

# Object flow

Anonymous ACCV 2012 submission

Paper ID \*\*\*

**Abstract.** Motion analysis in image sequences has undoubtedly shown good progress in terms of its two main research branches. Optical flow estimation and object visual tracking have been mostly studied as isolated problems, and high accuracy algorithms are available when needed as independent bricks. This paper presents a framework for combining object tracking techniques with optical flow methods aiming towards a precise motion description for objects in video sequences. Firstly, we introduce a method to extend max-flow min-cut based segmentation techniques to videos, without adding the computational load of performing a graph cut optimization approach for 3-dimensional graphs. This is done by exploiting the inherent foreground-background separation hints given by object trackers, and the novel concept of super pixel flow. Then, we show that long-motion awareness obtained from object tracking, together with a per frame object segmentation can improve the precision of the object motion description in comparison to several optical flow techniques. We may call the proposed approach Object flow as it offers a dense and semantic aware description of the current motion state of the studied object.

## 1 Introduction

Object tracking and optical flow are two of the main components in the Computer Vision toolbox, and have been focus of great research efforts, leading to significant progress in the last years [16][17]. The object tracking problem consist on estimating the position of the target in future frames, given an initialization. In the other hand, the optical flow problem consist on finding a motion vector for each pixel of interest in the initial image. Even though for several applications a full motion-field is needed, other applications like human-computer interaction, object editing in video or structure-from-motion, may only focus on an interest object and, thus, only motion vectors within its space may be of interest. In such scenarions combining optical flow and object tracking in a unified framework would become useful and the precision of the object motion description could be enhanced. For instance, even with modern optical flow approaches, the long motion problem remains a challenge. However, the problem is more bearable for object tracking techniques.

Among the state of the art segmentation methods for objects in video sequences, point trajectories based ones stand for its performance and reliability, even when only sparse trajectories are known because of computational reasons.

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

2 ACCV-12 submission ID \*\*\*

In the other hand, for the problem of extracting out a preselected object in still frames, max-flow min-cut based approaches have demonstrated to be a powerful tool. We propose to mix these two ideas for reliable object segmentation through video. Moreover, in this paper we introduce the superpixel flow method for tracking of background regions, rather than sparse points. We show how this extra information can be used to expand the max-flow min-cut techniques for an efficient foreground-background segmentation. Experiments on several video sequences demonstrate the characteristics of the proposed method.

Superpixels and over segmentation techniques became a widely used pre-processing stage for a large number of machine vision applications, after the original concept was introduced [1]. Superpixels are traditionally used as performance booster for several other techniques. However, it is still mostly related to single frame processing [1][10][11]. In the search for consistency in superpixel labeling through video, some authors have proposed different techniques, which go from simple extension to supervoxels[9][11], to more complicated approaches [8]. These approaches, nonetheless, usually require a global processing and knowledge of all (or several of) the video frames beforehand. One of the contributions of this work is a superpixel matching technique which assumes a flowlike behavior in the image sequences (natural video), and propose an application for improving object segmentation in videos.

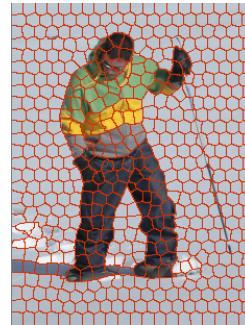
For some kind of video sequences and computer vision applications the superpixel labelling usually loses coherency between frames, even when computed as supervoxels. This means superpixel matching techniques offer some interest and can still be explored. Some work have been done towards a superpixel based image comparison using the Earth Mover's Distance, by taking super pixels as bins of a global histogram [2]. The label propagation or superpixel flow can be achieved with this technique by selecting the superpixel in the second frame that maximize the flow from each superpixel in the first frame. However, we look at the problem of superpixel propagation between a pair of frames, in terms of the indexes of the already computed superpixelization in both of them.

By taking into account superpixels computed separately in images we can open the advantages of streaming like approaches, so the video process can be performed by processing only two frames at a time, saving memory and moving towards a more time friendly approach. This propagation means finding the superpixel labels in the second frame that correspond to a given label in the first frame. This matching, however, has to comply with a set of constraints. Firstly, two correspondent superpixels should be similar in terms of some appearance feature, which most likely depends on the way the superpixelization was performed (color, texture, shape). Also, the superpixel propagation should maintain certain global homogeneity (at least for superpixels that belong to the same object). This is because the superpixel propagation or matching is not completely a one-to-one combinatorial problem, but more a superpixel flow. In this sense, it seems natural that the problem of superpixel propagation could be solved with a discrete energy minimization procedure. If the size compactness of the superpixels is maintained, it actually seems to share some of the properties

090 of the optical flow problem, with the difference that the smoothness is usually  
 091 a very strong constraint for the last one. The strength of this smoothness prior  
 092 relies not only in the nature of the problem, but also because it gives better cues  
 093 towards an easy-to- minimize global approach.

## 095 2 PROBLEM DEFINITION

097 The objective of the superpixel propagation is therefore to find the best labeling  
 098  $L$  for every superpixel  $p$  (with  $L_p \in 0, 1, \dots, N - 1$ ) between a pair of frames  
 099 ( $I_1, I_2$ ), but holding a flow-like behavior. Thus, the superpixelization should  
 100



113 **Fig. 1.** SLIC super segmentation in the Snow Shoes sequence.

115 maintain certain regularity in size within a single frame. Some super pixel tech-  
 116 niques can cope with this requirement [9][10]. For the experiments presented in  
 117 this work, we prefer the SLIC (Fig. 1) method [9], which usually gives good  
 118 results in homogeneity of the superpixelization. The proposed steps to solve the  
 119 propagation problem assume this requirement is hold. For other kind of the  
 120 techniques, other approaches should be followed.

### 123 2.1 Energy Formulation

124 Inspired by a large number of optical flow and stereo techniques [7][12][13], the  
 125 super-pixel propagation can be modeled with pairwise Markov Random Fields.  
 126 If the matching is performed with MAP inference, its posterior probability is:  
 127

$$128 P(L|I_0, I_1) = \prod_{p \in \Omega} e^{-D_p(L_p; I_0, I_1)} \prod_{p, q \in \mathcal{N}_p} e^{-S_{p,q}(L_p; L_q)} \quad (1)$$

131 With  $L$  the set of labels of the super pixels in  $I_0$ , that match with those in  $I_1$ .  
 132  $\mathcal{N}_p$  is a neighborhood of the superpixel  $p$ , which defines its adjacency. Given this  
 133 posterior probability, the equivalent energy function can be directly obtained by  
 134 extracting the negative logarithm of the posterior,

135

$$E(L) = \sum_{p \in \Omega} -D_p(L_p; I_0, I_1) + \sum_{p,q \in \mathcal{N}_p} -S_{p,q}(L_p, L_q) \quad (2)$$

138

The terms  $D$ , and  $S$  in 2 stand for data and spatial as they are popularly known in the MRF literature. The first one determines how accurate is the labeling in terms of consistency of the measured data (Color, Shape, etc.). In the usual optical flow counterpart of this equation, the data term corresponds to the pixel brightness [7][5]. However, as superpixels are a set of similar (or somehow homogenous) pixels, an adequate color based feature can be a low binned histogram or its average color. So it can be written more precisely as

145

$$D_p(L_p; I_0, I_1) = \rho(h(p), h(p')) \quad (3)$$

146

Where  $h(p)$  and  $h(p')$  are the histograms of the super-pixel  $p$  and its correspondent superpixel in the second frame ( $I_1$ ).  $\rho$  can be replaced by the Bhattacharyya distance. Note that the low binned histogram or average color gives certain robustness against noise, and slowly changing colors between frames. The spatial term is a penalty function for horizontal and vertical changes of the vectors that have origin in the centroid of the super-pixel of the first frame and end in the centroid of the super-pixel of the second frame.

155

$$S_{p,q}(L_p, L_q) = \lambda(p) \sqrt{\frac{|u_{p_c} - u_{q_c}|}{\|p_c - q_c\|} + \frac{|v_{p_c} - v_{q_c}|}{\|p_c - q_c\|}} \quad (4)$$

156

$$\text{where, } \lambda(p) = (\rho(h(p), h(p')))^2$$

157

In 4 the operator  $\rho$  has the same meaning as in the data term 3. The histograms distance is nonetheless computed between super pixels  $p$  and  $q$ , which belong to the same neighborhood. The super pixels centroids are noted as  $q_c$  and  $p_c$ , and  $u$  and  $v$  are the horizontal and vertical changes between centroids. This term is usual in the MRF formulation and has a smoothing effect in superpixels that belong to the same object. It has to be observed that when two superpixels are different, thus, more probable to be in different semantic groups within the image, the term  $\lambda$  allows them to have matches that do not hold the smoothness prior with the same strength. It has to be noted that the proposed energy function is highly non-convex and robust.

158

159

## 2.2 Energy Minimization

160

A fair amount of work had been dedicated to discrete optimization techniques in computer vision, leading to a couple of well-defined and widely tested approaches to solve the pairwise MRF of labels [3][4]. However, some of the approaches restrict the construction of the spatial term, and/or enforce limitations in the number of labels [3]. Because of the high amount of possible labels for each element in the proposed approach, the use of the Fusion Moves [7] technique seems

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180 to be well suited. This algorithm employs the Quadratic Pseudo- Boolean Opti-  
 181 mization (QPBO) graph-cut, to combine incremental sets of proposal labeling,  
 182 resulting a semi- globally-optimal solution [4]. Thus, the minimization starts by  
 183 proposing a set of possible solutions, and iteratively merge them with the QPBO  
 184 graph-cut technique.

185 The possible solutions that can be given depend on the kind of problem that is  
 186 intended to be solved. For example, in stereo super-pixel matching, some assump-  
 187 tions related to the cameras organization can be made to generate solutions. In a  
 188 more generic sense, other assumptions can be made towards option generation.  
 189 For instance, for a given super-pixel in the initial frame, in the second frame  
 190 the corresponding matching would be the most similar one in terms of color,  
 191 or the most similar un terms of shape, or the spatially closer super-pixel. More  
 192 proposal solutions can be added by defining a neighborhood in the second frame  
 193 and select random pairs from every neighborhood of every super-pixel in the first  
 194 frame. This is more suitable for problems where the images are extracted from  
 195 the same video sequence. It is interesting to notice that different assumptions  
 196 for this neighborhood can lead to a technique for generic image based retrieval,  
 197 where the total cost of matching can be used as metric.

198

### 199 3 EXPERIMENTAL RESULTS

200

201 The Fig. 2 shows some examples of superpixel matching with the presented  
 202 method. It can be seen that the matching performs well even in difficult cases,  
 203 like the hands in the top row. It has to be noted as well that even in superpixels  
 204 where there is a lack of texture, there is correct matching. This seems to be the  
 205 effect of enforcing the regularization between superpixels that are close, but are  
 206 also similar to each other.

207 Moreover, unlike most of the optical flow methods, superpixel flow extends  
 208 naturally for more distant frames. The Fig. 3 shows results for larger separations  
 209 between frames, without tweaking or adjusting any parameters. For this case,  
 210 however, the matches in the textureless part of the scene are mostly invalids.  
 211 Though this is expected because because of the aperture problem and heavy  
 212 occlusions.

213

#### 214 3.1 Superpixel propagation for object segmentation in videos

215 The GrabCut algorithm proposed in [14], offers a good deal in terms of background-  
 216 foreground separation from user interaction. A technique like this, however, per-  
 217 forms very well in still images, but it may not be well adapted for sequential  
 218 videos. The GrabCut works by implementing an iterative graph cut based min-  
 219 imization to separate regions according to appearance information, that can be  
 220 extracted from the user interaction. This interaction, however, could be mini-  
 221 mized in videos, given the extra information that offers the flow of the sequence.  
 222 Some authors had approached the GrabCut or similar graph based segmen-  
 223 tation techniques in sequential videos, to propagate a consistent segmentation [15].

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

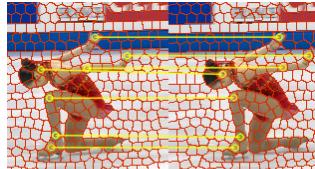
242

243

244

245

**Fig. 2.** The yellow lines show selected super-pixel matching between pairs of subsequent frames in a video with the proposed method. The video frames go from right to left.



225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

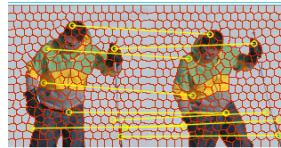
265

266

267

268

269



**Fig. 3.** The yellow lines show selected superpixel matching between a pair of distant frames in the Snow Shoes sequence.

270 However, some more work on reducing user interaction given the extra flow-like  
271 information that video sequences offer is still needed. We propose to combine  
272 the presented super pixel propagation as an automatic method to initialize the  
273 Grabcut (or similar) algorithm to perform object segmentation through frames  
274 in a video sequence.

275

276 The idea is to track (or more exactly, match) super pixels that are labeled  
277 as foreground, thanks to an object tracker initialization. Thus, the super pixels  
278 that are initially outside the ROI, can be propagated through the sequence,  
279 and if they fall into the ROI of the next frame, they can be safely labeled as  
280 foreground again. The process is repeated for any labeled superpixel through the  
281 video. Having several labeled superpixels can reduce widely the necessity for user  
282 interaction in subsequent frames. Thus, to perform object segmentation in a full  
283 video sequence, the required user interaction would only be the initial bounding  
284 box. Moreover, a fully automatic approach can be obtained if a reliable object  
285 detector is available.

286 Fig. 4 shows how the initialization of a tracker and a superpixelization pro-  
287 vides information for better background-foreground separation. The background-  
288 foreground models are updated as the frames go on, giving more robustness for  
289 sequential propagation of the segmentation. The method is tested in the Walk-  
290 ing Couple sequence, by allowing only a small amount of iterations in the graph  
291 based segmentation. Observe how the contour in the mans head is correctly delin-  
292 eated when another persons head occludes part of it. In this case, the super-pixels  
293 that belong to the womans face were correctly propagated and thus, labeled as  
294 background.

295

296 In order to understand the effect of including super- pixel propagation in a  
297 video sequence for object segmentation, some results are shown in the Fig. 5.  
298 For these experiments only one iteration is allowed in both grab-cuts initialized  
299 with only the tracker, and the one performed with the super-pixel propagation.  
300 Observe that in general, the contour delineated is usually better in terms of  
301 precision and stability for the later one.

302

## 303 References

304

## 305 References

306

- 307 1. J. Malik and X. Ren, Learning a classification model for segmentation, *Computer  
308 Vision, International Conference*, 2003.
- 309 2. S. Boltz; F. Nielsen and S. Soatto, Earth mover distance on superpix-  
310 els *International Conference on Image Processing*, 2010.
- 311 3. E. Boros; P. Hammer and G. Tavares, Preprocessing of unconstrained quadratic  
312 binary optimization *RUTCOR*, 2010
- 313 4. E. Boros and P. Hammer, Pseudo-boolean optimization *Discrete applied Mathematics*,  
314 2002
- 315 5. B. Horn and B. Schunck, Determining Optical Flow *Artificial Intelligence*, 1981



**Fig. 4.** Segmentation through the sequence Walking Couple (Yellow contour) initialized in the mans head.

354

352

355

353

356

354

357

355

358

356

359

357

360

358

361

359



**Fig. 5.** Face segmentation in the Amelie Retro and the Snow shoes sequences in three different frames. For each group, the Top Row: One-iteration window-based grabcut; and the Bottom Row: One-iteration grabcut with super pixel propagation.

- 360
- 361
- 362
- 363
- 364
- 365
- 366
- 367
- 368
- 369
- 370
- 371
- 372
- 373
- 374
- 375
- 376
- 377
- 378
- 379
- 380
- 381
- 382
- 383
- 384
- 385
- 386
- 387
- 388
- 389
- 390
- 391
- 392
- 393
- 394
- 395
- 396
- 397
- 398
- 399
- 400
- 401
- 402
- 403
- 404
- 360
- 361
- 362
- 363
- 364
- 365
- 366
- 367
- 368
- 369
- 370
- 371
- 372
- 373
- 374
- 375
- 376
- 377
- 378
- 379
- 380
- 381
- 382
- 383
- 384
- 385
- 386
- 387
- 388
- 389
- 390
- 391
- 392
- 393
- 394
- 395
- 396
- 397
- 398
- 399
- 400
- 401
- 402
- 403
- 404
- 6. H. Ishikawa and P. Bouthemy, Multimodal estimation of discontinuous optical flow using Markov random fields *TPAMI*, 1993
- 7. V. Lempitsky, S. Roth and C. Rother, Fusion Flow: Discrete-Continuos optimization for optical flow estimation *Computer Vision and Pattern Recognition*, 2008
- 8. M. Reso and J. Jachalsky, Temporally Consistent Superpixels *International Conference Computer Vision.*, 2011
- 9. R. Achanta; A. Shaji; K. Smith; Aurelien Lucchi; P. Fua and S. Susstrunk SLIC Superpixels compared to state of the art superpixel methods *Discrete applied Mathematics.*, 2002
- 10. F. Perbet and A. Maki, Homogeneous superpixels from random walks *MVA*, 2011
- 11. C. Xu and J.J. Corso. Evaluation of super-voxel methods for early video proccesing. *Computer Vision and Pattern Recognition*. 2012.
- 12. A. Shekhovtsov, I. Kovtun and V. Hlavac. Efficient MRF deformation model for non-rigid image matching. *Computer Vision and Pattern Recognition*. 2007.
- 13. J. Sun, N.N Shen and H.Y. Shum. Stereo matching using propagation belief. *TPAMI*. 2003.
- 14. C. Rother, V. Kolmogorov and A. Blake. Grabcut: Interactive foreground extraction using iterated graph cuts. *SIGGRAPH*. 2004.
- 15. L. Yang, Y. Guo, X. Wu and X. Wang. A new video segmentation approach: Grabcut in local window. *Soft Computing and Pattern Recognition*. 2011.
- 16. Y. Wu, J. Lim and M.-H. Yang. Online object tracking: A benchmark. *Computer Vision and Pattern Recognition*. 2013.

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

10 ACCV-12 submission ID \*\*\*

- 405 17. S. Baker, D. Scharstein, J.P. Lewis, S. Roth, M.J. Black and R. Szeliski. A Database  
406 and Evaluation Methodology for Optical Flow *International Journal Computer*  
407 *Vision*. 2013.

408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449

405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449