

Object flow

Anonymous ACCV 2014 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 superpixel 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 between a pair of frames 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 scenarios 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 term motion problem remains a challenge. However, the problem is more bearable for object tracking techniques. In contrast, object trackers are more global in motion description, and its information can be completed by optical flow sub-pixel precision. Moreover, even when object trackers and optical flow could give good hints for object segmentation in video, these

2 ACCV-14 submission ID ***

elements are not deeply studied in the literature as a unified problem. We introduce the object flow problem as the computation of dense motion flow fields of the set of pixels that belong to an interest object. In other words, the object flow by definition induces the segmentation of the target and its motion field.

We can define more precisely the object flow by starting with an image sequence and an initial position of the interest object in the first frame of this sequence, and letting \mathcal{R} be the region corresponding to the support of the object in 2D, such that $\mathcal{R} \subset \Omega$. If Ω is the set of all the possible grid positions, the object flow problem consist in finding the displacement vector $d_{0,t}(x)$ from the image I_0 to I_t , $\forall x \in \mathcal{R}$.

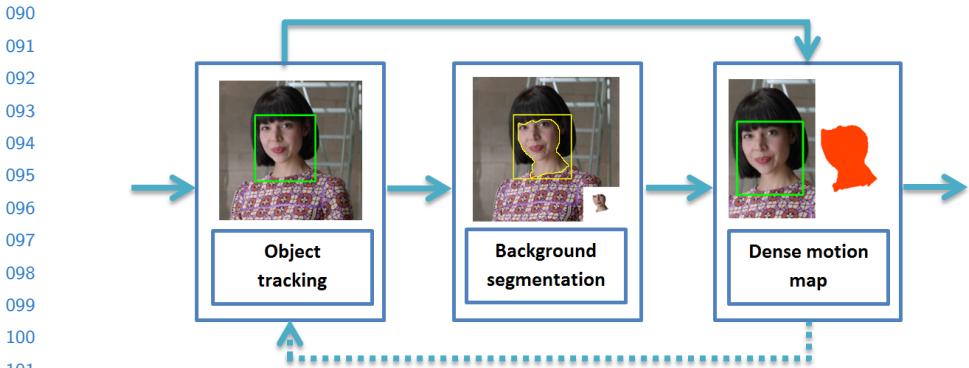
A straightforward solution to this problem would be to compute the optical flow motion field, and apply a segmentation mask to recover the desired motion vectors. Nevertheless, this approach carries several problems. For example, a globally computed optical flow method can affect small objects motion, because of the common use of heavy regularization priors. Moreover, even if the segmentation mask is extracted from a tracker position by a graph-cut based method, is likely that this mask is not going to be well suited for the interest object and some extra user interaction would be needed to refine this process. We propose an approach to reduce these problems.

The present paper is organized as follows. We describe our pipeline for object flow, including the novel concept of superpixel flow in Sec. 2 and its use in object segmentation in videos. In following sections some results showing how the object flow overpass state of the art optical flow methods for object motion flows estimation are discussed. Finally, some insights and conclusions are given.

2 Algorithm description

The Fig. 1 shows a simplified block diagram of the proposed system. Two details are important, the use of the tracker window to initialize a segmentation procedure, and the use of this segmentation over the tracked window to perform a more precise motion flow computation in the interest pixels. The dotted line represents the possible interaction between precise flow information with the next tracker state. For instance, the current object flow can work as direction hint, and the segmentation information can be used to improve the sampling process of the learning stage in several trackers by detection methods [22], and thus the tracker and motion flow algorithm can work for mutual enhancement.

The first step in the object flow pipeline can be selected according to specific need for a given application. We prefer, in general, tracking-by-detection methods like *Struck* [22] or *MIL* [23], but other approaches could be followed. In the second place, for the object segmentation in video we propose the use of labelled background regions through the concept of superpixel flow, which is explained in the next section.

**Fig. 1.** Block diagram of the proposed pipeline.

2.1 Superpixel flow

As a preprocessing step in the object flow pipeline, we propose a superpixel matching technique which assumes a flowlike behavior in the image sequences (natural video), which can be used to track superpixels. 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 flow should maintain certain global regularity (at least for superpixels that belong to the same object). If the size compactness of the superpixels is maintained, it actually seems to share some of the properties of the optical flow problem, with the difference that the smoothness is usually a very strong constraint for the last one. The strength of this smoothness prior relies not only in the nature of the problem, but also because it gives better cues towards an easier-to-minimize global approach.

The objective of the superpixel flow is therefore to find the best labeling l for every superpixel p (with $l_p \in 0, 1, \dots, N - 1$) between a pair of frames (I_0, I_1) , but holding a flow-like behavior.

Thus, the superpixelization should maintain certain size homogeneity within a single frame. Some super pixel techniques can cope with this requirement [9][10]. For the experiments presented in this work, we prefer the SLIC method [9], which usually gives good results in terms of homogeneity of the superpixelization across the sequence.

Inspired by a large number of optical flow and stereo techniques [7][12][13], the superpixel flow can be modeled with pairwise Markov Random Fields. If the matching is performed with MAP inference, its energy function extracted from the posterior probability is:

$$E(l) = \sum_{p \in \Omega} D_p(L_p; I_0, I_1) + \sum_{p, q \in \mathcal{N}} S_{p,q}(L_p, L_q) \quad (1)$$

135 With l the set of labels of the super pixels in I_0 , that match with those in I_1 .
 136 \mathcal{N} is a neighborhood of the superpixel p , which defines its adjacency. Given this
 137 posterior probability, the equivalent energy function can be directly obtained by
 138 extracting the negative logarithm of the posterior,

139 The terms D , and S in (1) stand for data term and spatial smoothness terms
 140 as they are popularly known in the MRF literature. The first one determines
 141 how accurate is the labeling in terms of consistency of the measured data (color,
 142 shape,etc.). In the classical optical flow formulation of this equation, the data
 143 term corresponds to the pixel brightness conservation[7][5]. However, as super-
 144 pixels are a set of similar (or somehow homogenous) pixels, an adequate color
 145 based feature can be a low dimensional color histogram. So D can be written
 146 more precisely as the Hellinger distance between the histograms:

$$147 \quad D_p(l_p; I_0, I_1) = \sqrt{1 - \frac{1}{\sqrt{h(p)h(p')}N^2} \sum_i \sqrt{h_i(p)h_i(p')}} \quad (2)$$

148 Where $h(p)$ and $h(p')$ are the histograms of the superpixel p and its cor-
 149 respondent superpixel in the second frame I_1 . Note that the low dimensional
 150 histogram gives certain robustness against noise, and slowly changing colors be-
 151 tween frames.

152 In the other hand, the spatial term is a penalty function for horizontal and
 153 vertical changes of the vectors that have origin in the centroid of the superpixel
 154 of the first frame and end in the centroid of the superpixel of the second frame.

$$155 \quad S_{p,q}(l_p, l_q) = \lambda(p) \sqrt{\frac{|u_{p_c} - u_{q_c}|}{\|p_c - q_c\|} + \frac{|v_{p_c} - v_{q_c}|}{\|p_c - q_c\|}} \quad (3)$$

156 where, $\lambda(p) = (1 + \rho(h(p), h(q)))^2$

157 In (3) the operator ρ is the Hellinger distance as used in the data term (2).
 158 The histogram distance is nonetheless computed between superpixels p and q ,
 159 which belong to the same neighborhood. The superpixels centroids are noted as
 160 q_c and p_c , and u and v are the horizontal and vertical changes between centroids.
 161 This term is usual in the MRF formulation and has a smoothing effect in su-
 162 perpixels that belong to the same object. It has to be observed that when two
 163 close superpixels are different, thus, more probable to belong to different objects
 164 within the image, the term λ allows them to have matches that do not hold the
 165 smoothness prior with the same strength. It has to be noted that the proposed
 166 energy function is highly non-convex.

167 The Quadratic Pseudo-Boolean Optimization (QPBO) [3][4] is used to min-
 168 imize the proposed energy function, by merging a set of candidate matches for
 169 every superpixel in the first frame. For instance, for a given superpixel in the
 170 initial frame, the corresponding matching would be the most similar one in terms
 171 of color, shape, or the spatial distance. More candidate solutions can be added
 172 by defining a neighborhood in the second frame and select random pairs from
 173 every neighborhood of every superpixel in the first frame.

180

181

182

183

184

185

186

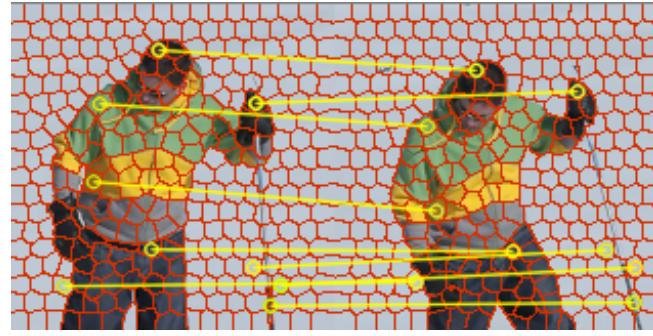
187

188

189

190

Fig. 2. The yellow lines show selected superpixel matching between a pair of distant frames in the Snow Shoes sequence.



191

192

193

194

The Fig. 2 shows results for large separations between frames. For this case, however, the matches in the textureless part of the scene are mostly invalids. Though this is expected because because of the aperture problem and heavy occlusions.

195

196

197

198

199

200

2.2 Background regions tracking for object segmentation

201

202

The main idea to perform object segmentation consist in tracking (or more exactly, match) superpixels that are labeled as background, thanks to an object tracker initialization.

Thus, the superpixels that are initially outside the tracker region of interest, can be propagated through the sequence, and if they fall into the window on a subsequent frame, they can be safely labeled as background

(Fig. 3).

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

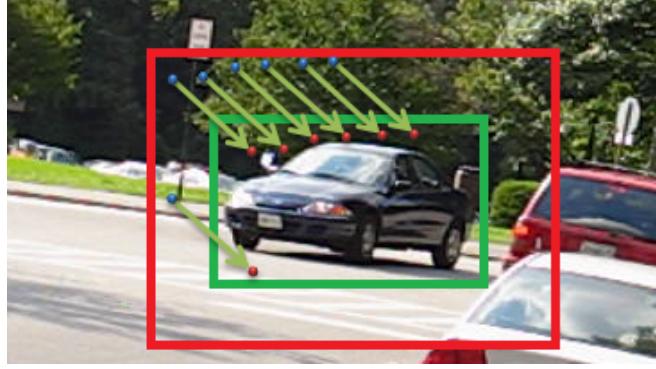


Fig. 3. Example image of points entering a tracking region (green) due to object motion in a video sequence.

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

To save computational power, the tracked superpixels are limited to the ones that fall inside a control region (red box in the Fig. 3). Usually, after several frames, the labeled superpixels will almost completely cover the unwanted areas in a dinamyc scene. We call this process background segments tracking. The Fig. 4 shows this idea in a real scenario. From left to right, initially the superpixels with elements outside the bounding box are labeled as background (green), then, as the sequence changes, the labeled superpixels flow inside the window, giving hints for the model initialization in the background-foreground separation algorithm. At this point, some generic segmentation technique can be connected to the pipeline to refine the segmentation (e.g. region growing). We prefer, however graph based segmentation methods ([18][15]) because the usual user interaction can be replaced by the tracked background regions.

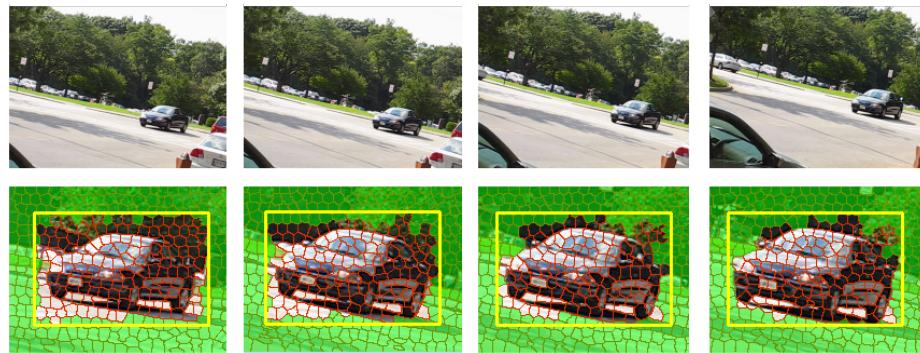


Fig. 4. Background segments automatic labeling and propagation, the flow goes from left to right.

2.3 Segmentation results

Fig. 5 shows the results for an image sequence where the interest object is the head of a person. The head tracker and the superpixel flow provide information for better background-foreground separation. The background-foreground models are updated as the frames go on, giving more robustness for sequential propagation of the segmentation. The method is tested in the Walking Couple sequence, by allowing only a small amount of iterations in the graph based segmentation. Observe how the contour in the man's head is correctly delineated when another person's head occludes part of it. In this case, the superpixels that belong to the womans face were correctly propagated and thus, labeled as background.

In order to understand the effect of including superpixel propagation in a video sequence for object segmentation, some results are shown in the Fig. 6. For these experiments only one iteration is allowed in the graph-cut based methods.

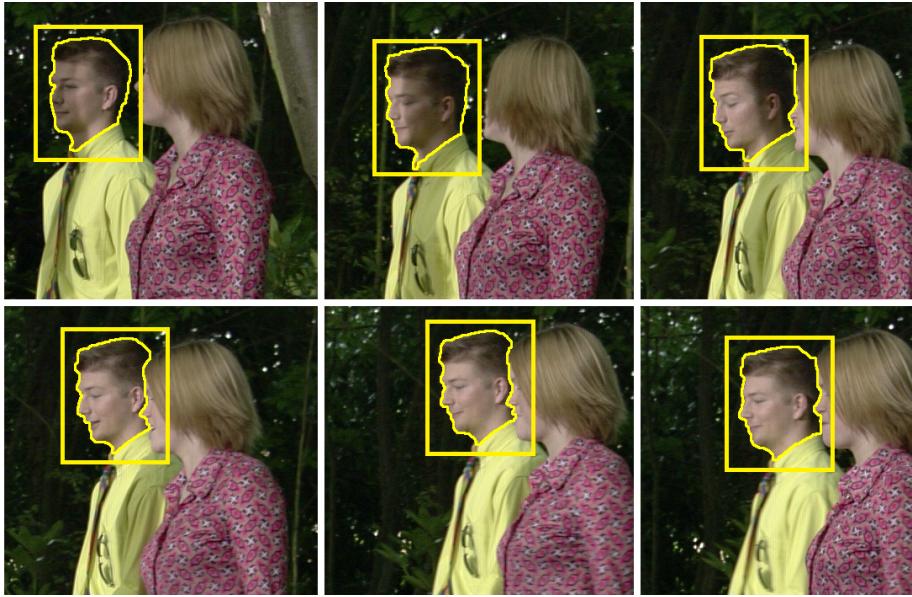


Fig. 5. Segmentation through the sequence Walking Couple (Yellow contour) initialized in the mans head. The yellow box correspond to the tracker output. The labeled background superpixel are not shown for clarity.

The top row frames (Fig. 6) were initialized only with the tracker, and the bottom row was initialized with the superpixel tracking technique. Observe that in general, the contour delineated is usually better in terms of precision and stability for the later one.

2.4 Flow estimation

The object flow consist on computing the motion field for an object of interest through an image sequence. The most usual approach to solve a problem like this is to implement some of the available optical flow techniques through the complete sequence and perform the flow integration. However, this process results in high levels of motion drift [18][19] and usually the motion of the interest object is affected by a global regularization. In some extreme cases, the interest object motion may be totally blurred and other techniques have to be incorporated. Moreover, the diversity of natural video sequences makes difficult the choice of one technique over another, even when specialized databases are at hand [17], because currently no single method can achieve a strong performance in every of the available datasets. Most of these methods consist in the minimization of an energy function with two terms (As was previously mentioned in the Sec. 2.1). The data term is mostly shared between different approaches, but the prior or spatial term is different, and basically states under what conditions the optical flow smoothness should be maintained or not. In a global approach,

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

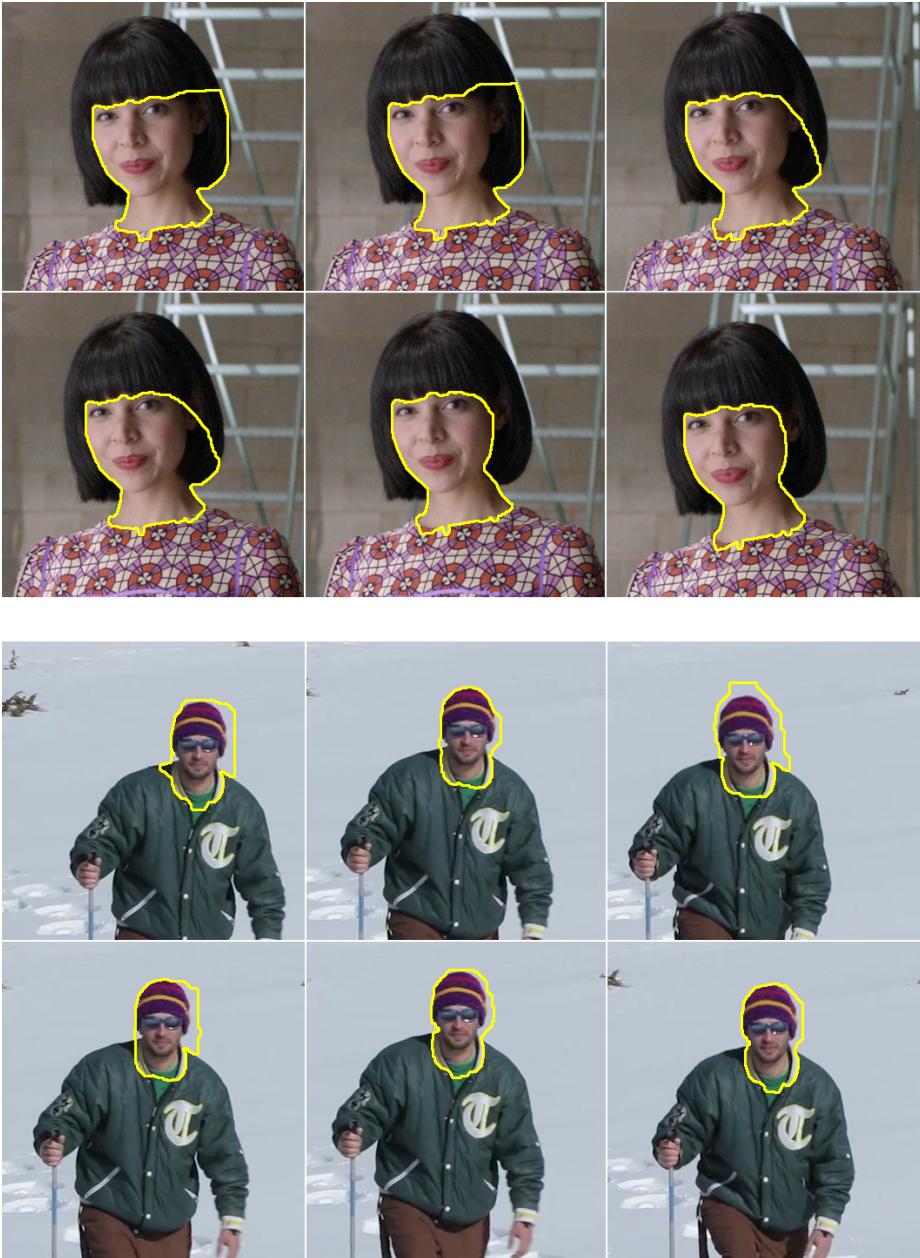


Fig. 6. Face segmentation in the Amelie Retro and the Snow shoes sequences in three different frames. For each group, the Top Row: One-iteration window-based graph-cuts; and the Bottom Row: One-iteration graph-cuts initialized with superpixel tracking.

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360 however, this is a difficult concept to define. Most of these smoothness terms
 361 rely in appearance differences or gradients. All these meaning that, unavoidably,
 362 some methods may be more reliable for some cases but weaker for others. It can
 363 be argued that this behaviour may be caused because most of the techniques do
 364 not count with a way to identify firmly where exactly this smoothness prior can
 365 be applied. The main idea behind the object flow is that given the availability
 366
 367



368
 369 **Fig. 7.** Object flow with the color code of [17] (bottom)
 370 sequence (up).
 371
 372

373
 374 of several robust tracking techniques, and the proposed segmentation method
 375 for video, the optical flow computation can be refined by computing it succes-
 376 sively between pairs of tracked windows. The basic proposal to perform this
 377 refinement consist on considering the segmentation limits as reliable smoothness
 378 boundaries. This is, of course, under the assumption that the motion is indeed
 379 smooth within the object region. This assumption is not far from reality in
 380 most scenes with an interest object. Naturally, as the object tracker is included,
 381 is expected that the object flow should be more robust to rapid motions than the
 382 optical flow. Thus, the full motion is split in two, the long range motion, given
 383 by the tracker window, and the precision part, given by the targeted optical flow.
 384 The Fig. 7 shows the object flow for a frame in the Puppy sequence. Observe
 385 the motion vectors are computed only inside the object of interest, preserving a
 386 strong smoothing prior, but also allowing internal variations in the flow.
 387
 388

389 As a first approximation to the object flow, the Simple Flow technique [21]
 390 is taken as core base. This is because of its scalability to higher resolutions and
 391 because its specialization to the concept of object flow is only natural. This is
 392 because in the Simple Flow pipeline the smoothness localization can be easily
 393 specified through computation masks. More specifically, the initial computa-
 394 tion mask is derived from the segmentation performed as prior step. The resulting
 395 flow is then filtered only inside the mask limits to enhance precision and fas-
 396 tening the implementation. However, direct modifications in other optical flow
 397 methods can be further studied. For instance, in graph-cut based minimization
 398 approaches, the regularity constraints can be precisely targeted by disconnecting
 399 foreground pixels from background ones.
 400
 401
 402
 403
 404

405 **3 Experimental results**

406

407

408

409

410

411 **REAL OBJECT**412 **OBJ FLOW BW**413 **OPT FLOW BW**414 **OBJ FLOW FW**415 **OPT FLOW FW**

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

Fig. 8. Extrapolation results from integrated flow in 4 sequences. In descending order: Amelia Retro, Boy, Walking, Puppy. From Left to Right: Annotated object, Backward object flow, Backward optical flow, Forward object flow, Backward optical flow.

432

433

434

435

To evaluate the performance of the object flow in comparison with optical flow techniques, we performed a number of experiments on several video sequences. We annotated an initial bounding box for the videos, and a segmentation contour of the interest object for every frame. The experiment measures the ability of the method to extrapolate an image from the initial frame and the integrated flow. For every pair of frames the PSNR between the annotated current state of the object and the extrapolated images is computed. The Fig. 8 is a sample of the performed experiment, each column is an image generated from the given flow.

436

437

438

439

440

441

442

443

444

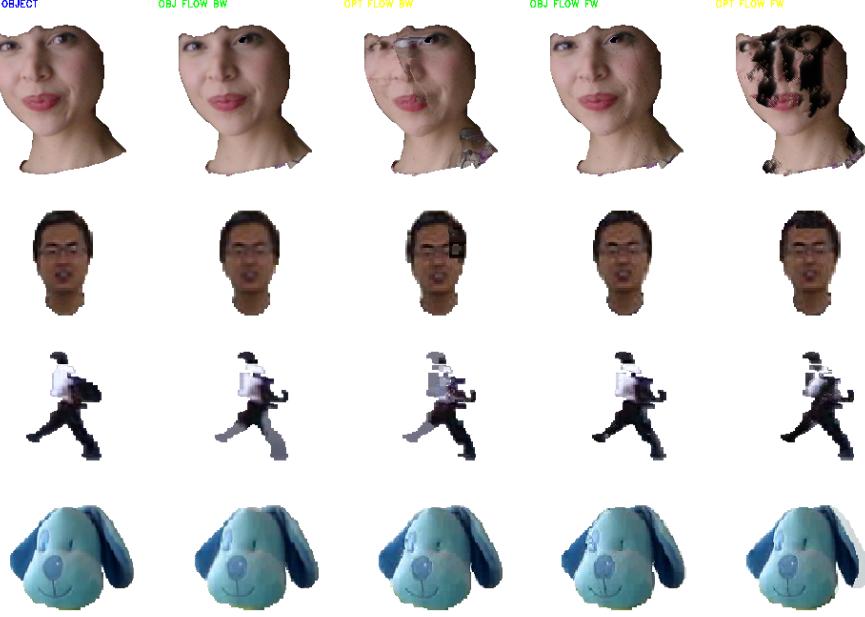
445

446

447

448

449



The Fig. 9 shows PSNR graphics for 4 different sequences. For every pair of frames an image is extrapolated, and the PSNR is computed. The measure is computed using Euler integration (Labelled as *forward* in the figs.) of the used flow (object or optical flows), and using the integration method described in [20], labeled as *backward* in the figures.

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

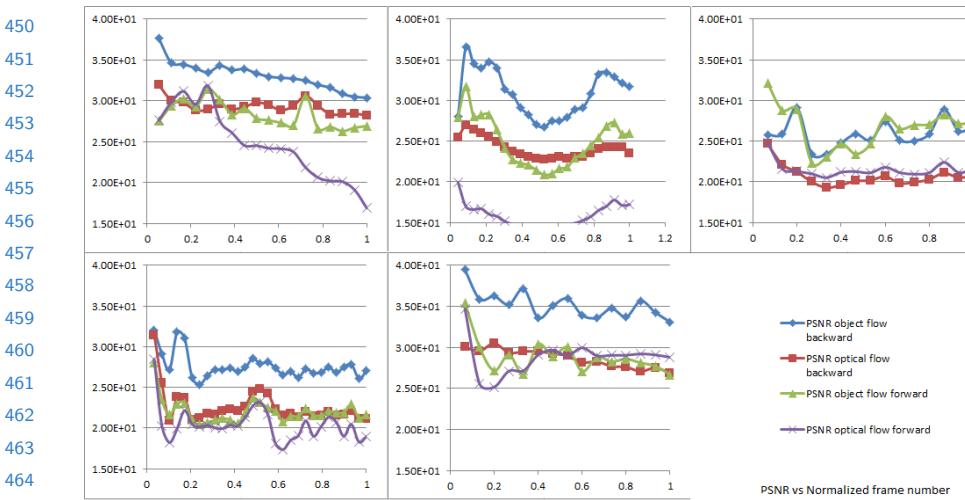


Fig. 9. PSNR graphs for extrapolated images using Object flow and the Simple Optical Flow for 4 sequences. In descendent order: Puppy Seq.; Amelie Retro Seq.; Boy Seq.; Walking Seq.

4 Conclusions

A framework to combine tracking and optical flow methods to improve object based dense motion description is presented. The pipeline is composed of three main steps, object tracking, segmentation and flow estimation. For the segmentation step a new video object segmentation algorithm was proposed and some results are shown. The last step, flow estimation is a modification of the simple-flow method to use the obtained segmentation mask. The experiments showed that the object flow improves the dense motion estimation in comparison to optical flow techniques. Future work includes the use of different optical flow techniques as base of the object flow. Also, applications of the object flow in the structure-from-motion pipeline, or for automatic video edition can be more deeply studied.

References

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

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

12 ACCV-14 submission ID ***

- 495 6. H. Ishikawa and P. Bouthemy, Multimodal estimation of discontinuous optical flow
496 using Markov random fields *TPAMI*, 1993
497 7. V. Lempitsky, S. Roth and C. Rother, Fusion Flow: Discrete-Continuos optimization
498 for optical flow estimation *Computer Vision and Pattern Recognition*, 2008
499 8. M. Reso and J. Jachalsky, Temporally Consistent Superpixels *International Conference Computer Vision.*, 2011
500 9. R. Achanta; A. Shaji; K. Smith; Aurelien Lucchi; P. Fua and S. Susstrunk SLIC
501 Superpixels compared to state of the art superpixel methods *Discrete applied Mathematics.*, 2002
502 10. F. Perbet and A. Maki, Homogeneous superpixels from random walks *MVA*, 2011
503 11. C. Xu and J.J. Corso. Evaluation of super-voxel methods for early video proccesing.
504 *Computer Vision and Pattern Recognition*. 2012.
505 12. A. Shekhovtsov, I. Kovtun and V. Hlavac. Efficient MRF deformation model for
506 non-rigid image matching. *Computer Vision and Pattern Recognition*. 2007.
507 13. J. Sun, N.N Shen and H.Y. Shum. Stereo matching using propagation belief.
508 *TPAMI*. 2003.
509 14. C. Rother, V. Kolmogorov and A. Blake. Grabcut: Interactive foreground extraction
510 using iterated graph cuts. *SIGGRAPH*. 2004.
511 15. L. Yang, Y. Guo, X. Wu and X. Wang. A new video segmentation approach:
512 Grabcut in local window. *Soft Computing and Pattern Recognition*. 2011.
513 16. Y. Wu, J. Lim and M.-H. Yang. Online object tracking: A benchmark. *Computer
514 Vision and Pattern Recognition*. 2013.
515 17. S. Baker, D. Scharstein, J.P. Lewis, S. Roth, M.J. Black and R. Szeliski. A Database
516 and Evaluation Methodology for Optical Flow *International Journal Computer
517 Vision*. 2013.
518 18. Y. Boykov, M-P. Jolly. Interactive Graph Cuts for Optimal Boundary & Region
519 Segmentation of Objects in N-D images *International Conference of Computer
520 Vision*. 2013.
521 19. W. Li, D. Cosker and M. Brown. An anchor patch based optimization framework for
522 reducing optical flow drift in long image sequences. *Asian Conference of Computer
523 Vision*. 2012.
524 20. T. Crivelli, P.-H. Conze, P. Robert, M. Fradet and P. Perez. Multi-step flow fusion:
525 towards accurate and dense correpondence in long video shots. *British Conference
526 Machine Vision*. 2012.
527 21. M. Tao, J. Bai, P. Kohli, and S. Paris. SimpleFlow: A Non-iterative, Sublinear
528 Optical Flow Algorithm. *Computer Graphics Forum, Eurographics*. 2012.
529 22. S. Hare, A. Saffari, P.H.S. Torr. Struck: Structured output tracking with kernels.
530 *International Conference of Computer Vision* . 2011.
531 23. B. Babenko, M.-H. Yang and S. Belongie. Robust Object Tracking with Online
532 Multiple Instance Learning. *Pattern analysis and Machine Learning* . 2010.
533
534
535
536
537
538
539