Camouflage Assessment: Machine and Human

Timothy N. Volonakis[1,3*], Olivia E. Matthews[1], Eric Liggins[4], Roland J. Baddeley[1], Nicholas E. Scott-Samuel[1], Innes C. Cuthill[2]

*Corresponding author

[1] School of Experimental Psychology, University of Bristol, 12a Priory Road, Bristol, BS8 1TU, UK

[2] School of Biological Sciences, University of Bristol, 24 Tyndall Avenue, Bristol, BS8 1TQ, UK

[3] Centre for Machine Vision, Bristol Robotics Laboratory, University of the West of England, Frenchay Campus, Coldharbour Lane, Bristol, BS16 1QY, UK

[4] QinetiQ Ltd, Cody Technology Park, Farnborough, Hampshire, GU14 0LX, UK

1   **Abstract**

2   A vision model is designed using low-level vision principles so that it can perform as a

3   human observer model for camouflage assessment.  In a camouflaged-object assessment

4   task, using military patterns in an outdoor environment, human performance at detection

5   and recognition is compared with the human observer model.   This involved field data

6   acquisition and subsequent image calibration, a human experiment, and the design of the

7   vision model.  Human and machine performance, at recognition and detection, of military

8   patterns in two environments was found to correlate highly.  Our model offers an

9   inexpensive, automated, and objective method for the assessment of camouflage where it is

10  impractical, or too expensive, to use human observers to evaluate the conspicuity of a large

11  number of candidate patterns. Furthermore, the method should generalize to the

12  assessment of visual conspicuity in non-military contexts.

13

**Key Words: Camouflage Assessment, Observer Modelling, Visual Search**

14  **Declarations of interest**

17

18

19

20

21

22

28

29   **1.  Introduction**

30   Military personnel and equipment need protection from detection during conflict.

31   Camouflage is the primary method to achieve this, frequently through coloured textures

32   that match the background and/or disrupt the object's outline (Hartcup 2008; Merilaita et

33   al. 2017; Talas et al. 2017). Assessment of effectiveness can be carried out in a number of

34   ways. The most intuitive method is to use human participants as observers. Such an

35   apparently straightforward procedure, however, is not only limited by uncontrollable

36   conditions, such as the weather: it is also impractical given the large variety of

37   objects/patterns that one might want to evaluate and the range of environments one might

38   want them to be assessed in. Field trials are also expensive and, in some circumstances, may

39   not even be possible. They also do not lend themselves to precise isolation of exactly what

40   leads to the failure of camouflage, something that a paired comparison of otherwise

41   identical target-present and target-absent scenes would allow. Photo-simulation attempts

42   to overcome weather constraints and problems with inaccessible environment-types by

43   using photographic or synthetic imagery. Recent advances in synthetic rendering are

44   impressive; however, current methods are still computationally expensive and the images

45   are unrealistic at small spatial scales due to the current limitations of simulating realistic ray

46    scattering. Furthermore, human experiments are necessarily subjective and do not readily

47    allow evaluation of camouflage against autonomous systems perhaps operating using

48    different spectral bandwidths than the human vision. A computational approach is therefore

49    helpful in overcoming the limitations of assessing camouflage when using human observers.

50    Such a computational model should be ideally designed, in the first instance, in accordance

51    with the human visual system, since it will be performing the task of a human observer and,

52    if it is to replace subjective assessment, needs to be compared with human performance.

53    More generally, however, such a system could be adapted to have a different 'front end'

54    (e.g. infra-red sensor, hyperspectral sensor). Therefore it is surprising that a biologically

55    motivated design for the assessment of camouflage has not been implemented.

56    This omission means that the confidence and extendibility of current models and metrics

57    are low, falling short in their ability to cope with high dynamic range (i.e. natural) (Bhajantri

58    and Nagabhushan, 2006; Hecker, 1992; Sengottuvelan et al., 2008), semi-automatic labelling

59    or tracking of the target (Chandesa et al., 2009), non-probabilistic and non-scalable distance

60    metrics to high dimensional data or multiple observations given many images (Birkemark,

61    1999; Heinrich and Selj, 2015; Kiltie et al., 1995). Human behavioural data need to be

62    recorded to assess the coherence between human and model observers. This requires

63    tasking human and model observers with the same experiment, based on a stimulus set

64    from the real world: outdoor environments and militarily relevant objects.

65

66    **2. Method**

67    An experiment was devised so that human participants and a model observer could both be

68    tasked with it, allowing for direct comparison.  This method section is broken down into the

69    three components that comprise this study: (i) images of objects placed in real world scenes

70     were photographed and calibrated; (ii) a human experiment, using a protocol from

71     psychophysics, recorded unbiased performance for recognition and detection of these

72     objects; and (iii) the design of the visual observer model, and modelling the discrimination

73     task.

74

75     **2.1 Stimuli**

76     Targets were photographed in two outdoor environments in the UK: Leigh Woods National

77     Nature Reserve in North Somerset (2°38.6' W, 51°27.8' N), which is mixed deciduous

78     woodland, and Woodbury Common in Devon (3°22' W, 50°40' N), a heathland used for

79     Royal Marine training. A replica military PASGT helmet (Personnel Armor System for Ground

80     Troops, the US Army's combat helmet until the mid-2000's) was the chosen object used in

81     the experiment and visibility was manipulated by changes in helmet covers varying in both

82     colour and textural appearance (Figure 1). The camouflage patterns worn by the helmet

83     were United Nations Peacekeeper Blue (UN PKB), Olive Drab, Multi-Terrain Pattern (MTP, as

84     used by the British Army since 2012), Disruptive Material Pattern (DPM, the dominant

85     British Army pattern prior to the adoption of MTP), US Marine Pattern (MarPat) and, for the

86     Woodbury Common experiment, Flecktarn (as used by the Bundeswehr, the German Army).

87     These patterns were chosen not for the purpose of evaluation per se, but to reflect a range

88     of styles (e.g. unpatterned Olive Drab, DPM as a subjective human design, MTP and MarPat

89     based on spatio-chromatic analysis of natural scenes, but MarPat being 'digital' or

90     pixellated), with UN PKB as a high visibility control.

91     For the computational approach to be useful, the spectrum of visibility across the patterns

92     should be highly correlated in the model and human observers. Scene locations were

93     selected on a meandering transect through the habitats, at 20 m intervals and alternating

94   left and right. If the predetermined side was inaccessible or inappropriate due to occlusions

95   then the opposite side of the transect path was used, and if neither side was accessible the

96   interval was ignored and the next location in the transect was used. At each location the

97   object was placed in a 3 × 3 grid resulting in nine images. The distance of each row of the

98   grid was 3.5, 5 and 7.5 metres. The scene was also divided into 3 arcs: left, middle and right.

99   The combination of distance and left-right positioning mean that, in the subsequent tests on

100  humans, the location of the target within the scene was unpredictable. This resulted in nine

101  images of each helmet per location for analysis, plus a scene including a Gretag-Macbeth

102  Color Checker chart (X-Rite Inc., Grand Rapids, Michigan, USA) for calibration. The

103  orientation of the helmet in each photograph was set an angle drawn randomly from the

104  uniform distribution {0, 45, 90, 135, 180, 225, 270, 315°}. For efficiency of implementation,

105  the list of random angles was generated before going into the field. Each scene was also

106  photographed without a helmet present. Photographs were taken using a Nikon D80 digital

107  SLR (Nikon Ltd., Tokyo, Japan) with focal length 35mm, exposure 1/30 and F-Number 16.

108  RAW images (Nikon NEF format) were captured and these were subsequently converted to

109  uncompressed 8-bit TIFF and calibrated. Images were calibrated by recording luminance and

110  chromatic spectral values of the Gretag-Macbeth colour chart in the field using a Konica

111  Minolta Chroma Meter CS - 100A colour and luminance meter (Konika, Tokyo, Japan). This

112  process was repeated three times to average over the natural variation in lighting from

113  moment to moment. The spectral values were transformed to the CIE sRGB colour space

114  after first converting them to the CIE XYZ colour space. The process was then repeated in

115  the lab from a projected image from the projector. A cubic polynomial approximated the

116  relationship between the two sets of RGB measurements. Images were then calibrated

117  using the coefficients of the polynomial for each RGB channel. Not only does this procedure

118　avoid having a colour chart in every single image, but also it calibrates the entire pipeline in

119　a single step: calibrating the camera, projector and images individually could result in over-

120　fitting or multiplicative errors.

121



122

123　**Figure 1. Example cropped helmet images from real world scenes**
124　 An example of each camouflaged helmet cropped for recognition purposes. From left to
125　right the patterns that the helmet wears are DPM, MarPat, MTP, UN PKB, Olive drab and
126　Flecktarn. The top row are the helmets from Leigh Woods and the bottom row are helmets
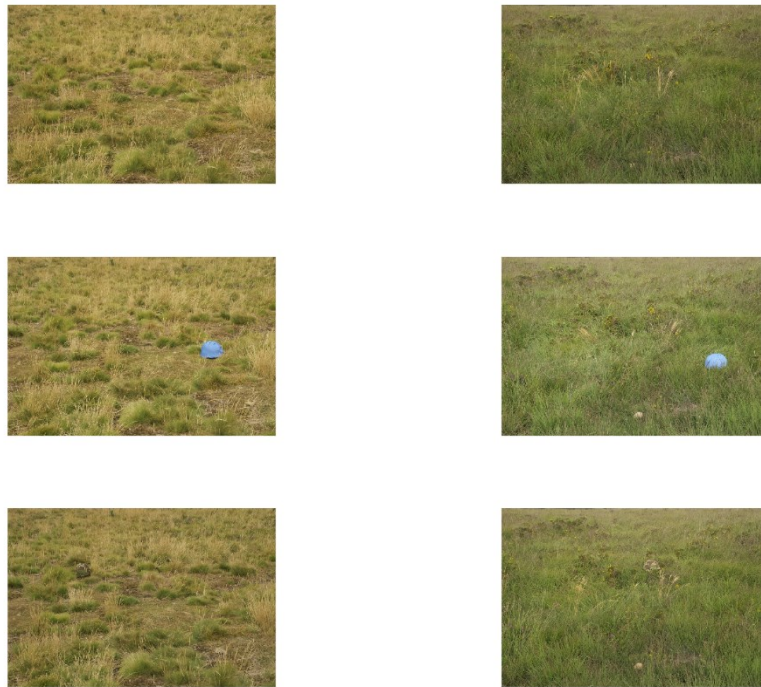127　from Woodbury Common. Flecktarn was only used in Woodbury Common.
128
129

130
131
132 **Figure 2. Example Leigh Woods scenes**
133 Two example scenes from the Leigh Woods environment. The left column and the
134 right column are two different scenes. The top two scenes do not contain a helmet. The
135 middle two contain a UNPKB helmet. The bottom two contain the DPM helmet.
136
137

Figure 3.  Example Woodbury Common scenes
Two example scenes from the Woodbury Common environment. The left column and the right column are two different scenes. The top two scenes do not contain a helmet. The middle two contain a UNPKB helmet. The bottom two contain the DPM helmet.

**2.2 Human Experiment**

**2.2.1 Participants and Materials**

A human experiment using 22 participants for the Leigh Woods dataset and another 20

participants for the Woodbury Common dataset was conducted.  Each of the two

experiments had an equal proportion of each gender. Images were projected onto a 190 ×

107cm screen (Euroscreen, Halmstad, Sweden) from 310cm using a 1920 × 1080 pixel HD

(contrast ratio 300,000:1) LCD Projector (PT-AE7000U; Panasonic Corporation, Kadoma,

154    Japan). Participants were seated at a distance of 255 cm from the screen and therefore

155    images subtended 41° horizontally and 24° vertically.

156

157    **2.2.2 Procedure**

158    At the start of each block participants were informed which helmet to search for by

159    presenting an image of the helmet; only one camouflage type was present in any one block.

160    There were 27 and 22 trials per block, respectively, for Leigh Woods and Woodbury

161    Common, and the order of patterns across blocks and replicated within blocks were

162    separately randomised for each participant. A trial consisted of sequentially presenting two

163    scenes for 250 ms with a 250ms blank screen, of luminance and chromaticity equal to the

164    mean of all the test images, immediately followed by a 250 ms cue screen, prior to each

165    scene. One of the scenes presented contained a helmet and the other did not, the order

166    being randomised. The participant's task was a two alternative force choice, reporting which

167    of the two scenes contained the helmet. Responses were given using the number keys one

168    and two on the keyboard, reporting the first scene or the second scene respectively during a

169    1000 ms response period after each pair of scenes.   A presentation time of 250 ms is

170    enough time to ensure a single saccade to the cued location, but not long enough to allow

171    more complicated scan patterns. Since these scan patterns will be possibly highly variable,

172    they will introduce variability into the responses above and beyond that due to the stimuli,

173    and hence decrease the power of the study. One thousand milliseconds were allowed for

174    subject's responses to allow more than adequate time to respond, but not so long as to
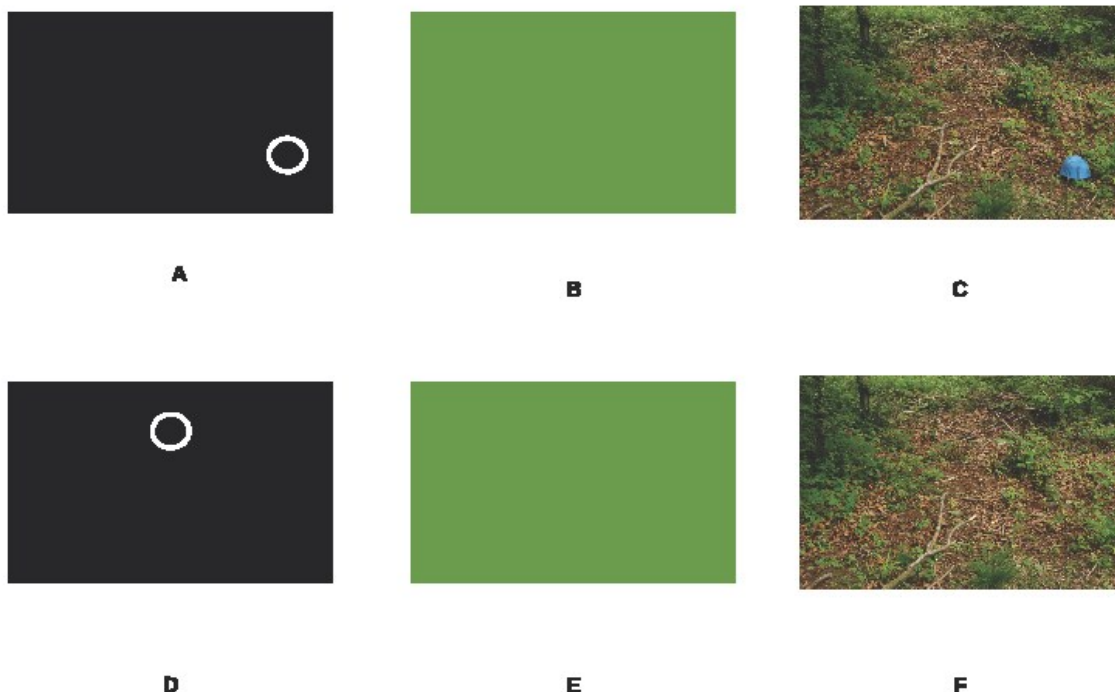
175    increase the time of the total experiment.  There were four general conditions of viewing,

176    the factorial combination of two levels of colour information and two levels of location

177    cueing. Cueing was of interest to separate effects of pattern recognition from detection,

178    because the model was initially designed for recognition. Colour was of interest because it

179    has been suggested that camouflage is more effective when there is chromatic as well as

180    spatial noise (Melin et al., 2007; Morgan et al., 1992). In the first cueing condition, ('cued'),

181    participants were cued to the location of the helmet.  In the scene that did not contain the

182    helmet, this cue's location was a random selection of one of nine possible pre-determined

183    target locations. In the second condition, ('uncued'), the cue was presented in the centre of

184    the screen for both scenes. The spatial cue was a white circle, 50 pixel diameter, 5 pixel line

185    width, circle that was presented for 250ms. The whole experiment was repeated in

186    greyscale and colour. As with pattern, the order of conditions for each participant was

187    randomised.



Human Experiment Story Board

188
189

190    Figure 4.  Human experiment storyboard

11

191 Storyboard for one trial in the experiment. Sequence is in alphabetical order.  Duration of
192 each interval was 250msec.  Either **C** or **F** contains the helmet.  Intervals **A** and **D** cue the
193 participant to the spatial location of the helmet.   Intervals **B** and **E** present a blank interval of
194 average chromaticity across all scenes.  At the end of the sequence, participants are asked
195 which scene the helmet was in and are given 1000msec to respond.  The procedure is
196 identical for the uncued condition however the spatial cue in **A** and **D** are uninformative.

197

198

199

200

201 **2.3 The Human Observer Model**

202

203 **2.3.1 The Model Framework**

204 The model is a four-stage process as outlined below. By modelling low level visual

205 processing, a side effect of the features chosen produces Gaussian variation from small

206 metric distortions.  The resultant Gaussian variation can then be approximated using a

207 mixture of multivariate Gaussian distributions. The centre of each Gaussian distribution

208 stores a familiar view. Probabilistic principal components (Tipping and Bishop, 1999b)

209 describe the variability in an interpretable way to recognise unseen and unfamiliar views.

210 Estimating the density and evaluating the maximum posterior probability determines the

211 object class. This method turns the difficult problem of learning a complex invariant

212 representation of an object into the simple problem of estimating parameters of a mixture

213 of multivariate Gaussian distributions.

214

215 **Stage 1. Filter Images with a Log Gabor Filter Bank**

216 Grey scale images are cropped to a square and resized to 128×128 pixels, preserving the

217    aspect ratio of the object. They are then filtered by a log Gabor wavelet filter bank,

218    comprising three spatial scales (wavelengths of 16, 32 and 64 pixels), and four orientations

219    (0, 45, 90 and 135°) (Kovesi, 2000). This first stage captures the early linear properties of the

220    visual system. Whilst 2D Gabors can be used to approximate simple cells (Daugman, 1985;

221    Jones and Palmer, 1987), we know that (i) simple cells are tuned to spatial frequency with a

222    Gaussian bell-shaped tuning curve on a log frequency scale (De Valois et al., 1982; Field,

223    1987) and (ii) the Gabor filter has a D.C. component. The power in natural images is

224    dominated by the D.C. component (Field, 1987), and, given that the cosine Gabor is

225    sensitive to it and the sine Gabor is not, it will corrupt any computation of phase

226    information in the next stage. The solution to both these problems is to employ log Gabors

227    instead, which do not have a D.C. component (Kovesi, 1999).

228

229    **Stage 2. Process the Filtered Output**

230    Next we compute local energy and phase from the filtered output in stage 1. Stage 2

231    accounts for two non-linear properties of the visual system, illumination invariance and shift

232    invariance. The energy is logged and the effect is two-fold: (i) the energy is positive, and not

233    symmetrical for Gaussian approximation in the fourth stage; and (ii) introducing logarithms

234    will turn differences in illumination into additive offsets. Denoting the response of the real

235    and imaginary filters as $R(x,y)$ and $I(x,y)$, where x and y indicate the index in the image and

236    atan2 computes the four quadrant arc tangent, log local energy and phase can be computed

237    as Energy = $\ln|R(x,y)+I(x,y)|+c$ and Phase = $\text{atan2}(I(x,y),R(x,y))$, where c is a small

238    constant, 0.05, to avoid the undefined logarithm of zero and || is the absolute. The absolute

239    is the magnitude of the real (cosine log Gabor) and imaginary (sine log Gabor) filters. The

240   sum of the squared filter responses is the magnitude, since $sin^2 + cos^2 = 1$. The energy

241   loses local position, but confers some translational invariance and therefore small shifts are

242   turned into small variations.  Local energy represents lines as symmetrical Gaussians.

243   Therefore the variance of these features is Gaussian through small metric distortions such

244   as shift and object pose.

245

246   Phase angles will cycle from π to −π as the distortion moves through sampling locations,

247   resulting in correlated variation. Phase information is a polar, circular variable; in order to

248   use this feature for Gaussian approximation one must convert this feature into Cartesian

249   space. Therefore the sine and cosine of the phase are computed, doubling the number of

250   dimensions required for phase information. Concatenating this sampled local logged energy,

251   sine and cosine phase information creates the feature vector.

252

253   **Stage 3. Sample the Local Energy and Phase.**

254   A hexagonal lattice, of equal size to the image, is placed over the image and the local energy

255   and phase is sampled at the centres of each hexagon. A hexagonal lattice provides optimal

256   sampling where samples are equidistant from each other (Yfantis et al., 1987). Phase angles

257   vary less at larger spatial scales and therefore, to avoid over complete and redundant

258   sampling, hexagonal lattices at larger spatial scales have fewer hexagons.

259

260   **Stage 4. Evaluate Recognition Decision Using Bayes' Rule**

261   The Gaussian variation computed in stage 2 can now be approximated. A unimodal

262   distribution can represent a single view of an object. A mixture of Gaussians can model a

263  multimodal distribution where multiple views of an object are learnt. The dimensions of

264  each Gaussian component should represent the local variation of that view. The

265  concatenation of the local energy and phase results in a high-dimensional feature vector

266  and therefore a mixture of probabilistic components (Tipping and Bishop, 1999a,b) or a

267  mixture of factor analysers (Ghahramani and Hinton, 1996) provides a local subspace for

268  each Gaussian component and approximates the high dimensional covariance structure. To

269  evaluate the recognition of an object, a model is created explicitly for each class. Likelihoods

270  are computed for each explicit class and the posterior probability that an unseen object

271  came from each object class is then evaluated using Bayes' rule, $P(A|B) = P(A)P(B|A)$. Where

272  $P(A|B)$ is the posterior probability that the data A is from the object class B and $P(B|A)$ is the

273  likelihood of data A under the object class B. The prior probability $P(A)$ is equal for each

274  object class and this therefore cancels out.

275

276  **2.3.2 Modelling the 2AFC Recognition Task**

277  Human participants were tasked with recognising a helmet given two different images.  One

278  of the images contained a helmet and the other did not.  For a direct comparison, both

279  observers need to be tasked in a similar way.  Ten-fold cross-validation was used to assess

280  the model's accuracy.  However, instead of evaluating a single image at a time, two images,

281  one with a helmet and one without, were both evaluated under both background and

282  helmet models.  Therefore each image needs to be evaluated under both models, producing

283  four likelihoods (Fig. 5). There are two scenarios: either the helmet is in image A or it is in

284  image B. In the first scenario the helmet is in image A, where there is a high likelihood that it

285  came from the helmet model and so the likelihood that image B came from the background

286  class will therefore have a high likelihood. Bayes' rule will integrate over the mutually

287 exclusive probabilities as shown in the diagram by incorporating the four likelihoods

288 P(A|Helmet), P(A|Background), P(B|Helmet) and P(B|Background).  Using Bayes' rule, the

289 probability that image A is a helmet is simply:

290

291   1.  $P\left(Helmet|A\right) = \frac{P(A|Helmet)\times P(B|Background)}{P(A|Helmet)\times P(B|Background)+P(B|Helmet)\times P(A|Background)}$ .
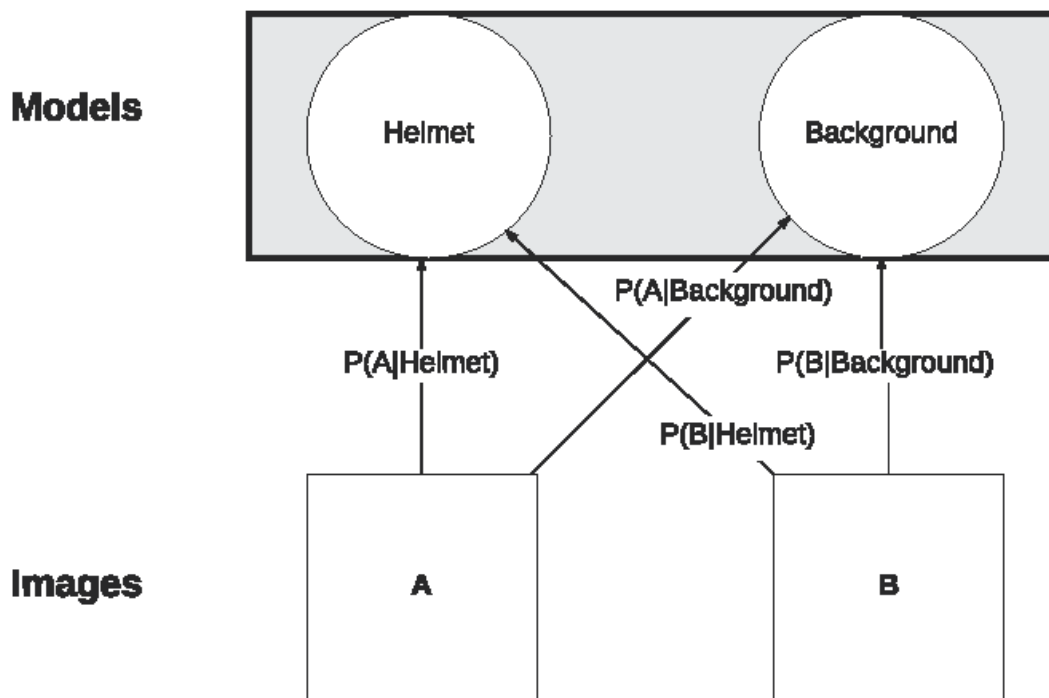
292

293

294

295

## Modelling the Two Alternative Force Choice Task



296
297

298  **Figure 5.**  Graphical illustration at modelling the 2AFC procedure
299  To model the 2AFC task that humans were given, likelihoods under both models are

300    computed for both images.

301

302    **2.3.3 Modelling the Detection Task**

303    The model is trained on a series of cropped images, where the object fills the crop. If the

304    model is presented with an image of the target at a different spatial scale, i.e. the object

305    does not fill the crop, it would be unable to recognise the object. To accommodate scale,

306    likelihoods are computed for both the helmet and background classes at different spatial

307    scales, at intervals of 10 ranging from the smallest helmet to the largest helmet across all

308    images. Weightings are computed for each scale using Bayes' rule by evaluating which scale

309    is most probable from the helmet class whilst evaluating that the other spatial scales belong

310    to the background class. The weightings are multiplied with the likelihoods from each scale

311    and summed. In short this procedure integrates probabilities over all spatial scales into a

312    single likelihood for classification. This probabilistic approach, graphically demonstrated

313    below where A and B denote two different sized crops at location in an image, is superior

314    over simply taking the maximum, because the maximum only considers one model and if

315    two scales are likely under the probabilistic approach the maximum would be too brittle and

316    would ignore one of the likely scales. Equations below 2 - 7, show how Bayes' rule integrates

317    the likelihoods over all the spatial scales, denoting two spatial scales A and B.  Detection was

318    modelled using leave-one-out cross-validation instead of the 2AFC approach. This was

319    because there were too few scenes to compare the helmet scenes with. Problematically, if

320    one were to compare likely peaks between two scenes, one scene would always have the

321    same area of interest and this would be compared to many helmets. Leave-one-out cross-

322    validation also provides a straightforward way to manipulate the training data so that the

323    model did not see any of the scene whilst detecting the helmet.

324

325

$$2. \quad P(Helmet|A,B) = \frac{P(A|Helmet) \times P(B|Background)}{P(A|Helmet) \times P(B|Background) + P(B|Helmet) \times P(A|Background)}$$

326

$$3. \quad P(Helmet|B,A) = \frac{P(B|Helmet) \times P(A|Background)}{P(A|Helmet) \times P(B|Background) + P(B|Helmet) \times P(A|Background)}$$

327

$$4. \quad L1 = P(A|Helmet) \times P(Helmet|A,B) + P(B|Helmet) \times P(Helmet|B,A)$$

$$5. \quad L2 = P(A|Background) \times P(Helmet|A,B) + P(B|Background) \times P(Helmet|B,A)$$

328

$$6. \quad Posterior\ probability\ that\ helmet\ is\ at\ (x,y) = \frac{L1}{L1+L2}$$

329  $$7. \quad Posterior\ probability\ that\ helmet\ is\ absent\ at\ (x,y) = \frac{L2}{L1+L2}$$

330

331 Equations 2-7 elaborate an example of how the model evaluates over spatial scale, where **A**

332 and **B** denote two images each at a different spatial scale.

333

334 **2.3.4 Colour**

335 There are three main issues to consider when including colour: i) colour in the periphery, ii)

336 efficient feature combination of texture and colour and iii) appropriate choice of colour

337 space for measuring the distance between colours. The representation of short, medium

338 and long wavelength receptors on its own is insufficient because computed distances in the

339 colour space do not correlate with human perception (Tkaclic and Tasic, 2003; Wyszecki and

340 Stiles, 1982). Projections in the CIE L*a*b* colour space are consistent with the judgements

341 of human observers and are appropriate for discrimination purposes (Renoult et al., 2015).

342 The model is a human observer model. Whilst recognition accuracy should be high, similar

343   to human observers, it should not be able to recognise camouflaged objects all the time.

344   The aim of the model is not to break camouflage and achieve perfect recognition. Therefore,

345   instead of opting to use the CIE L*a*b* colour space, the MacLeod-Boynton chromaticity

346   diagram is used. The MacLeod-Boynton chromaticity diagram (MacLeod and Boynton, 1979)

347   is an isoluminant cone excitation space that is particularly good at discriminating large

348   chromatic differences (Renoult et al., 2015). Modelling the detection of camouflaged

349   helmets therefore is being treated as evaluating saliency, which this colour space has been

350   shown to be successful at (Tatler et al., 2005). Colour is perceived differently in the

351   periphery, because there are fewer cone receptors outside of the fovea (Hubel, 1995).  The

352   receptive field sizes in the periphery increase with eccentricity (Abramov et al., 1991), and

353   therefore for objects to appear chromatically similar as if they were in the fovea, they must

354   be spatially larger (Hansen et al., 2009; Vakrou et al., 2005). Given that an object is big

355   enough to be scaled, the upper bound of eccentricity has been found to be 40° to 50°

356   (Abramov et al., 1991; Hansen et al., 2009), after which it has not been found to be possible

357   to simulate chromaticity as if it were in the fovea. An object that subtends 2° of visual angle

358   has been found to appear approximately chromatically similar as if it were in the fovea up to

359   20° away. Therefore colour patterns can be simulated by low-pass-filtering the image

360   (Mullen, 1985). Given the approximate appearance of foveal chromaticity with eccentricity

361   up to 20° (half of the display), of objects that subtend 2° of visual angle, the scene was

362   convolved with a Gaussian, whose standard deviation was measured to be 1° of visual angle,

363   which was chosen so that it was comfortably smaller than 2°. It must be noted that the

364   Gaussian blur is only an approximation and does not accommodate larger receptive fields as

365   objects are more distant. The brightness varies the most across an image. Without

366   processing the luminance, the mixture of Gaussians will have to explain this large variation,

367 which will result in noisy likelihoods. The luminance information across all images could be

368 normalised between one and zero, however that would no longer be Gaussian and, because

369 we are only interested in chromaticity and not luminance at this point, the luminance

370 channel was excluded and was therefore not modelled. Excluding the luminance channel is

371 straightforward to do using some colour spaces such as hue, saturation and value (HSV),

372 where luminance is represented in the channel named value, or opponency colour spaces

373 such as the Macleod and Boynton or L*a*b*, where again the luminance is represented in

374 its own channel. Removing the luminance channel is a standard method to avoid the large

375 variance of brightness in images (Cai and Goshtasby, 1999; Shadeed et al., 2003). Instead of

376 concatenating colour onto the feature vector of energy and phase, another Gaussian

377 mixture model was trained for colour, allowing probabilities of colour and texture to be

378 independent and a full covariance structure of colour to be modelled rather than a mixture

379 of factor analysers. For each posterior map, the probabilities in the region where the target

380 was located were logged and the maximum was taken. The summed log probabilities were

381 plotted against human performance to visualise the correlation.

382

383

## 3. Results

385 Human data were not normally distributed and therefore a Generalised Linear Mixed

386 (Effects) Model with binomial error and logit link function was used to generate

387 interpretable means and error for analysis.   Figures 6 - 9 compare the model accuracy with

388 that of human accuracy and below in table 1 are the correlation coefficients between the

389    model and human observers for each condition.  Correlations coefficients are very high, all

390    above 0.85 with the exception of detection in Woodbury Common in colour.
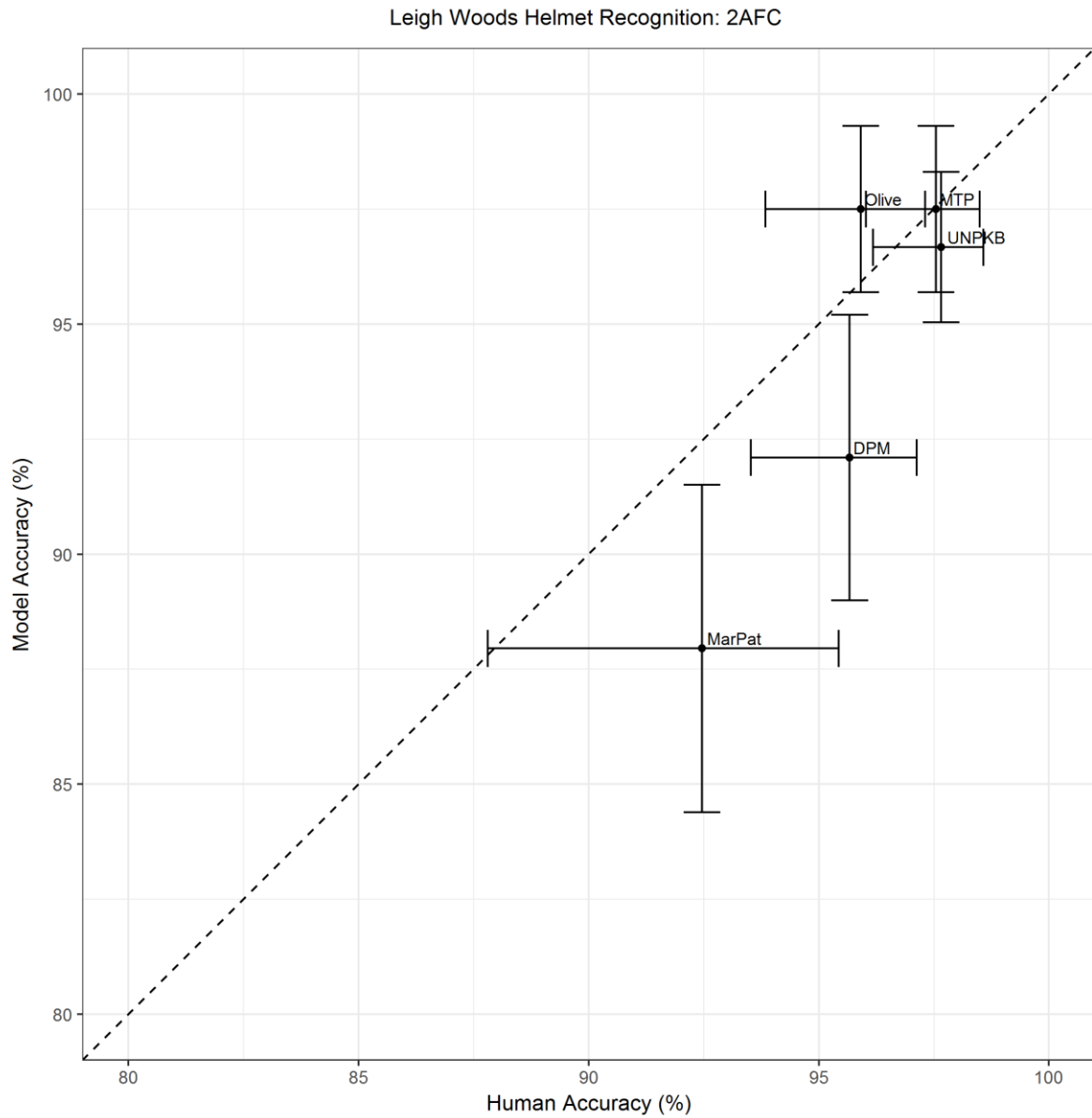
391

392

393

394

395

396

397

398

399

| Condition | Correlation |
|---|---|
| **Leigh Woods** | |
| Recognition | 0.90 |
| Detection Greyscale | 0.93 |
| Detection Colour | 0.89 |
| **Woodbury Common** | |
| Recognition | 0.91 |
| Detection Greyscale | 0.87 |
| Detection Colour | 0.68 |

400    Table 1. The correlation coefficients between the model and human participants at 3 different
401    conditions in two different environments, Leigh Woods and Woodbury Common
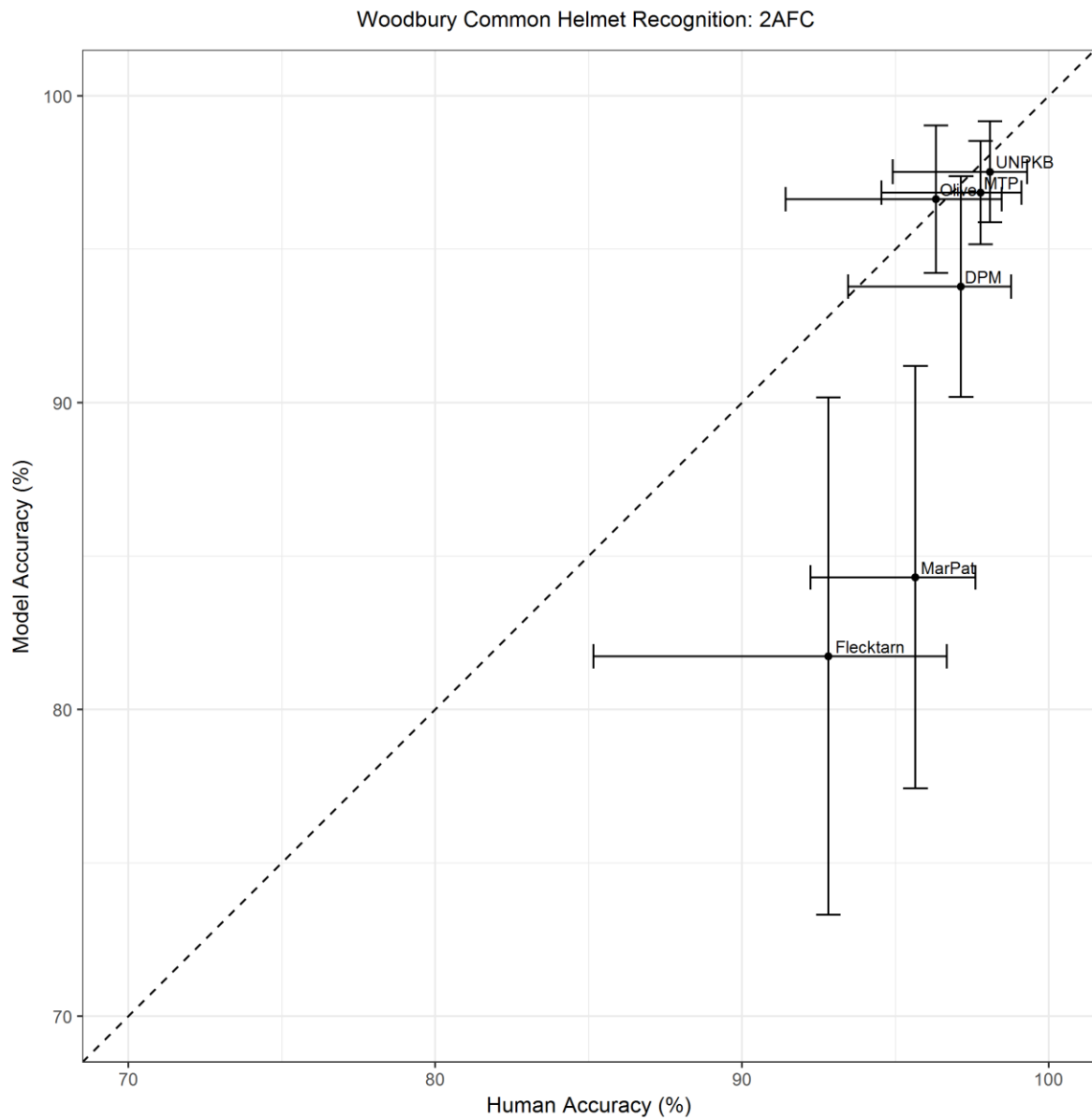402

Figure 6.  Human and model recognition accuracy: Leigh Woods

Leigh Woods model accuracy at recognition in greyscale plotted against human
accuracy at recognition in greyscale. Correlation coefficient: 0.937. Error bars are 95%
confidence intervals.

Figure 7. Human and model recognition accuracy: Woodbury Common
Woodbury Common model accuracy at recognition in greyscale plotted against
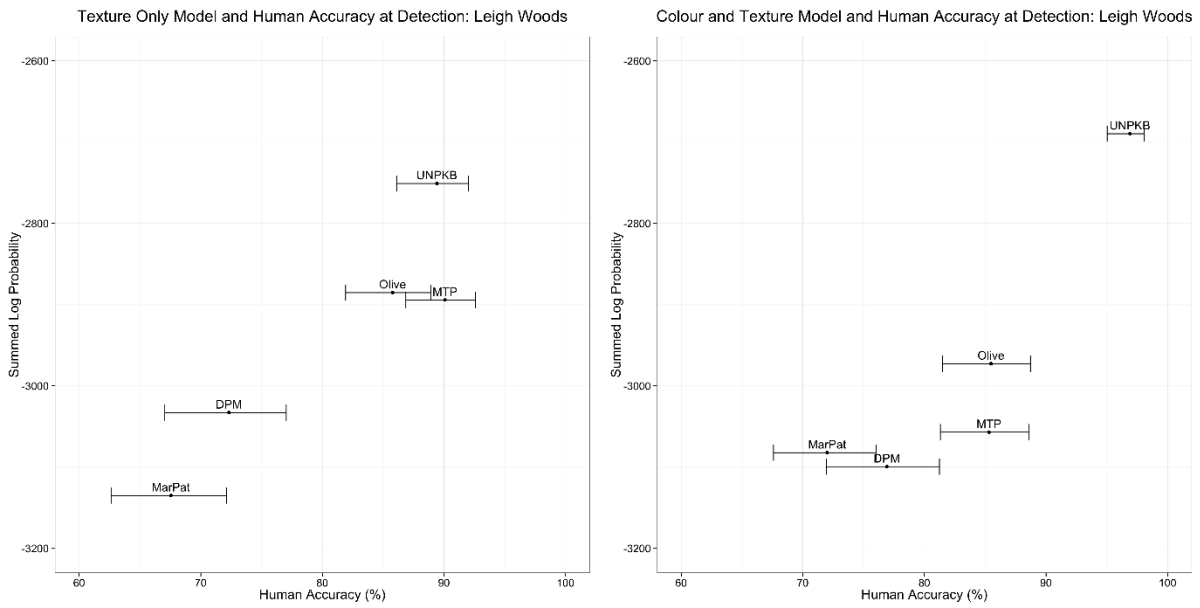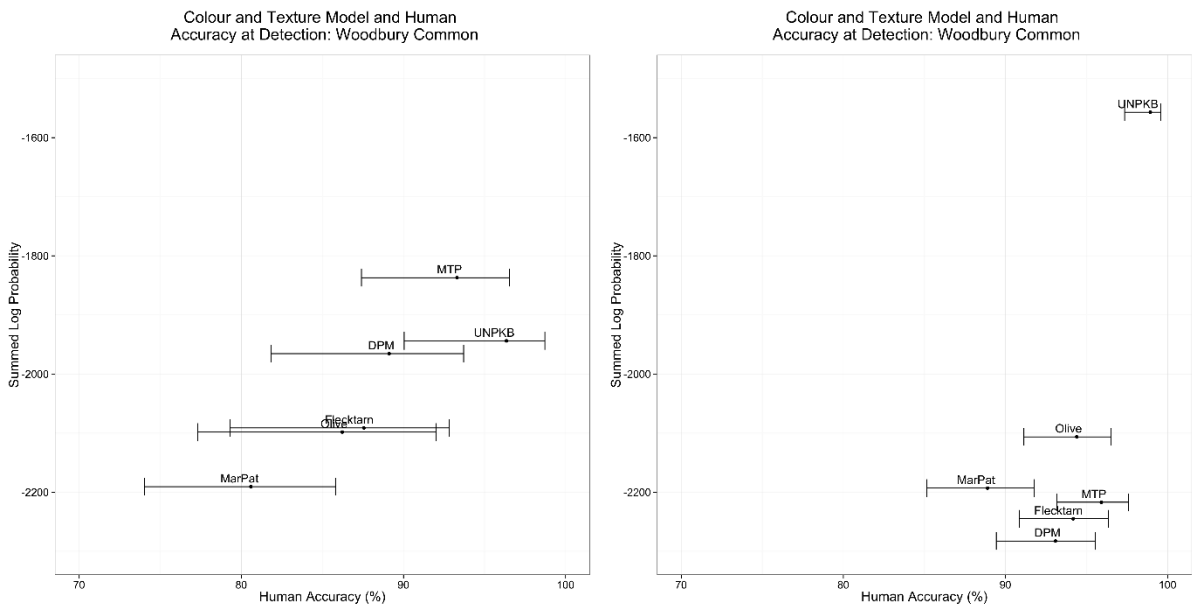human accuracy at recognition in greyscale. Correlation coefficient: 0.859.

**420**

Figure 8.  Human and model detection accuracy: Leigh Woods

Model and Human Accuracy at Detection in Leigh Woods. Left: Texture Only, Right: Colour
and texture. Error bars are 95% confidence intervals.



**426**

Figure 9.  Human and model detection accuracy: Woodbury Common

Model and Human Accuracy at Detection in Woodbury Common. Left: Texture Only, Right:
Colour and texture. Error bars are 95% confidence intervals.

24

## 4. Discussion

This paper has described and validated a visual recognition system that is designed to behave in a similar way to humans. The principles of its design are based upon low-level visual processing in the primary visual cortex. Although it is well-known that Gabor filters can approximate simple cells found in the primary visual cortex, and simple models using Gabor filters can achieve high recognition accuracy on simple datasets (Pinto et al., 2008), we present physiological evidence and a computational argument for the use of log Gabor filters. Such applicability of a human observer model is high, because using human participants is impractical given a variety of viewpoints, environments and objects. This paper also defined a task, a judgement of whether a target is present or absent in a scene, that would allow a direct comparison between the biologically motivated visual observer and human participants. The analysis of the behavior from both observers provides the necessary evidence to assess whether the model is an adequate surrogate for a human observer. The task was to estimate the accuracy with which camouflaged objects, military helmets with different coverings, could be detected and recognised. The selection of a single object class with different colour patterns, rather than an array of different objects, avoided the problem of object choice and allowed visibility to be easily controlled through only colouration and textural properties. The visibilities of the objects were unknown prior to the experiment because, to our knowledge, they had never been evaluated in the two environments nor directly compared. However, a priori, the UN PKB helmet was expected to be easy to detect, the Olive Drab harder to detect and the three (Leigh Woods) or four (Woodbury Common) patterned camouflages hardest to detect. It was essential that the visibility of the patterns varied. If human recognition and detection for all camouflaged objects was at ceiling performance, or all the patterns were equally visible, then we would

456 lack any evidence that the model reflects what human subjects find difficult and what they

457 find effortless.

458 There were clear differences in detectability of the patterns to human subjects (Figs. 6 and

459 7) and the patterns do indeed provide a spectrum of conspicuousness that is sufficient to

460 draw conclusions from. The two different environments did not contain bright blue

461 elements and the texture of the pattern was smooth and therefore UN PKB was, as

462 predicted, very visible and the motivation for its inclusion as a control was vindicated. Olive

463 Drab is also texturally smooth and its colouration is perceptually much closer to the

464 environments used than UN PKB. The cost of pattern design is expensive and if simple olive

465 drab were effective this would have implications for the design of camouflage; in fact this

466 was not the case, with the patterned Flecktarn, Marpat and DPM performing better in most

467 contexts. These patterns' visibilities could not be as easily predicted as UN PKB, because

468 they have never previously been compared in the two environments. We should not over-

469 interpret their relative effectiveness in our experiment, as the experiment was not designed

470 with this goal. Multiple replicates of each pattern type, and habitat class, would be needed

471 before we could conclude that, say, Marpat was better than MTP for these environments.

472 Similarly, we cannot be sure that tendency of humans to outperform the model for

473 Flecktarn, Marpat and DPM, but not MTP or the untextured patterns, is due to specifics of

474 the textures or colours involved.

475 The PASGT helmet, the standard issue for the US Armed Forces from the 1980s to 2000s,

476 was chosen as a typical item of camouflaged military equipment but unvarying in shape

477 (unlike a soldier or combat uniform) and easily portable.  It is difficult to predict how the

478 model might perform with larger objects such as vehicles because these objects would have

479 to be placed much further away from the camera and so the spatial scale of the background

480    textures relative to the object would change. However, given the success of the model in

481    this task and the multiresolution nature of log Gabor filters, there are grounds for thinking it

482    has general applicability.  The primary function of camouflage is to avoid detection in plain

483    sight by enemies.  But it is also the case that friendly personnel need to identify peers, and

484    therefore there is a trade off in visibility and identification such that one needs, not to be

485    easily visible (to avoid attack) and yet remain identifiable (to avoid friendly fire) (Talas et al.

486    2017).  The framework elaborated here, where classification was evaluated in a paired

487    manner, helmet versus background, can be easily extended for this problem as a multi-class

488    classification task.

489

490    **5. Conclusion**

491    A human observer model has been designed, and its detection and recognition behavior

492    was compared with human participants.   Its behavior correlated highly with human

493    participants.  There is large applicability for such a human observer model, where it is

494    impractical to use human participants.  We have shown that an inexpensive and automated

495    objective assessment of camouflage effectiveness is possible in a real-world setting.

## References

Abramov, I., Gordon, J., and Chan, H. (1991). Color appearance in the peripheral retina: effects of stimulus size. Journal of the Optical Society of America, A, 8(2):404–414.

Bhajantri, N. U. and Nagabhushan, P. (2006). Camouflage defect identification: a novel approach. ICIT'06 9th International Conference on Information Technology, pages 145–148.

Birkemark, C. M. (1999). Cameva: a methodology for computerized evaluation of camouflage effectiveness and estimation of target detectability. International Society for Optics and Photonics. In AeroSense 1999, pages 229–238.

Cai, J. & Goshtasby, A. (1999). Detecting human faces in color images. Image and Vision Computing, 18(1):63–75.

Chandesa, T., Pridmore, T., and Bargiela, A. (2009). Detecting occlusion and camouflage during visual tracking. IEEE International Conference on Signal and Image Processing Applications (ICSIPA), pages 468–473.

Daugman, J. G. (1985). Uncertainty relation for resolution in space, spatial frequency, and orientation optimized by two-dimensional visual cortical filters. Optical Society of America, 2(7):1160–1169.

De Valois, R. L., Albrecht, D. G., and Thorell, L. G. (1982). Spatial frequency selectivity of cells in macaque visual cortex. Vision Research, 22(5):545–559.

Field, D. J. (1987). Relations between the statistics of natural images and the response properties of cortical cells. Journal of the Optical Society of America, 4(12):2397–2394.

Ghahramani, Z. and Hinton, G. E. (1996). The EM algorithm for mixtures of factor analyzers. Technical report, University of Toronto.

Hartcup G, 2008. Camouflage: The History of Concealment and Deception in War. Barnsley, UK: Pen and Sword.

Hecker, R. (1992). Camaeleon – camouflage assessment by evaluation of local energy, Spatial frequency, and orientation. In Aerospace Sensing. International Society for Optics and Photonics, pages 343–349.

Hansen, T., Pracejus, L., and Gegenfurtner, K. R. (2009). Color perception in the intermediate periphery of the visual field. Journal of Vision, 9(4):26–26.

Heinrich, D. H. and Selj, G. K. (2015). The effect of contrast in camouflage patterns on detectability by human observers and camaeleon. In SPIE Defense and Security. International Society for Optics and Photonics, pages 947604–947604.

542  Hubel, D. H. (1995). Eye, Brain, and Vision. Scientific American Library/Scientific American
543  Books.

544

545  Jones, J. P. and Palmer, L. A. (1987). An evaluation of the two-dimensional gabor filter
546  model of simple receptive fields in cat striate cortex. Journal of Neurophysiology,
547  58(6):1233–1258.

548

549  Kiltie, R. A., Fan, J., and Laine, A. F. (1995). A wavelet-based metric for visual texture
550  discrimination with applications in evolutionary ecology. Mathematical Biosciences,
551  126(1):21–39.

552

553  Kovesi, P. (1999). Phase preserving denoising of images. Signal, 4(3):1.

554

555  Kovesi, P. D. (2000). MATLAB and Octave functions for computer vision and image
556  processing. Centre for Exploration Targeting, School of Earth and Environment, The
557  University of Western Australia.

558

559  MacLeod, D.I.A. and Boynton, R.M. (1979).  A chromaticity diagram showing cone excitation
560  by stimuli of equal luminance. Journal of the Optical Society of America, A, 69:1183-1186.

561

562  Melin, A. D., Fedigan, L. M., Hiramatsu, C., Sendall, C. L., and Kawamura, S. (2007).
563  Effects of colour vision phenotype on insect capture by a free-ranging population of
564  white-faced capuchins, *Cebus capucinus*. Animal Behaviour, 73(1):205–214.

565

566  Merilaita, S., Scott-Samuel, N. E., and Cuthill, I. C. (2017). How camouflage works.
567  Philosophical Transactions of the Royal Society B 372:20160341.

568

569  Morgan, M., Adam, A., and Mollon, J. (1992). Dichromats detect colour-camouflaged
570  objects that are not detected by trichromats. Proceedings of the Royal Society of London
571  B248:291–295.

572

573  Mullen, K. T. (1985). The contrast sensitivity of human colour vision to red-green and
574  blue-yellow chromatic gratings. The Journal of Physiology, 359(1):381–400.

575

576  Pinto, N., Cox, D. D., and DiCarlo, J. J. (2008). Why is real-world visual object recognition
577  hard? PLoS Computational Biology, 4(1):e27.

578

579  Renoult, J. P., Kelber, A., and Schaefer, H. M. (2015). Colour spaces in ecology and
580  evolutionary biology. Biological Reviews 92: 292–315.

581

582  Shadeed, W., Abu-Al-Nadi, D. I., and Mismar, M. J. (2003). Road traffic sign detection in
583  color images. ICECS 2003. Proceedings of the 2003 10th IEEE International Conference
584  on Electronics, Circuits and Systems, 2003., 2:890–893.

585

586  Sengottuvelan, P., Wahi, A., & Shanmugam, A. (2008). Performance of decamouflaging
587  through exploratory image analysis. First International Conference on Emerging Trends in
588  Engineering and Technology, 2008. ICETET'08 (pp. 6-10).

589

590 Talas L, Baddeley R, Cuthill IC, (2017). Cultural evolution of military camouflage. Phil Trans R
591 Soc B 372, 20160351.

592

593 Tatler, B. W., Baddeley, R. J., and Gilchrist, I. D. (2005). Visual correlates of fixation
594 selection: effects of scale and time. Vision Research, 45(5):643–659.

595

596 Tkaclic, M. and Tasic, J. F. (2003). Colour spaces: perceptual, historical and application.
597 EUROCON 2008. Computer as a Tool. The IEEE Region, 8 (1), 304–308.

598

599 Tipping, M. E. and Bishop, C. M. (1999a). Mixtures of probabilistic principal component
600 analyzers. Neural Computation, 11(2):443–482.

601

602 Tipping, M. E. and Bishop, C. M. (1999b). Probabilistic principal component analysis.
603 Journal of the Royal Statistical Society: Series B (Statistical Methodology), 61(3):611–622.

604

605 Vakrou, C., Whitaker, D., McGraw, P. V., and McKeefry, D. (2005). Functional evidence for
606 cone-specific connectivity in the human retina. The Journal of Physiology, 566(1):93–102.

607

608 Wyszecki, G. and Stiles, W. S. (1982). Color Science: Concepts and Methods, Quantitative
609 Data and Formulae. 2nd edition. New York: Wiley.

610

611 Yfantis, E. A., Flatman, G. T., and Behar, J. V. (1987). Efficiency of kriging estimation for
612 square, triangular, and hexagonal grids. Mathematical Geology, 19(3):183–205.

613
614
615
616
617