between sensory views and the  
57 pose of an agent. In this study we primarily focus on sensory pre-processing for place recognition, our  
58 interest being in the performance of PCNs to transform samples from co-localised but disparate sensory  
59 modalities into a representation suitable for efficiently determining proximity between locations in an  
60 environment. This is interesting for two reasons, firstly, PCN is a mechanistic implementation of the  
61 parsimonious theory that perception arises from generative, inferential representations of what causes  
62 sensory inputs to arise Helmholtz (1867); Gregory (1980); Mumford (1992); Pennartz (2015); Friston  
63 (2010). This framework describes many of the phenomenological and anatomical observations from animal  
64 behaviour and neurophysiology. By incorporating this model into robots we have the opportunity to  
65 reproduce these observations in an autonomous agent and thus better understand the principles at work  
66 in the brain from an algorithmic level of abstraction. Secondly, combining information from different

67 sensory modalities in mobile robots overcomes unimodal aliasing and sensor drop-out but introduces new  
68 challenges such as dimensionality mismatch and registration Khaleghi et al. (2013). Sensor fusion is a  
69 well-established field of research with numerous approaches developed and successfully incorporated  
70 into widespread use. The three predominant approaches of sensor fusion are probabilistic, evidential,  
71 and model-free/neural networks, with Kalman filtering Roumeliotis and Bekey (1997), Dempster-Shafer  
72 Theory Murphy (1998), and Variational Auto-Encoders (VAEs) Kingma and Welling (2013); Suzuki et al.  
73 (2017); Korthals et al. (2019) being prominent examples of the respective approaches. Auto-Encoders are  
74 well established tools in machine learning for approximating a higher dimensional input space using a  
75 lower dimensional representation space. A VAE is a generative modelling approach that uses variational  
76 inference methods for training with large-scale and high dimensional data sets Kingma and Welling (2013).  
77 More recently, this has been extended for learning bi-directional, joint distributions between different  
78 sensory modalities Suzuki et al. (2017); Korthals et al. (2019). This allows inferences in one sensory  
79 modality based on evidence in another modality via a jointly trained generative model. Suzuki et al.  
80 Suzuki et al. (2017) demonstrated that visual images and textual labels could be associated using a Joint  
81 Multimodal VAE (JMVAE) such that either modality could be used to reconstruct meaningful inferences  
82 about the other modality. Korthals et al. Korthals et al. (2019) used a modified version of JMVAE to  
83 jointly infer coloured geometric objects from visual and LiDAR data gathered from a simulated robot.  
84 Here we introduce a Multimodal Predictive Coding Network model (MultiPredNet) which is rooted in  
85 neurobiology and psychology, and utilise JMVAEs with the visual tactile datasets gathered in this study  
86 as a contemporary machine learning approach for comparison. The MultiPredNet presented here also fits  
87 the *model-free* learning category of sensory fusion as, similar to VAEs, it requires a period of training  
88 before it can be reliably used. Both JMVAEs and MultiPredNet learn the structural properties of the  
89 sensory information pertaining to the environment in which they are trained, i.e., they are computationally  
90 equivalent but differ in their algorithmic approach, including the learning rules. In this study we train both  
91 a contemporary JMVAE and the novel MultiPredNet using visual and tactile sensory data sets sampled  
92 from a biomimetic robot as it explores the world. The robot head has been physically designed to mimic  
93 a whiskered rat, with an array of individually actuated tactile whisker sensors and wide angle cameras  
94 in place of eyes. The dimensionality, timing and registration, or reference frame, of these two sensory  
95 modalities are different, with the salient information available from each being dependent on the current  
96 pose of the robot. Intuitively, combining cues from both sensory systems should reduce ambiguity in  
97 place recognition which ultimately will result in less frequent incorrect re-localisations from a robot pose  
98 mapping module. We extend this further by testing the ability of JMVAE and MultiPredNet to generalise  
99 between poses through the representation space itself; in other words, subsuming some of the functionality  
100 of a mapping module by the sensory preprocessing. To demonstrate the extensibility and applicability of  
101 MultiPredNet more explicitly to the robot localisation problem we use a simulation of the robot within a  
102 larger scale environment to reveal examples of loop-closure detection through its integration with a simple  
103 associative memory. Finally, we compare separate data sets that cover similar regions of the pose space  
104 to apply more conventional precision-recall curve analysis as a second measure of performance for place  
105 recognition by MultiPredNet.

## 2 RESULTS

### 106 2.1 Experimental procedure

#### 107 2.1.1 Data capture

108 The study was conducted using a custom built robotic platform called “WhiskEye”, modelled on previous  
109 whiskered mobile robots developed in collaboration with biological scientists Prescott et al. (2009);

110 Pearson et al. (2013). The body of WhiskEye is a Robotino<sup>TM</sup> mobile platform from Festo Didactic with  
111 an additional embedded computer and a 3 degree of freedom neck installed as shown in figure 1. The  
112 head, which was mounted as the end-effector of the neck, has 24 individually actuated artificial tactile  
113 whisker sensors and two forward facing cameras (eyes). The embedded computer collected all sensory  
114 data and coordinated the motor action of the platform using the Robot Operating System (ROS) execution  
115 framework Stanford Artificial Intelligence Laboratory et al. (2018). The actions of the robot were directed  
116 and controlled using a model of tactile attention distributed across functional models of distinct regions  
117 of the brain. In brief, the collective behaviour of this model was to direct the nose of WhiskEye toward  
118 the most salient region of a head centred map of space representing the volume surrounding WhiskEye's  
119 head Mitchinson and Prescott (2013). WhiskEye's whiskers, which occupy this space, can be actively  
120 rotated around their base mimicking the cyclic whisking behaviour expressed by many small mammals  
121 such as mice and rats Gao et al. (2001). If a whisker makes contact with a solid object during a whisk  
122 cycle then the sensory consequences of that collision will be interpreted as a more salient location in the  
123 head centred map, thus increasing the likelihood of the robot orienting its nose toward the point of contact.  
124 An orient is enacted when the saliency of a point in the map exceeds a certain level; this can be excited  
125 by whisker contacts or through a random background noise pattern that increases in relation to the time  
126 since a previous orient. In this way the robot moves through the world through a sequence of regular  
127 orients whilst preferentially attending to objects that it encounters with its whiskers. To prevent repetitive  
128 orienting behaviour a mechanism based on visual inhibition of return was included to temporarily suppress  
129 the salience of regions in the map that have recently been explored. In addition, there was a low level  
130 reflexive behaviour built into the whisker motion controllers such that the drive force to each whisker was  
131 inhibited by deflection of the shaft. This reduces the likelihood of damage and constrains the magnitude of  
132 deflections measured by the whisker, effectively normalising the sensor range Pearson et al. (2013). The  
133 data sets that were collected for this study were composed of samples taken from the whisker array and the  
134 forward facing cameras at the point of peak protraction of the whisker array. Alongside these data were  
135 stored the robot odometry, motor commands, and ground truth 2D location and orientation of WhiskEye's  
136 head as determined by an overhead camera and associated robot mounted markers. For longer duration and  
137 larger scale experiments a simulation of the WhiskEye platform was instantiated into an on-line robotics  
138 simulator called the NeuroRobotics Platform (NRP) Falotico et al. (2017). The interface with the NRP is  
139 based on ROS which enabled the simulated WhiskEye robot to use the same control software and capture  
140 the same format of sensory data as the physical platform. Gaussian noise was added to the simulated tactile  
141 sensory responses to match statistics from the physical whiskers and the resolution of the visual frames  
142 captured from the simulated cameras were scaled appropriately to match.

143

144 The physical environment was bounded by a 600mm high ellipsoid shaped perimeter measuring 3m  
145 by 5m which in turn was bounded by 1.5m high blue partition boards on a smooth grey painted concrete  
146 floor to minimise distinct visual cues. Within the bounded arena we placed black coloured 600mm high  
147 boxes and cylinders in various configurations to delineate different environments for gathering training  
148 and test data sets. The training set was gathered in batches as the robot explored the training arena  
149 which was then concatenated together into 1270 visual-tactile data points representing 30 minutes of real  
150 time exploration (nominal whisk rate 1Hz). Test sets were gathered from arenas composed of differently  
151 configured geometric shapes (Figure 1b) in batches of 73 samples for each set. The trajectory of the robot  
152 was governed by the attention driven model of control and therefore not repeatable between test runs.  
153 However, the robot did adopt similar poses at multiple points between the test data sets and training set as  
154 shown in the quiver plots of figure 2. The simulated environment consisted of 4 interconnected quadrants

155 each the same size as the physical arena but with alternate black and white walls and different coloured  
156 and configured cylinders and boxes in each quadrant. A training set of 2400 visual-tactile samples and  
157 6 test sets of 400 samples each were captured as the simulated robot explored in different regions of the  
158 environment (see figure 3). This simulated arena serves as a controlled intermediate step toward larger  
159 scale unstructured environments to systematically test the efficacy of MultiPredNet for robot localisation.  
160 Specifically, we used the simulator here to perform longer duration experiments in order to capture loop  
161 closure events between data sets as is clear in the quiver plots of figure 3.

162

### 163 2.1.2 Comparative network model architecture

164 The network structure and learning rules for the proposed MultiPredNet and the VAE are presented in  
165 the materials and methods section. Briefly, the MultiPredNet consists of 3 modules such that one module  
166 (called the visual module) receives visual data as input, the second module (called the tactile module)  
167 receives tactile data, and the third module (called multisensory module) receives the concatenated higher  
168 order representations inferred by the visual and tactile modules (see figure 4). The synaptic weights of the  
169 three modules are learned using the same Hebbian-like learning rule. The representation inferred from  
170 the last layer of the multisensory module denotes the joint representation inferred using MultiPredNet.  
171 We compared the place recognition performance of MultiPredNet with existing VAE approaches for  
172 inferring multisensory representations, namely Joint Multimodal VAEs (JMVAEs) or more specifically  
173 a JMVAE-zero and JMVAE-kl Suzuki et al. (2017) as shown in figure 14. JMVAE-zero consists of two  
174 VAEs for handling visual and tactile inputs respectively. The last layer of the encoders in both VAEs is  
175 connected to a common layer whose activities are used for multisensory place recognition. JMVAE-kl  
176 uses the same network architecture as used in JMVAE-zero with two additional VAE encoders that infer  
177 unisensory representations based on visual and tactile inputs respectively. The multisensory and unisensory  
178 representations in JMVAE-kl are optimised together to be similar to each other using a Kullback-Leibler  
179 divergence component in the objective function. This allows JMVAE-kl to generate better crossmodal  
180 reconstruction in case of sensor drop-out. For a fair comparison, the dimensions of the multisensory  
181 representations obtained from MultiPredNet, JMVAE-zero and JMVAE-kl were fixed at 100.

182

### 183 2.1.3 Performance metrics

184 To quantitatively measure the performance of a place recognition system we need to relate the ground truth  
185 3D pose of the robot head ( $(x, y, \theta_{head})$  relative to a global reference frame) to the representations generated  
186 by the sensory pre-processing modules. It leads, therefore, that to perform efficient place recognition the  
187 similarity between representations should correlate with similarity between robot poses. Here we adopt  
188 a technique from computational neuroscience called Representational Similarity Analysis (RSA) that  
189 was originally developed to compare measurements from brain activity, behavioural measurements and  
190 computational modelling Kriegeskorte et al. (2008). For the current study, we computed a Representational  
191 Dissimilarity Matrix (RDM) for both pose and the generated representations from candidate systems for  
192 each run of the robot in the testing arena and compared their rank order using Spearman's rank correlation.  
193 Briefly, an RDM is a symmetric matrix around a diagonal of zeros with each element encoding the dis-  
194 similarity of the row sample to the column sample, i.e., the distance from each sample in the data set  
195 to all the others. Comparing the rank order of the RDMs for ground truth pose against representations  
196 provides an intuitive measure of performance for place recognition. A second measure of performance  
197 in this study was the error in inferring the tactile (visual) modality based on the representations inferred  
198 from the MultiPredNet and JMVAEs in presence of visual (tactile) modality (i.e., by sensor drop-out). This  
199 experiment was carried out using the physical test data only. To evaluate the extensibility of MultiPredNet

200 and the validity of RSA as a measure of performance for place recognition in larger scale, on-field settings  
201 we use the simulated data sets and a simple mapping system to reveal loop closure recognition as a  
202 qualitative demonstration of its performance in robot localisation. This is evaluated further using more  
203 conventional ROC curve analysis between 2 of the simulated testsets that have similar but not identical  
204 ground truth trajectories.

205

## 206 **2.2 Model performance at place recognition**

207 The three models, namely MultiPredNet, JMVAE-zero and JMVAE-kl, were trained on the same physical  
208 data set which was shuffled and decomposed into mini-batches of 10 samples. Each model was trained for  
209 200 epochs. Once trained, the physical test sets were presented to each model one sample at a time to infer  
210 corresponding sets of joint latent representations. These were used to estimate RDMs for corresponding  
211 models which were used to compute the rank order with respect to the RDM of corresponding poses of the  
212 robot from each of the 4 test sets (see Figure 2 for ground truth poses). Figure 5 contains example RDMs  
213 displayed as heatmaps from typical instances of each model validated against test set 1. The boxplots  
214 summarise the statistics of the Spearman’s rank correlation coefficient ( $\rho$ ) calculated for each sample in  
215 the test set, with ( $p < 0.001$ ,  $N = 73$ ) indicated by the horizontal green line. For control, a random set of  
216 representations was also compared to reveal the structural relationship that each model has found between  
217 pose and multisensory view.

218 Using the same analysis across all test data sets,  $N = 292$ , the average  $\rho$  and the percentage of samples  
219 that scored above statistical significance ( $p < 0.001$ ,  $N = 292$ ) were (0.289, 69.17%) for MultiPredNet,  
220 (0.141, 47.26%) for JMVAE-zero, and (0.140, 49.31%) for JMVAE-kl. Applied to place recognition, a true  
221 positive correlation between representation distance and pose distance will result in a correct re-localisation.  
222 Therefore, we can expect the frequency of true positive re-localisations generated by the MultiPredNet to  
223 be 20 – 22% higher than JMVAE.

224

225 We next trained a MultiPredNet model on the simulated training set (2400 samples) using the same  
226 network topology, batch size, learning rates and epochs as adopted for the physical data set model. For  
227 visual clarity, the RDMs for only the first 100 samples from each of the 6 simulated test sets are shown  
228 in figure 6, whilst the box plots summarise the statistics of  $\rho$  for all samples in each set ( $n = 400$ ). As  
229 in the physical tests the above significance positive correlation is clear (mean  $\rho$  beneath each boxplot  
230  $p < 0.001$ ,  $n = 400$ ), suggesting that MultiPredNet can infer structural relationships in the simulated  
231 sensory modalities appropriate for place recognition as in the physical demonstration.

## 232 **2.3 Model performance during sensor drop-out**

233 The models were also evaluated for place recognition in a sensor drop-out scenario. For this purpose, we  
234 evaluated the place recognition performance of the three models using either visual or tactile input with  
235 the other sensory modality set to zeros. All three models performed well at place recognition (average  
236  $p < 0.001$ ) across each physical test set when only visual sensory information was available (see Figure 7).  
237 However, with only tactile sensory information available the JMVAE-zero model could not maintain the  
238 positive correlation between representations and ground truth pose data above the significance threshold  
239 in the majority of cases. These results are presented in figure 7, which summarises the Spearman’s rank  
240 correlations for each model across all 4 test sets with both sensory modes available, only vision, and only  
241 tactile available. The line plots in the bottom panels track the cumulative number of samples that returned  
242 an above chance positive correlation ( $p < 0.001$ ) between pose and representation distance implying  
243 a positive contribution towards place recognition. In summary, the mean  $\rho$  and  $p < 0.001$  percentage

244 scores for each model in the two drop-out conditions, only vision available and only tactile available, were  
245 (0.294, 72.95%) & (0.279, 69.52%) for MultiPredNet, (0.138, 46.23%) & (0.036, 6.51%) for JMVAE-zero,  
246 and (0.131, 48.97%) & (0.126, 44.18%) for JMVAE-kl. These results indicate that the MultiPredNet model  
247 has the potential to correctly re-localise on average 25% more often than both the JMVAE based models in  
248 the absence of either tactile or visual cues ( $p < 0.001$ ).

249 As these models are generative in nature we can compare their ability to reconstruct the missing modality  
250 inferred from the conditioned evidence derived from the other. Indeed, the loss function used during  
251 training of the three models is computed using the reconstruction error. During training the JMVAE models  
252 generate sensory reconstructions by propagating the joint latent representation through a decoder network  
253 which is trained end-to-end with the encoder network using back-propagation. In contrast, all layers of  
254 the MultiPredNet model generate predictions about the activity of neurons in the previous layer of the  
255 network and the error in these predictions is used to update the weights during training according to a  
256 Hebbian-like learning rule Dora et al. (2018). In the absence of input in a given sensory modality, the  
257 joint latent representation inferred using a single modality is propagated backwards, by way of feedback  
258 projections to the input layer, to reconstruct the sensory input in the missing modality.

259 The Mean Squared Errors (MSE) for the tactile and visual reconstructions from each of the network  
260 models in the three sensory conditions are plotted in figure 8. The results from the JMVAE-kl model  
261 revealed that it had successfully accommodated a systematic positive off-set in the tactile reconstruction  
262 which both the JMVAE-zero and MultiPredNet had failed to (See figure 9). With this off-set removed, the  
263 tactile reconstruction errors from the JMVAE-kl model were still significantly lower than the other two  
264 ( $p < 0.001$ ). Another interesting observation was that the JMVAE-kl model performed worse (relative to  
265 the others) at visual reconstruction when only tactile sensory input was available. However, as with the  
266 MultiPredNet, it performed consistently at place recognition under this condition on which JMVAE-zero  
267 failed to maintain as shown in figure 7. This suggests that performance on place recognition (measured by  
268  $\rho$ ) and sensory reconstruction (as measured by MSE) are not correlated.

269

## 270 **2.4 MultiPredNet performance at robot localisation in simulated field trials**

271 To evaluate place recognition by MultiPredNet more explicitly, a simple memory module based on the  
272 view-cell memory of RatSLAM Milford et al. (2004) was adopted to associate poses with the joint modal  
273 representations inferred by the MultiPredNet at each sample step. The distance ( $1 - \text{Pearson correlation}$ )  
274 between the current representation and others already stored in the memory was calculated, if this was  
275 above a certain threshold (discussed below) then it is considered novel and a new view-cell is added to the  
276 memory containing the representation and associated pose. If the distance was below the threshold, i.e.,  
277 they were deemed similar, then a re-localisation event was registered and the representation not stored into  
278 memory. All 6 simulated test sets were concatenated together into a continuous run of 2400 samples and  
279 presented to the view-cell memory in sequence. The results shown in figure 12 demonstrate that similar  
280 poses are recalled from the view-cell memory triggered by similarity in representation. The black asterisks  
281 in the quiver plot highlight the sample points at which re-localisation events were detected by the view-cell  
282 memory, note that these occur during loop-closures within and between test sets. This is more clearly  
283 shown in the plot of the lower panel of figure 12 where the horizontal coloured panes indicate the regions  
284 of the view-cell memory that are composed of view-cells created during each test set in sequence (blue,  
285 red, green, magenta, cyan and yellow) with the coloured vertical lines marking the start of each region.  
286 The black dots indicate the view-cell index address associated with each sample in the concatenated data  
287 set, re-localisation events, therefore, are indicated by a sharp decrease in view-cell index between samples.

288 This is most evident during sets 2 and 3 (red and green) which include re-localisation events occurring  
289 during set 3 that reference view-cells created during set 2. Referring back to the quiver plot we can see  
290 that these relate to the loop-closures that occur in the shared pose space adopted by the robot during the  
291 acquisition of these 2 test sets. The same phenomenon is seen between sets 5 and 6 (cyan and yellow) and,  
292 to a lesser extent, between sets 5 and 1 (cyan and blue). Test set 4 (magenta) has multiple re-localisation  
293 events, however, these are confined to its own region of the view-cell memory which corresponds to the  
294 unique region of pose space that it represents.  
295

296 A more quantitative measure for the performance of a system at place recognition can be obtained  
297 through the analysis of the precision-recall rate Kazmi and Mertsching (2016), Flach and Kull (2015). To  
298 calculate this we selected test sets 2 and 3 from the simulator as we have seen that they approximately  
299 share the pose space within the arena during their independent runs with loop-closure events evident from  
300 the qualitative analysis described above. The distance between representation of each sample in each test  
301 set to each sample in the other was calculated again using  $1 - \text{Pearson correlation}$  as the metric of distance.  
302 The distance in pose between each sample in a set against the pose of all samples in the other set was  
303 also calculated (Euclidean). These inter test set distance matrices are displayed as heatmaps in figure 13  
304 allowing us to visualise which regions of both the pose and representation space are similar between the  
305 2 test sets. Intuitively, regions of low distance in representation space should correspond to an equally  
306 low distance in pose space to perform place recognition. Putting this into the context of the view-cell  
307 memory demonstration, a below threshold representation distance should trigger a re-localisation event  
308 from one test set to the other which should correspond to a similar pose. Therefore, for each sample in  
309 test set 2 (columns in heatmaps of figure 13) we select the sample from test set 3 (row) with the lowest  
310 distance in representation space as the candidate classification. If the distance of a classification is below a  
311 representation threshold we label it as *Positive*, if higher then it is *Negative*. The [column, row] coordinates  
312 of the candidate classifications are relayed to the pose space distance matrix to determine whether they  
313 were *True* or *False* classifications by comparing the pose distance to a threshold which we fixed arbitrarily  
314 at 0.2. To determine an appropriate representation threshold to maximise performance from this system  
315 we calculated the peak geometric mean of the Receiver Operating Characteristic (ROC) curve generated  
316 through a sweep of 1000 threshold values from 0.001 to 1 registering the classifications generated from  
317 each as *True Positive* (TP), *False Positive* (FP), *False Negative* (FN) or *True Negative* (TN) accordingly.  
318 The area under the ROC curve was found to be 0.836 with a peak geometric mean found at iteration 290  
319 indicating an optimal representation threshold of 0.29 to maximise the classification performance of the  
320 system. With this threshold the Precision (70.7%), Recall (91.4%) and F1-Score (79.7%) for the classifier  
321 was calculated.

### 3 DISCUSSION

322 The potential for networks trained using predictive coding and Variational Auto Encoders for learning  
323 joint latent representations of multimodal real-world sensory scenes to perform place recognition has been  
324 demonstrated. The MultiPredNet model proposed here consistently outperformed the JMVAE-zero and  
325 JMVAE-kl models in place recognition as evident from the RSA in all 3 test conditions using the physical  
326 platform (figure 5). Importantly, each model was composed of the same number of layers and nodes,  
327 trained and tested using the same data sets, and their weight spaces learnt through the same number of  
328 training epochs. The analysis used to compare performance at place recognition between models serves  
329 as a proxy to more direct measures of performance at place recognition through navigation. To clarify  
330 this we have demonstrated that coupling the MultiPredNet to a simple associative memory system and

331 capturing longer duration data sets, enables more conventional metrics for quantifying place recognition to  
332 be derived. What now remains to be demonstrated is how these metrics compare to other model-free or  
333 model-based place recognition systems that combine visual and tactile sensory data. Toward this we are  
334 unaware of any suitable model for comparison other than the ViTa-SLAM system Struckmeier et al. (2019)  
335 introduced by co-authors in a previous study that inspired this work and as such would be uninformative  
336 to compare against. What we have shown is that RSA does enable an empirical comparison of complex  
337 representation spaces to low dimensional pose spaces in an intuitive manner to guide in the evaluation  
338 of candidate models and to adjust network parameters prior to full integration with a SLAM back-end.  
339 The MultiPredNet returned the lowest visual sensory reconstruction errors whilst the JMVAE-kl model  
340 performed best at tactile reconstruction (figure 8, 10 and 11). The Kullback Leibler divergence term  
341 included in the JMVAE-kl loss function during training was introduced to bring the representation spaces  
342 of the disparate sensory modalities closer to enable bi-directional multimodal reconstruction. This appears  
343 to be the case for tactile reconstruction from visual input, however, it did not result in an improved  
344 performance in place recognition nor did it improve visual reconstruction from tactile as evidenced in the  
345 lower panel of figure 8 and example reconstructions shown in figure 10. The large offset apparent in the  
346 tactile reconstruction errors from the JMVAE-zero and MultiPredNet models suggest that these models  
347 did not accommodate this disparity. However, the MultiPredNet model maintained an above significance  
348 correlation in place recognition when only tactile information was available, which JMVAE-zero could  
349 not. Interpreting the representation space of MultiPredNet is, therefore, subject to further investigation,  
350 which reinforces the position that VAEs are certainly better understood machine learning tools and as  
351 such are the obvious choice for adoption by robotics engineers. However, PCNs stand as an algorithmic  
352 level solution to learning that more closely approximates the physiology of a “cortical compute unit”.  
353 The base compute unit in a PCN is the same throughout the network, referred to by Rao and Ballard as a  
354 Predictive Estimator Rao and Ballard (1999), wherein only local computation and updates are performed  
355 during training and inference. By contrast the JMVAE approaches require separate decoder networks for  
356 training, which are then disregarded during inference if sensory reconstructions are not required. In the  
357 case of the JMVAE-kl network which learns unimodal representations in parallel to, and in support of, the  
358 joint modal distribution during training, the additional encoder-decoder network pairs are also disregarded  
359 during inference. In a purely software-based system this inefficiency is not an issue, however, as we look  
360 toward the future of embedded machine learning, particularly within robotics and edge based applications,  
361 we anticipate the increasing adoption of energy efficient hardware platforms, such as neuromorphic devices  
362 Krichmar et al. (2019). The immediate practical advantage in adopting the PCN approach, therefore, is  
363 that the algorithm is highly amenable to hardware optimisation through parallel distributed learning and  
364 processing. Unlike VAE networks, the local learning rule applied at each layer of a PCN requires no global  
365 back propagation of error in agreement with the physiology of mammalian cortex Roelfsema and Holtmaat  
366 (2018). Moreover, the basic feedforward-feedback structure of PCNs resemble the core architecture of the  
367 sensory neocortex Felleman and Van Essen (1991); Bastos et al. (2012); Pennartz et al. (2019). Further,  
368 the modular nature of the PCN compute unit that encapsulates the encoder-decoder pairing of the VAE  
369 but at a local level, allows a graceful scaling of the algorithm through simple duplication of the basic unit.  
370 This principle extends to including additional sensory modalities for joint representation learning, whereby  
371 any additional modality specific networks could be integrated into the multisensory network with a linear  
372 increase in complexity. By contrast, the JMVAE-kl approach would require a combinatorial increase in  
373 encoder-decoder pairs to correctly integrate additional modalities into the joint space. In conclusion, PCNs  
374 not only offer computational advantages to autonomously learning robots in terms of place recognition,

375 but also convey a considerable neurobiological plausibility and better scalability as compared to VAE  
376 approaches.

## 4 MATERIALS AND METHODS

### 377 4.1 Multi-modal predictive coding network algorithm

#### 378 4.1.1 Multimodal Predictive Coding Network Architecture

379 The network consists of three modules, namely the visual module, tactile module and multisensory  
380 module as shown in figure 4a. The visual module processes visual information and consists of a neural  
381 network with  $N_V$  layers. Activity of the neurons in the  $l^{th}$  layer for the  $i^{th}$  input is denoted by a  $n_{V(l)}$   
382 dimensional vector,  $\mathbf{y}_i^{V(l)}$  where  $n_{V(l)}$  denotes the number of neurons in the  $l^{th}$  layer of the visual module.  
383 Each layer in the network predicts the activity ( $\hat{\mathbf{y}}$ ) of the preceding layer according to

$$\hat{\mathbf{y}}_i^{V(l-1)} = \phi \left( \left( \mathbf{y}_i^{V(l)} \right)^T \mathbf{W}_{l(l-1)}^V \right)^T \quad (1)$$

384 where  $\mathbf{W}_{l(l-1)}^V$  denotes the synaptic weights of the projections between the  $l^{th}$  and  $(l-1)^{th}$  layer  
385 in the visual module and  $\phi$  is the activation function of the neurons (ReLU). The lowest layer in the  
386 network predicts the visual input ( $\mathbf{X}_i^V$ ) and other layers predict the activities of neurons in the preceding  
387 layer. All layers in the network generate these predictions in parallel using Eq. 1. This aspect of the  
388 network is different from commonly employed feedforward networks in deep learning, where information  
389 is sequentially propagated from the first to last layer of the network. The tactile module consists of a  
390 similar neural network with  $N_T$  layers that processes tactile information. The multisensory module consists  
391 of a single layer which predicts the activities of neurons in the last layers of both the visual and tactile  
392 modules. The activity patterns of neurons in this layer are denoted by  $y_i^D$  for the  $i^{th}$  input and serve as the  
393 representations used for place recognition.

#### 394 4.1.2 Learning Algorithm

395 Predictive coding is used to update the synaptic weights and infer neuronal activities in the network. A  
396 graphical depiction of the inter-layer connectivity is shown in figure 4b. The  $l^{th}$  layer in the visual module  
397 generates a prediction about the neuronal activities in the  $(l-1)^{th}$  layer and also receives a prediction of  
398 its own neuronal activity from the  $(l+1)^{th}$  layer. The goal of the learning algorithm is to infer  $l^{th}$  layer  
399 neuronal activity ( $\mathbf{y}_i^{V(l)}$ ) for the  $i^{th}$  input that generates better predictions about neuronal activity in the  
400  $(l-1)^{th}$  layer and is predictable by the  $(l+1)^{th}$  layer. For this purpose,  $\mathbf{y}_i^{V(l)}$  is updated by performing  
401 gradient descent on the error function

$$\mathbf{e}_i^{V(l)} = \left( \hat{\mathbf{y}}_i^{V(l-1)} - \mathbf{y}_i^{V(l-1)} \right)^2 + \left( \hat{\mathbf{y}}_i^{V(l)} - \mathbf{y}_i^{V(l)} \right)^2 \quad (2)$$

402 which results in the following update rule for  $\mathbf{y}_i^{V(l)}$

$$\Delta \mathbf{y}_i^{V(l)} = \eta_y \left( \mathbf{W}_{l(l-1)}^V \phi'(\hat{\mathbf{y}}_i^{V(l-1)}) \left( \mathbf{y}_i^{V(l-1)} - \hat{\mathbf{y}}_i^{V(l-1)} \right) + \left( \mathbf{y}_i^{V(l)} - \hat{\mathbf{y}}_i^{V(l)} \right) \right) \quad (3)$$

403 where  $\eta_y$  is the rate for updating neuronal activities and  $\phi^{prime}$  is the derivative of the activation function  
404 used in the predictive coding network. The update rule in Eq. 3 is used to infer neuronal activity in all layers  
405 of the visual module for all inputs. Weights ( $\mathbf{W}_{l(l-1)}^V$ ) between  $l_{th}$  and  $(l-1)_{th}$  layers in the network are

406 updated by performing gradient descent on the error in the prediction generated by the  $l^{th}$  layer neurons  
 407 which results in the update rule for weights:

$$\Delta \mathbf{W}_{l(l-1)}^V = \eta_w \mathbf{y}_i^{V(l)} \phi'(\hat{y}_i^{V(l-1)}) \left( \mathbf{y}_i^{V(l-1)} - \hat{\mathbf{y}}_i^{V(l-1)} \right)^T \quad (4)$$

408 where  $\eta_w$  is the learning rate for updating weights. Note that the learning rule is Hebbian-like in the sense  
 409 that weight changes depend on the pre- and post-synaptic activity (pre:  $y^{V(l)}$ ; post:  $y^{V(l-1)} - \hat{y}^{V(l-1)}$ ).  
 410 The learning approach for the tactile module is identical to the visual module. In case of the multisensory  
 411 module, the representations are inferred based on prediction errors of topmost layers in both the visual and  
 412 tactile modules.

## 413 4.2 Multi-modal Variational Auto-Encoder algorithm

414 Both JMVAE-zero and JMVAE-kl extend Variational Autoencoders (VAE) to handle multisensory inputs.  
 415 Therefore, this section first provides a description of the VAE and then presents extensions pertaining to  
 416 JMVAE-zero and JMVAE-kl.

### 417 4.2.1 Variational Autoencoders

419 A VAE is an autoencoder with an encoder-decoder architecture that allows estimating a latent distribution  
 420 which can be used to sample data from the input space. Given input data  $x$  with a distribution of  $p(x)$   
 421 and a prior distribution  $p(z)$ , the encoder in a VAE estimates an approximate posterior distribution  
 422  $q_\phi(z|x)$  for the actual posterior  $p(z|x)$ . Here,  $\phi$  represents the parameters associated with the encoder. The  
 423 decoder maximizes the likelihood of the data  $p_\theta(x|z)$  given this approximate posterior distribution where  $\theta$   
 424 represents the parameters associated with the decoder. To overcome the intractable problem of computing  
 425 the marginal distribution, VAEs are trained to maximize a lower bound for the input data distribution  $p(x)$   
 426 by maximizing the following objective function

$$\mathcal{L}_{VAE} = -D_{KL}(q_\phi(z|x)||p(z)) + E_{q_\phi(z|x)}(\log p_\theta(x|z)) \quad (5)$$

427 where the first term  $D_{KL}(q_\phi(z|x)||p(z))$  represents the Kullback-Leibler divergence between the  
 428 approximate posterior and the prior distribution  $p(z)$ . The second term represents the reconstruction error  
 429 in the output of the decoder. VAE's represent both  $q_\phi(z|x)$  and  $p(z)$  using Gaussian distributions. The  
 430 mean and variance of  $q_\phi(z|x)$  are determined by the output of the encoder.  $p(z)$  is assumed to be a standard  
 431 normal distribution  $\mathcal{N}(0, I)$  where  $I$  denotes the identity matrix. This assumption allows VAEs to estimate  
 432 an approximate posterior that is closer to the standard normal distribution. This enables sampling from the  
 433 learned latent distribution to generate samples from the input space.

434

435 JMVAE builds upon VAE by enabling inference of joint representations based on input in multiple  
 436 modalities. In this paper, multiple modalities constitute the input from tactile (denoted by  $w$  for whisker)  
 437 and vision (denoted by  $v$ ) sensors on the WhiskEye robot. A straightforward approach for inferring  
 438 multimodal representations using a VAE is to learn a joint approximate posterior distribution  $q_\phi(z|w, v)$   
 439 using a network as shown in figure 14a. In this approach, a VAE maximizes the following objective function  
 440 to achieve a maximal lower bound on the marginal joint distribution

$$\mathcal{L}_{JM} = -D_{KL}(q_{\phi}(z|w, v)||p(z)) + E_{q_{\phi}(z|w, v)}(\log p_{\theta}(w|z)) + E_{q_{\phi}(z|w, v)}(\log p_{\theta}(v|z)) \quad (6)$$

441 Equation 6 represents the objective function for JMVAE-zero. It has been shown that JMVAE-zero is  
 442 not able to generate good crossmodal reconstructions when there are large structural differences between  
 443 different modalities Suzuki et al. (2017). To overcome this issue, a better VAE was developed in Suzuki et al.  
 444 (2017) called JMVAE-kl. JMVAE-kl employs a JMVAE-zero with two additional encoders for inferring the  
 445 approximate posteriors for the individual modalities  $w$  (denoted by  $q_{\phi_w}(z|w)$ ) and  $v$  (denoted by  $q_{\phi_v}(z|v)$ )  
 446 as shown in figure 14b. It is trained using an objective function that minimizes the Kullback-Leibler  
 447 divergence between the joint approximate posterior distribution and the approximate posterior distributions  
 448 for individual modalities, given by

$$\mathcal{L}_{JM(KL)} = \mathcal{L}_{JM} - \alpha \left( D_{KL}(q_{\phi}(z|w, v)||q_{\phi_v}(z|v)) + D_{KL}(q_{\phi}(z|w, v)||q_{\phi_w}(z|w)) \right) \quad (7)$$

449 where  $\alpha$  controls the strength of regularization due to the KL-divergence between the different posterior  
 450 distributions,  $\mathcal{L}_{JM(KL)}$  encourages inference of similar multimodal and unimodal approximate posterior  
 451 distributions thereby resulting in better crossmodal reconstructions.

### 452 4.3 RSA and statistical measures

#### 453 4.3.1 Representational Dissimilarity matrices (RDM)

454 To transform the 100-dimensional representations inferred by each network model in response to each  
 455 sample of a test set composed of  $n$  samples into an RDM, the dissimilarity distance between each  
 456 representation to all others was calculated. In this case we use 1 - Pearson correlation coefficient as  
 457 preferred by Kriegeskorte Kriegeskorte et al. (2008):

$$d_{x,y} = 1 - \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} \quad (8)$$

458 where  $d_{x,y}$  is the dissimilarity between representation  $x$  and  $y$ , which in turn will be the index address  
 459 into the  $n \times n$  symmetric RDM.

460

461 The RDMs for pose were constructed using the Euclidean distance between each 3D pose sample to all  
 462 others in the test set. For correctness the orientation and position components of the poses ( $x_{rot}$  and  $x_{trans}$ )  
 463 were independently scaled Bregier et al. (2018):

$$d_{x,y} = a||x_{rot} - y_{rot}|| + b||x_{trans} - y_{trans}|| \quad (9)$$

464 For the scaling factors  $a$  and  $b$  for rotation and translation respectively, we found that  $a = 0.3$  and  $b = 1$   
 465 to be appropriate in all experiments.

#### 466 4.3.2 Representational Similarity Analysis (RSA)

467 The vector of representational dissimilarity distances and accompanying vector of pose distances for  
 468 each sample in the test set were sorted into rank order and compared using Spearman's rank correlation  
 469 coefficient ( $\rho$ ):

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)} \quad (10)$$

470 where  $d_i$  is the difference between the rank order in pose against representation distance, and  $n$  being the  
 471 number of samples in the test set. The significance tests (p-values) were taken from  $t = \rho \sqrt{\frac{n-2}{1-\rho^2}}$  which  
 472 approximately follows Student's t with  $n - 2$  degrees of freedom under the null hypothesis.

473

#### 474 4.3.3 Precision-Recall analysis

475 To build the ROC curve the True Positive Rate (TPR) and False Positive Rates (FPR) were calculated at  
 476 each iteration of representation threshold to be tested as follows:

$$TPR = \frac{TP}{TP + FN} \quad FPR = \frac{FP}{FP + TN} \quad (11)$$

477 Where the cumulative True Positive (TP), False Positive (FP), True Negative (TN) and False Negative (FN)  
 478 scores from each iteration were used .

479 The geometric mean was calculated at each iteration as follows:

$$\sqrt{TPR * (1 - FPR)} \quad (12)$$

## CONFLICT OF INTEREST STATEMENT

480 The authors declare that the research was conducted in the absence of any commercial or financial  
 481 relationships that could be construed as a potential conflict of interest.

## AUTHOR CONTRIBUTIONS

482 MJP authored, conducted MultiPredNet training and RSA; MJP & BM developed robotic hardware  
 483 and captured data; SD, SB & CMAP developed MultiPredNet algorithm; OS, KT & VK developed  
 484 robotic experimental procedure; OS & TCK trained JMVAE; TCK and MJP generated simulated data and  
 485 conducted ROC curve analysis.

## ACKNOWLEDGMENTS

486 This research has received funding from the European Union's Horizon 2020 Framework Programme for  
 487 Research and Innovation under the Specific Grant Agreement No. 945539 (Human Brain Project SGA3).

## REFERENCES

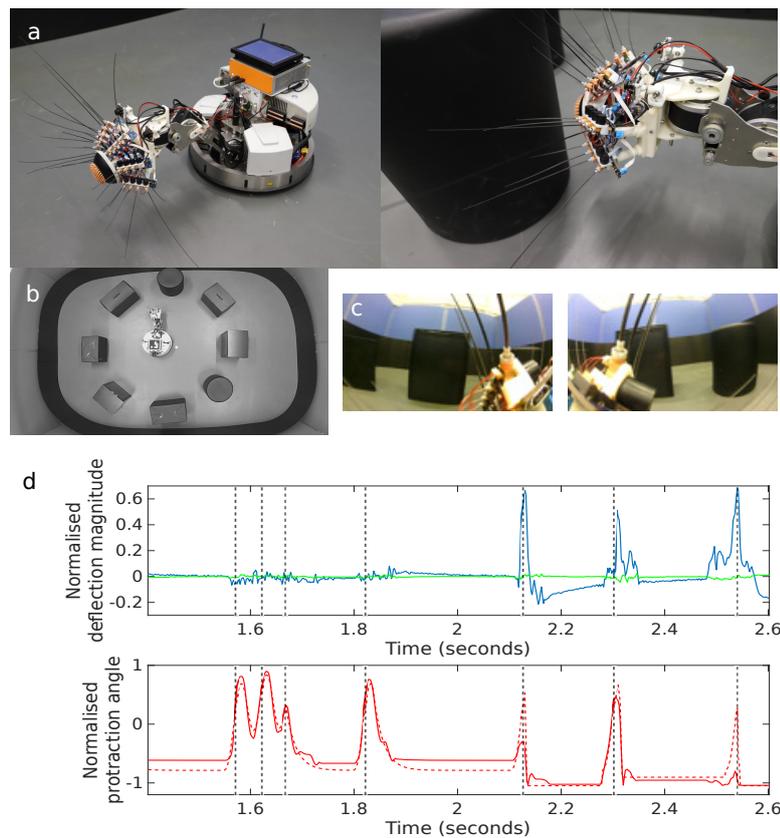
- 488 Badrinarayanan, V., Kendall, A., and Cipolla, R. (2017). Segnet: A deep convolutional encoder-decoder  
 489 architecture for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*  
 490 39, 2481–2495. doi:10.1109/TPAMI.2016.2644615
- 491 Bastos, A., Usrey, W., Adams, R., Mangun, G., Fries, P., and Friston, K. (2012). Canonical microcircuits  
 492 for predictive coding. *Neuron* 76, 695–771. doi:10.1016/j.neuron.2012.10.038
- 493 Bregier, R., Devernay, F., Leyrit, L., and Crowley, J. (2018). Defining the pose of any 3d rigid object and  
 494 an associated distance. *Int J Comput Vis* 126, 571–596. doi:https://doi.org/10.1007/s11263-017-1052-4
- 495 Choi, M. and Tani, T. (2018). Predictive coding for dynamic visual processing: Development of functional  
 496 hierarchy in a multiple spatiotemporal scales RNN model. *Neural Computation* 30, 237–270. doi:10.  
 497 1162/neco\_a\_01026

- 498 Dora, S., Pennartz, C., and Bohte, S. (2018). A deep predictive coding network for inferring hierarchical  
499 causes underlying sensory inputs. In *Artificial Neural Networks and Machine Learning – ICANN 2018.*,  
500 eds. K. V., M. Y., H. B., I. L., and M. I. (Springer), vol. 11141. doi:10.1007/978-3-030-01424-7\_45
- 501 Eitel, A., Springenberg, J. T., Spinello, L., Riedmiller, M., and Burgard, W. (2015). Multimodal deep  
502 learning for robust rgb-d object recognition. In *2015 IEEE/RSJ International Conference on Intelligent  
503 Robots and Systems (IROS)*. 681–687. doi:10.1109/IROS.2015.7353446
- 504 Eslami, S. M. A., Jimenez Rezende, D., Besse, F., Viola, F., Morcos, A. S., Garnelo, M., et al. (2018).  
505 Neural scene representation and rendering. *Science* 360, 1204–1210. doi:10.1126/science.aar6170
- 506 Falotico, E., Vannucci, L., Ambrosano, A., Albanese, U., Ulbrich, S., Vasquez Tieck, J. C., et al. (2017).  
507 Connecting artificial brains to robots in a comprehensive simulation framework: The neurorobotics  
508 platform. *Frontiers in Neurobotics* 11, 2. doi:10.3389/fnbot.2017.00002
- 509 Felleman, D. and Van Essen, D. (1991). Distributed hierarchical processing in the primate cerebral cortex.  
510 *Cereb Cortex* 1, 1–47. doi:10.1093/cercor/1.1.1
- 511 Flach, P. and Kull, M. (2015). Precision-recall-gain curves: Pr analysis done right. In *Advances in Neural  
512 Information Processing Systems*, eds. C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett  
513 (Curran Associates, Inc.), vol. 28
- 514 Friston, K. (2010). The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience* 11,  
515 127–138
- 516 Gao, P., Bermejo, R., and Zeigler, H. (2001). Whisker deafferentation and rodent whisking patterns:  
517 behavioural evidence for a central pattern generator. *The Journal of Neuroscience* 21, 5374–5380
- 518 Giusti, A., Guzzi, J., Cireşan, D. C., He, F., Rodríguez, J. P., Fontana, F., et al. (2016). A machine learning  
519 approach to visual perception of forest trails for mobile robots. *IEEE Robotics and Automation Letters* 1,  
520 661–667. doi:10.1109/LRA.2015.2509024
- 521 Gregory, R. (1980). Perceptions as hypotheses. *Philos Trans R Soc Lond B Biol Sci* 290, 181–97.  
522 doi:10.1098/rstb.1980.0090
- 523 Gu, Y., Wang, Y., and Li, Y. (2019). A survey on deep learning-driven remote sensing image scene  
524 understanding: Scene classification, scene retrieval and scene-guided object detection. *Applied Sciences*  
525 9, 2110. doi:10.3390/app9102110
- 526 Helmholtz, H. V. (1867). *Treatise on Physiological Optics Vol. III* (Dover Publications)
- 527 Hwang, J., Kim, J., Ahmadi, A., Choi, M., and Tani, J. (2017). Predictive coding-based deep dynamic  
528 neural network for visuomotor learning. In *2017 Joint IEEE International Conference on Development  
529 and Learning and Epigenetic Robotics (ICDL-EpiRob)*. 132–139. doi:10.1109/DEVLRN.2017.8329798
- 530 Kazmi, S. M. A. M. and Mertsching, B. (2016). Gist+ratslam: An incremental bio-inspired place  
531 recognition front-end for ratslam. In *Proceedings of the 9th EAI International Conference on Bio-  
532 Inspired Information and Communications Technologies (Formerly BIONETICS)* (ICST (Institute for  
533 Computer Sciences, Social-Informatics and Telecommunications Engineering)), 27–34. doi:10.4108/eai.  
534 3-12-2015.2262532
- 535 Khaleghi, B., Khamis, A., Karray, F., and Razavi, S. (2013). Multisensor data fusion: A review of the state  
536 of the art. *Information Fusion* 14, 28–44. doi:10.1016/j.inffus.2011.08.001
- 537 Kingma, D. P. and Welling, M. (2013). Auto-encoding variational bayes. In *ICLR 2014 Workshop*
- 538 Korthals, T., Hesse, M., Leitner, J., Melnik, A., and Rückert, U. (2019). Jointly trained variational  
539 autoencoder for multi-modal sensor fusion. In *2019 22th International Conference on Information  
540 Fusion (FUSION)*. 1–8
- 541 Krichmar, J. L., Severa, W., Khan, M. S., and Olds, J. L. (2019). Making bread: Biomimetic strategies for  
542 artificial intelligence now and in the future. *Frontiers in Neuroscience* 13, 666. doi:10.3389/fnins.2019.

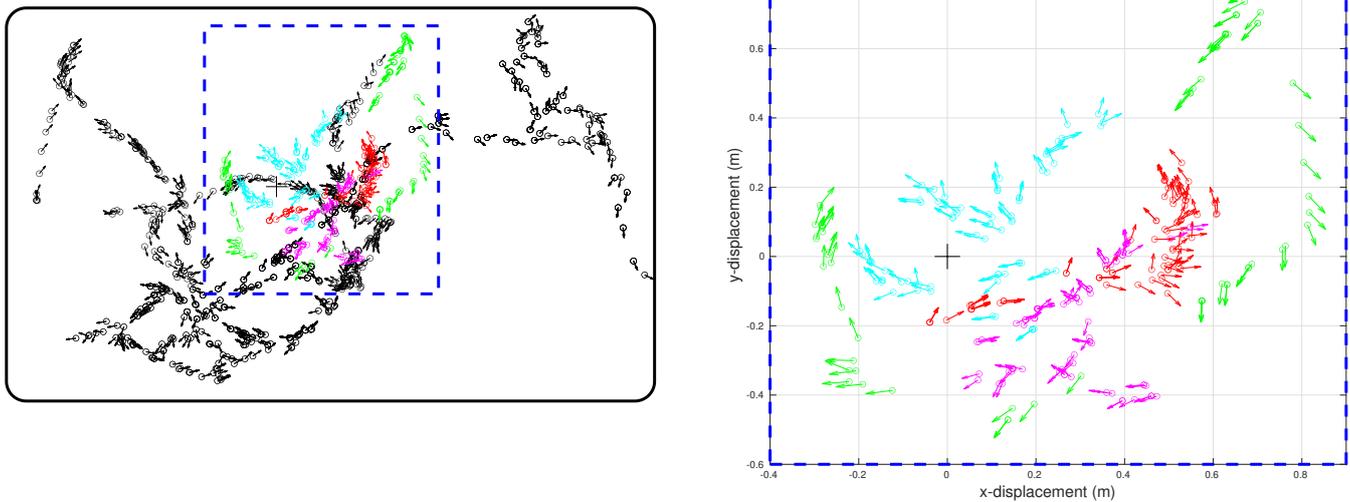
- 543 00666
- 544 Kriegeskorte, N., Mur, M., and Bandettini, P. (2008). Representational similarity analysis - connecting the  
545 branches of systems neuroscience. *Frontiers in systems neuroscience* 2. doi:10.3389/neuro.06.004.2008
- 546 Lanillos, P. and Cheng, G. (2018). Adaptive robot body learning and estimation through predictive coding.  
547 In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 4083–4090
- 548 Lowry, S., Sünderhauf, N., Newman, P., Leonard, J. J., Cox, D., Corke, P., et al. (2016). Visual place  
549 recognition: A survey. *IEEE Transactions on Robotics* 32, 1–19. doi:10.1109/TRO.2015.2496823
- 550 Milford, M., Wyeth, G., and Prasser, D. (2004). Ratslam: a hippocampal model for simultaneous  
551 localization and mapping. In *IEEE International Conference on Robotics and Automation, 2004. Proceedings. ICRA '04. 2004*. vol. 1, 403–408 Vol.1. doi:10.1109/ROBOT.2004.1307183
- 552
- 553 Mitchinson, B. and Prescott, T. J. (2013). Whisker movements reveal spatial attention: A unified  
554 computational model of active sensing control in the rat. *PLOS Computational Biology* 9, 1–16.  
555 doi:10.1371/journal.pcbi.1003236
- 556 Mnih, V., Kavukcuoglu, K., Silver, D., A.A.Rusu, Veness, J., Bellemare, M., et al. (2015). Human-level  
557 control through deep reinforcement learning. *Nature* 518, 527–533. doi:10.1038/nature14236
- 558 Mumford, D. (1992). On the computational architecture of the neocortex. *Biological Cybernetics* 66,  
559 241–251. doi:10.1007/BF00198477
- 560 Murphy, R. R. (1998). Dempster-shafer theory for sensor fusion in autonomous mobile robots. *IEEE*  
561 *Transactions on Robotics and Automation* 14, 197–206. doi:10.1109/70.681240
- 562 Nagai, Y. (2019). Predictive learning: its key role in early cognitive development. *Philosophical*  
563 *Transactions of the Royal Society B: Biological Sciences* 374, 20180030. doi:10.1098/rstb.2018.0030
- 564 Park, J.-C., Lim, J. H., Choi, H., and Kim, D.-S. (2012). Predictive coding strategies for developmental  
565 neurorobotics. *Frontiers in Psychology* 3, 134. doi:10.3389/fpsyg.2012.00134
- 566 Pearson, M. J., Fox, C., Sullivan, J. C., Prescott, T. J., Pipe, T., and Mitchinson, B. (2013). Simultaneous  
567 localisation and mapping on a multi-degree of freedom biomimetic whiskered robot. In *Robotics and*  
568 *Automation (ICRA), 2013 IEEE International Conference on (IEEE)*, 586–592
- 569 Pennartz, C., Dora, S., Muckli, L., and Lorteije, J. (2019). Towards a unified view on pathways and  
570 functions of neural recurrent processing. *Trends Neurosci.* 42, 589–603. doi:10.1016/j.tins.2019.07.005
- 571 Pennartz, C. M. (2015). *The Brain's Representational Power: On Consciousness and the Integration of*  
572 *Modalities* (Cambridge: The MIT Press)
- 573 Prescott, T., Pearson, M., Mitchinson, B., and Pipe, T. (2009). Whisking with robots: From rat vibrissae to  
574 biomimetic technology for active touch. *IEEE Robotics and Automation Magazine* 16, 42–50
- 575 Rao, R. and Ballard, D. (1999). Predictive coding in the visual cortex: a functional interpretation of some  
576 extra-classical receptive-field effects. *Nature Neuroscience* 2, 79–87. doi:https://doi.org/10.1038/4580
- 577 Roelfsema, P. and Holtmaat, A. (2018). Control of synaptic plasticity in deep cortical networks. *Nature*  
578 *reviews Neuroscience* 19, 166–180. doi:10.1038/nrn.2018.6
- 579 Roumeliotis, S. I. and Bekey, G. A. (1997). Extended Kalman filter for frequent local and infrequent global  
580 sensor data fusion. In *Sensor Fusion and Decentralized Control in Autonomous Robotic Systems*, eds.  
581 P. S. Schenker and G. T. McKee. International Society for Optics and Photonics (SPIE), vol. 3209, 11 –  
582 22. doi:10.1117/12.287638
- 583 Sheppard, C. and Rahnemoonfar, M. (2017). Real-time scene understanding for uav imagery based on deep  
584 convolutional neural networks. In *2017 IEEE International Geoscience and Remote Sensing Symposium*  
585 *(IGARSS)*. 2243–2246. doi:10.1109/IGARSS.2017.8127435
- 586 Spratling, M. (2017). A review of predictive coding algorithms. *Brain and cognition* 112, 92–97.  
587 doi:10.1016/j.bandc.2015.11.003

- 588 [Dataset] Stanford Artificial Intelligence Laboratory et al. (2018). Robotic operating system  
589 Struckmeier, O., Tiwari, K., Salman, M., Pearson, M., and Kyrki, V. (2019). ViTa-SLAM: A bio-inspired  
590 visuo-tactile slam for navigation while interacting with aliased environments. In *IEEE International*  
591 *Conference on Cyborg and Bionic Systems CBS2019* (Curran Associates, Inc.), 97–103  
592 Sunderhauf, N., Brock, O., Scheirer, W., Hadsell, R., Fox, D., Leitner, J., et al. (2018). The limits and  
593 potentials of deep learning for robotics. *The International Journal of Robotics Research* 37, 405–420.  
594 doi:10.1177/0278364918770733  
595 Suzuki, M., Nakayama, K., and Matsuo, Y. (2017). Joint multimodal learning with deep generative models.  
596 In *ICLR 2017 Workshop*

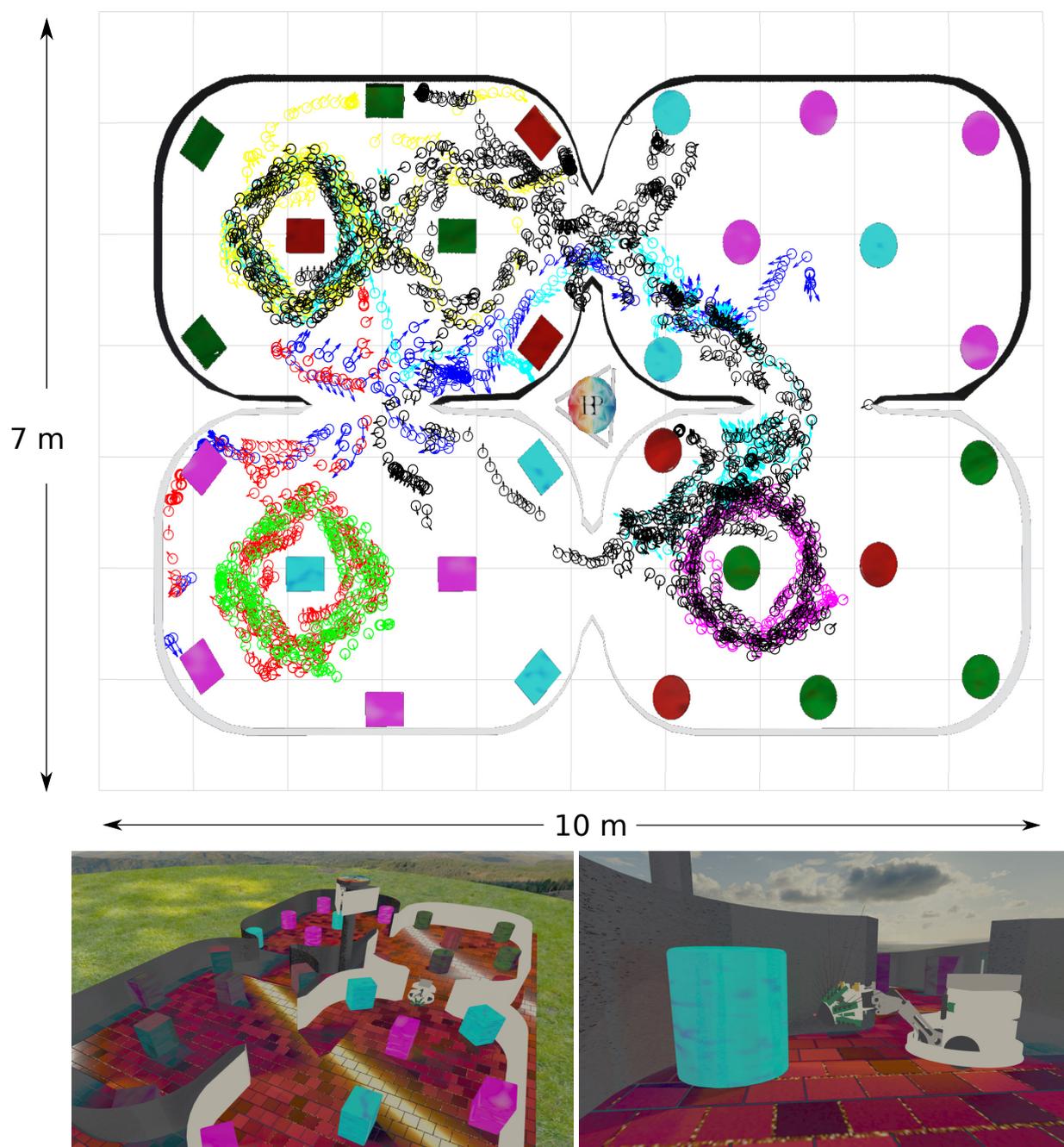
## FIGURE CAPTIONS



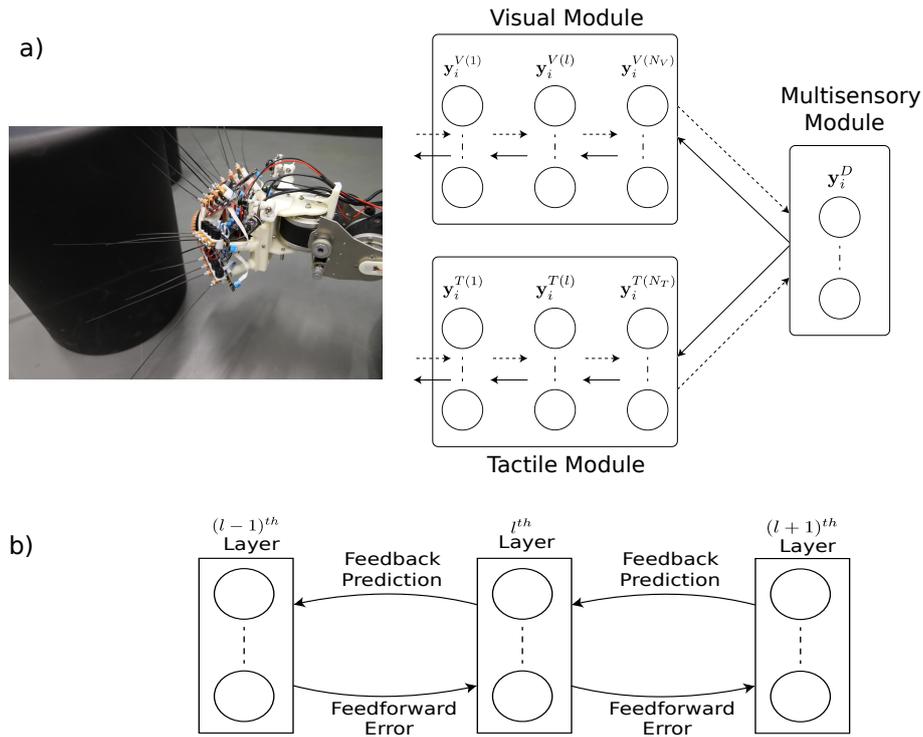
**Figure 1.** The WhiskEye robot (a) has 24 actuated tactile whiskers and camera eyes on its head mounted at the end a 3 dof neck and omnidrive Robotino body. Simultaneous visual frames (panel c, left and right eye cameras) and tactile samples are taken at the point of peak whisker protraction as indicated by the black vertical dashed lines in the plots of panel d. The example time series data shown in panel d are taken from a single whisker through 7 whisk cycles with only the final 3 whisks making contact with an object. The red dashed trace in the lower plot is the drive or desired protraction angle of the whisker scaled to  $\pm 1$  of the full whisk angle range of  $\pm 80$  degrees of rotation. The solid red line is the measured protraction angle of the whisker ( $\theta_{whisk}$ ) which can be inhibited by contacts as is clear on the 5th whisk. The blue trace in the upper plot is the  $x$ - and the green the  $y$ -deflection of the whisker scaled to  $\pm 1$  of maximum deflection magnitude. The three positive whisker contacts are clear in the final three whisks of this sample. The  $x$ ,  $y$  and  $\theta_{whisk}$  samples taken at the point of peak protraction for all 24 whiskers constitute the tactile ‘view’ of the robot at that instance. The experimental arena shown in panel b, was populated with matt black cylinders and boxes, the configuration was changed between collecting data to train the networks and for testing.



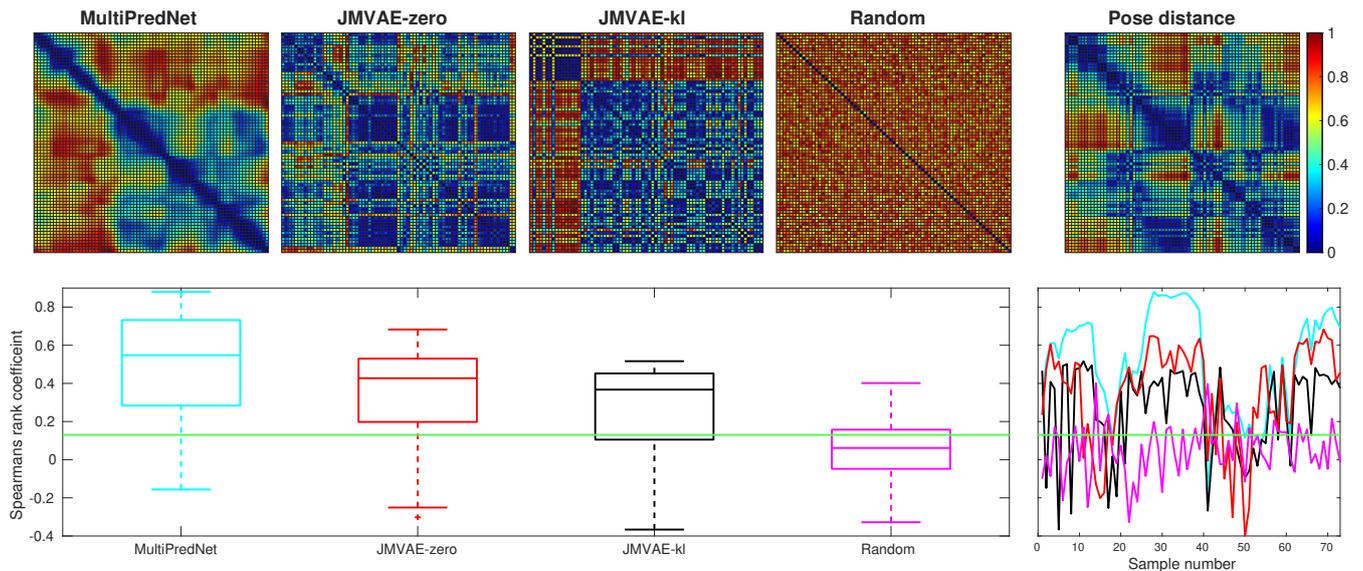
**Figure 2.** Quiver plots of the poses of WhiskEye’s head (defined as x,y position and head direction) at each sample point taken for the training set (black) and for each of the test sets (red, cyan, magenta, green) used to evaluate the models. Each test set was recorded from a different initial pose of WhiskEye and in the arena populated with different object configurations to test for generalisation. The right panel is a scaled view of the region indicated by the blue dashed rectangle in the left panel. The bold black rectangle enclosing the sample points in the left panel indicates the boundary walls of the arena. The arena measured 3.5m by 2.25m



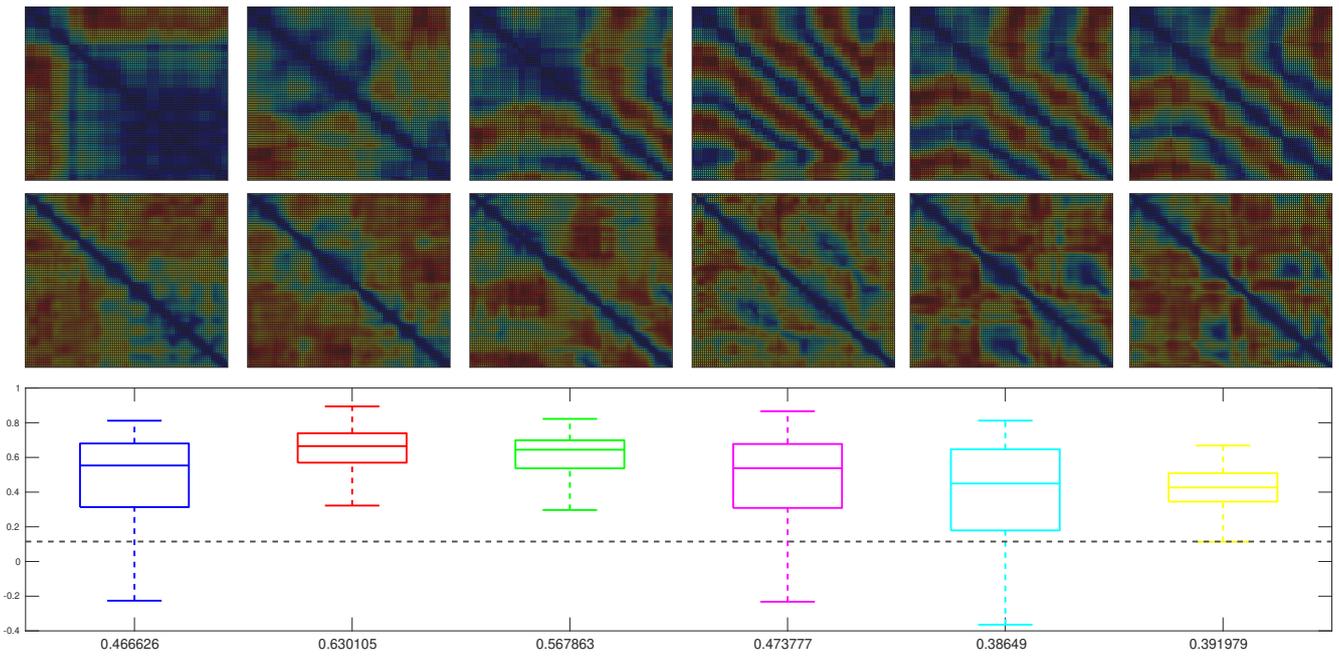
**Figure 3.** Simulated WhiskEye in the NeuroRobotics Platform. Top panel) Quiver plot of head poses ( $x, y, \theta$ ) at sample points taken as the training set (black) and 6 test sets (blue:1, red:2, green:3, magenta:4, cyan:5, yellow:6) The arena walls and coloured objects have been superimposed onto the quiver plot for reference. Lower panels) Screen shots taken from the simulator showing the arena and simulated WhiskEye robot as it explores the arena. The tactile attention model used to control the physical platform is the same as used in the simulator.



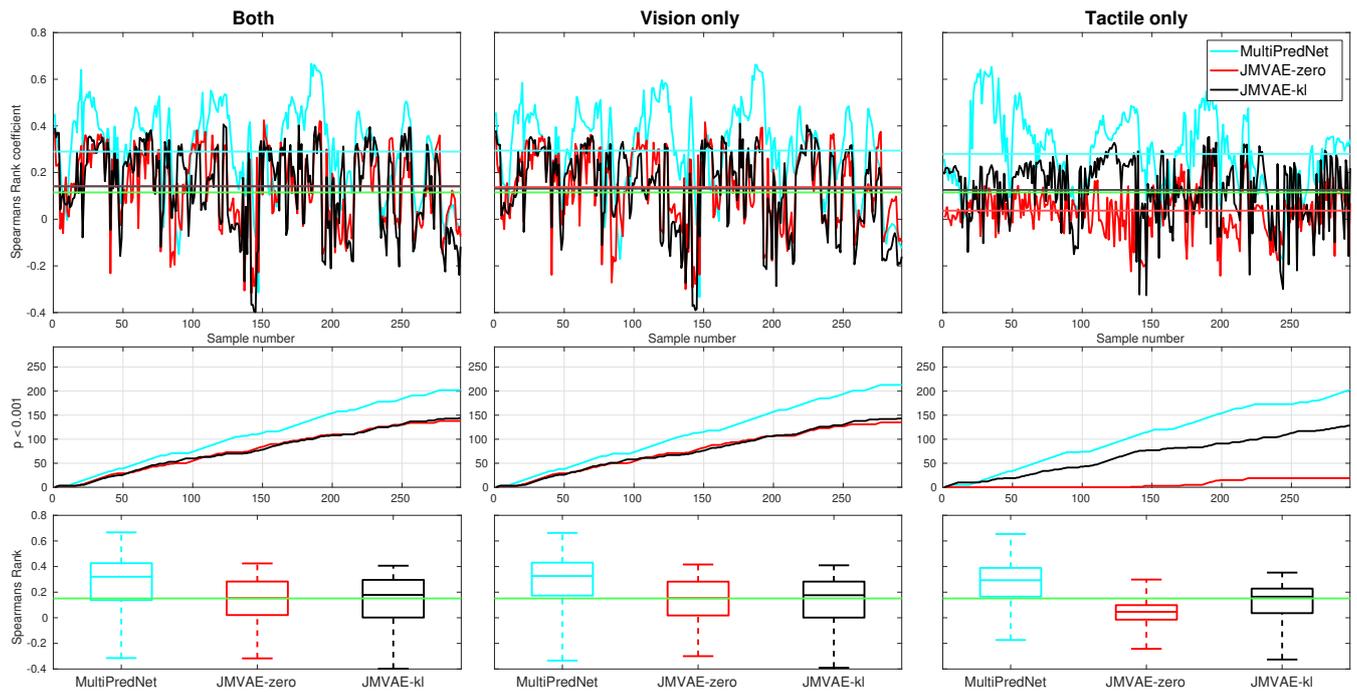
**Figure 4.** Multimodal predictive coding network architecture. a) The network consists of 3 neural network modules; the Visual and Tactile modules taking visual and tactile input from the WhiskEye robot with  $N_V$  and  $N_T$  layers respectively; and the Multisensory module consisting of a single layer that predicts the activities of the neurons in the last layers of both the visual and tactile modules. Note that each layer in the Visual and Tactile modules are also predicting the activities of neurons in their preceding layers as shown in panel b, passing prediction errors forward through the network



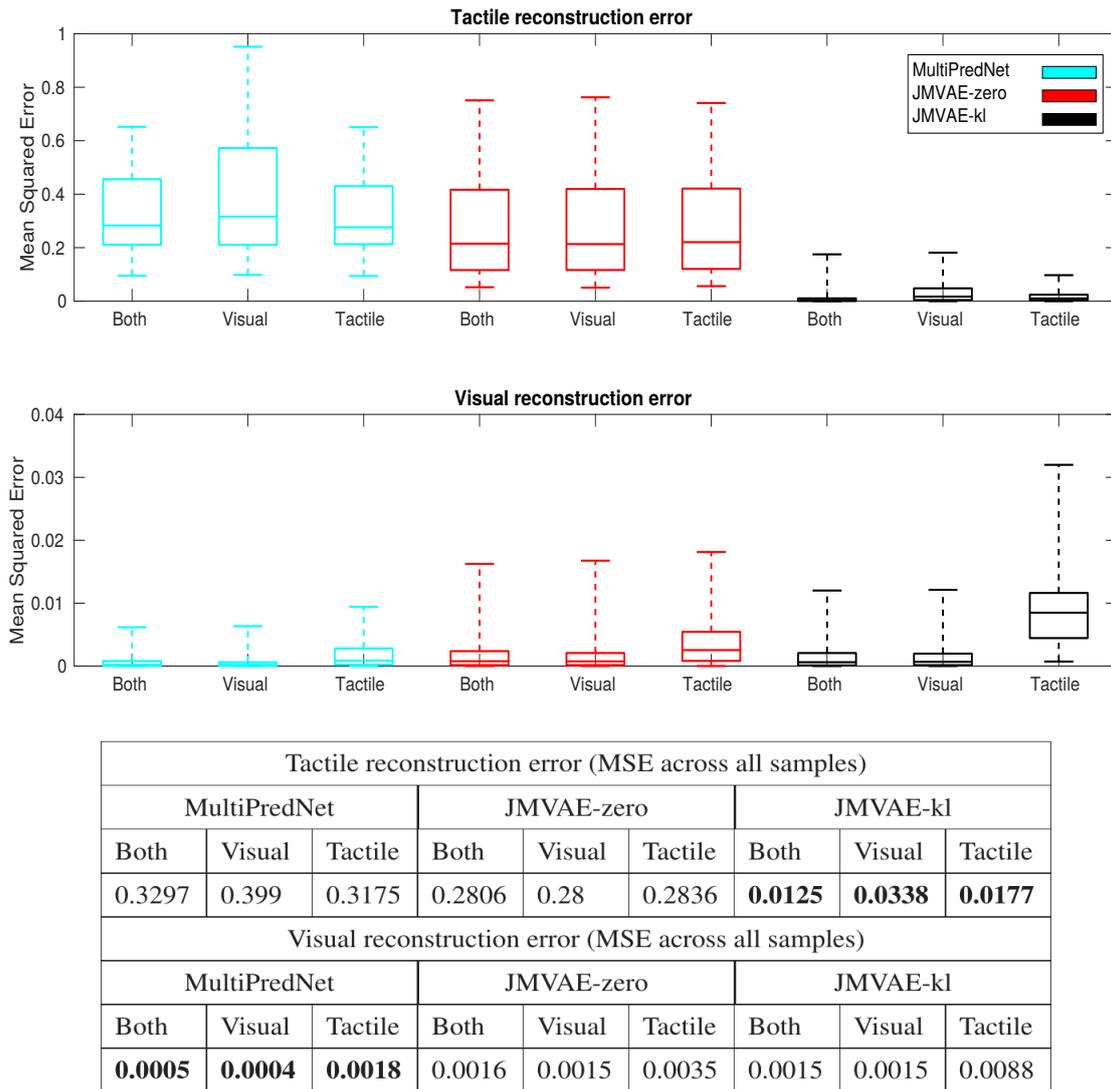
**Figure 5.** Representational Similarity Analysis of three trained network models to measure performance at place recognition across a small example test set. The top row presents the Representational Dissimilarity Matrices (RDM) generated from (left to right) MultiPredNet, JMVAE-zero, JMVAE-kl and random representations in response to visual and tactile samples taken from physical test set 1 (73 samples). The RDM on the far right was generated from the associated ground truth 3D poses of WhiskEye's head ( $x, y, \theta_{head}$ ) for each sample in the test set. The boxplots in the lower panel summarise statistics of the Spearman's rank correlation coefficient ( $\rho$ ) calculated between the pose and representation distances across the test set for each model as shown in the coloured line plots in the panel to the right (cyan: MultiPredNet, red: JMVAE-zero, black: JMVAE-kl, and magenta: Random). The Green horizontal line indicating 99% significance above chance ( $p < 0.001, N = 73$ )



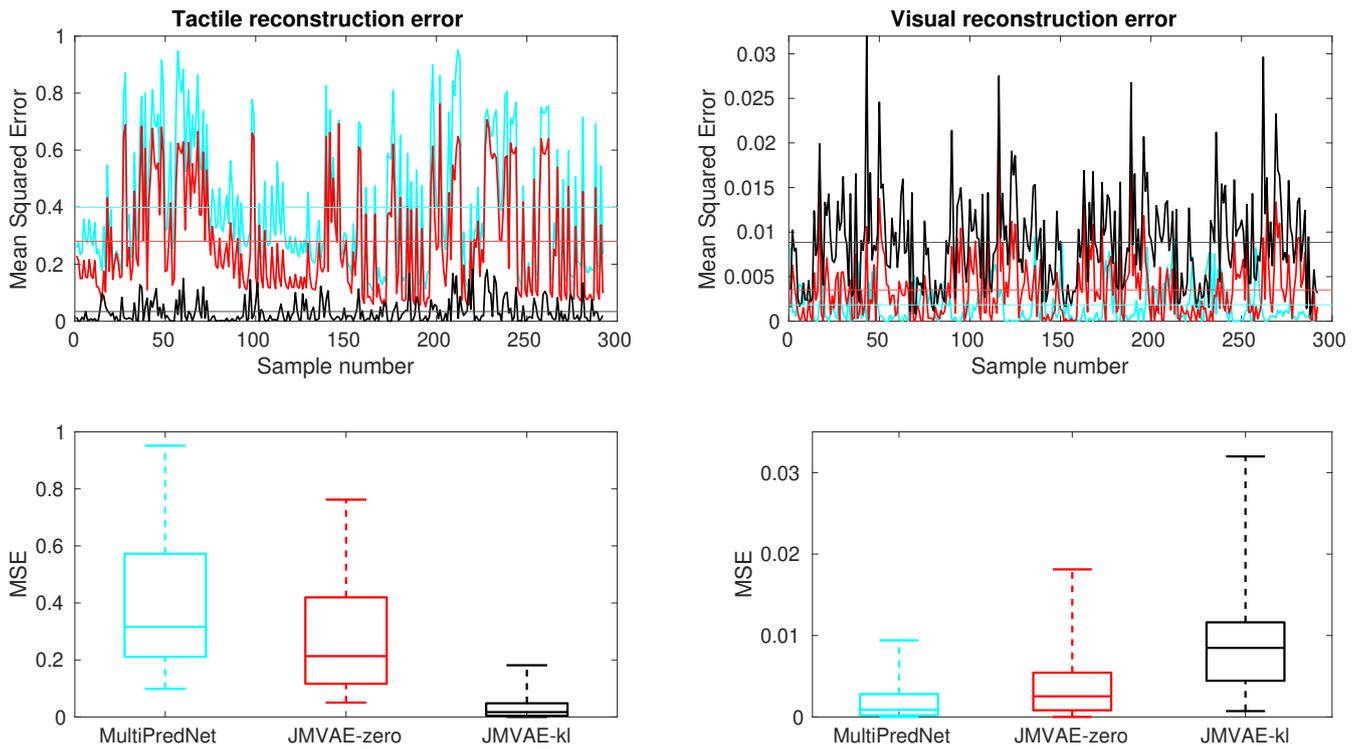
**Figure 6.** Representation Similarity Analysis applied to each of the 6 test sets sampled from the simulator and inferred by a trained MultiPredNet Model using same network topology as for physical data sets. Top row) Representation Dissimilarity Matrices (RDMs) for the first 100 samples of pose from each set (1 to 6 left to right). Middle row) RDMs for the first 100 inferred joint latent representations from each set. Lower panel) Box plots summarising the Spearman's Rank coefficient ( $\rho$ ) calculated for the full 400 samples of each test set. The colour of each box plot is the same colour as quiver plots for each test set shown in figure 3, with the black dashed line indicating significance ( $p < 0.001, n = 400$ ). The mean value of  $\rho$  for each set is printed beneath each box plot for clarity.



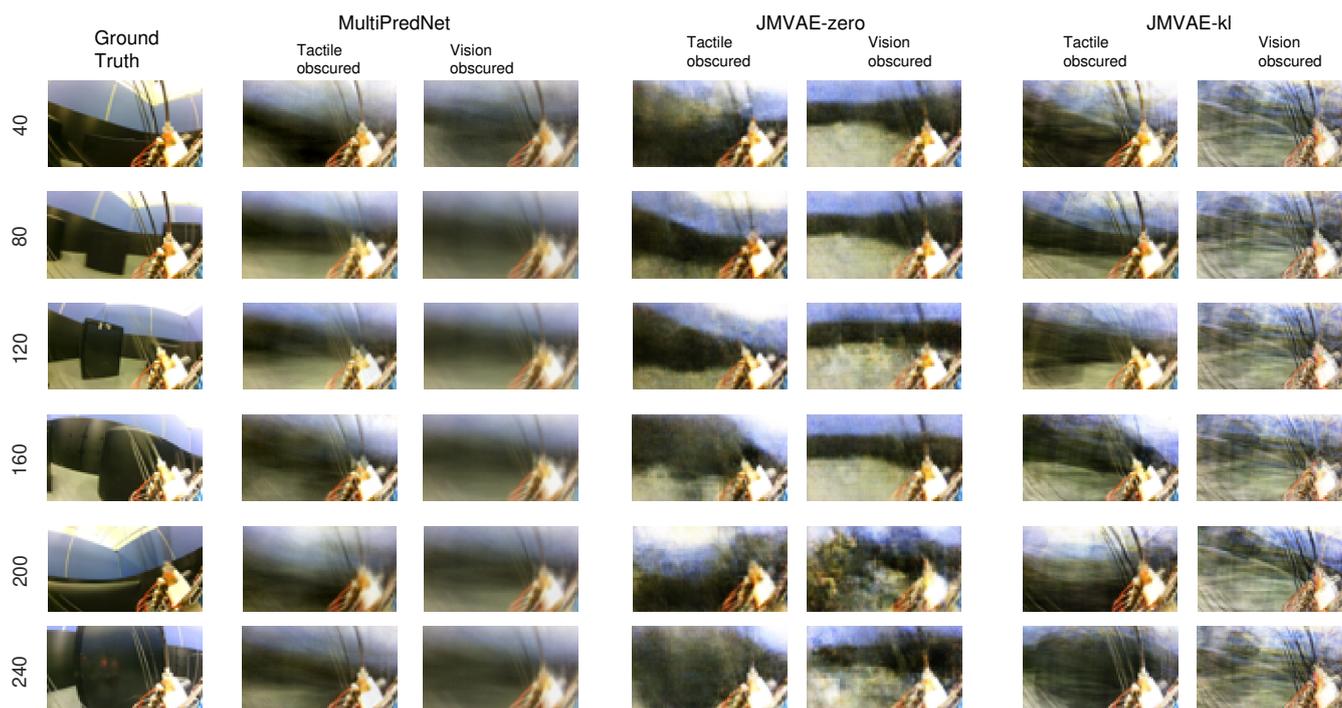
**Figure 7.** Spearman's Rank correlation coefficient ( $\rho$ ) to measure performance at place recognition by 3 trained models across all 4 physical test sets when both sensory modalities are available (left column), only vision available (middle column) and only tactile available (right column). The top row of panels trace  $\rho$  with the coloured horizontal lines indicating the mean value for each model and green line indicating significance ( $p < 0.001$ ,  $N = 292$ ). The middle row of traces show the cumulative number of samples across the test set that scored an above chance positive correlation in each of the sensory conditions. The statistics for  $\rho$  across each model is distilled into boxplots in the bottom row, again with the green line indicating significance.



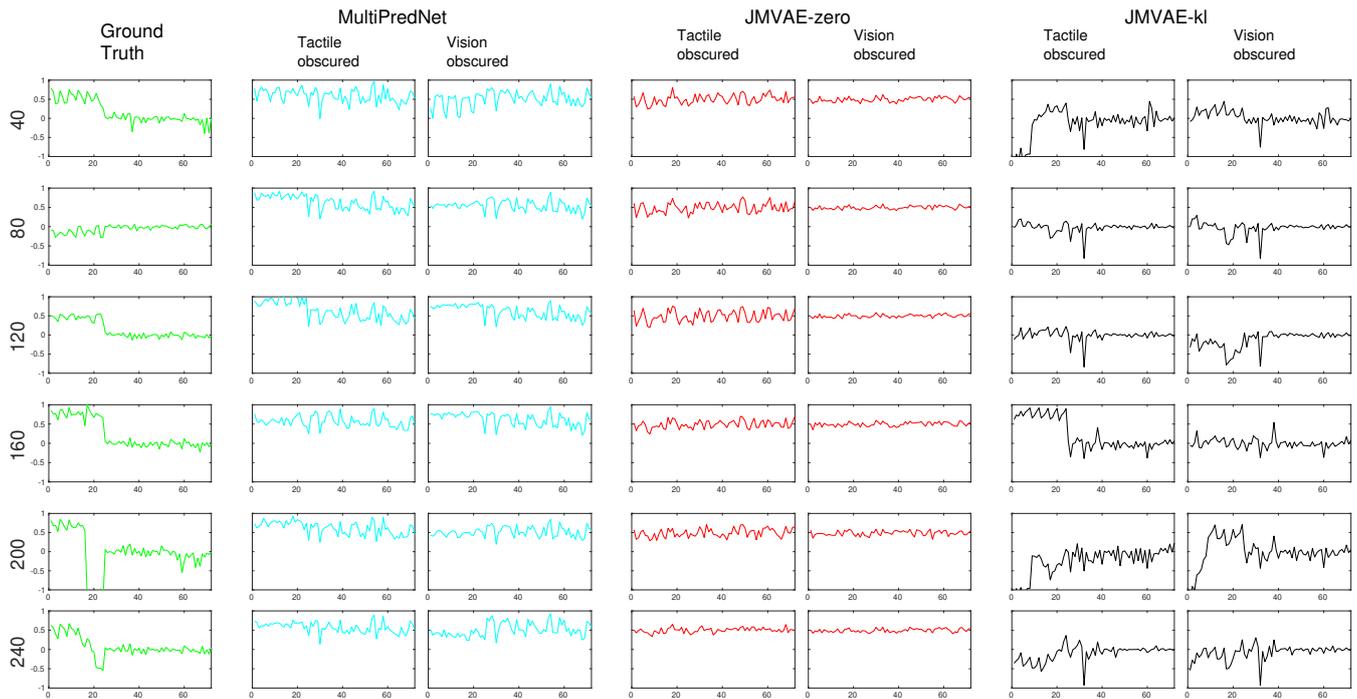
**Figure 8.** Sensory reconstruction errors from the 3 trained model networks in response to all 4 concatenated test sets under 3 test conditions; *Both* visual and tactile sensory data available; *Visual* data available and tactile masked; and *Tactile* data available with visual data masked. The top panel box plots summarise statistics of the mean squared error between actual tactile sensory data and the reconstructed tactile impressions generated by each of the networks. The middle panel applied the same analysis to the visual reconstructions with the average MSE for each model in each of the conditions summarised in the table beneath (bold highlighting lowest error condition)



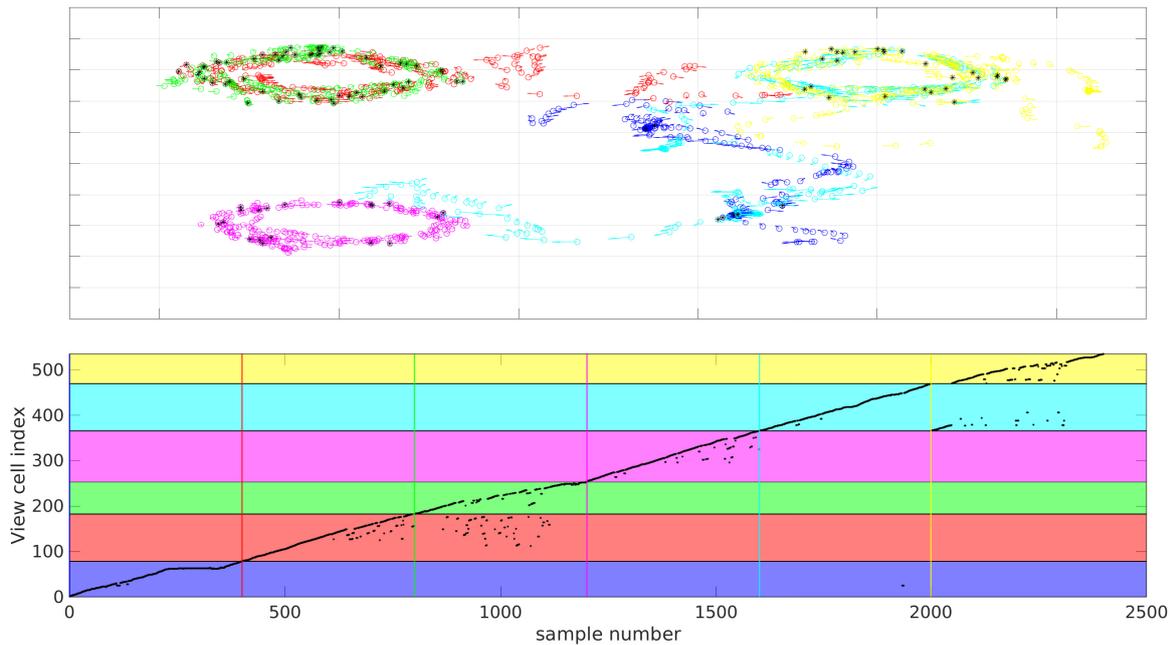
**Figure 9.** Tactile and Visual reconstruction errors from each sample in all test sets inferred by the MultiPredNet, JMVAE-zero and JMVAE-kl models with the tactile (and visual respectively) sensory input obscured. Note the systematic offset in tactile reconstruction error from the MultiPredNet and JMVAE-zero model as indicated by the mean of the Mean Squared Error (MSE) across all samples (horizontal coloured lines).



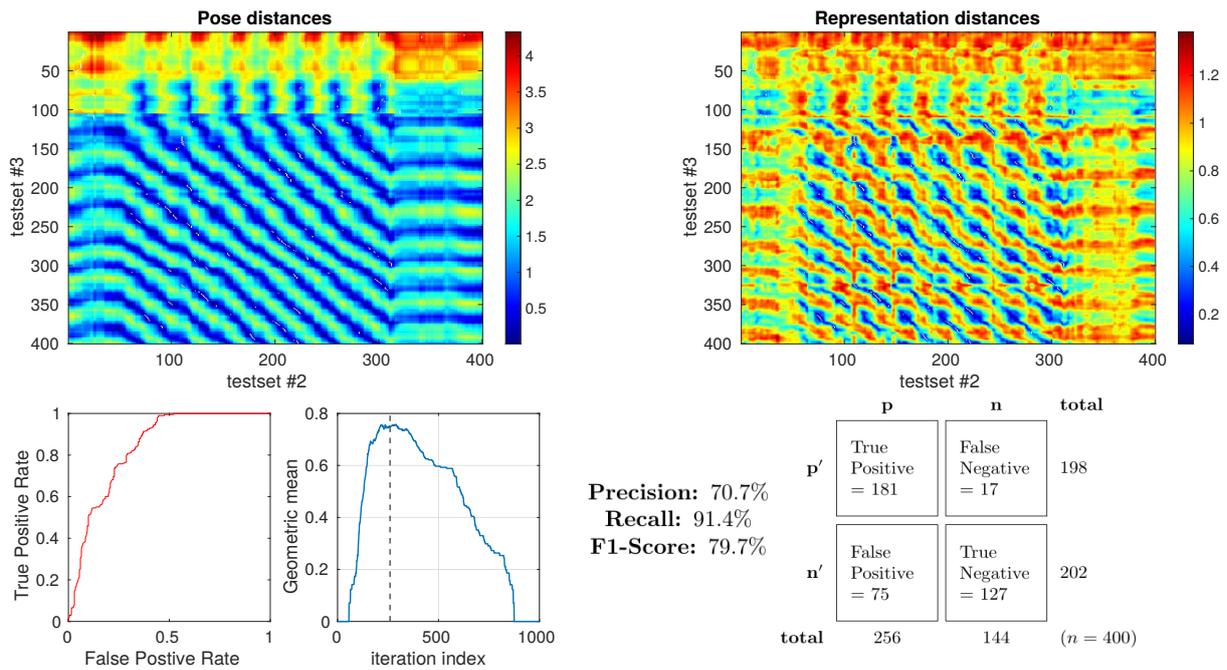
**Figure 10.** Example frames captured by the eye camera that are presented as visual input to the trained models with their subsequent reconstructions in the 2 sensory drop-out conditions. The Ground truth images in the left column were taken at the sample number indicated to the left of each panel with each reconstruction from the models presented in that row. The reconstructions are qualitatively similar in their quality across all 3 models, however, MultiPredNet did return the lowest mean squared reconstruction errors in all conditions as shown in figure 8.



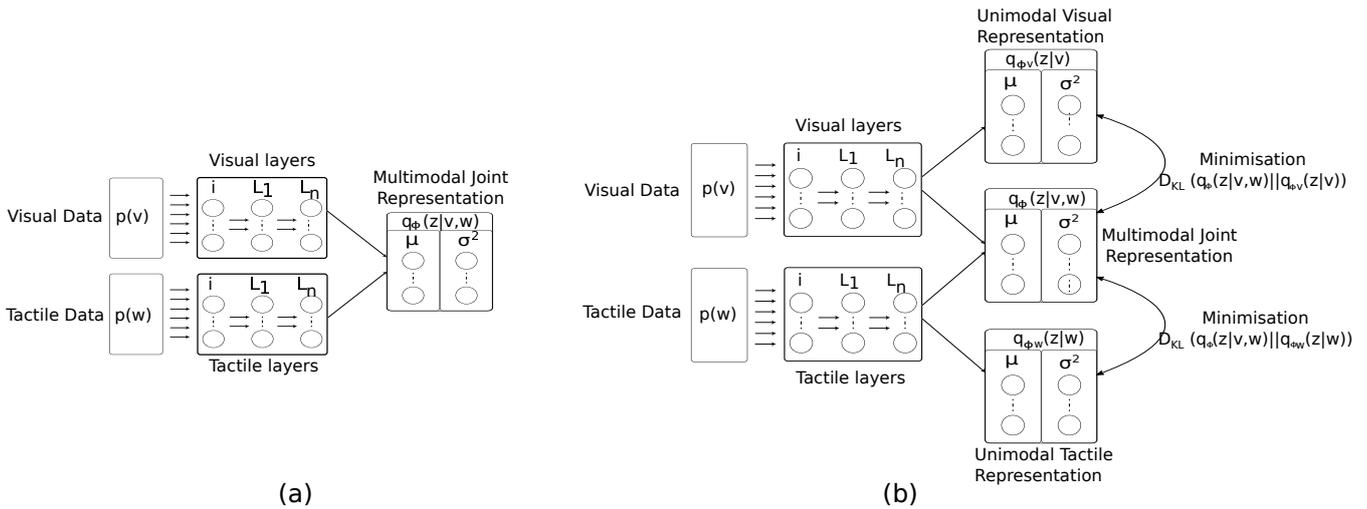
**Figure 11.** Example tactile information captured from the whisker array at point of peak protraction at the sample points in the test set as indicated by the number on the left of each row. The first 24 values in each sample is the protraction angle of each whisker ( $\theta_{whisk}$ ) in the array with the remaining 48 values indicating the magnitude of deflection experienced by each whisker in the  $x$  and  $y$  dimensions (refer to figure 1 for description). As in figure 10, the plots in each row to the right of the ground truth were reconstructed by the different trained models under test during the 2 drop-out conditions. The JMVAE-kl model performed best at tactile reconstruction in all conditions (see figure 8) as is clear from these indicative examples. Both MultiPredNet and JMVAE-zero failed to accommodate the systematic off-set in deflection angle which JMVAE-kl has, i.e., nominally zero in the last 48 values of reconstructed vector.



**Figure 12.** View-cell memory sequentially presented with representations generated from all 6 simulated test sets to associate poses with proximal representations. Top panel) Colour coded quiver plots of ground truth poses of the simulated WhiskEye head during each of the 6 test sets (1 =blue, 2 =red, 3 =green, 4 =magenta, 5 =cyan & 6 =yellow). The black asterisks indicate samples which triggered a re-localisation event in the view-cell memory, i.e., the representation of that sample was close to a previous representation stored in the view-cell memory. Lower panel) Graph of view-cell index against sample number from the concatenated 6 test sets ( $n = 2400$ ). The coloured horizontal bars highlight the region within the view-cell memory that store representations encountered during a particular test set (colour matched to test sets). The vertical coloured lines set the start of each new region in the view-cell memory. Re-localisation events are marked by a step decrease in view cell index between samples, significantly, samples from test set 3 are triggering re-localisation events that reference to view-cells created in test set 2.



**Figure 13.** Precision-recall curve analysis applied to simulated test sets 2 and 3 (see figure 3) classifying for place recognition through the representation space of the trained MultiPredNet. Top left) Heatmap summarising the distance in pose space from each sample in test set 2 against each sample in test set 3. Top right) Heatmap of distances in representation space between samples in test set 2 against all samples in test set 3. White dots in the heatmap indicate classification points for each sample ( $n = 400$ ) selected as the lowest distance in representation space between the 2 test sets. These points have been translated into the pose distance heatmap to determine a true or false classification. Lower left plot) Receiver Operating Characteristic (ROC) curve summarising impact of representation threshold for determining positive versus negative classifications. Lower right plot) Geometric mean of the ROC curve at each threshold iteration with the peak highlighted by the vertical dashed line. The confusion matrix contains the summed classification classes when using the optimal representation threshold determined from ROC curve with the Precision, Recall and F1-Score for the classifier calculated from them.



**Figure 14.** Architectural diagrams of the encoder networks of the two JMVAE models; **(a)** JMVAE-zero, and **(b)** JMVAE-kl. Both networks attempt to represent the two input modalities (visual  $p(v)$  and tactile  $p(w)$ ) as a joint multimodal latent representation  $q_\phi(z|v, w)$ . This is done by minimising both the reconstruction error for each modality, and the KL-divergence ( $D_{KL}$ ) between a standard normal distribution and the joint multimodal distribution. The resulting continuous distribution is encoded by the activity of two parallel layers of nodes representing the mean and variance ( $\mu, \sigma^2$ ) of each latent dimension. The JMVAE-kl model trains 2 further encoders, one using only visual input  $q_{\phi_v}(z|v)$  and the other only tactile  $q_{\phi_w}(z|w)$ , such that the KL-divergence measures between the unimodal and multimodal approximate distributions can be included into the loss function during training.