

VIEW-BASED MODELLING OF HUMAN VISUAL NAVIGATION ERRORS

Lyndsey Pickup^A, Andrew Fitzgibbon^B, Stuart Gilson^C, Andrew Glennerster^A

^ASchool of Psychology and Clinical Language Sciences, University of Reading

^BMicrosoft Research Ltd., UK

^CDepartment of Physiology, Anatomy and Genetics, University of Oxford

ABSTRACT

View-based and Cartesian representations provide rival accounts of visual navigation in humans, and here we explore possible models for the view-based case. A visual “homing” experiment was undertaken by human participants in immersive virtual reality. The distributions of end-point errors on the ground plane differed significantly in shape and extent depending on visual landmark configuration and relative goal location. A model based on simple visual cues captures important characteristics of these distributions. Augmenting visual features to include 3D elements such as stereo and motion parallax result in a set of models that describe the data accurately, demonstrating the effectiveness of a view-based approach.

Index Terms— navigation, visual perception, virtual reality

1. INTRODUCTION

When we view a scene with two eyes and move our heads to and fro we get a powerful sense of the 3D structure of the scene and our location within it. Is the brain really constructing a model of the scene in any 3D frame? Rival accounts of how humans navigate support either a view-based (*i.e.* no 3D reconstruction) or Cartesian representation of the environment [1, 2, 3, 4, 5].

To our knowledge, there are not yet any workable computational models implementing the dominant biological model of visual representation (*i.e.* involving transformations from retinal to egocentric and then world-based reference frames). Our aim in this work is to test the hypothesis that human navigation is instead based on view-based principles such as snap-shot recognition [6], or view-graph navigation [1]. Unlike earlier work on large-scale navigation [7] or online control of movement [8, 9], here we apply the view-based framework to peri-personal space.

Specifically, we show that the distribution of errors in navigation is strongly influenced by the scene geometry in ways that can be modelled using only simple view-based features. This is done by running a homing experiment in immersive virtual reality, where the participant’s view and position can be recorded accurately, and then testing a set of possible view-based models for their ability to describe the types of navigation errors observed. An example of some homing data is shown in Figure 1; the top row shows data gathered from our experiment for two different conditions, and the bottom row shows the same data, this time overlaid on a likelihood map created from our model. In the longer term, our aim is to make a detailed comparison of the predictions of view-based and 3D reconstruction models for a range of behaviours.

This paper begins with an overview of the biological background to our work, then in Section 3 we describe the experimental setup

Thanks to the Wellcome Trust for funding.

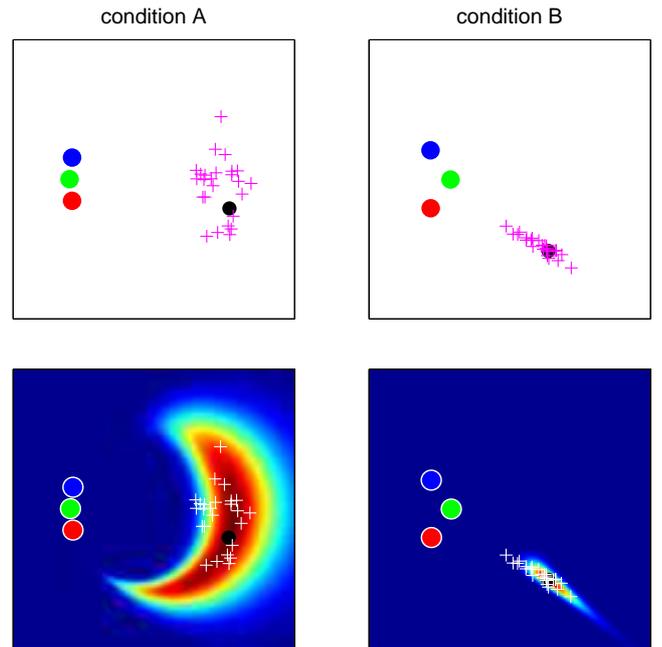


Fig. 1. Two example conditions (left and right), each showing a plan view of a $4m \times 4m$ room in which participants were asked to navigate with only three coloured poles visible as landmarks. Top row: raw data showing three coloured poles, black goal location (not visible to participants), and 25 points (magenta pluses) recorded from 4 different subjects when they thought they had returned to the goal location. Bottom row: likelihood map for these points, based on a model trained using 22 different experimental conditions.

and the data obtained from human participants. Details of the view-based modelling and model fitting are given in Section 4, followed by results and discussion in Section 5.

2. BIOLOGICAL BACKGROUND

There is considerable evidence that insects such as ants, bees and wasps use view-based strategies in order to navigate [10, 11, 12]. They can store a visual “snap-shot” of their goal to guide their return, as is clear from the fact that manipulating the configuration of landmarks around the goal causes the insect to search in a region where the view is similar to the original snap-shot [10, 11].

Navigation in mammals is more flexible and robust than that of insects. Here, it is commonly suggested that “cognitive maps” [5,

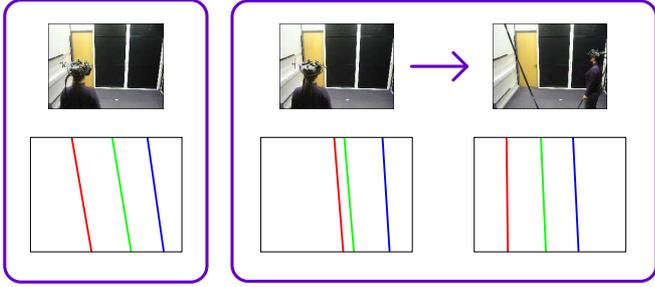


Fig. 2. Left: In interval one, the participant sees a view of the poles. Right: in interval two, the participant starts in a different place relative to the poles, and has to walk to a point where the view matches that of interval one.

13, 14] are used and that the hippocampus and entorhinal cortex may contribute to such a map, representing space in an allocentric (world-based) reference frame [15, 16]. There is some support from behavioural experiments in humans for this view [17], although counter-arguments have also been made [2, 18, 19].

Superficially, the proposed hippocampal representations are similar to the 3D, world-based reconstructions of a scene that are generated in computer vision from multiple views. Photogrammetry (*i.e.* reconstructing scene geometry from photographs) from two, three or more views has been shown to be accurate and robust when applied to a pre-recorded sequence of images [20, 21, 22] and even to simultaneous localisation and mapping (SLAM) using real-time data from a moving camera [23]. However, the principles underlying photogrammetry are very different from those assumed to take place in biological visual systems. For example, in photogrammetry there is no attempt to build a retinotopic depth map of the scene as is observed in the primary visual cortex nor an ego-centric representation of space as is posited in mammalian parietal cortex as an intermediate stage on the way to generating an allocentric map of visual space [3, 4, 24].

3. EXPERIMENTAL SETUP

The experiment was designed to show up the different patterns of errors made by human subjects in a simple visual navigation task. To this end, a very sparse visual world was created in a fully immersive 3D virtual reality environment, and participants were asked to find their way back to a location from which they had viewed the scene previously, based on visual landmarks. Patterns of errors depend on what objects are visible in the VR world, the configuration of those objects, and where the goal point is in relation to them.

The landmarks we used were three differently-coloured infinitely long vertical poles whose angular width was fixed at one pixel irrespective of viewing distance. Figure 2 illustrates the two intervals of the experiment, which proceeded as follows:

Interval 1: Participant sees the set of three poles from a particular viewing point, but limited to an $20cm \times 80cm$ (depth \times width) axis-aligned viewing box centred on viewing point, outside of which the stimulus blanks out. In all cases, the midpoint of the red and blue poles was directly in front of the participant down the z -axis of the room. The participant is encouraged to move around within this box to gain motion parallax information. When the participant wishes to proceed, he or she presses a button on a hand-held pointer. A $0.5s$ blank *inter-stimulus interval* follows interval one, during which nothing appears in the field of view.

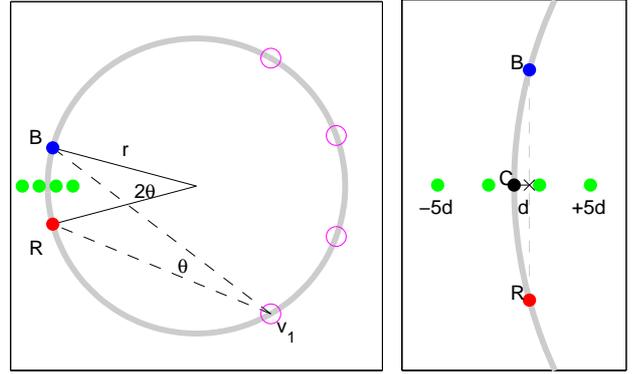


Fig. 3. The overall layout of the experiment, in plan view. The red and blue poles, along with the goal-point, were arranged on a circle of radius r , such that the RB angle is θ anywhere on the circle. The four possible goal points (magenta circles) were positioned at $\pm 20^\circ$ and $\pm 60^\circ$. The spacing of the four possible green pole positions is determined by the distance d from the midpoint of RB to the closest point on the circle.

Interval 2: The poles re-appear in the VR space at a different location relative to the participant, and the participant’s task is then to walk in the VR room until his or her relative position to the poles matches what it was when they pressed the button at the end of the first interval. When he or she believes the goal point has been reached, the participant presses the button to signal the end of the trial.

Interval two is followed by a reset period, where the participant is helped back to the “home” location in the physical room by a plan view displayed in the headset. Note that participants obtain no feedback on how accurate their performance on the homing task was.

If the experiment is considered in a frame of reference in which the poles are static, then the observer first views the poles from the “goal point”, then is transferred to another location (“start point”), and moves along some trajectory to a final “last point” at which the button is pressed again. When a given condition is repeated a number of times, the last-point locations that are associated with the condition’s goal point form samples from some underlying distribution, and it is this distribution we are interested in modelling.

3.1. Further experimental details

The basic layout of the experiment is shown in Figure 3, with four possible goal point positions, and four possible positions of the green pole. To determine the complete set of pole and goal locations, a radius r and viewing angle θ must be specified. We used three pairs for the study: $(r = 0.8m, \theta = 15^\circ)$, $(r = 1.2m, \theta = 15^\circ)$, and $(r = 1.2m, \theta = 20^\circ)$. This generated a total of 48 layouts.

Each of these 48 layouts was coupled with either a full stereo view in both intervals, with the ability to move around the start zone in interval one, or with a synoptic view (*i.e.* one in which the headset was configured to show views as if viewed from a common optic centre by both eyes), and with only a static image available in the first interval. This brings the total number of conditions to 96. Conditions were partitioned at random into three 32-trial blocks, which were deemed comfortable for most participants to work through in a single run. Within each block, trials were presented in a different random order each time.

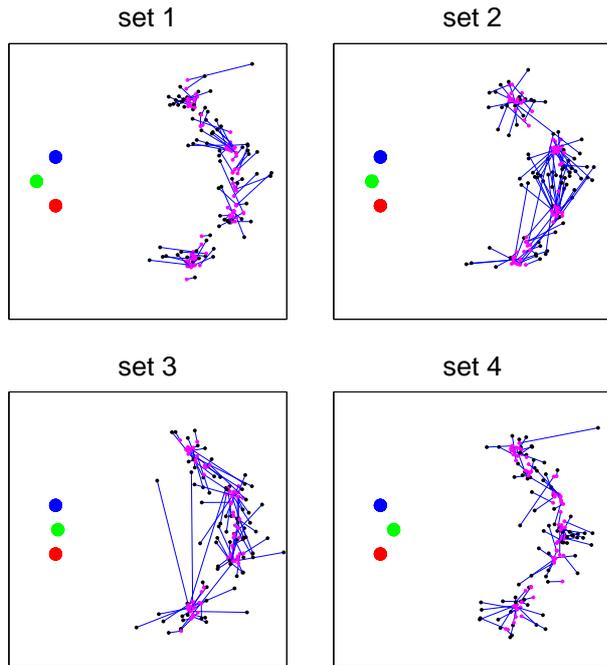


Fig. 4. Data gathered from two of the participants, with goal points in magenta and endpoints in black, and goal-endpoint pairs marked in blue. These data-points all come from the cases where a full stereo view with motion parallax was available to the participants. Four conditions are shown in each plot (pole locations are constant within a plot; goal locations change).

We used a head mounted display (SX111 from nVis) with a wide field of view (108°), a real time head tracker and a computer generating appropriate binocular images according to the observer’s pose and head position. The system had a total latency of less than two frames. Further details can be found in [25].

3.2. Participant data

Two naïve subjects (*i.e.* with no knowledge of the experiment or its goals) completed 10 sets each. The components of their data with $r = 1.2$ and $\theta = 15^\circ$ given full stereo stimuli are shown in composite plots in Figure 4, where the four conditions in which the poles are in any given configuration are shown collapsed onto the same axes.

Notice that there is some spread in the goal points for each different pole configuration. This arises because participants are free to move within the start box, and the goal point they are instructed to return to in interval two is the one which matches the *final* view they had in interval one, at the moment they pressed the button to proceed onto interval two.

In order to highlight the different shapes and extents of the distributions for a subset of the data only, we augmented the experiment to include an intermediate interval. Participants familiarised themselves with the three poles from within the starting box as usual, but were then given a static stereo view of the poles from the exact goal location to which they had to return, then they pressed the button again to proceed to interval two as normal. Some of the data collected in this way are shown in Figure 5, which makes the differences in shape readily apparent (in these cases, $r = 1.2$ and $\theta = 15^\circ$).

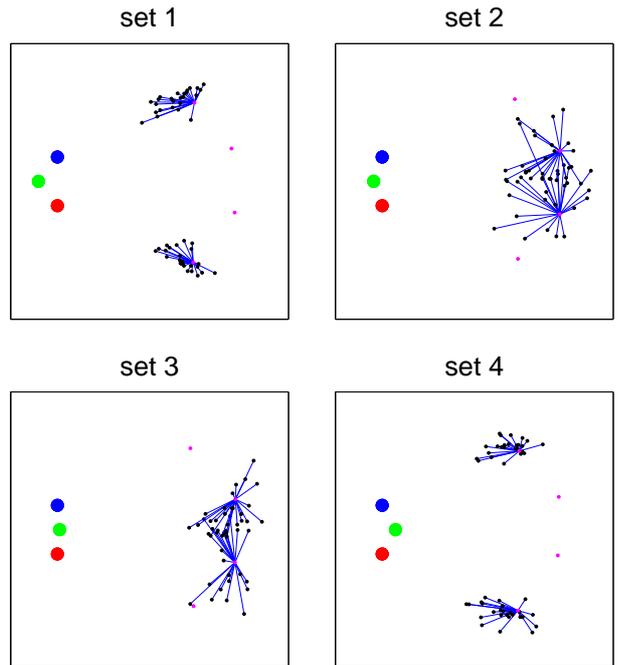


Fig. 5. Extra data for illustrative purposes only, gathered under a variation in experimental protocol which forced goal points for each condition to be aligned exactly. The differing shapes of the distributions is clear here, with large spreads for the less certain cases where the green pole is closer to the circle marked out in Figure 3.

The distributions of end-point errors on the ground plane differed significantly in shape and extent depending on pole configuration and goal location. Where the three poles and the goal point almost lie on a single circle, the errors tend to be distributed around this arc. Where there is a relatively small visual angle between two of the poles, as in set 1 of Figure 5, the error distribution is elongated on the ground plane in such a way as to preserve the ratio of angles between poles from these viewpoints, even if the overall scale (*i.e.* the red-blue angle) is not always accurately reproduced.

4. DATA MODELLING

The purpose of our modelling of these data was to test whether a simple view-based navigation model – *i.e.* a model which assumed *no* 3D reconstruction of the layout of the poles – could provide an accurate account of the errors obtained in the experiment.

We selected a set of simple visual features to describe the views available to the participants. Some were *monocular* single-view features, such as angles between the poles as measured from the cyclopean point (directly between the optic centres for the two eyes), or ratios of these angles. Others were inherently two-view features (stereo or motion) such as disparity, relative disparity, and disparity gradient. The full set of features is defined in detail in Section 4.1.

Feature vectors are calculated at the goal and endpoint locations for each trial. Errors in endpoint location then transform into errors with respect to goal-point features in this feature space, though the mapping between the two is nonlinear. If suitable features are chosen, the error distribution in feature space can be analysed for the full set of trials at once – even though this data includes many dif-

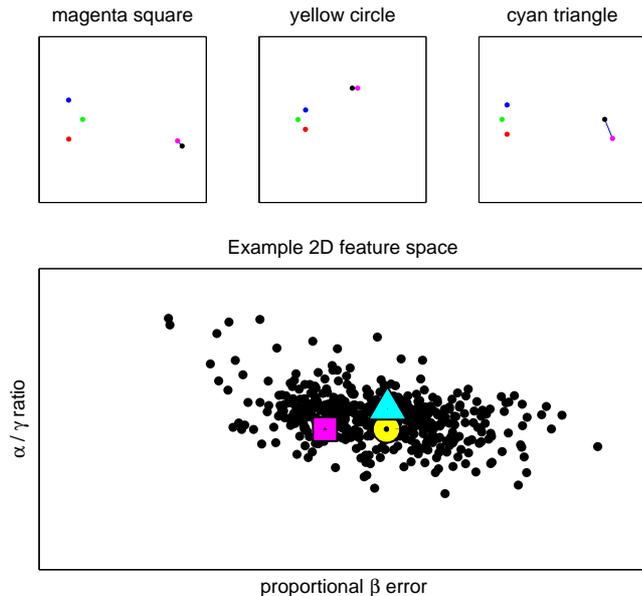


Fig. 6. Three individual datapoints (top row) from different conditions, projected in feature-error space (bottom row). All 480 datapoints (48 stereo conditions, 10 complete runs) are plotted, with the three example cases highlighted in red, green and blue respectively. Note that errors which seem different in “room space” still share the same distribution in feature space.

ferent pole and goal-point configurations. This is illustrated for one pair of features in Figure 6, where spatial errors in three different conditions are mapped into feature-error space (*i.e.* end-point error – goal-point error) along with every other trial carried out by this participant.

The distribution in feature space is very stable, *i.e.* when a Gaussian is fitted to features from one set of conditions, its parameters are much the same as those for a disjoint set of conditions. This means that one single Gaussian in feature space describes *all* the different error shapes in the virtual reality room.

To make a likelihood map of endpoint locations in room space, we simply evaluated the feature-error vectors for a grid of room points, and found those vectors’ likelihoods using the Gaussian mean and covariance from the fitted model. Figure 7 shows likelihood maps for the same three example trials as Figure 6, along with the simple 2D Gaussian fitted to the overall set of points in feature space. Notice how well the shapes of the likelihood distributions match up with the different endpoint patterns seen in Figures 4 and 5.

4.1. Features in detail

We give specifications here for the 20 different features we proposed to describe the views in the view-based model. In Section 4.2 we will explain how we picked a subset of these that best described the experimental data.

The full set of features gives a high-dimensional representation of the data, though not all dimensions are linearly independent.

The first nine features used are *monocular*, because they are calculated from a single viewing location, taken to be the cyclopean centre, *i.e.* the point half-way between the optic centres of the two eyes. The remainder rely on stereo information obtained by compar-

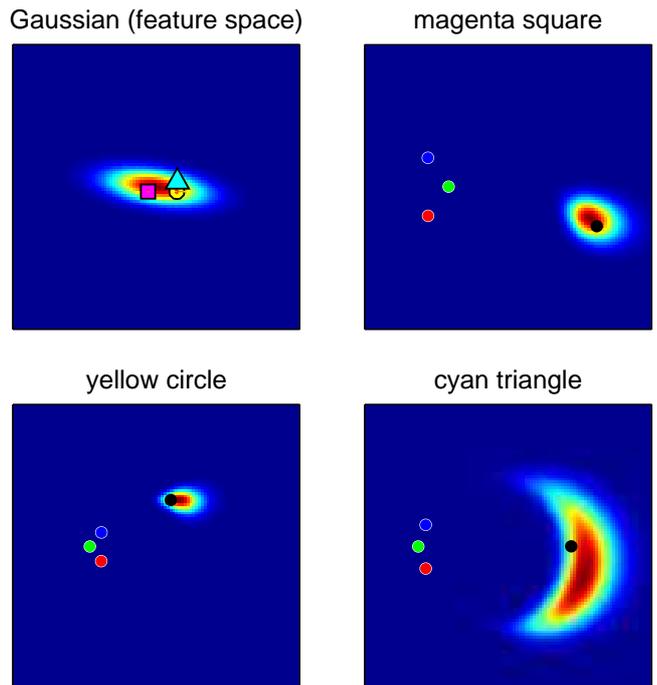


Fig. 7. Top left: the 2D Gaussian fitted to the cloud of points in feature-error-space shown in Figure 6. Other plots: the same three example cases in the $4m \times 4m$ room, showing the different shapes of endpoint likelihoods given by this single Gaussian model. The differences come about because of the variation in pole and goal-point locations.

ing left and right views, or alternatively, views taken from different points within the width of the start box.

The first three features are simply the set of angles between the three poles as viewed from a point: $\{\alpha, \beta, \gamma\}$. These are labeled according to the relative sizes of the three angles as seen from a given viewing point, as illustrated in Figure 8.

If α_G is the alpha angle viewed from the goal location, then α_X is the same angle viewed from some point X in the room measured between the same two *colours* of poles as α_G (even if the ordering on angle size is not the same). The features (feature-error vectors) reported for each trial are $\{\alpha_G - \alpha_X, \beta_G - \beta_X, \gamma_G - \gamma_X\}$. A second trio of monocular features was taken to be simply the proportional error in these angles, since a small absolute error in a small angle is much more significant than a small absolute error in a larger angle: $\left\{ \frac{\alpha_G - \alpha_X}{\alpha_G}, \frac{\beta_G - \beta_X}{\beta_G}, \frac{\gamma_G - \gamma_X}{\gamma_G} \right\}$.

The final three monocular features are ratios of the three simple angles, as participants frequently report using ratios or proportions of angles to guide themselves in the task, for instance “the green pole was a third of the way between the red and blue poles”. As before, the features reported for a trial are $\left\{ \frac{\alpha_G}{\gamma_G} - \frac{\alpha_X}{\gamma_X}, \frac{\alpha_G}{\beta_G} - \frac{\alpha_X}{\beta_X}, \frac{\beta_G}{\gamma_G} - \frac{\beta_X}{\gamma_X} \right\}$. Technically, there is no need to include both $\frac{\alpha}{\gamma}$ and $\frac{\beta}{\gamma}$, since $\frac{\alpha}{\gamma} = 1 - \frac{\beta}{\gamma}$, so only one or other cue were used at once, to avoid linearly dependent sets of features in the models.

Another 11 cues were constructed using information only available when considering binocular stereo or motion parallax. The ver-

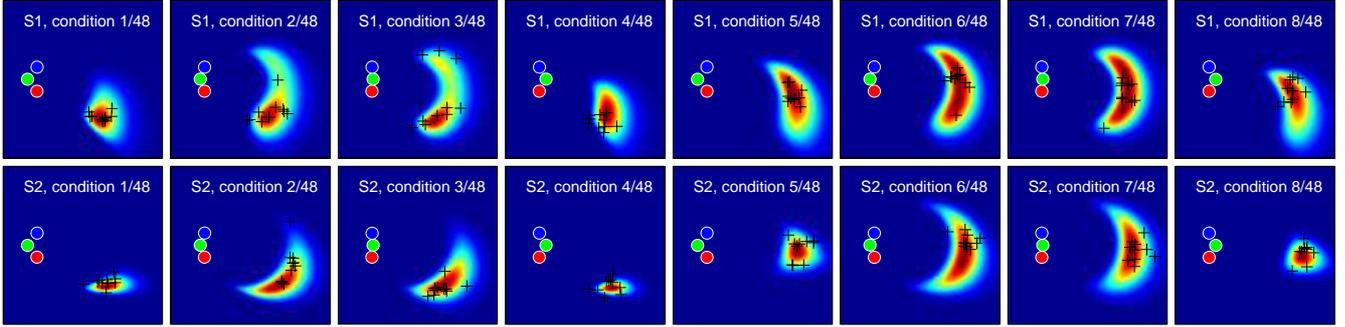


Fig. 9. The first eight (out of 48) conditions under the best four-feature models for two different naïve participants (top row and bottom row). All ten datapoints per participant are plotted, and the distributions shown were calculated for the mean of the ten goal points.

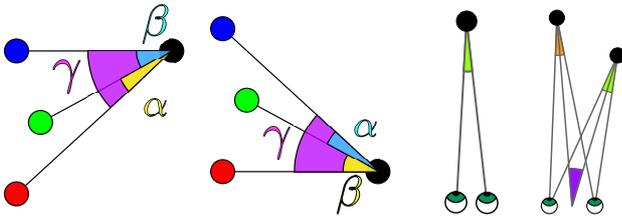


Fig. 8. Left to Right: (a-b) The labelling of these angles depend on their relative sizes: α is always the smallest, and γ is always the largest of the three; (c) The vergence angle for a single coloured pole, as measured from the cyclopean point; (d) *Disparity gradient*: change in vergence angle between two poles, and divided by the angle between them.

gence angle, shown middle-right of Figure 8 gives some indication of the distance an observer is from an object. The first three stereo cues are therefore the vergence angles from the three coloured poles. The next three cues are the three pair-wise differences between the vergence angles, termed the *relative disparity* of the poles, since it gives information about how much farther (or nearer) one object is to an observer than another.

Disparity Gradient is the term given to the relative disparity (difference in vergence angles) divided by the angle between the objects, measured from the cyclopean point. This can be viewed as a finite first-order difference approximation to the slant of a surface on which the two objects sit. As for the first two types of stereo cue, there are three possible ways to compute this pairwise cue from the set of three poles.

There are two stereo cues computed from all three poles at once. First, analogous to the first-order difference cue, we constructed a second-order difference cue which approximates the *curvature* of an underlying surface around the green pole. To do this, the disparity gradient was calculated for the red-green pole pair, subtracted from the disparity gradient of the green-blue pole pair, and finally divided by the overall angle between red and blue (e.g. γ from Figure 8).

The final cue is calculated using an imagined pole. In 2D, an extra point g' is defined as the intersection of the red-blue line with the viewer-green line. This means g' lies in the red-blue plane, projected so that it lines up exactly with the green pole as seen from the goal point. The cue itself is then the relative disparity of g (the green pole) and g' (the green pole's location if it were really in the plane of the other poles). This cue is inspired by work showing sensitivity

to depth relief on slanted planes [26].

4.2. Learning and evaluating models

We are interested in testing how well a view-based model is able to describe the experimental data. To quantify how well a model describes the observations, we turned the the distributions of Figure 7 into probabilities by finding the normalizing constant for each datapoint. This allows us first of all to compare different view-based models to one another, and secondly to quantify how well in general a view-based model can do, for the purpose of making comparisons to other families of models (e.g. explicit 3D reconstruction models) in the future.

The step of finding the normalization constant for every datapoint is computationally expensive. We calculated the 2D integral over the ground plane numerically using MATLAB's "dblquad" routine. Note, however, that this step is only required for model *comparison*, rather than for actual model *use*, e.g. in some visual navigation strategy.

In order to find good view-based models, all possible models using 1–4 linearly independent features were evaluated, and ranked according to data likelihood. Excluding feature combinations with linear dependencies, this gave almost 6000 combinations from the 20 available visual features. For each feature combination, each datapoint x in turn was excluded from the dataset, and the Gaussian \mathcal{G} was fitted in feature-error space on the remainder of the points (i.e. points used in testing were not part of the training set for each individual trial). The integral over the ground plane was evaluated out to $\pm 5m$, and used to normalize the likelihood of x under \mathcal{G} . The total likelihood of the dataset for the chosen feature combination was then the product of all these normalized likelihoods (assuming independence of trials).

5. RESULTS

Many of the view-based models we tested were very successful at describing the shapes of the end-point distributions observed in homing experiment. The results of the comparison between sets of visual features varied between participants, as did the maximum data likelihood obtained. For the single-view data, models based on the γ angle and the $\frac{\alpha}{\beta}$ ratio describe errors successfully. Adding in two-view features allows for better modelling of the stereo-view half of the dataset for most participants, and while the best features seem to be subject-dependent, vergence angles featured in most of the best models.

Of the two naïve participants who had each completed 10 runs, one appeared to make extensive use of these stereo features, while the other's "best" features came entirely from the monocular set. The set of best-performing features for each of these two participants (according to the procedure described in Section 4.2) was used to create two separate 4D feature-error-space models (one per participant). These two distributions were used to draw the corresponding rows of plots in Figure 9, representing the first eight of the full 48-condition set. The top row shows the results for S1 (who used stereo features), while the bottom shows the equivalent likelihood maps and end-points for S2. Note that average end-points were used for these plots, so some mis-match of endpoints and distribution centres is due to the fact that the actual goal locations varied.

6. CONCLUSIONS

In this paper, we have explored a number of ways in which a view-based model might describe datasets gathered from human subjects performing a simple homing task. The results fit the different shapes of the data distributions, which are due to different configurations of visual landmarks, very convincingly, and appear to support to the hypothesis that humans employ a view-based strategy when faced with this simple homing task in virtual reality.

This now puts us into a position to compare view-based models to equivalent 3D models in a bid to probe possible mechanisms described in Section 2. Preliminary modelling using a 3D reconstruction algorithm (not shown here) has shown a quite different pattern of predicted behaviour. Future experiments will also explore the pattern of navigation errors in a richer or more natural visual environment.

7. REFERENCES

- [1] H. A. Mallot and S. Gillner, "Route navigation without place recognition: what is recognized in recognition-triggered responses?," *Perception*, vol. 29, pp. 43–55, 2000.
- [2] P. Foo, W. H. Warren, A. Duchon, and M. J. Tarr, "Do humans integrate routes into a cognitive map? Map versus landmark-based navigation of novel shortcuts," *J. Exp. Psych. : Learning, Memory and Cognition*, vol. 31, pp. 195–215, 2005.
- [3] R. A. Andersen, L. H. Snyder, D. C. Bradley, and J. Xing, "Multi-modal representation of space in the posterior parietal cortex and its use in planning movements," *Ann. Rev. Neuroscience*, vol. 20, pp. 303–330, 1997.
- [4] N. Burgess, K. J. Jeffery, and J. O'Keefe, *The hippocampal and parietal foundations of spatial cognition*, Oxford: OUP, 1999.
- [5] E. A. Maguire, N. Burgess, and J. O'Keefe, "Human spatial navigation: cognitive maps, sexual dimorphism, and neural substrates," *Current Opinion in Neurobiology*, vol. 9, no. 2, pp. 171–177, 1999.
- [6] M. O. Franz, B. Scölkopf, H. A. Mallot, and H. H. Büthoff, "Where did I take that snapshot? Scene-based homing by image matching," *Biological Cybernetics*, vol. 79, pp. 191–202, 1998.
- [7] S. Gillner and H. A. Mallot, "Navigation and acquisition of spatial knowledge in a virtual maze," *Journal of Cognitive Neuroscience*, vol. 10, pp. 445 – 463, 1998.
- [8] J. J. Gibson, *The ecological approach to visual perception*, Boston: Houghton Mifflin, 1979.
- [9] B. R. Fajen, W.H. Warren, S. Temizer, and L. P. Kaelbling, "A dynamical model of visually-guided steering, obstacle avoidance, and route selection," *International Journal of Computer Vision*, vol. 54, pp. 13–34, 2003.
- [10] B. A. Cartwright and T. S. Collett, "Landmark learning in bees: experiments and models," *Journal of Comparative Physiology*, vol. 151, pp. 521–543, 1983.
- [11] R. Wehner and F. Rüber, "Visual spatial memory in desert ants, *cataglyphis bicolor* (Hymenoptera: Formicidae)," *Cellular and Molecular Life Sciences*, vol. 35, pp. 1569–1571, 1979, 10.1007/BF01953197.
- [12] P. Graham and T. S. Collett, "View-based navigation in insects: how wood ants (*formica rufa* L.) look at and are guided by extended landmarks," *J Exp Biol*, vol. 205, no. 16, pp. 2499–2509, 2002.
- [13] E. C. Tolman, "Cognitive maps in rats and men," *Psychological Review*, vol. 55, no. 4, pp. 189–208, 1948.
- [14] C. R. Gallistel, "Animal cognition, the representation of space, time and number," *Ann. Rev. Psychology*, vol. 40, pp. 155–189, 1989.
- [15] J. O'Keefe and L. Nadel, *The Hippocampus as a Cognitive Map*, Oxford University Press, Oxford, UK, 1978.
- [16] B. L. McNaughton, F. P. Battaglia, O. Jensen, E. I Moser, and M-B Moser, "Path integration and the neural basis of the 'cognitive map'," *Nature Reviews Neuroscience*, vol. 7, pp. 663–678, 2006.
- [17] T. Hartley, I. Trinkler, and N. Burgess, "Geometric determinants of human spatial memory," *Cognition*, vol. 94, no. 1, pp. 39 – 75, 2004.
- [18] A.T. Bennett, "Do animals have cognitive maps?," *Journal of Experimental Biology*, vol. 199, pp. 219224, 1996.
- [19] R. F. Wang and E. S. Spelke, "Human spatial representation: insights from animals," *Trends in Cognitive Sciences*, vol. 6, pp. 376382, 2002.
- [20] A. W. Fitzgibbon and A. Zisserman, "Automatic camera recovery for closed or open image sequences," in *LNCS 1406: Computer Vision—ECCV '98*. 1998, pp. 311–326, Springer.
- [21] R. I. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, Cambridge, UK: Cambridge University Press, second edition, 2004.
- [22] 2d3 Ltd, "Boujou 2," 2003, <http://www.2d3.com>.
- [23] A. J. Davison, "Real-time simultaneous localisation and mapping with a single camera," in *Proceedings. Ninth IEEE International Conference on computer vision*, 2003, pp. 1403–1410.
- [24] L. H. Snyder, K. L. Grieve, P. Brotchie, and R. A. Andersen, "Separate body- and world-referenced representations of visual space in parietal cortex.," *Nature*, vol. 394, pp. 887–891, 1998.
- [25] S. J. Gilson, A. W. Fitzgibbon, and A. Glennerster, "Spatial calibration of an optical see-through head mounted display," *Journal of Neuroscience Methods*, vol. 173, pp. 140–146, 2008.
- [26] A. Glennerster and S. P. McKee, "Bias and sensitivity of stereo judgements in the presence of a slanted reference plane," *Vision Research*, vol. 39, pp. 3057–3069, 1999.