

# *Biologically-inspired robust motion segmentation using mutual information*

Article

Accepted Version

Ellis, A.-L. and Ferryman, J. (2014) Biologically-inspired robust motion segmentation using mutual information. *Computer Vision and Image Understanding*, 122. 47 - 64. ISSN 1077-3142 doi: 10.1016/j.cviu.2014.01.009 Available at <https://centaur.reading.ac.uk/36796/>

It is advisable to refer to the publisher's version if you intend to cite from the work. See [Guidance on citing](#).

Published version at: <http://www.sciencedirect.com.idproxy.reading.ac.uk/science/article/pii/S1077314214000228>

To link to this article DOI: <http://dx.doi.org/10.1016/j.cviu.2014.01.009>

Publisher: Elsevier

All outputs in CentAUR are protected by Intellectual Property Rights law, including copyright law. Copyright and IPR is retained by the creators or other copyright holders. Terms and conditions for use of this material are defined in the [End User Agreement](#).

[www.reading.ac.uk/centaur](http://www.reading.ac.uk/centaur)

**CentAUR**

Central Archive at the University of Reading

Reading's research outputs online

1  
2

3

4  
5

## 6

7  
8  
9  
10  
11  
12  
13  
14  
15  
16

17

18

19

## 20

21

22

community. The capacity to provide real-time segmentations - silhouettes and bounding boxes - of objects (especially pedestrian) assists in the tracking and reasoning of the behaviour. Surveillance scenes often contain change that may be inaccurately detected as object motion such as changes in lighting, periodic motion, moving shadows and reflections. In addition the quality of surveillance footage is often poor, and at a low resolution resulting in noisy motion and ghosts. An example of these challenges is shown in Figure 1. The extraction of objects of interest is frequently tackled by removing all irrelevant pixels in each frame. This is referred to as motion segmentation. To date no segmentation algorithm is robust under all these conditions.

In this paper, we propose a new formulation of pixel-based foreground segmentation which is motivated by recent results in biological vision which exploit the mutual information between multiple segmentation channels. The paper is divided as follows. Firstly, Section 2 details the biological motivation and mapping to a combination of parametric background modelling approaches. This is followed in Section 3 by approaches to fusing the outputs of multiple segmentation algorithms and introduces the multivariate mutual information formulation adopted in this work. In Section 4 the datasets, evaluation methodology and the results of experiments are presented before concluding in Section 5 with conclusions and recommendations for future research.

## 2. Biologically-Inspired Segmentation

The ability of primates to recognise objects of interest, regardless of illumination and background, drives much of the biologically inspired computa-



Figure 1: PETS 2009 dataset original frame annotated with automated visual surveillance challenges.

47 tional vision systems. A new biologically inspired vision system is introduced  
 48 in this section that models current vision research which has not previously  
 49 been examined by the computational vision community.

50 In Section 2.1 the model of primate vision conventionally accepted by the  
 51 computer vision community is presented. Section 2.2 provides descriptions  
 52 of state of the art biologically inspired computational vision systems that  
 53 refer to this model. Section 2.3 progresses on to accounts of current pub-  
 54 lished neuro-biological, physiological and psychological vision research and  
 55 highlights descriptions of retinal functions, inputs to the ventral and dorsal  
 56 streams, and ventral and dorsal stream behaviour that have not been consid-  
 57 ered in modelling primate visual systems in the computer vision community.  
 58 Based on this, a new model of understanding is presented and the behaviours

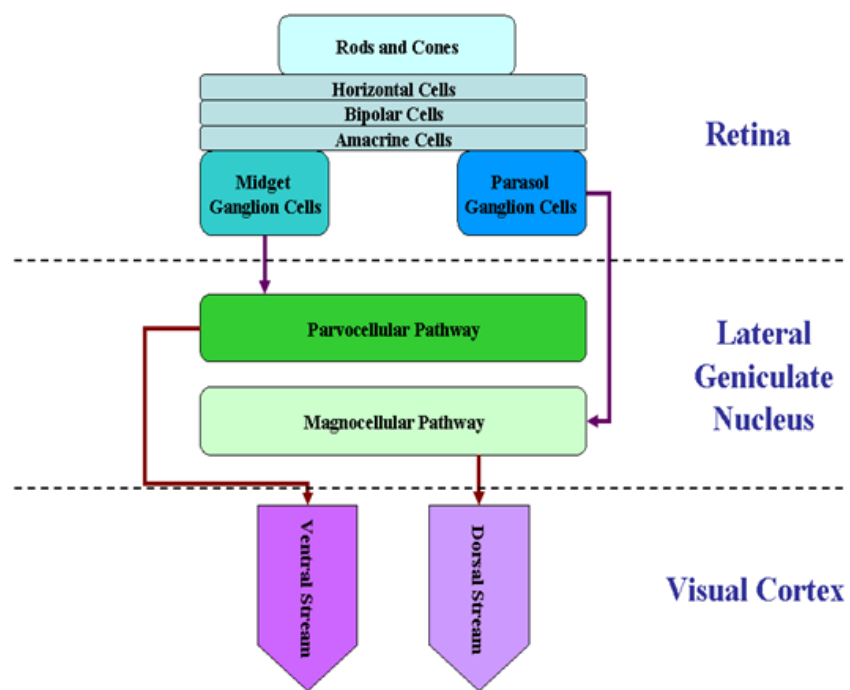


Figure 2: Model of traditional computational vision process

59 of these retinal functions are summarised.

## 60 *2.1. Conventional Model of Primate Vision*

61 It is widely acknowledged that the rods and cones (photoreceptors) of the  
62 primate retina detect light and cells of the inner retina providing the initial  
63 stages of the visual processing. The retinal ganglion cells convey this infor-  
64 mation, via pathways in the lateral geniculate nucleus, to the ventral and  
65 dorsal streams in visual cortex. Figure 2 represents a model of these tradi-  
66 tionally accepted components, frequently referred to in biologically inspired  
67 computational vision systems.

68 Within the retina, shown in Figure 2 as the blue area, the photorecep-  
69 tor rod cells respond to achromatic brightness and the photoreceptor cone  
70 cells respond to short (blue), medium (green) and long (red) chromatic wave-  
71 lengths. These nerve impulses are passed on to the network of horizontal,  
72 amacrine and bipolar cells, which provide cumulative information to retinal  
73 ganglion cells, shown in Figure 2 as the midget and parasol ganglion cells.  
74 The midget ganglion cells have been associated with providing chromatic  
75 information and parasol ganglion cells with luminance and contrast.

76 The lateral geniculate nucleus (LGN), illustrated as the green area in  
77 Figure 2, receives the assembled information from the ganglion cells, in the  
78 form of pathways. The parvocellular pathway is conventionally understood  
79 to receive information from the midget ganglion cells, and as such provides  
80 a means to direct colour information to the visual cortex. It is customary  
81 to describe the magnocellular pathway as a swiftly responsive structure, pre-  
82 senting the visual cortex with luminance and contrast information.

83 Finally, the visual cortex (VC), emphasised as the purple area in Figure 2,

84 includes two different streams: the ventral stream, associated with form, and  
85 the dorsal stream associated with motion.

## 86 *2.2. Existing Bio-Inspired Computational Models*

87 (Mota et al., 2006) state that because bio-inspired vision models based  
88 on a vertebrates visual system are limited and require high computational  
89 cost, real-time applications are seldom addressed. As flies are capable of  
90 exploiting optical flow, which modelled by calculating the local image mo-  
91 tion with Reichardt motion detectors (and referred to as Elementary Motion  
92 Detectors), they use this as inspiration and employ EMD as the first ex-  
93 traction primitive to characterise motion in a scene. Sequences are initially  
94 pre-processed by extracting edges within each frame using a Sobel edge ex-  
95 traction procedure. The Reichardt motion detector is then used to extract  
96 sideways moving features. Noise is removed from the resulting saliency map  
97 with a neural structure that allows the emergence of rigid bodies (indep-  
98 dent moving objects in the scene) using “velocity channels”. The technique  
99 is limited to greyscale images and suffers from being unable to identify to  
100 objects moving in parallel at the same speed. The system proposed by (Serre  
101 et al., 2007) follows on from their own theory of a feed forward path of object  
102 recognition that accounts for the first 100-200 milliseconds of processing in  
103 the ventral stream of primate visual cortex. It is based on Hubel and Wiesel’s  
104 findings in 1962 of a cat’s visual cortex (Hubel and Wiesel, 1985). Unlike the  
105 conventionally accepted chromatic input to the primate ventral stream, the  
106 approach takes a grey scale input and uses a set of scale and position-tolerant  
107 feature detectors, to simulate the properties of V1 and V4 (Figure 2 shows  
108 V1 and V4 within the ventral stream). A major limitation of the system



109 for real-time application is the processing speed which is limited by some of  
 110 its modules that typically take tens of seconds, depending on the size of the  
 111 input image. The authors have yet to address whether the recognition re-  
 112 sults obtained can be extended to the analysis of video. (Huang et al., 2011)  
 113 offer an improvement on the system proposed by (Serre et al., 2007) focusing  
 114 on improving the biological Standard Model Feature (SMF) for scene clas-  
 115 sification in a video surveillance environment. They develop a new energy  
 116 computation component to improve SMF in occlusion and disorder cases as  
 117 basic SMF models can only handle shift and invariance. An energy function  
 118 is used in order that patches for saliency are not chosen randomly. An earlier  
 119 analysis of energy density is used to conduct a local energy measurement after  
 120 the initial basic feature extraction stage. Again the technique is limited to  
 121 greyscale images. Using accounts of the primate visual cortex (Bayerl et al.,  
 122 2007) have developed a neurodynamical computational vision model of mo-  
 123 tion segregation in the dorsal stream, as described in (Mishkin et al., 1983).  
 124 The model includes two modules, corresponding to the primate visual cortex  
 125 (highlighted as the purple area in Figure 2): V1 represents a motion hypoth-  
 126 esis on the same scale of resolution on which it was detected, and V5 uses a  
 127 coarser spatial resolution, where the accuracy of both location and velocity  
 128 is reduced by a factor of five in accordance with physiological findings of Al-  
 129 bright and Destmone in 1987 (Albright et al., 1987). The authors conclude  
 130 that it is a step towards producing a biologically inspired model which may  
 131 be capable of real-time computation. (Thriault et al., 2013) use a principle  
 132 referred to as Slow Features Analysis (SFA) which bears foundations in neu-  
 133 roscience. SFA extract slowly varying features from a quickly varying input

134 signal. These features have been shown by (Thriault et al., 2013) to reveal  
 135 sensible motion components correlated with specific semantic classes such as  
 136 complex flame motion, waterfalls and fountains. As perceptions vary on a  
 137 slower timescale compared to input signals from the environment, the SFA  
 138 model learns to generate a slower, more invariant output signal. Temporal  
 139 variations created by motion are minimised to in order to learn the stable  
 140 representations of objects in motion. Motion features are defined by thread-  
 141 ing together short temporal sequences of SFA outputs. The motion features  
 142 can be interpreted as spatio-temporal atoms describing the stable motion  
 143 components inside a small space time window. Again this model relies on  
 144 grey scale video as an input. The authors state that employing it for motion  
 145 segmentation is a direction for future work. In (Yuen et al., 2009) features  
 146 of objects are extracted “in a way similar to that of the ventral stream pro-  
 147 cessing”, referring to Diddays two visual stream model (Didday et al., 1975)  
 148 published in 1975 and Mishkins slightly earlier publication than previously  
 149 mentioned, with Ungerleider, in 1982 (Ungerleider et al., 1982). They use an  
 150 RGB image input and proceed with a cortex-like centre surround operation  
 151 in the spatiotemporal domain, by sub-sampling the image data into various  
 152 spatial scales resulting in a set of images with horizontal and vertical scale re-  
 153 ductions. Sets of features are extracted from the spatiotemporal stream and  
 154 manipulated across various scales to detect those which locally stand out  
 155 from their surround, similar to that of an edge detector. The authors state  
 156 that due to the lack of a full understanding about the object recognition pro-  
 157 cess in the visual cortex, the recognition mechanism that was implemented  
 158 was a statistical classifier (SVM). In contrast Benoit et al. (Benoit et al.,

159 2010) recognise that consideration must be taken of the processing of the  
 160 retinal signals that occur in primate vision, in order to assist further pro-  
 161 cessing of that input, in a primate biologically inspired manner, in the visual  
 162 cortex. They base their retinal architecture on Meads silicon model (Mead et  
 163 al., 1988) albeit improved in terms of spatial and temporal properties. Their  
 164 system contains two processing modules, one based on the retina for motion  
 165 information extraction and the second representing a model of the V1 cortex  
 166 area providing motion event detection. Their focus on the retinal processing  
 167 includes passing information to their parvocellular channel model and mag-  
 168 nocellular channel model from the midget ganglion cells model and parasol  
 169 ganglion cells model respectively. These are shown in Figure 2 in green. This  
 170 transformed information then is presented to their V1 model of the visual  
 171 cortex. The system concentrates on using grey level image processing as the  
 172 authors state the cell actions at the retinal level are unknown and further  
 173 investigation is required to produce a better model.

### 174 *2.3. Current Primate Vision Research*

175 Current neurobiology, visual neuroscience, physiology and psychology re-  
 176 search provide descriptions of the input to the ventral and dorsal streams that  
 177 have not been considered in computational vision systems modelling primate  
 178 visual systems. Ganglion cell types other than midget and parasol cells also  
 179 project to the LGN (Nieuwenhys et al., 2008; Dacey et al., 2000; Chatterjee  
 180 and Callaway, 2003). (Dacey et al., 2000) provides a detailed description of  
 181 these cell types, referred to as bistratified ganglion cells. They project their  
 182 information to a further pathway in the lateral geniculate nucleus which is  
 183 referred to as the koniocellular pathway (Nieuwenhys et al., 2008; Dacey

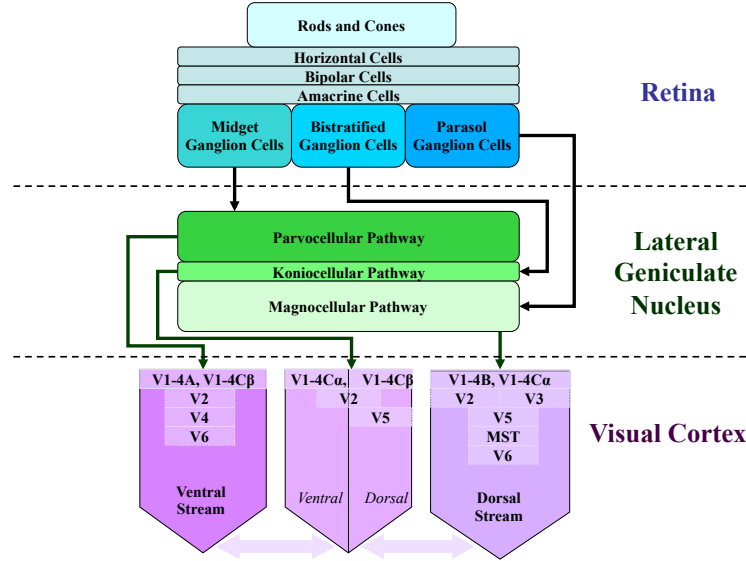


Figure 3: Model of recognised primate vision processes

et al., 2000; Chatterjee and Callaway, 2003; Hendry, 2000; Morand et al., 2000; Briggs and Usrey, 2011). A new illustration representing these recognised processes, including the bistratified ganglion cells and the koniocellular pathway is shown in Figure 3.

The retinal ganglion cells function in a distinct manner. The received wavelength signals can be used in the course of perceiving form or motion, independent of their role in the subjective experience of colour. Contrastively to the traditional accepted processes, the networked routing provides the midget cells with some contrast information (Kentridge et al., 2002), alongside the bistratified and parasol cells and therefore contrast information is present within both the ventral and dorsal streams. In addition prominent computation has been found to occur in the retina: the detection

196 of object motion while rejecting background motion (resulting from subtle eye  
 197 movements) (Baccus et al., 2008) through specific interactions of amacrine  
 198 and bipolar cells and presented to the ganglion cells. The koniocellular layer  
 199 has been found to project to both the ventral and dorsal streams (Hendry,  
 200 2000). Finally recent primate vision research suggests there is communica-  
 201 tion between the dorsal and ventral streams, contrary to the traditionally  
 202 accepted definitions used by the computer vision community of independent  
 203 luminance motion information and colour object information occurring in  
 204 the dorsal and ventral streams respectively. (McKeefry et al., 2010) ascer-  
 205 tain that both luminance and chromatically defined motion is analysed in  
 206 the dorsal stream and (Farivar et al., 2009) provide evidence that the dorsal  
 207 stream participates in object recognition and some dorsal-ventral integration  
 208 may be considered. Furthermore the study by (Zanon et al., 2010) states that  
 209 the continuous interchange of information between the two streams is nec-  
 210 essary and provides evidence that interaction is present in order to produce  
 211 adaptive behaviour, for example, in order to elaborate the position in space  
 212 and the shape of a 3D object. In effect the individual streams of information  
 213 are weaved back together.

### 214 *2.3.1. Ganglion Cells and the Lateral Geniculate Nucleus Pathways*

215 The current understanding of the individual behaviours of the three types  
 216 of ganglion cells is described in detail in a vast array of vision research liter-  
 217 ature. These components in turn project this information to their respective  
 218 lateral geniculate nucleus (LGN) streams, and these three streams have been  
 219 ascertained by the neuroscience vision research community to have distinct  
 220 behaviours and output. In this section brief descriptions of these components

221 and their respective LGN streams and behaviours are presented.

222     Parasol retinal ganglion cells receive many inputs and are responsively  
223 fast. They react to achromatic information and low contrast stimuli from  
224 the rods, and medium and long wavelength cones. They are unable to trans-  
225 mit information about wavelength independent of intensity and as such are  
226 not very sensitive to changes in colour. These cells are more sensitive to light  
227 since they are three times larger in diameter to the midget retinal ganglion  
228 cells. This information is relayed to the magnocellular pathway which is a  
229 fast system which contributes to the perception of luminance and motion  
230 derived from both achromatic and chromatic wavelengths, though it is un-  
231 able to transmit any chromatic wavelength signals (Nieuwenhys et al., 2008;  
232 Kentridge et al., 2002; Dacey et al., 2000; Chatterjee and Callaway, 2003;  
233 Briggs and Usrey, 2011).

234     Midget retinal ganglion cells are involved in colour encoding. They react  
235 to chromatic information from the rods, and medium and long wavelength  
236 cones (green and red cones respectively) in the retina. They have low sen-  
237 sitivity because of their small receptive fields, but because of that they are  
238 densely packed and their resolution ability is higher. They respond weakly  
239 to changes in contrast unless that change is great. However, though these  
240 cells are found predominantly in the fovea of the retina, those located in the  
241 periphery show a non-opponent luminance response, indistinguishable from  
242 the parasol cells. The red/green colour opponent information and achromatic  
243 contrast detection information, provided by the synergy of the medium and  
244 long wavelength cones in the fovea, and those of the periphery able to dis-  
245 tinguish brightness only, are relayed through the slow parvocellular pathway.

246 This pathway transmits information about long and medium wavelengths  
 247 and fine detail. Motion perception information is presented but is far weaker  
 248 than that of the magnocellular pathway and is dependent on the available  
 249 chromatic contrast (Nieuwenhys et al., 2008; Kentridge et al., 2002; Dacey  
 250 et al., 2000; Chatterjee and Callaway, 2003; Briggs and Usrey, 2011).

251 Bistratified retinal ganglion cells are involved in colour perception. They  
 252 receive inputs from all rods and cone types but respond to rods and small  
 253 wavelength cones (blue cones) 23 only. They have the lowest resolution abil-  
 254 ity, their density is extremely low and they have very large receptive fields.  
 255 They have moderate to low spatial resolution and react to moderate changes  
 256 in contrast. This information is projected to the koniocellular pathway which  
 257 contributes to colour perception dependant on the small wavelength cone out-  
 258 put and contributes to motion perception (Nieuwenhys et al., 2008; Kentridge  
 259 et al., 2002; Dacey et al., 2000; Chatterjee and Callaway, 2003; Morand et al.,  
 260 2000; Briggs and Usrey, 2011). Table 1 summarises the functions of the Mag-  
 261 nocellular, Parvocellular and Koniocellular streams in the Lateral Geniculate  
 262 Nucleus.

	Magnocellular	Parvocellular	Koniocellular
Ganglion Cell	Parasol	Midget	Bistratified
Colour	No	Yes (R, G cones)	Yes (B cones)
Sensitivity to Contrast	High	Low	Moderate
Spatial Resolution	Low	High	Low
Temporal Resolution	Fast	Slow	Slow

Table 1: Magnocellular, Parvocellular and Koniocellular Functions

## 263 2.4. *Modelling the Lateral Geniculate Nucleus Pathways*

264 Recent research in (Zanon et al., 2010; Briggs and Usrey, 2011) have  
265 shown that the output of the magnocellular, koniocellular and parvocellular  
266 pathways provide mutual information to both ventral and dorsal streams, in  
267 order to supply the visual cortex with robust data about objects of interest  
268 and their location. Modelling this behaviour a form of multivariate mutual  
269 information is employed to enable the quantification of the amount of mu-  
270 tual information provided by the foreground segmentations of the modelling  
271 approaches described in this section. Background models may be seen to be  
272 analogous with the retinal suppression of global image motion as described  
273 by (Baccus et al., 2008). Using RGB colour space video sequences as input,  
274 the function of each of the parvocellular, magnocellular and koniocellular  
275 streams may each be modelled in a similar statistical manner. This sec-  
276 tion provides details of how these streams may be mapped to computational  
277 vision pixel-based background models.

### 278 2.4.1. *Parvocellular*

279 A background statistical model, which approximates behaviour of the  
280 parvocellular stream function (Kentridge et al., 2002), is able to distinguish  
281 between the brightness and its chromaticity of any one pixel, over time. This  
282 relates most closely to the method of (Horprasert et al., 1999). It is able to  
283 separate its wavelength (colour) information to include pixels with changes  
284 in luminance and contrast within its background model. The remaining  
285 pixels, with changes in colour and a limited amount of motion information.  
286 Figure 4 represents a graphical representation of the brightness distortion  
287 and chromaticity distortion in three dimensional RGB colour space.  $E_i$  is the



288 initial (background) colour value for pixel  $i$ , and  $I_i$  is the current colour value  
 289 of the image. The line OE from the origin to  $E_i$  represents the chromaticity  
 290 line. Brightness distortion is a scalar value  $\alpha$  and scales the point along OE  
 291 where the orthogonal line from  $I_i$  intersects OE. Chromaticity distortion  $CD_i$   
 292 is the orthogonal distance between the observed colour and the line OE. The  
 293 values for  $\alpha$  and  $CD$  are calculated for each of  $N$  background frames

$$\alpha_i = \frac{\left( \frac{I_R(i)\mu_R(i)}{\sigma_R^2(i)} + \frac{I_G(i)\mu_G(i)}{\sigma_G^2(i)} + \frac{I_B(i)\mu_B(i)}{\sigma_B^2(i)} \right)}{\left( \left[ \frac{\mu_R(i)}{\sigma_R(i)} \right]^2 + \left[ \frac{\mu_G(i)}{\sigma_G(i)} \right]^2 + \left[ \frac{\mu_B(i)}{\sigma_B(i)} \right]^2 \right)}$$

294 where  $\sigma_R(i)$ ,  $\sigma_G(i)$  and  $\sigma_B(i)$  are the standard deviation and  $\mu_R(i)$ ,  $\mu_G(i)$   
 295 and  $\mu_B(i)$  are the means of the  $i^{th}$  pixel's red green and blue values computed  
 296 over  $N$  background frames

$$CD_i = \sqrt{\left( \frac{I_R(i) - \alpha_i \mu_R(i)}{\sigma_R(i)} \right)^2 + \left( \frac{I_G(i) - \alpha_i \mu_G(i)}{\sigma_G(i)} \right)^2 + \left( \frac{I_B(i) - \alpha_i \mu_B(i)}{\sigma_B(i)} \right)^2}$$

297 and then normalised to find a single threshold for all pixels

$$a_i = \sqrt{\frac{\sum_{i=0}^N (\alpha_i - 1)^2}{N}}$$

$$\hat{\alpha}_i = \frac{\alpha_i - 1}{a_i}$$

$$b_i = \sqrt{\frac{\sum_{i=0}^N (CD_i)^2}{N}}$$

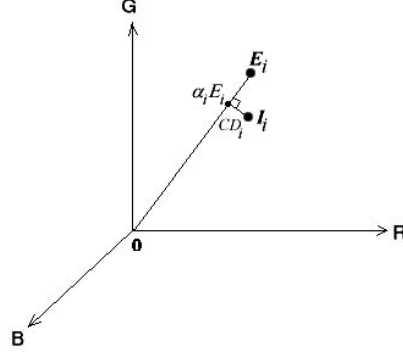


Figure 4: Graphical representation of the brightness distortion and chromaticity distortion in 3D RGB colour space.

$$\widehat{CD}_i = \frac{CD_i}{b_i}$$

298 The method constructs histograms of the normalised  $\hat{\alpha}$  and  $\widehat{CD}$  values  
 299 and takes a detection rate as input to automatically select thresholds. For  
 300 segmentation, incoming pixels are used to calculate  $\hat{\alpha}_i$  and  $\widehat{CD}_i$  values which  
 301 are compared to those of the background model. The pixel classification for  
 302 the  $i$ th pixel as defined by (Horprasert et al., 1999) is:

- 303 1. Original background if both  $\hat{\alpha}_i$  and  $\widehat{CD}_i$  are within a threshold of those  
 304 in the background model
- 305 2. Shadows or shaded background if the chromaticity  $\widehat{CD}_i$  is within the  
 306 threshold, but the brightness  $\hat{\alpha}_i$  is below
- 307 3. Highlighted background if the chromaticity  $\widehat{CD}_i$  is within the threshold,  
 308 but the brightness  $\hat{\alpha}_i$  is above

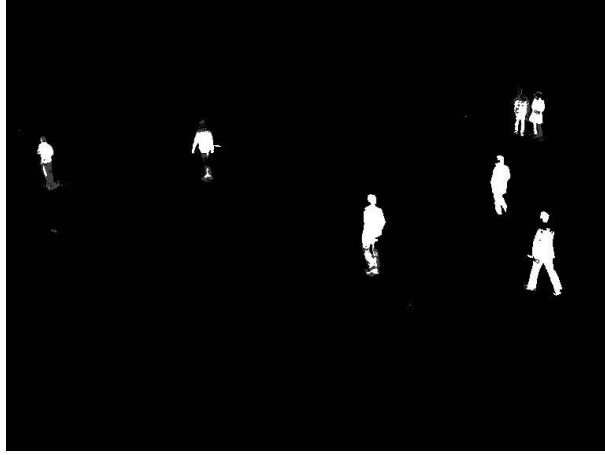


Figure 5: PETS 2009 dataset frame - BC algorithm approximating Parvocellular behaviour.

309 4. Moving foreground object if the chromaticity  $\widehat{CD}_i$  is outside of the  
 310 threshold

311 The resulting motion segmentation (Figure 5) from the original frame  
 312 (Figure 1) show the model is able distinguish subtle differences in colour due  
 313 to its motion sensitivity, but because of its motion sensitivity (due to both  
 314 the temporal resolution and contrast sensitivity) parts of fluttering tape in  
 315 the wind appear as foreground. Both the illumination and motion sensitivity  
 316 provide the foreground segmentation with shadows.

#### 317 2.4.2. *Magnocellular*

318 A statistical model that presents foreground segmentation approximating  
 319 behaviour of the magnocellular stream function is one that is able to provide  
 320 high contrast information but does not distinguish between colour and its  
 321 intensity. It must be sensitive to changes in luminance and motion (Ken-  
 322 tridge et al., 2002). This most closely relates to the mixture model approach

323 of Stauffer and Grimson (Stauffer et al., 1999). Gaussian mixture models  
 324 (GMM)s are able to model each component distribution as a soft classifica-  
 325 tion; that is they are able to produce a distribution without specifying exactly  
 326 what each cluster must represent. Yet as a whole, the mixture model covers  
 327 the entire set of features (colour, brightness, intensity and luminance) that  
 328 the data represents. The clusters formed represent more than one feature  
 329 of information, and in this way the model becomes sensitive to contrast and  
 330 motion. The resulting motion segmentations show that the model is able  
 331 distinguish subtle differences in colour due to its motion sensitivity. Both  
 332 the illumination and motion sensitivity provide the foreground segmentation  
 333 with shadows. The recent history of a pixel is modelled by a mixture of  $K$   
 334 Gaussians ( $K$  usually varies from 3 - 5). The mixture is weighted by the  
 335 frequency with which each of the Gaussians explains the background. The  
 336 probability of observing a foreground pixel  $x$  is:

$$P(x) = \sum_{j=1}^K w_j N(x, \mu_j, \Sigma_j) \quad (1)$$

337 where  $w$  is the weight of the  $K$ th Gaussian distribution,  $\mu$  is the mean,  $\Sigma$   
 338 is the covariance matrix and  $N$  is a multivariate Gaussian density function.

339 The resulting motion segmentation (Figure 6) from the original frame  
 340 (Figure 1) show the model is able distinguish subtle differences in colour due  
 341 to its motion sensitivity, but because of its motion sensitivity (due to both  
 342 the temporal resolution and contrast sensitivity) parts of fluttering tape in  
 343 the wind appear as foreground. Both the illumination and motion sensitivity  
 344 provide the foreground segmentation with shadows.



Figure 6: PETS 2009 dataset frame - GMM algorithm approximating Magnocellular behaviour.

#### 345 2.4.3. *Koniocellular*

346 Similar to that of the Gaussian Mixture Model, the Colour Mean and  
 347 Variance (CMV) algorithm, described in (Wren et al., 1997) captures the  
 348 brightness, motion and colour information but only for a single colour chan-  
 349 nel. In this way the algorithm is able to provide foreground segmentation,  
 350 similar to the behaviour of the koniocellular pathway (Kentridge et al., 2002).  
 351 Encapsulating features in distinct distributions, using one independent chan-  
 352 nel value, removes the ability to capture some of the colour contrast infor-  
 353 mation in the model, enabling any subtle changes to appear as foreground.  
 354 The changes in the objective luminance of a pixel provide additional nec-  
 355 essary motion information, but it is not as precise a measure as perceived  
 356 brightness change and as such the motion sensitivity is coarser. The result-  
 357 ing motion segmentations show the model is able distinguish between some  
 358 subtle differences in colour, however is of lower resolution and provides low

359 resolution shadow information from its motion sensitivity. CMV builds a  
 360 statistical background model to represent an independent Gaussian distribu-  
 361 tion for each normalised colour channel (R,G,B) and a Gaussian distribution  
 362 of the luminance (A) of each normalised pixel colour:

$$n(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp^{-(x-\mu)^2/2\sigma^2} \quad (2)$$

363 where  $x$  is the value of a single channel R, G, or B, or luminance (A),  
 364  $\mu$  is the mean and  $\sigma$  is the standard deviation of that channel. A pixel is  
 365 classified as foreground if it is found to be more than 3 standard deviations  
 366 of the R, G, B or A distributions.

367 The resulting motion segmentation (Figure 7) from the original frame  
 368 (Figure 1) show the model is able to distinguish between some subtle differ-  
 369 ences in colour, but is of lower resolution (shown by the merging of moving  
 370 objects in close proximity in Figure 7 and provides low resolution shadow  
 371 information from its motion sensitivity.

### 372 **3. Combining Algorithms**

373 A number of approaches have been adopted in the literature for com-  
 374 bining or fusing the outputs of multiple motion segmentation algorithms.  
 375 (Martin et al., 2006) exploit optimal algorithm selection and key parameters  
 376 tuning. A library of segmentation algorithms are fine tuned against predeter-  
 377 mined ground truth images. The features extracted, alongside the optimal  
 378 algorithm parameters, are saved as a case. They are ranked by a number of  
 379 criteria. For each image a new case is created composed of a vector of image  
 380 features, the chosen algorithm, and its optimised parameters. A multilayer



Figure 7: PETS 2009 dataset frame - CMV algorithm approximating Koniocellular behaviour.

381 perceptron (MLP) neural network is trained with this stored knowledge for  
 382 algorithm selection. As the technique relies on predetermined ground truth  
 383 this rules out generality. A Support Vector Machine (SVM), used by (Avidan et al., 2004), views the feature information as two sets of vectors in  
 384 an n-dimensional space. It constructs a separate hyper-plane in that space  
 385 which maximizes the margin between the two data sets. (Farmer et al., 2006)  
 386 employ Expectation Maximisation (EM) as a fusion engine. Principal Component Analysis (PCA) is first applied to perform dimensionality reduction  
 387 to improve the performance of EM and reduce the computational load. It is  
 388 claimed that the approach applied to fusion of three popular optical flow algorithms (where the U and V component images are treated as image planes  
 389 and EM applied to them) reduces the percentage of missing target pixels by  
 390 33%, although only one outdoor driving sequence has been used for evaluation.  
 391 Boosting is an alternative. In (Zhou et al., 2004) each base classifier

395 must be trained, sequentially, using feature points that are weighted. The  
 396 weight of a feature point is increased if a previous classifier misclassifies it.  
 397 Once all of the classifiers are trained, their decisions can be combined through  
 398 a weighted majority vote method or others. Popular boosting methods Ad-  
 399 aboost and LogitBoost both have structural space, a cost function, and a  
 400 selection algorithm. The AdaBoost algorithm minimises an upper bound of  
 401 the target misclassification error, and LogitBoost minimises a negative bi-  
 402 nomial log-likelihood, as cost functions. Serre, Wolf, Bileschi, Riensenhuber  
 403 and Poggio model a neurobiological design of a primate cortex (Serre et al.,  
 404 2007). It is designed using hierarchical alternating layers of simple units and  
 405 complex units. Simple units (16 Gabor filters for each layer) combine their  
 406 inputs with a (bell shaped) tuning function to increase selectivity. Complex  
 407 units pool their inputs (from the output of the previous Simple unit layer)  
 408 through a MAX function. The image (grey scale only) is propagated through  
 409 the hierarchical architecture. Standard Model Features (SMFs) are extracted  
 410 from the complex units and classified using SVM or boosting (Gentle boost-  
 411 ing providing the best performance). It was discovered that because there  
 412 are variations in the amount of clutter and in the 2D transformations, it  
 413 is beneficial to allow the classifier to choose the optimal features extracted  
 414 from either the high or low level SMFs at a point in time, to improve the  
 415 performance. A major limitation of the system in the use of real world  
 416 applications remains its processing speed which is typically tens of seconds  
 417 per image. (Jodoin and Mignotte, 2005) fusion of motion segmentation ap-  
 418 proach is based on a K-nearest-neighbour-based fusion procedure that mixes  
 419 spatial and temporal data taken from two input label fields. The first one



420 is a spatial segmentation of a frame at time  $t$  which contains regions of uni-  
421 form brightness while the second label field is an estimated version of the  
422 motion partition. The two segmentation maps are estimated separately with  
423 an unsupervised Markovian segmentation routine. The fusion occurs with  
424 an iterative optimization algorithm called Iterative Conditional Mode whose  
425 maximum local energy for each site, at each iteration, is obtained with a  
426 K-nearest neighbour algorithm.

427 Mazeed, Nixon and Gunn (Al-Mazeed et al., 2004), whose work is closest  
428 to the work described in this paper, employ Bayes. Two background models  
429 are produced using a Mixture of Gaussians algorithm and a brightness and  
430 chromaticity algorithm referred to as Statistical Background Disturbance  
431 Technique (SBD). When the classifiers agree (pixel is foreground or back-  
432 ground) a decision is set accordingly. When classifiers disagree, conditional  
433 probability for the chosen class by each class is calculated. The product of  
434 each class of conditional probabilities provide the parameters for the final  
435 decision

$$\arg \max_{i \in \{1,2\}} p(x|w_{CLSF_i})P(w_{CLSF_i}) \quad (3)$$

436 where  $w$  is a class of either a background (BG) or a foreground (FG) for  
437 the classifier  $CLSF_i$ . The maximum conditional probability for each classifier  
438 is used with the classifier's confidence measure  $P(w_{CLSF_i})$  to find the decision  
439 for the algorithm. The main limitation of the approach is that it is limited  
440 to combination of two classifiers and that the priors are calculated using an  
441 exhaustive search method based on the training data to obtain the optimal  
442 values giving minimum classification errors.

443 While Bayesian inference, as well as other methods details above, have  
 444 been exploited for classification in motion segmentation, application of mu-  
 445 tual information to fuse multiple motion segmentation outputs has not been  
 446 studied. The approach taken here in selecting mutual information as a  
 447 method to combine multiple classifiers (the output from the LGN pathways)  
 448 is threefold: Firstly, in the same way the recognised behaviours of the LGN  
 449 pathways influenced the modelling of such, the identified interactions be-  
 450 tween these channels of visual information that occur in the visual cortex  
 451 influenced the choice of mathematical approach we use to model such find-  
 452 ings. Recent neurophysiological and vision research highlight that the output  
 453 of all three LGN pathways is shared within the visual cortex (McKeefry et  
 454 al., 2010; Farivar et al., 2009; Zanon et al., 2010; Briggs and Usrey, 2011).  
 455 Indeed (Clery et al., 2013) state that when considering the encoding of visual  
 456 information in the brain, the statistical independence between luminance and  
 457 chromatic edges in natural scenes vary depending on the dataset of natural  
 458 images used and “mutual information” may be found. These findings rule  
 459 out choosing methods of combining classifiers where the classifiers are com-  
 460 peting and a single classifier is found to be the “expert” at each instance for  
 461 example Behaviour Knowledge Space (Raudys et al., 2003) and those such  
 462 as the majority vote and K-nearest neighbour algorithm. As the information  
 463 theory principle of mutual information measures the amount of information  
 464 one random variable contains about another it is seemingly a sensible map-  
 465 ping to choose to model the neurophysiological and vision findings. Secondly,  
 466 consideration is taken regarding the data used from a statistical view point.  
 467 Multiple classifiers that produce probabilities as an output may be combined

468 using the product or average of the probabilities or the “Naïve Bayes” rule  
 469 however these combiners require that the individual classifiers use mutually  
 470 independent subsets of features (Kuncheva, 2001). This is not the case with  
 471 the output from the LGN pathways as each pathway produces an interpreta-  
 472 tion of identical data that each is presented with. Mutual information may  
 473 also be described as a technique that measures the mutual dependency of  
 474 one random variable with another and it is certainly the case with the LGN  
 475 outputs that there will be some commonality. In addition mutual informa-  
 476 tion classifiers have been found to provide an objective solution (Hu, 2012).  
 477 Finally, as the LGN pathways are modelled using real-time computational vi-  
 478 sion techniques, it is pertinent to choose a combining method such as mutual  
 479 information which, unlike techniques such as boosting, requires no additional  
 480 training on the data presented and may provide a fused result “on-the-fly”.

### 481 *3.1. Mutual Information*

482 In information theory the entropy of a discrete random variable  $X$  is  
 483 the measure of the amount of uncertainty associated with the value of  $X$ .  
 484 Shannon entropy, denoted by  $H$ , of a discrete random variable  $X$ , includes  
 485 a probability measure. If  $p$  represents a probability mass function of  $X$  then  
 486 Shannon entropy can be described in terms of a discrete set of probabilities

$$H(X) = - \sum_{i=1} p(x_i) \log p(x_i) \quad (4)$$

487 Mutual information  $I$  measures the amount of information that can be  
 488 obtained about one random variable by observing another. Mutual informa-  
 489 tion can be expressed as

$$\begin{aligned}
I(X;Y) &= H(X) - H(X|Y) \\
&= H(Y) - H(Y|X) \\
&= H(X,Y) - H(X|Y) - H(Y|X) \\
&= H(X) + H(Y) - H(X,Y)
\end{aligned} \tag{5}$$

490 where  $H(X)$  and  $H(Y)$  are the marginal entropies,  $H(X|Y)$  and  $H(Y|X)$   
 491 are the conditional entropies, and  $H(Y|X)$  is a measure of what  $Y$  does not  
 492 say about  $X$ .  $I(X;Y)$  is non-negative. Mutual information is a well estab-  
 493 lished technique for medical image registration of several modalities (Pluim  
 494 et al., 2003; Cheah, 2012) due to its insensitivity to changes in lighting condi-  
 495 tions ability to address a wide range of non-linear image transformations. It  
 496 has also been shown to be well suited to registration of images of the same  
 497 modality (Pluim et al., 2003).

498 Trivariate mutual information is described in various ways by authors  
 499 of research literature with reference to both the definition and in the use  
 500 of notation. Figure 8 provides examples of the assorted ways that (Pluim,  
 501 2003) discovered it had been defined and used in his survey of multivariate  
 502 mutual information in terms of entropies. The darker shaded areas represent  
 503 the mutual information in each case. (Pluim, 2003) asserts that a property  
 504 of the definition of Figure 8a. is that it is not necessarily nonnegative. In  
 505 Figure 8b. the deeper shaded middle section denotes that this area is counted  
 506 twice.

507 Figure 9 provides examples of how the notation varies between authors.  
 508 The diagrams labelled Figure 9a., Figure 9b. and Figure 9c. depict a bi-

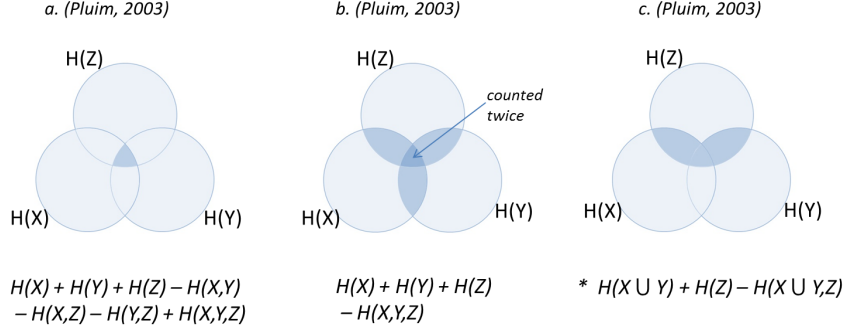


Figure 8: Different definitions of trivariate mutual information in terms of Shannon entropies. Each circles denote the entropy of an image. \*Definition from (Pluim, 2003) text.

509 variate and two trivariate examples respectively and the notation to describe  
 510 them given by (Studholme, 1996). He uses a ‘;’ to separate the arguments  
 511 for mutual information, while a ‘,’ denotes a union of two variables. The  
 512 notation used by (Pluim, 2003) differs in that to describe the same examples  
 513 in the diagrams labelled Figure 9d., Figure 9e. and Figure 9f. ‘,’ is used  
 514 as the separator between the arguments and is not a union. Further to the  
 515 differences found in literature in the notation, (MacKay, 2003) states that  
 516 the term  $I(X;Y;Z)$  is illegal. For clarity in this work the notation used  
 517 throughout is that of (MacKay, 2003) which is consistent with (Studholme,  
 518 1996) and later authors (Escolano et al., 2009).

519 In this work the variables  $X$ ,  $Y$  and  $Z$  are the probability in each LGN  
 520 stream (parvocellular, magnocellular, and koniocellular) that a pixel is fore-  
 521 ground. Here mutual information is used as a measure of the information  
 522 or interaction between any two or all three LGN streams. To this end,  
 523 CMI (Combined Mutual Informations) is defined as a linear combination of

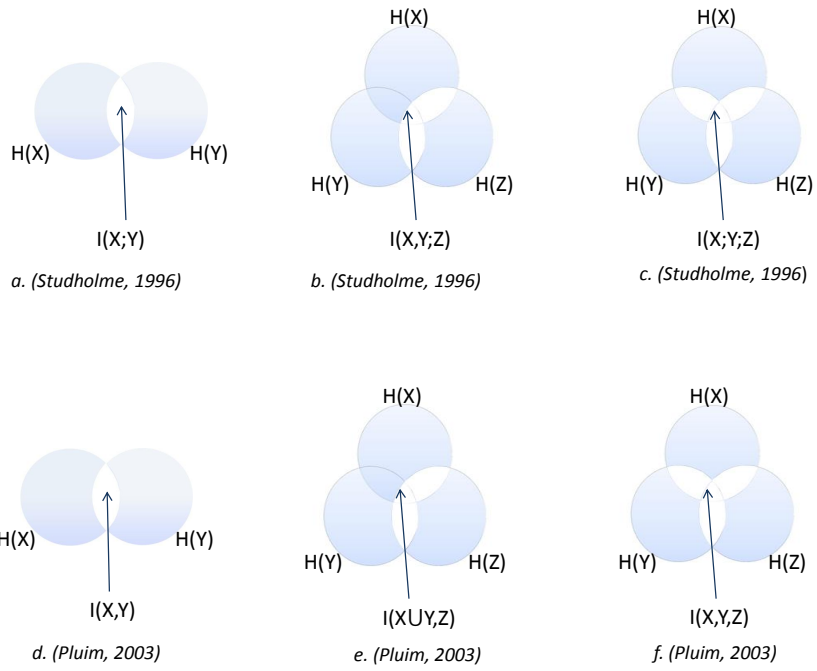


Figure 9: Differing notations describing the same mutual information examples

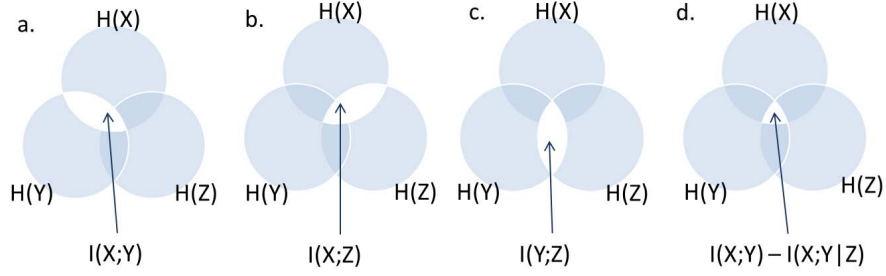


Figure 10: Bivariate and trivariate mutual information in terms of Shannon entropies

trivariate mutual information for all three LGN streams and bivariate mutual information for each pair of LGN streams such that none of the constituent entropies are counted twice. To avoid the use of any terms which could be considered illegal, the only trivariate mutual information used here will be of the form  $I(X;Y|Z)$  which is the mutual information between  $X$  and ( $Y$  given  $Z$ ) and is considered a legal term (MacKay, 2003).

Bivariate mutual informations are  $I(X;Y)$ ,  $I(X;Z)$  and  $I(Y;Z)$  (Figure 10 a., b. and c. respectively) and are expressed in terms of Shannon entropies as

$$\begin{aligned}
 I(X;Y) &= H(X) + H(Y) - H(X,Y) \\
 I(X;Z) &= H(X) + H(Z) - H(X,Z) \\
 I(Y;Z) &= H(Y) + H(Z) - H(Y,Z)
 \end{aligned}
 \tag{6}$$

Trivariate mutual informations are  $I(X;Y|Z)$ ,  $I(X;Z|Y)$  and  $I(Y;Z|X)$ . In

534 terms of Shannon entropies  $I(X; Y|Z)$  is defined as

$$I(X; Y|Z) = -H(Z) + H(X, Z) + H(Y, Z) - H(X, Y, Z) \quad (7)$$

535 The quantity  $I(X; Y) - I(X; Y|Z)$  is shown in Figure 10d. and may also be  
536 defined as

$$\begin{aligned} I(X; Y) - I(X; Y|Z) &= I(X; Z) - I(X; Z|Y) \\ &= I(Y; Z) - I(Y; Z|X) \end{aligned} \quad (8)$$

537 Therefore a consistent quantity  $CMI$ , with no overlapping entropies may be  
538 defined as

$$\begin{aligned} CMI &= I(X; Y) + I(X; Z) + I(Y; Z) \\ &\quad - 2[I(X; Y) - I(X; Y|Z)] \end{aligned} \quad (9)$$

539  $CMI$  can thus be expanded to give

$$\begin{aligned} CMI &= I(X; Y) + I(X; Z) + I(Y; Z) - 2[I(X; Y)] \\ &\quad + 2[I(X; Y|Z)] \\ &= -I(X; Y) + I(X; Z) + I(Y; Z) \\ &\quad + 2[I(X; Y|Z)] \end{aligned} \quad (10)$$



540 which can be expressed in terms of Shannon entropies as

$$\begin{aligned}
CMI &= -H(X) - H(Y) + H(X, Y) \\
&+ H(X) + H(Z) - H(X, Z) \\
&+ H(Y) + H(Z) - H(Y, Z) \\
&+ 2[H(X, Z) + H(Y, Z) - H(X, Y, Z) - H(Z)]
\end{aligned} \tag{11}$$

541 and can be simplified as

$$CMI = H(X, Y) + H(X, Z) + H(Y, Z) - 2H(X, Y, Z) \tag{12}$$

542 Since

$$H(X) = - \sum_{i=1} p(x_i) \log p(x_i) \tag{13}$$

543  $CMI$  may be rewritten as

$$\begin{aligned}
CMI &= - \sum_{x,y} p(x, y) \log p(x, y) - \sum_{z,y} p(y, z) \log p(y, z) \\
&- \sum_{x,z} p(x, z) \log p(x, z) + 2 \sum_{x,y,z} p(x, y, z) \log p(x, y, z)
\end{aligned} \tag{14}$$

544 and yields an expected value over all possible instances of  $X, Y$  and  $Z$ .

545 The quantities given below, that are summed to find CMI, exist at all

546 points  $x, y, z$ .

$$\begin{aligned}
& p(x, y) \log p(x, y) \\
& p(x, z) \log p(x, z) \\
& p(y, z) \log p(y, z) \\
& p(x, y, z) \log p(x, y, z)
\end{aligned}
\tag{15}$$

547 The two variable quantities are each defined on a 2D grid and the three vari-  
548 able quantity is defined on the 3D space  $(x, y, z)$ . Hence  $p(x, y, z) \log p(x, y, z)$   
549 may have a different value at all points  $(x, y, z)$  where as  $p(x, y) \log p(x, y)$  is  
550 only defined on the  $x, y$  grid and values at any point  $(x, y)$  are the same for  
551 all  $z$ . It is therefore possible to define a quantity  $pVC$  at each point based  
552 on the point wise constituents of CMI.

$$\begin{aligned}
pVC = & -p(x, y) \log p(x, y) - p(y, z) \log p(y, z) - p(x, z) \log p(x, z) \\
& + 2p(x, y, z) \log p(x, y, z)
\end{aligned}
\tag{16}$$

553 This provides a nonnegative result and is referred to as the Visual Cortex  
554 (VC) model in the following text.

555 The approximated probability mass functions produced by respectively  
556 the GMM, Brightness and Chromaticity, and Colour, Mean and Variance  
557 algorithms provide the mutual information required to produce silhouettes  
558 of objects of interest. For Brightness and Chromaticity, the probability that a  
559 pixel is foreground (FP) may be computed as (see Section 2.4.1 for notation)

$$FP = \frac{p(1 - p(\widehat{CD}_i))p(\widehat{\alpha}_i)}{p(\widehat{\alpha}_i)}
\tag{17}$$

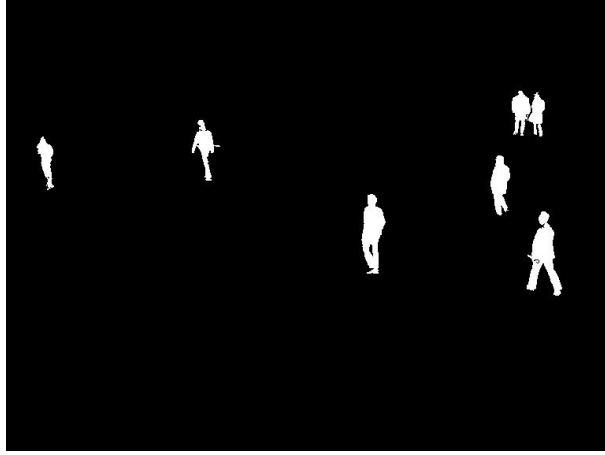


Figure 11: PETS 2009 dataset frame - resulting segmentation using the VC model

560 For Colour, Mean and Variance, the probability that a pixel is foreground  
 561 (FP) may be computed as follows:

$$FP = p(R_i \cup G_i \cup B_i \cup A_i) \quad (18)$$

562 The probability for the Gaussian Mixture Model may be computed as  
 563 given in equation 1.

564 Figure 7 represents the classification by the VC model of foreground pixels  
 565 (white) from the original frame in Figure 1

## 566 4. Experimental Results

### 567 4.1. Ground Truth

#### 568 4.1.1. Silhouettes

569 The binary silhouettes of both the MuHAVi and PAMELA datas were  
 570 hand labelled for all frames. For MuHAVi, Manually Annotated Silhouette  
 571 Data (MAS) consists of annotated footage of 5 action classes. They include

two different actors and two separate camera views. In this case the annotation consists of white silhouettes of the actors performing their actions on a black background.

#### 4.1.2. *Objects*

Each of the PETS2009 seven independent 2D camera views (views 1,3,4,5,6,7,8) and CAVIAR “Walk” and “Walk 2” sequences were ground truthed frame by frame using the Video Performance Evaluation Resource (ViPER-GT) ground truth tool (Mariano et al., 2002). The ground truth consists of bounding boxes that are created around the objects and the coordinate positions of these boxes within the scene are given in a ground truth XML file.

#### 4.2. *Background Learning*

Each of the three motion segmentation methods used to model the LGN pathways require an initial “learning” phase, where the algorithms produce a statistical interpretation of the initial scene. Visual surveillance scenes are frequently dynamic in nature and whilst lengthy “background learning” sequences may produce a better motion segmentation from each of the algorithms this is mostly not practical due to rapidly changing scenes. To capture a scene or “background” where there is little of interest happening it is prudent to use as short a number of frames as is possible when initialising each of the motion segmentation algorithms. With this in mind for all datasets and sequences the following initialisations to the algorithms were given. The BC algorithm was set to a “background run length” of 100 frames, the initial  $a_i$  and  $b_i$  calculations used 50 frames and the initial histograms were created with just 10 frames. The GMM in this case was set to three gaussians, had

596 a “background run length” of 100 frames and calculated Expected Maximisation (EM) from just 20 frames. The CMV algorithm initialised with 10  
597 background frames. For all algorithms a weight of 0.0001 was set for the  
598 learning rate.

### 600 4.3. Datasets

601 Four different datasets are used to test the performance of the proposed  
602 Visual Cortex model, the publicly available MuHAVi (Singh et al., 2010),  
603 CAVIAR, PETS2009 (Ferryman and Ellis, 2009), and the datasets produced  
604 for the Background Models Comparison (BMC) challenge (Vacavant et al.,  
605 2012).

606 The first dataset, MuHAVi (Singh et al., 2010), introduces the challenge  
607 of real night-time street lighting, street paving (reflective) and real high street  
608 surveillance camera footage (with glare and large prominent shadows) to the  
609 motion segmentation algorithms. There is also some camouflage of individuals  
610 present, where the clothing and the background are similar in colour.

611 CAVIAR Walk 1 and Walk 2 indoor datasets include sunlight shining  
612 through large glass panels and producing variable lighting within an indoor  
613 scene, alongside intermittent and unpredictable shadows of the panel frames  
614 on the floor. Reflections appear intermittently on additional glass panels that  
615 reside inside the building, and sunlight reflects from these panels. Shadows  
616 are present when individuals walk through the scene and some camouflage is  
617 present with the clothing of certain individuals and the background.

618 The third dataset, (Ferryman and Ellis, 2009), comprises multi-sensor sequences  
619 containing crowd scenarios with increasing scene complexity. Dataset  
620 S2, used in this evaluation, addresses people detection and tracking. Spe-

621 cific challenges include occluding ,moving objects encompassing whole scenes;  
 622 moving vegetation; vehicles; motion behind translucent windows; reflective  
 623 surfaces; objects appearing both very large and close to the camera and small  
 624 and in the far distance; lack of natural lighting to entire footage.

625 Finally, the BMC dataset consists of both synthetic and real world videos.  
 626 The synthetic videos present a variety of cloudy, sunny, foggy and windy  
 627 scenes with and without acquisition noise. The real world videos contain  
 628 challenges such as outdoor scenes, lengthy videos, varying ground types,  
 629 presence of vegetation, casted shadows and the presence of continuous flow  
 630 of objects.

#### 631 4.4. *Evaluation Metrics*

632 Performance evaluation was based on Precision and F1 Score Metrics  
 633 and the framework by (Kasturi et al., 2009), a well established protocol for  
 634 performance evaluation of object detection and tracking in video sequences.  
 635 These metrics are formally used by the Video Analysis and Content Extrac-  
 636 tion (VACE) programme and the CClassification of Events, Activities, and  
 637 Relationships (CLEAR) consortium.(Vacavant et al., 2012) provides details  
 638 for the F-score and SSIM metric used for the Background Model Challenge  
 639 dataset.

#### 640 *Notation.*

- 641 •  $G_i^t$  denotes  $i^{th}$  ground-truth object in frame  $t$ ;  $G_i$  denotes the  $i^{th}$  ground-  
 642 truth object at the sequence level;  $N_{frames}$  is the number of frames in  
 643 the sequence

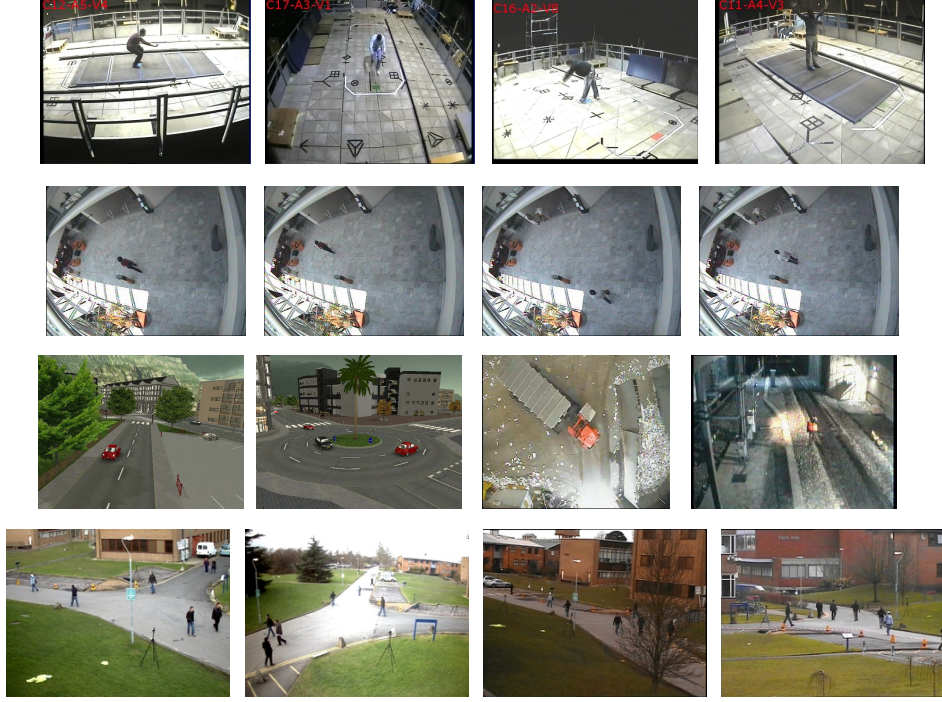


Figure 12: Datasets used. Top row: Four views from MuHAVi which contains sequences with realistic street scenes. Second row: Four example frames from CAVIAR Walk 1 (left two images) and Walk 2 (right two images) sequences. Third row: Four example frames from Background Model Challenge dataset which contains both synthetic and real videos. Fourth row: Four views from the PETS2009 dataset which contains a range of crowd-based scenarios.

- 644 •  $D_i^t$  denotes the  $i^{th}$  detected object in frame  $t$ ;  $D_i$  denotes the  $i$ th de-  
645 tected object at the sequence level
- 646 •  $N_G^t$  and  $N_D^t$  denote the number of ground-truth objects and the num-  
647 ber of detected objects in frame  $t$ , respectively;  $N_G$  and  $N_D$  denote  
648 the number of unique ground-truth objects and the number of unique  
649 detected objects in the given sequence, respectively

- 650 •  $N_{frames}^i$  refers to the number of frames where either ground-truth object  
651 ( $G_i$ ) or the detected object ( $D_i$ ) existed in the sequence
- 652 •  $N_{mapped}$  refers to sequence level detected object and ground truth pairs,  
653  $N_{mapped}^t$  refers to frame  $t$  mapped ground truth and detected object  
654 pairs
- 655 •  $m_t$  represents the missed detection count, ( $fp_t$ ) is the false positive  
656 count,  $c_m$  and  $c_f$  represent respectively the cost functions for missed  
657 detects and false positives, and  $c_s = \log_{10}ID - SWITCHES_t$

#### 658 4.4.1. Precision and F1 Score

659 Pixel based metrics are computed from pixel counts that may be classified  
660 as true positives (TP), false positives (FP), false negatives (FN), and true  
661 negatives (TN). FP and FN refer to those that are misclassified as pixels  
662 belonging to the objects of interest (FP) or the background (FN) while TP  
663 and TN account for accurately classified pixels.

664 The precision of a silhouette is an important factor for the reasoning of  
665 behaviour using pose and gait techniques, and is found by:

$$Precision = 100 - \left[ \left( \frac{FN + FP}{TP + FN} \right) \times 100 \right] \quad (19)$$

666 The F1 score is a popular metric for evaluation of segmentation and  
667 represents a measure of the accuracy of an algorithm and is found by:

$$F1Score = \frac{2TP}{((TP + FN) + (TP + FP))} \quad (20)$$



#### 668 4.4.2. Sequence Frame Detection Accuracy (SFDA)

669 SFDA uses the number of objects detected, the number of missed de-  
 670 tections, the number of falsely identified objects, and the calculation of the  
 671 spatial alignment between the algorithm’s output for detected objects and  
 672 that of the ground truthed objects. It is derived from a Frame Detection  
 673 Accuracy (FDA) measure. The FDA is calculated using a ratio of the spa-  
 674 tial intersection and union of an output object and mapped ground truthed  
 675 objects

$$OverLapRatio = \sum_{i=1}^{N_{mapped}^t} \frac{|G_i^t \cap D_i^t|}{|G_i^t \cup D_i^t|} \quad (21)$$

$$FDA(t) = \frac{OverLapRatio}{\left\lceil \frac{N_G^t + N_D^t}{2} \right\rceil} \quad (22)$$

$$SFDA = \frac{\sum_{t=1}^{N_{frames}} FDA(t)}{\sum_{t=1}^{N_{frames}} \exists (N_G^t \vee N_D^t)} \quad (23)$$

676 For this study although the annotation of the ground truth was challeng-  
 677 ing, an overlap threshold of 100 percent for the intersection over union scores,  
 678 was used.

679 For both detection and tracking metrics in the following descriptions the  
 680 accuracy metrics provide a measure of the correctness of the detections or  
 681 tracks. The precision metrics provide the measure of, in the instance where  
 682 there has been a correct detection or track, how close to the ground truth  
 683 that detection or track may be.

#### 684 4.4.3. Multiple Object Detection Accuracy (MODA)

685 MODA is an accuracy measure that uses the number of missed detections  
686 and the number of falsely identified objects. Cost functions to allow weighting  
687 to either of these errors are included, however for the sake of both PETS 2009  
688 evaluations they were equally set to 1.

$$MODA = 1 - \frac{c_m(m_t) + c_f(f_{p_t})}{N_G^t} \quad (24)$$

#### 689 4.4.4. Multiple Object Detection Precision (MODP)

690 MODP gives the precision of the detection in a given frame. Again, with  
691 this metric, an overlap ratio is calculated as previously defined in (1), and, in  
692 addition to a count of the number of mapped objects, the MODP is defined  
693 as:

$$MODP(t) = \frac{OverLapRatio}{N_{mapped}^t} \quad (25)$$

### 694 4.5. Results

#### 695 4.5.1. MuHAVi

696 The three individual segmentation algorithms and Visual Cortex algo-  
697 rithm were evaluated on the MuHAVi dataset against ground truth using  
698 the Precision and F1 Metrics. Comparisons are then made frame by frame  
699 between the algorithms resulting silhouette and the ground truth. True posi-  
700 tive, false positive, true negative and false negative pixels are counted for each  
701 frame. Figure 13 shows the robust nature of the Visual Cortex model, respec-  
702 tively for F1 score (14) and Precision (13), using the mutual information of  
703 the three LGN pathways, in comparison to their independent performances.

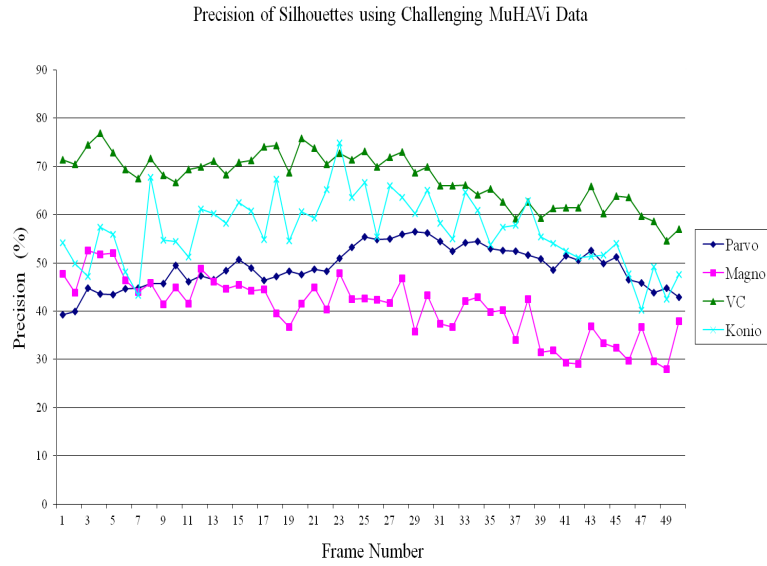
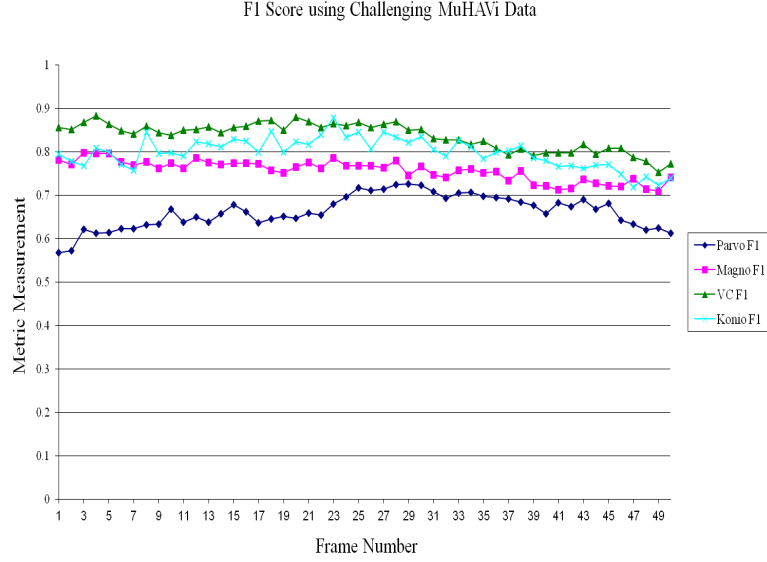


Figure 13: Accuracy (top) and precision (bottom) of the silhouettes produced by the independent LGN pathways versus the mutual information of the VC model on the challenging MuHAVi dataset.

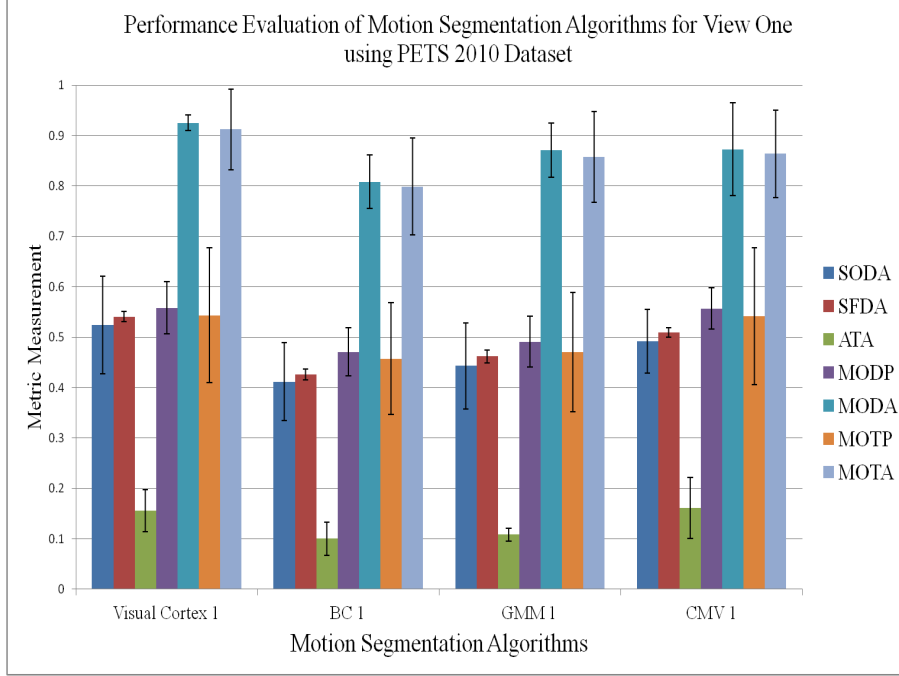


Figure 14: Performance of Visual Cortex and individual motion segmentation algorithms for view one of PETS2009 dataset.

#### 4.5.2. PETS2009

The next set of evaluations show comparisons of the performance of individual motion segmentation algorithms against the Visual Cortex model for the PETS2009 dataset. Figure 14 represents the evaluation results for sequence S2.L1, at time sequence 12.34, for the first camera view. and illustrates the superior performance of the Visual Cortex model, in comparison to the established motion segmentation algorithms, for the detection of objects within the surveillance scene. Every object detection metric, SODA, SFDA, MODA and MODP evaluates the Visual Cortex model (VC) as the best in

713 performance for its criteria, with the detection precision (MODP) metric  
 714 proving the performance of the CMV algorithm as equal to that of the Vi-  
 715 sual Cortex model. Referring to the MOTA tracking metric, further analysis  
 716 of Figure 14 demonstrates the increase in performance in tracking accuracy  
 717 using the Visual Cortex model as the motion segmentation algorithm base  
 718 for the tracker.

719 Next, to assess robustness in real world scenarios the Kanade-Lucas-  
 720 Tomasi (KLT) tracking algorithm (Tomasi and Kanade, 1991) was used with  
 721 individual sets of motion segmentation silhouette results using the PETS2009  
 722 dataset to produce tracking results, and in turn 2D bounding box coordi-  
 723 nate positions and unique identifiers for each object for view one of the  
 724 PETS2009 dataset. The performance evaluation results of the PETS 2009  
 725 and PETS2010 workshops (Ellis et al., 2010) were used to enable the com-  
 726 parisons. The SODA, SFDA, MODA and MODP metrics are relevant to the  
 727 evaluation of the motion segmentation algorithms of the workshop’s partic-  
 728 ipating authors systems in addition to that of the Visual Cortex model. A  
 729 summary of their motion segmentation/object detection techniques follow in  
 730 order that comparisons may be drawn:

731 (Arsic et al., 2009) employ a multi-layer homography, which is capable  
 732 of creating a three dimensional representation of the scene. Homography  
 733 frameworks rely on the fusion of previously segmented foreground regions  
 734 visible from multiple views. In the case of (Arsic et al., 2009) system, these  
 735 foreground segmentations are produced by finding the median of pixel values  
 736 and composing a reference image for simple background subtraction. Bright-  
 737 ness invariance is achieved by normalised cross covariance when compared

738 with the reference image and contrast invariance is achieved using normalised  
739 cross-correlation. A graph cut optimisation algorithm is then optionally car-  
740 ried out to fill in small holes in foreground silhouettes.

741 (Breitenstein et al., 2009) presents a HOG object detector producing the  
742 input for the observation model of a particle filter, which includes not only the  
743 objects detected, but their confidence density of that detection (rep-resented  
744 as a colour heat map). Each object has its own particle filter initialised which  
745 includes its position and velocity. Bounding boxes are created by a boosted  
746 ensemble of weak classifiers employing colour histograms.

747 (Yang et al., 2009) utilises dynamic appearance models, using single Gaus-  
748 sians for foreground descriptions, and a Gaussian background model.

749 (Alahi et al., 2009) creates degraded foreground silhouettes from some  
750 binary silhouette image and its approximation, using rectangular and ellipse  
751 shapes. These then help form the input to a Multi-Silhouette Dictionary  
752 which is made up of atoms modelling the presence of individuals at give  
753 locations on an occupancy grid. The atoms are generated using homogra-  
754 phies mapping points in a three dimensional scene to their two dimensional  
755 coordinates in the planar view.

756 (Bolme et al., 2009) approaches the challenge with the object detection  
757 filtering method Average of Synthetic Exact Filters which considers the entire  
758 output of the filter un-der a full convolution operation. He also uses a Viola  
759 and Jones cascade classifier with both visual and motion features used for  
760 detection. The third detector he uses is based on the deformable parts model  
761 system.

762 (Ge et al., 2009) regard people in a crowd scene as a realisation of a

763 Marked Point Process. Each person is associated with a random mark that  
764 specifies their location and size within the frame. A binary foreground mask  
765 is obtained by an adaptive background subtraction method and is subjected  
766 to further morphological processing. This then becomes the input to the  
767 detector.

768 (Conte et al., 2010) utilise an adaptive background image difference al-  
769 gorithm to detect moving objects. In order to make the system robust in  
770 realistic environments this has been extended to included processes that han-  
771 dle illumination, camouflage detection, noise filtering, shadow filtering and  
772 reflection removal.

773 (Berclaz et al., 2009) employ an object detector that produces a proba-  
774 bilistic occupancy grid, using a set of prob-abilities of the presence of objects,  
775 at a discrete set of locations, at each time step. These objects are represented  
776 as cylinders that project to rectangles in the frame sequences.

777 Figure 15 shows that the Visual Cortex model outperforms the evaluation  
778 of the individual algorithms with respect to the accuracy of both the detection  
779 of the objects and the tracking, using view one of the PETS 2009 datasets  
780 and the SODA, SFDA, MODP and MODA metrics.

781 It should be noted that the accuracy of the tracking algorithm used im-  
782 proves with the accuracy of the segmentation. The precision of any single  
783 detected object in this case refers to the precision of the location of its bound-  
784 ing box enclosing the object, that the tracker has produced, and not the pre-  
785 cision of the silhouettes previously measured. Note that the standard error  
786 of mean (SEM) error bars have been added to the performance evaluation  
787 results charts. These quantify how precisely the true mean is known, taking

788 into account both the standard deviation and the sample size. Looking at  
789 whether the error bars overlap, therefore enables comparison of the difference  
790 between the mean with the precision of those means. It is very important to  
791 note that if two SEM error bars do overlap, and the sample sizes are equal  
792 the difference is not statistically significant, however if two SEM error bars  
793 do not overlap no conclusions may be made about statistical significance.

794 It is clear that for this sequence, the systems described by (Breitenstein  
795 et al., 2009) performed strongly at multiple object detection and tracking,  
796 with (Yang et al., 2009) outperforming all others. However the Visual Cor-  
797 tex model provides a strong performance in object detection and outperforms  
798 Breitenstein’s system for detection accuracy (MODA) using the Visual Cor-  
799 tex model motion segmentation algorithm alone. Most detection and track-  
800 ing systems employ further processing filters after any initial segmentation  
801 to improve the motion segmentation quality. This is not the case with the  
802 Visual Cortex model. The tracking accuracy (MOTA) gained from using the  
803 Visual Cortex model is second only to the system produced by Yang. As  
804 both Breitenstein and Yang did not provide results for views 5,6, and 8 no  
805 further comparisons or analysis of robustness using these systems may be  
806 drawn. (Ge et al., 2009), (Berclaz et al., 2009) and (Conte et al., 2010) de-  
807 tection accuracy measures (MODA) also suggested a good performance for  
808 these particular areas, as do (Berclaz et al., 2009), (Conte et al., 2010), and  
809 AlahiOlasso (Alahi et al., 2009) for tracking accuracy (MOTA).

#### 810 4.5.3. CAVIAR

811 Two “Walk” sequences from CAVIAR were evaluated against using the  
812 SODA, SFDA, MODP and MODA metrics. The Visual Cortex model again



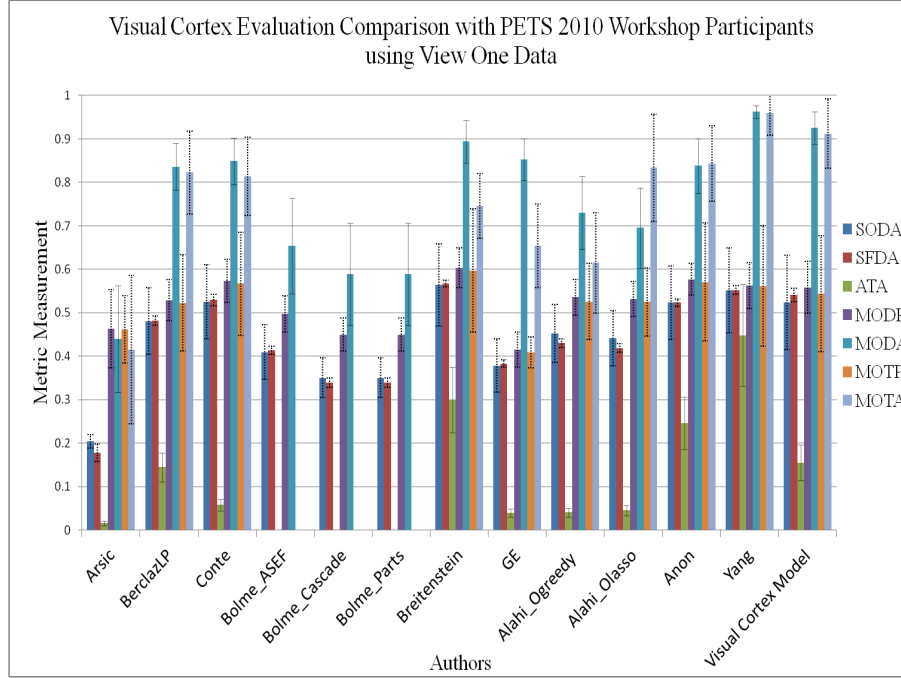


Figure 15: Performance of participating authors' systems, using CLEAR and VACE metrics for view one of PETS2009 dataset, mean SEM, N=109.

813 outperforms all three motion segmentation algorithms for each metric cate-  
814 gory despite the datasets being of a completely different nature to MuHAVi  
815 and PETS2009.

#### 816 4.5.4. BMC dataset

817 Finally, the synthetic and real datasets provided for this BMC special  
818 issue were evaluated and are shown in Figure 17. You can see from these that  
819 the VC model generally performs more robustly to the variety of sequences  
820 than published algorithms BC, GMM and CMV, in both synthetic and real  
821 world scenarios. The results for the synthetic videos show improvement on

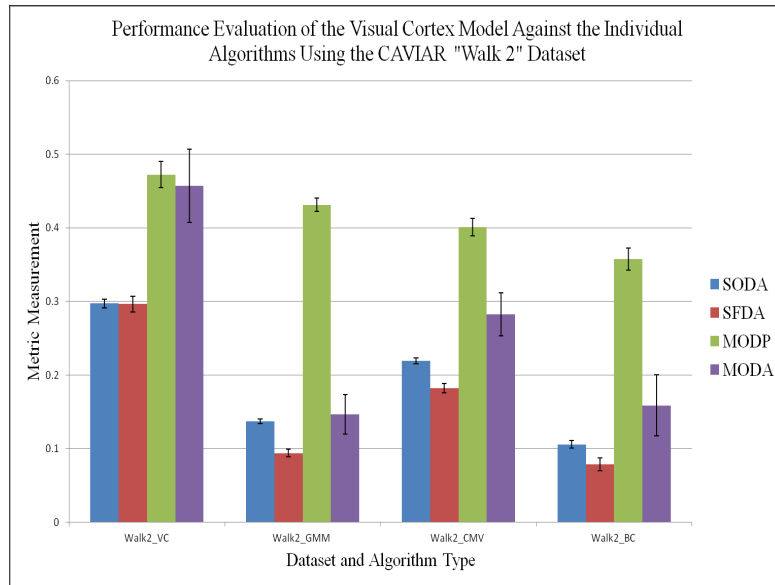
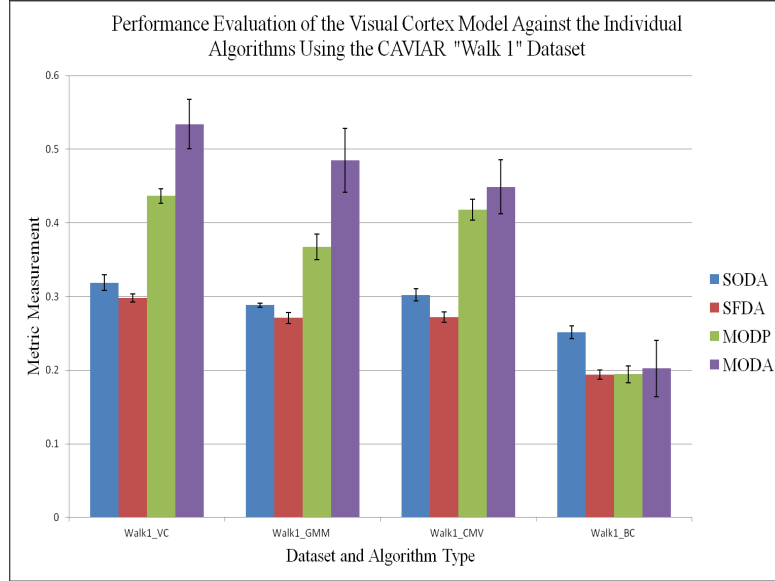


Figure 16: Comparing the Performance Evaluation of the Visual Cortex model with established motion segmentation algorithms using the CAVIAR (top) "Walk 1" and (bottom) "Walk 2" dataset, mean SEM, N=610.

822 the CMV, BC and GMM algorithms by employing the VC model, using both  
823 the F-Score and the SSIM metric as a measure, for all cases of videos tested.  
824 The individual algorithms however do not include any form of additional  
825 object recognition processing (and this is outside the scope of the biological  
826 model presented) that would distinguish between the cars travelling on the  
827 road and moving ground-truthed objects in the car park within the real  
828 world Video 1 scenario. In addition the VC model attempts to create a better  
829 silhouette of both the cars on the road and the ground-truthed cars in the car  
830 park than ones presented by the individual CMV, GMM and BC algorithms  
831 and as such is penalised by the pixel-based F Score metric for doing so. This  
832 is also the case for Video 8 where there is an additional flow of traffic to that  
833 which has been ground-truthed. It should be noted that pixel based metrics  
834 such as the F score can be heavily biased towards the larger moving objects  
835 within a frame when a video sequence contains more than one object and/or  
836 perspective plays a part. This bias is inherent in the results. The SSIM  
837 metric measures, for each real video sequence, highlight the visual structural  
838 (silhouettes) improvement gain made using the VC model, as opposed to the  
839 individual CMV, BC, and GMM algorithms.

840 The performance evaluation results of the Background Models Challenge  
841 workshop (Vacavant et al., 2012) participating authors' systems are shown in  
842 Figure 18. The VC model represents the results of motion segmentation only  
843 and does not include any additional processing techniques that may be added  
844 to assist in the elicitation of objects from the background. The VC model  
845 shows a noticeable comparison to all participating authors' background model  
846 systems with regard to the SSIM metric. The F-score metric highlights the

847 difficulty in producing a robust background model system for all scenarios,  
848 where generally the performances of each individual system appears to vary  
849 depending on the scenario it is presented with. A summary of the workshop’s  
850 participating authors’ techniques follow:

851 (Yoshinga et al., 2013) use illumination invariant local features and de-  
852 scribe their distribution by Gaussian Mixture Models. The local feature has  
853 the ability to tolerate the effects of illumination changes, and the GMM can  
854 learn the variety of motion changes. Radial distances control the local feature  
855 and the localized regions focused by each pixel.

856 For (Shah et al., 2013) A Gaussian mixture model is used as a background  
857 basis and a new match function is used by computing separate variances for  
858 colour and intensity channels. For every foreground blob SURF features are  
859 matched and irrelevant features are removed using RANSAC sampling. The  
860 weight of winning Gaussian is increased a little for foreground blobs detected  
861 as paused objects. Automatic parameter adaptation is achieved using a fixed  
862 length sliding window to keep the most recent N frames in order to capture  
863 continuing statistical changes.

864 (Glazer et al., 2013) use one-class SVM classifiers to model the distribu-  
865 tion of the background. Three levels of resolution are used: block, region and  
866 frame. Images are divided in to equal-sized blocks of pixels and the one-class  
867 SVMs are independently trained on each block to model its background dis-  
868 tribution. Inter block relationships are used to refine the classification results  
869 at region level and at frame level an adaptive background method is used to  
870 re-initialise the model with regions considered to be part of the background.

871 (Tavakoli et al., 2013) introduce a method of estimating motion saliency

872 based on temporal cues obtained using frame de-correlation. Temporal salience  
873 maps are computed, presenting the amount of motion in a frame. Salient mo-  
874 tion is assumed steady and the focus is on the detection of firm movements.  
875 Principal components analysis is applied for reconstruction whilst suppress-  
876 ing background clutter and noise.

877 (Guyon et al., 2013) use Robust Principal Components Analysis (RPCA)  
878 to separate moving objects from the background. The background sequence is  
879 then modelled by a low rank subspace, using a low-rank matrix factorization  
880 with iteratively reweighted least squares that can gradually change over time.  
881 The moving foreground objects constitute the correlated sparse outliers.

## 882 5. Conclusions and Future Work

883 This paper has presented a novel neuroscience inspired information the-  
884 oretic approach to motion segmentation. In applying current neurological  
885 and physiological research in primate vision, a system has been created to  
886 improve the robustness of a multidimensional motion segmentation system.  
887 The major result found in this investigation is in using the current under-  
888 standing of the primate visual system as inspiration and guidance for choos-  
889 ing both feature sets (the LGN pathways), and the means of fusing them  
890 (the Visual Cortex model), considerably improves the appearance of the ob-  
891 tained silhouettes, without the need for subjective parameter adjustments, or  
892 the use of arbitrary thresholds. This presents an advantage over established  
893 multidimensional models which frequently rely on decisions, based on some  
894 weighting, whether a feature set provides the correct segmentation. These  
895 techniques are burdened with adjusting parameters, which do not necessarily

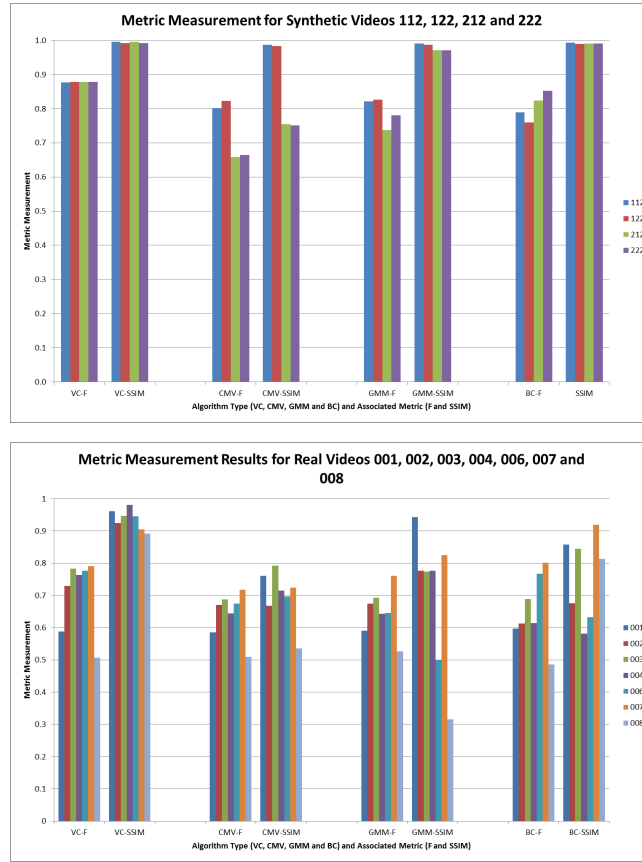


Figure 17: Comparing the performance of the Visual Cortex model with established motion segmentation algorithms using the BMC (top) synthetic and (bottom) real videos.

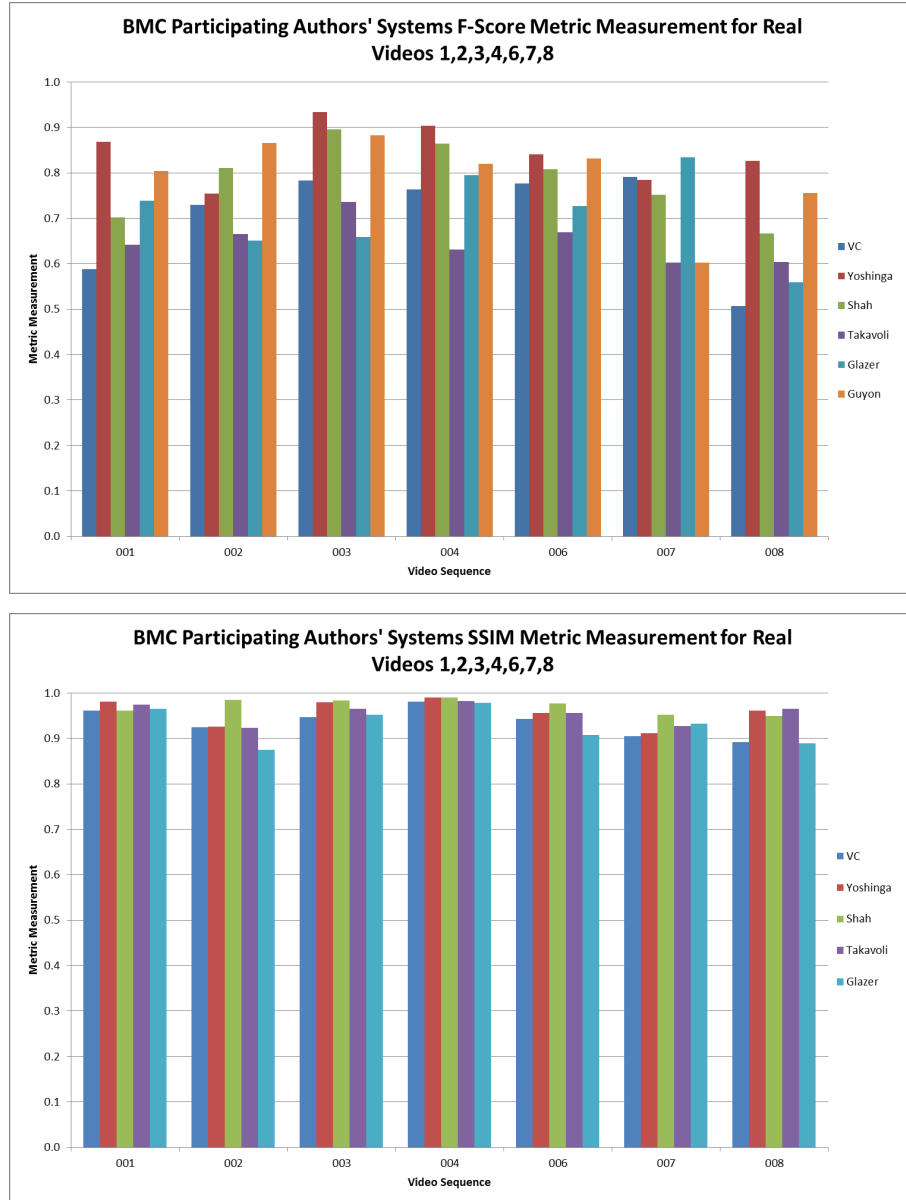


Figure 18: Comparing the performance of the Visual Cortex model with those of the participating authors' systems in the BMC challenge with the real videos dataset and F-Score(top) and SSIM (bottom) metrics.

896 provide the correct decision for all cases. This work has presented the perfor-  
897 mance evaluation of the biologically inspired motion segmentation system in  
898 challenging and diverse scenarios using a variety of evaluation metrics. In ad-  
899 dition the evaluation results of state of the art automated visual surveillance  
900 systems have been presented to enable comparisons to be drawn. It shows  
901 that biologically inspired automated visual surveillance detection systems  
902 may be considered comparable to the current state of the art surveillance  
903 systems in detection and tracking. Existing real-time computational vision  
904 techniques have been exploited in the production of feature sets similar to  
905 that which the primate retina produces with a view towards real-time bio-  
906 logically inspired visual surveillance systems. The “reasoning” made within  
907 the visual cortex model employs a technique already well-established in the  
908 registration of medical images. It is envisaged that refining the LGN pathway  
909 approximations to closer representations of the biological system may result  
910 in robust performance beyond that of the current model. Further research  
911 into biologically guided object detection may provide a further processing  
912 model with a view to presenting robust object detection in addition to mo-  
913 tion segmentation.

## 914 **Acknowledgements**

915 This work was supported by the EC project ARENA Grant Agreement  
916 No. 261658. Any opinions expressed in this paper do not necessarily reect  
917 the views of the European Community. The Community is not liable for any  
918 use that may be made of the information contained herein.

919 The authors would like to thank M. J. Lally, School of Mathematical and



920 Physical Sciences, University of Reading, UK

921 A. Alahi, L. Jacques, Y. Boursier and P. Vandergheynst, Sparsity-Driven  
922 People Localization Algorithm: Evaluation in Crowded Scenes Envi-  
923 ronments, Proceedings of the Twelfth IEEE International Workshop on  
924 Evaluation of Tracking and Surveillance (PETS-Winter), 2009, DOI:  
925 10.1109/PETSWINTER.2009.5399487

926 T. D. Albright, R. Desimone, Local Precision of Visuotopic Organization  
927 in the Middle Temporal Area (MT) of the Macaque, Experimental Brain  
928 Research, vol. 65(3), pp. 582-592, 1987.

929 A. Al-Mazeed, M. Nixon and S. Gunn, Classifiers Combination for Improved  
930 Motion Segmentation, Proceedings of International Conference on Image  
931 Analysis and Recognition, vol. 3212, pp. 363-371, 2004, ISBN: 3-540-23240-  
932 0

933 D. Arsic, A. Lyutskanov, G. Rigoll and B. Kwolek, Multi Camera Person  
934 Tracking Applying a Graph-Cuts Based Foreground Segmentation in a Ho-  
935 mography Framework, Proceedings of Twelfth IEEE International Work-  
936 shop Performance Evaluation of Tracking and Surveillance (PETSWinter),  
937 2009, DOI: 10.1109/PETS-WINTER.2009.5399723

938 S. Avidan, Support Vector Tracking, IEEE Transactions on Pattern Anal-  
939 ysis and Machine Intelligence, vol. 26(8), pp. 1064-1072, 2004, DOI:  
940 10.1109/TPAMI.2004.53

941 S. A. Baccus, B. P. Iveczky, M. Manu and M. Meister, A Retinal Circuit

- 942 That Computes Object Motion, The Journal of Neuroscience, vol. 28, pp.  
943 6807-6817, 2008, DOI: 10.1523/JNEUROSCI.4206-07.2008.
- 944 P. Bayerl and H. Neumann, A Fast Biologically Inspired Algorithm for  
945 Recurrent Motion Estimation, IEEE Transactions on Pattern Anal-  
946 ysis and Machine Intelligence, vol. 29(2), pp. 246-260, 2007, DOI:  
947 10.1109/TPAMI.2007.24
- 948 A. Benoit, A. Caplier, B. Durette and J. Herault, Using Human Visual Sys-  
949 tem modeling for Bio-Inspired Low Level Image Processing, Computer  
950 Vision and Image Understanding, vol. 114, pp. 758-773, 2010.
- 951 J. Berclaz, F. Fleuret and P. Fua, Multiple Object Tracking using Flow Linear  
952 Programming, Proceedings of the Twelfth IEEE International Workshop  
953 on Performance Evaluation of Tracking and Surveillance (PETSWinter),  
954 2009, DOI: 10.1109/PETS-WINTER.2009.5399488
- 955 D. Bolme, Y.M. Lui, B. Draper and J. Beveridge, Simple Real-Time Human  
956 Detection using a Single Correlation Filter, Proceedings Twelfth IEEE In-  
957 ternational Workshop on Performance Evaluation of Tracking and Surveil-  
958 lance (PETSWinter), 2009, DOI:10.1109/PETS-WINTER.2009.5399555
- 959 M. D. Breitenstein, F. Reichlin, B. Leibe, E. Koller-Meier and L. van Gool,  
960 L., Markovian Tracking-by-Detection from a Single, Uncalibrated Camera,  
961 Proceedings of the Eleventh IEEE International Workshop on Performance  
962 Evaluation of Tracking and Surveillance, pp. 7178, 2009.
- 963 F. Briggs and W. M. Usrey, Corticogeniculate Feedback and Visual Pro-

964     cessing in the Primate, *The Journal of Physiology*, vol. 589(1), pp.33-40,  
965     2011.

966     S. Chatterjee and E. M. Callaway, Parallel Colour-Opponent Pathways to  
967     Primary Visual Cortex, *Nature*, vol. 426, pp. 668-671, 2003.

968     T. C. Cheah, Medical Image Registration by Maximizing Mutual Information  
969     Based on Combination of Intensity and Gradient Information, *Proceedings*  
970     International Conference on Biomedical Engineering, pp. 368-372, 2012.

971     S. Cheng, L. Xingzhi and S. M. Bhandarkar, A Multiscale, Parametric Back-  
972     ground Model for Stationary Foreground Object Detection, *IEEE Work-*  
973     shop on Motion and Video Computing, pp. 18, 2007.

974     D. Conte, P. Foggia, G. Percannella and M. Vento, Performance Evaluation of  
975     a People Tracking System on the PETS Video Database, *Proceedings of the*  
976     Thirteenth IEEE International Workshop on Performance Evaluation of  
977     Tracking and Surveillance, pp. 119-126, 2010, DOI: 10.1109/AVSS.2010.87

978     D. Dacey, Parallel Pathways for Spectral Coding in Primate Retina, *An-*  
979     nual Review of Neuroscience, vol. 23, pp.743-775, 2000, DOI: 10.1146/an-  
980     nurev.neuro.23.1.743.

981     R. L. Didday and M. A. Arbib, Eye Movements and Visual Perception: A  
982     Two Visual Stream Model, *International Journal of Man-Machine Studies*,  
983     vol. 7, pp. 499-508, 1975.

984     R. P. W. Duin, The Combining Classifier: To Train or Not to Train?, *Pro-*  
985     ceedings of the Sixteenth International Conference on Pattern Recognition,  
986     pp. 765-770, 2002.

- 987 A. Elgammal, D. Harwood D, L. Davis, Non-Parametric Model for Back-  
 988 ground Subtraction, Proceedings of the Sixth European Conference on  
 989 Computer Vision, Part II, pp. 751767, 2000.
- 990 A. Ellis and J. Ferryman, PETS2010 and PETS2009 Evaluation of Results  
 991 Using Individual Ground Truthed Single Views, Proceedings of the Sev-  
 992 enth IEEE International Conference on Advanced Video and Signal Based  
 993 Surveillance, pp.135-142, 2010, DOI: 10.1109/AVSS.2010.89
- 994 F. Escolano, P. Suau and B. Bonev, Information Theory in Computer Vision  
 995 and Pattern Recognition, Springer, 2009, ISBN: 978-1-84882-296-2, DOI:  
 996 10.1007/978-1-84882-297-9
- 997 R. Farivar, O. Blanke and A. Chaudhuri, DorsalVentral Integration in the  
 998 Recognition of Motion-Defined Unfamiliar Faces, The Journal of Neuro-  
 999 science, vol. 29(16), pp. 5336 5342, 2009.
- 1000 M. E. Farmer, X. Lu, H. Chen and A. K. Jain, Robust Motion-  
 1001 Based Image Segmentation using Fusion, Proceedings of International  
 1002 Conference on Image Processing, vol. 5, pp. 3375-3378, 2004, DOI:  
 1003 10.1109/ICIP.2004.1421838
- 1004 J. Ferryman and A. Ellis, A., PETS2009: Dataset and Challenge, Proceed-  
 1005 ings of the Twelfth IEEE International Workshop on Performance Eval-  
 1006 uation of Tracking and Surveillance, pp. 1-6, 2009, DOI:10.1109/PETS-  
 1007 WINTER.2009.5399556
- 1008 W. Ge and R. Collins, Evaluation of Sampling-Based Pedestrian Detection  
 1009 for Crowd Counting, Proceedings of the Twelfth IEEE International Work-

shop on Performance Evaluation of Tracking and Surveillance (PETSWinter), 2009, DOI: 10.1109/PETS-WINTER.2009.5399553

A. Glazer, M. Lindenbaum and S. Markovitch, One-Class Background Model, Proceedings of the 11th International Conference on Computer Vision, vol. I, pp. 301-307, 2013, DOI: 10.1007/978-3-642-37410-4\_26

C. Guyon, T. Bouwmans and E. Zahzah, Foreground Detection Via Robust Low Rank Matrix Decomposition Including Spatio-Temporal Constraint, Proceedings of the 11th International Conference on Computer Vision, vol. I, pp.315-320, DOI: 10.1007/978-3-642-37410-4\_28

K. Huang, D. Tao, Y. Yuan, X. Li and T. Tan, Biologically Inspired Features for Scene Classification in Video Surveillance, IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics, vol. 41(1), pp. 307-313, 2011, DOI: 10.1109/TSMCB.2009.2037923

S. H. C. Hendry, The Koniocellular Pathway in Primate Vision, Annual Review of Neuroscience, vol. 23(1), pp 127, 2000, DOI: 10.1146/annurev.neuro.23.1.127.

T. Horprasert, D. Harwood, D. and L. S. Davis, A Statistical Approach for Real-Time Robust Background Subtraction and Shadow Detection, Proceedings of the Seventh IEEE ICCV Frame-rate Workshop, pp. 1-19, 1999.

B-G. Hu, What are the Differences between Bayesian Classifiers and Mutual-Information Classifiers?, CoRR, vol. abs/1105.0051v2, 2011.

D. H. Hubel and T. N. Wiesel, Receptive Fields and Functional Architecture

1032 of Monkey Striate Cortex, Journal of Physiology, vol. 196, pp.117-151,  
1033 1985.

1034 P. Jodoin and M. Mignotte, Motion Segmentation Using a K-Nearest-  
1035 Neighbor-Based Fusion Procedure of Spatial and Temporal Label Cues,  
1036 Proceedings of the Second international conference on Image Analysis and  
1037 Recognition, pp.778-788, 2005, DOI:10.1007/11559573-95

1038 P. KaewTraKulPong and R. Bowden, An Improved Adaptive Background  
1039 Mixture Model for Real-Time Tracking with Shadow Detection, Proceed-  
1040 ings of the Second European Workshop on Advanced Video Based Surveil-  
1041 lance Systems, pp. 149-158, 2001.

1042 K. Kasturi, D Goldgof, P. Soundararajan, V. Manohar, J. Garofolo, R. Bow-  
1043 ers, M. Boonstra, V. Korzhova and J. Zhang, Framework for Performance  
1044 Evaluation of Face, Text, and Vehicle Detection and Tracking in Video:  
1045 Data, Metrics, and Protocol, IEEE Transactions on Pattern Analysis and  
1046 Machine Intelligence, vol.31(2), pp. 319 336, 2009.

1047 R. Kentridge, C. Heywood, and J. Davidoff, Color Perception, The Handbook  
1048 of Brain Theory and Neural Networks, Second Edition, Part III: Articles,  
1049 pp. 230 233, 2002, ISBN-10: 0-262-01197-2, ISBN-13:978-0-262-01197-6

1050 D. J. C. MacKay, Information Theory, Inference and Learning Algorithms,  
1051 First Edition, Cambridge Press, 2003, ISBN-10: 0521642981, ISBN-13:  
1052 978-0521642989

1053 V.Y. Mariano, J. Min, J.H. Park, R. Kasturi, D. Mihalcik, D. Doermann  
1054 and T. Drayer, Performance Evaluation of Object Detection Algorithms,

- 1055 Proceedings of International Conference on Pattern Recognition, pp. 965-  
1056 969, 2002.
- 1057 V. Martin, M. Thomnat and N. Mailliot, A Learning Approach for Adap-  
1058 tive Image Segmentation, Proceedings of the Fourth IEEE International  
1059 Conference on Computer Vision Systems, pp. 40 40, 2006.
- 1060 D. J. McKeefry, M. P. Burton and A. B. Morland, The Contribution of Hu-  
1061 man Cortical Area V3A to the Perception of Chromatic Motion: A Tran-  
1062 scranial Magnetic Stimulation Study, European Journal of Neuroscience,  
1063 vol. 31, pp.575584, 2010, DOI:10.1111/j.1460-9568.2010.07095.x
- 1064 C. A. Mead and M. A. Mahhowald, A Silicon Model of Early Visual Process-  
1065 ing, Neural Networks, vol. 1, pp.91-97, 1988.
- 1066 M. Mishkin, L. G. Ungerleider, K. A. Macko, Object Vision and Spatial  
1067 Vision: Two Central Pathways, Trends in Neuroscience, vol. 6, pp 414-  
1068 417, 1983.
- 1069 S. Morand, G. Thut, R. Grave de Peralta, S. Clarke, A. Khateb, T. Landis  
1070 and C.M. Michel, Electrophysiological Evidence for Fast Visual Processing  
1071 through the Human Koniocellular Pathway When Stimuli Move, Cerebral  
1072 Cortex, vol.10(8), pp. 817-825, 2000.
- 1073 S. Mota, E. Ros, J. Díaz, R. Agis and F. de Toro, Bio-inspired  
1074 Motion-Based Object Segmentation, Proceedings of the Third Interna-  
1075 tional Conference on Image Analysis and Recognition, pp196-205, 2006,  
1076 DOI:10.1007/11867586\_19.

- 1077 R. Nieuwenhuys, J. Voogd and C. van Huijzen, The Human Central Nervous  
1078 System, Fourth Edition, Springer, 2008, ISBN: 978-3-540-346864-5
- 1079 J. P. W. Pluim, J. B. A. Maintz and M. A. Viergever, Mutual-Information  
1080 Based Registration of Medical Images: A Survey, In IEEE Transactions  
1081 Medical Imaging, vol. 22(8), pp. 986-1004, DOI:10.1109/TMI.2003.815867,  
1082 2003.
- 1083 T. Serre, L. Wolf, S. Bileschi, M. Riesenhuber and T. Poggio, Robust Object  
1084 Recognition with Cortex-Like Mechanisms, IEEE Transactions on Pattern  
1085 Analysis and Machine Intelligence, vol. 29(3), pp. 411-426, 2007.
- 1086 M. Shah, J. Deng and B. Woodford, Illumination Invariant Background  
1087 Model using Mixture of Gaussians and SURF Features, Proceedings of the  
1088 11th International Conference on Computer Vision , vol. I, pp. 308-314,  
1089 2013, DOI: 10.1007/978-3-642-37410-4\_27
- 1090 A. Shimada, D. Arita and R. Taniguchi, Dynamic Control of Adaptive  
1091 Mixture-of-Gaussians Background Model, IEEE International Conference  
1092 on Video and Signal Based Surveillance, pp. 5, 2006.
- 1093 S. Singh, S. Velastin and R. Hossein, Muhavi: A Multicamera Human Action  
1094 Video Dataset for the Evaluaton of Action Recognition Methods, Proceed-  
1095 ings Seventh IEEE International Conference on Advanced Video and Signal  
1096 Based Surveillance, pp. 48-55, 2010, DOI:10.1109/AVSS.2010.63, 2010
- 1097 C. Stauffer and W. Grimson, Adaptive Background Mixture Models for Real-  
1098 Time Tracking, Proceedings of IEEE Computer Society Conference on  
1099 Computer Vision and Pattern Recognition, vol. 2, pp. 23-25, 1999.



- 1100 C. Studholme, D. L. G. Hill, and D. J. Hawkes, Incorporating connected  
1101 region labelling into automated image registration using mutual informa-  
1102 tion, *Mathematical Methods in Biomedical Image Analysis*, A. A. Amini,  
1103 F. L Bookstein, and D. C. Wilson, Eds. 1996, pp. 2331, IEEE Computer  
1104 Society Press, Los Alamitos, CA.
- 1105 H. R. Tavakoli, E. Rahtu and J. Heikkilä, Temporal Saliency for Fast Motion  
1106 Detection, *Proceedings of the 11th International Conference on Computer*  
1107 *Vision* , vol. I, pp. 321-326, 2013, 10.1007/978-3-642-37410-4\_29
- 1108 C. Thriault, N. Thome, M. Cord, Dynamic Scene Classification: Learning  
1109 Motion Descriptors with Slow Features Analysis, *Proceedings of IEEE*  
1110 *ComputerVision and Pattern Recognition*, 2013, To be issued.
- 1111 C. Tomasi and T. Kanade, Detection and Tracking of Point Features,  
1112 Carnegie Mellon University Technical Report CMU-CS-91-132, 1991.
- 1113 L. C. Ungerleider and M. Mishkin, Two Cortical Visual Systems, *Analysis of*  
1114 *Visual Behavior*, pp. 549-586, Cambridge, MIT Press, 1982.
- 1115 A. Vacavant, T. Chateau, A. Wilhelm, and L. Lequivre, A Benchmark  
1116 Dataset for Outdoor Foreground/Background Extraction, *Proceedings of*  
1117 *the 11th International Conference on Computer Vision*, vol. I, pp. 291-300,  
1118 2012, DOI=10.1007/978-3-642-37410-4\_25
- 1119 C. R. Wren, A. Azarbayejani, T. Darrell and A. Pentland, Pfunder:  
1120 Real-Time Tracking of the Human Body, *International Conference*  
1121 *on Automatic Face and Gesture Recognition*, pp. 51-56, 1997, DOI:  
1122 10.1109/AFGR.1996.557243

- 1123 J. Yang, Z. Shi, P. Vela and J. Teizer, J., Probabilistic Multiple People Track-  
1124 ing through Complex Situations, Proceedings of the Eleventh IEEE Inter-  
1125 national Workshop on Performance Evaluation of Tracking and Surveil-  
1126 lance, pp. 7986, 2009.
- 1127 S. Yoshinaga, A. Shimada, H. Nagahara and R. Taniguchi, Background Model  
1128 Based on Statistical Local Difference Pattern, Proceedings of the 11th  
1129 International Conference on Computer Vision, vol. I, pp. 327-332, 2013,  
1130 DOI: 10.1007/978-3-642-37410-4\_30
- 1131 P. Yuen, A. Tsitiridis, K. Hong, T. Chen, F. Kam, J. Jackman, D. James  
1132 and M. Richardson, A Cortex Like Neuromorphic Target Recognition  
1133 and Tracking in Cluttered Background, Third International Conference  
1134 on Imaging for Crime Detection and Prevention, IET Seminar Digests,  
1135 vol. 2, pp.27, 2009, DOI:10.1049/ic.2009.0255
- 1136 S. K. Zhou, B. Georgescu, D. Comaniciu, S. Jie, BoostMotion: Boosting  
1137 a Discriminative Similarity Function for Motion Estimation, Proceedings  
1138 of IEEE Computer Society Conference on Computer Vision and Pattern  
1139 Recognition, vol. 2, pp. 1761-1768, 2006.
- 1140 M. Zanon, P. Busan, F. Monti, G. Pizzolato and P. Battaglini, Cortical Con-  
1141 nections Between Dorsal and Ventral Visual Streams in Humans: Evidence  
1142 By, Brain Topography, vol. 22(4), pp. 307-317, 2010.