

# 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 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 motion problem remains a challenge. However, the problem is more bearable for object tracking techniques. Moreover, even when object trackers and optical flow give good hints for object segmentation in video, these 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

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

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

045 of the set of pixels that belong to an interest object. In other words, the object  
 046 flow by definition copes with segmentation of the target.

047 Among the state of the art segmentation methods for objects in video se-  
 048 quences, point trajectories based ones stand for its performance and reliability,  
 049 even when only sparse trajectories are known because of computational reasons.  
 050 In the other hand, for the problem of extracting out a preselected object in still  
 051 frames, max-flow min-cut based approaches have demonstrated to be a powerful  
 052 tool. We propose to mix these two ideas together with the tracking of backgorund  
 053 regions via the novel concept of superpixel flow for reliable object segmentation  
 054 through video. We show how this extra information can be used to complement  
 055 the graph-cuts based techniques for an efficient foregorund-background segmen-  
 056 tation.

057 The present paper is organized as follows. We introduce the concept of super-  
 058 pixel flow in Sec. 2. Then, a method for object segmentation in video which uses  
 059 object tracker information and the background segments tracking is presented  
 060 in Sec. 3. In following sections the object flow basic approach is explained. Fi-  
 061 nally, results showing how the object flow overpass state of the art optical flow  
 062 methods for object motion flows computation.

## 063 2 Superpixel flow

### 064 2.1 Problem definition

065 Superpixels and over segmentation techniques became a widely used pre-processing  
 066 stage for a large number of machine vision applications, after the original concept  
 067 was introduced [1]. Superpixels are traditionally used as performance booster for  
 068 several other techniques. However, it is still mostly related to single frame pro-  
 069 cessing [1][10][11]. In the search for consistency in superpixel labeling through  
 070 video, some authors have proposed different techniques, which go from simple  
 071 extension to supervoxels[9][11], to more complicated approaches [8]. These ap-  
 072 proaches, nonetheless, usually require a global processing and knowledge of all  
 073 (or several of) the video frames beforehand. We propose a superpixel matching  
 074 technique which assumes a flowlike behavior in the image sequences (natural  
 075 video). Some previous work have been done towards a superpixel based image  
 076 comparison using the Earth Mover's Distance, by taking superpixels as bins of  
 077 a global histogram [2]. The label propagation or superpixel flow can be achieved  
 078 with this technique as a byproduct, by selecting the superpixel in the second  
 079 frame that maximize the flow from each superpixel in the first frame.

080 By taking into account superpixels computed separately in images, so the  
 081 video process can be performed with only two frames at a time, we move to-  
 082 wards a more time efficient approach. This matching, however, has to comply  
 083 with a set of constraints. Firstly, two correspondent superpixels should be simi-  
 084 lar in terms of some appearance feature, which most likely depends on the way  
 085 the superpixelization was performed (color, texture, shape). Also, the superpixel  
 086 flow should maintain certain global regularity (at least for superpixels that be-  
 087 long to the same object). In this sense, it seems natural that the problem of

superpixel flow 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 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 homegenity 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 homegenity of the superpixelization across the sequence. The proposed steps to solve the propagation problem assume this requirement is hold. For other kind of the techniques, other approaches should be followed.

## 2.2 Energy Formulation

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 posterior probability is:

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

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

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

The terms  $D$ , and  $S$  in 2 stand for data term and spatial smoothness terms 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 classical optical flow formulation of this equation, the data term corresponds to the pixel brightness conservation[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

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

Where  $h(p)$  and  $h(p')$  are the histograms of the superpixel  $p$  and its correspondent superpixel in the second frame ( $I_1$ ). The distance  $\rho$  can be replaced

135 by the Bhattacharyya distance. Note that the low binned histogram or average  
 136 color gives certain robustness against noise, and slowly changing colors between  
 137 frames. The spatial term is a penalty function for horizontal and vertical changes  
 138 of the vectors that have origin in the centroid of the superpixel of the first frame  
 139 and end in the centroid of the super-pixel of the second frame.

$$140 \quad S_{p,q}(l_p, l_q) = \lambda(p) \sqrt{\frac{|u_{p_c} - u_{q_c}|}{\|p_c - q_c\|} + \frac{|v_{p_c} - v_{q_c}|}{\|p_c - q_c\|}} \quad (4)$$

$$143 \quad \text{where, } \lambda(p) = (1 + \rho(h(p), h(p')))^2$$

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

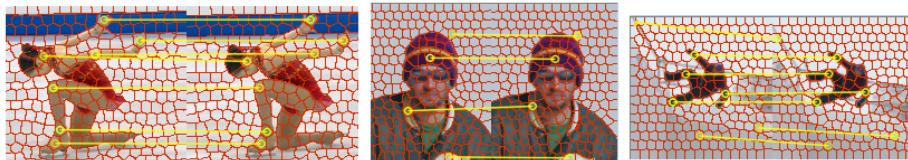
### 156 2.3 Energy Minimization

158 A fair amount of work had been dedicated to discrete optimization techniques in  
 159 computer vision, leading to a couple of well-defined and widely tested approaches  
 160 to solve the pairwise MRF[3][4]. However, some of the approaches restrict the  
 161 construction of the spatial term, and/or enforce limitations in the number of  
 162 labels [3]. Because of the high amount of possible labels for each element in  
 163 the proposed approach, the use of the Fusion Moves [7] technique seems to be  
 164 well suited. This algorithm employs the Quadratic Pseudo-Boolean Optimiza-  
 165 tion (QPBO), to combine incremental sets of proposal labelings, resulting a  
 166 semi-globally-optimal solution [4]. Thus, the minimization starts by proposing a  
 167 set of possible solutions, and iteratively merge them with the QPBO technique.  
 168 The possible solutions that can be given depend on the kind of problem that is  
 169 intended to be solved. For example, in stereo superpixel matching, some assump-  
 170 tions related to the cameras organization can be made to generate solutions. In  
 171 a more generic sense, other assumptions can be made towards option generation.  
 172 For instance, for a given superpixel in the initial frame, in the second frame the  
 173 corresponding matching would be the most similar one in terms of color, or the  
 174 most similar in terms of shape, or the spatially closer superpixel. More proposal  
 175 solutions can be added by defining a neighborhood in the second frame and  
 176 select random pairs from every neighborhood of every super-pixel in the first  
 177 frame. This is more suitable for problems where the images are extracted from  
 178 the same video sequence. To speed-up the minimization procedure, the QBPO  
 179 properties can be exploited. For instance, the fusion of the proposed solutions is

180 always guaranteed of lowest or equal energy than the two proposals. Thus, one  
 181 could split the fusion procedure in several cores and build a hierachichal chain  
 182 as fusions of proposal are subsequently fused.

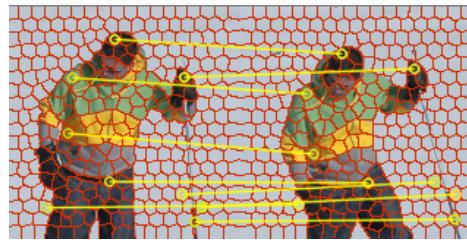
183  
 184 **2.4 Matching results**  
 185

186 The Fig. 1 shows some examples of superpixel matching with the presented  
 187 method. It can be seen that the matching performs well even in difficult cases,  
 188 like the hands in the top row. It has to be noted as well that even in superpixels  
 189 where there is a lack of texture, there is correct matching. This seems to be the  
 190 effect of enforcing the regularization between superpixels that are close, but are  
 191 also similar to each other.  
 192



200  
 201 **Fig. 1.** The yellow lines show selected superpixel matching between pairs of consecutive  
 202 frames in a video with the proposed method. The video frames go from right to left.  
 203

204  
 205 Moreover, unlike most of the optical flow methods, superpixel flow extends  
 206 naturally for more distant frames. The Fig. 2 shows results for larger separations  
 207 between frames, without tweaking or adjusting any parameters. For this case,  
 208 however, the matches in the textureless part of the scene are mostly invalids.  
 209 Though this is expected because because of the aperture problem and heavy  
 210 occlusions.  
 211



212  
 213 **Fig. 2.** The yellow lines show selected superpixel matching between a pair of distant  
 214 frames in the Snow Shoes sequence.  
 215

225        **3 Superpixel propagation for object segmentation in**

226        **videos**

227

228        The algorithm proposed in [18], offers a good deal in terms of background-  
 229        foreground separation from user interaction. A technique like this, however, per-  
 230        forms very well in still images, but it may not be well adapted for sequential  
 231        videos. Extensions to this method, like, GrabCut ([14]) work by implement-  
 232        ing an iterative graph cut based minimization to separate regions according to  
 233        appearance information, that can be extracted from the user interaction. This  
 234        interaction, however, could be minimized in videos, given the extra information  
 235        that offers the flow of the sequence. Some authors had approached the GrabCut  
 236        or similar graph based segmentation techniques in sequential videos, to propa-  
 237        gate a consistent segmentation [15]. However, some more work on reducing user  
 238        interaction given the extra flow-like information that video sequences offer is still  
 239        needed. We propose to combine the presented superpixel flow as an automatic  
 240        method to initialize the desired min-cut max flow based algorithm to perform  
 241        object segmentation through frames in a video sequence.

242        The main idea consist in tracking (or more exactly, match) superpixels that  
 243        are labeled as background, thanks to an object tracker initialization. Thus, the  
 244        superpixels that are initially outside the ROI, can be propagated through the  
 245        sequence, and if they fall into the ROI of the next frame, they can be safely  
 246        labeled as background again. We call this process background segments tarcking.  
 247        The process is repeated for any labeled superpixel through the video. Having  
 248        several labeled superpixels can reduce widely or totally the necessity for user  
 249        interaction in subsequent frames. Thus, to perform object segmentation in a full  
 250        video sequence, the required user interaction would only be the initial bounding  
 251        box. Moreover, a fully automatic approach can be obtained if a reliable object  
 252        detector is available.

253        Fig. 3 shows the results for an image sequence where the interest object is  
 254        the head of a person. The head tracker and the superpixel flow provide informa-  
 255        tion for better background-foreground separation. The background-foreground  
 256        models are updated as the frames go on, giving more robustness for sequential  
 257        propagation of the segmentation. The method is tested in the Walking Couple  
 258        sequence, by allowing only a small amount of iterations in the graph based seg-  
 259        mentation. Observe how the contour in the man's head is correctly delineated  
 260        when another person's head occludes part of it. In this case, the superpixels  
 261        that belong to the womans face were correctly propagated and thus, labeled as  
 262        background.

263        In order to understand the effect of including superpixel propagation in a  
 264        video sequence for object segmentation, some results are shown in the Fig. 4. For  
 265        these experiments only one iteration is allowed in two grab-cut based methods.  
 266        One initialized only with the tracker, and the other complemented with the  
 267        superpixel propagation. Observe that in general, the contour delineated is usually  
 268        better in terms of precision and stability for the later one.

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

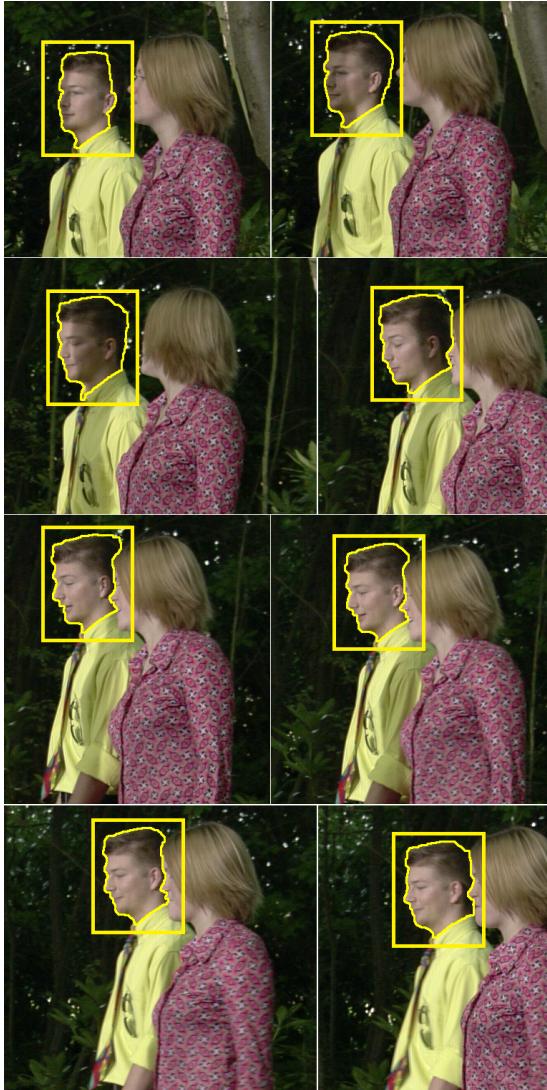
305

306

307

308

**Fig. 3.** Segmentation through the sequence Walking Couple (Yellow contour) initialized in the mans head. The yellow box correspond to the tracker output.



270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309  
310  
311  
312  
313  
314



**Fig. 4.** 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.

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

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      **4 Object flow**

361

362      The object flow consist on computing the motion field for an object of interest  
 363      through an image sequence. The most usual approach to solve a problem like this  
 364      is to implement some of the available optical flow techniques through the com-  
 365      plete sequence and perform the flow integration. However, this process results in  
 366      high levels of motion drift [18][19] and usually the motion of the interest object  
 367      is affected by a global regularization. In some extreme cases, the interest object  
 368      motion may be totally blurred and other techniques have to be incorporated.  
 369      Moreover, the diversity of natural video sequences makes difficult the choice of  
 370      one technique over another, even when specialized databases are at hand [17],  
 371      because currently no single method can achieve a strong performance in every  
 372      of the available datasets. Most of these methods consist in the minimization of  
 373      an energy function with two terms (As was previously mentioned in the Sec. 2).  
 374      The data term is mostly shared between different approaches, but the prior or  
 375      spatial term is different, and basically states under what conditions the optical  
 376      flow smoothness should be maintained or not. In a global approach, however, this  
 377      is a difficult concept to define. Most of these smoothness terms rely in appearance  
 378      differences or gradients. All these meaning that, unavoidably, some methods may  
 379      be more reliable for some cases but weaker for others. It can be argued that this  
 380      behaviour may be caused because most of the techniques do not count with a  
 381      way to identify firmly where exactly this smoothness prior can be applied. The

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

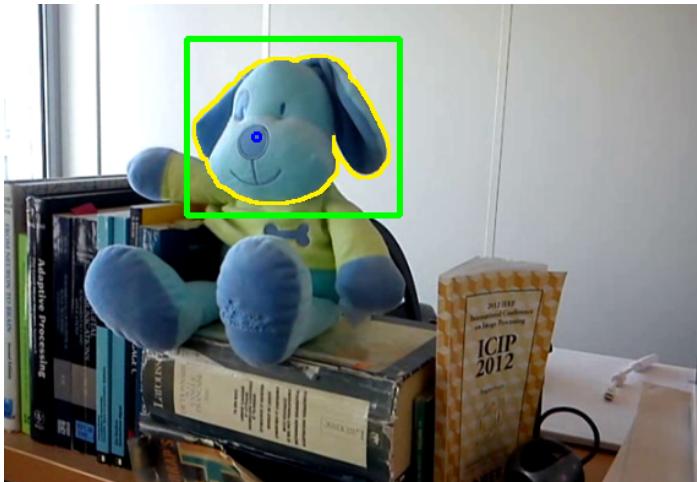
400

401

402

403

404



**Fig. 5.** Object flow with the color code of [17] (Right) for one frame in the Puppy sequence (Left). Green Box: Current tracker state. Yellow: Segmentation contour.

main idea behid the object flow is that given the availability of several robust tracking techniques, and the proposed segmentation method for video, the opti-

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

405 cal flow computation can be refined by computing it successively between pairs  
 406 of tracked windows. The basic proposal to perform this refinement consist on  
 407 considering the segmentation limits as reliable smoothness boundaries. This is,  
 408 of course, under the assumption that the motion is indeed smooth within the  
 409 object region. This is assumption is not far from reality in most scenes with an  
 410 interest object. Of course, as the object tracker is included, is expected that the  
 411 object flow should be more robust to rapid motion than the optical flow. Thus,  
 412 the full motion is split in two, the long range motion, given by the tracker win-  
 413 dows, and the precision part, given by the targeted optical flow. The Fig. 5 shows  
 414 the object flow for a frame in the Puppy sequence. Observe the motion vectors  
 415 are computed only inside the object of interest, preserving a strong smoothing  
 416 prior, but also allowing internal variations in the flow.  
 417

#### 418 4.1 Implementation details and results

419  
 420 As a first approximation to the object flow, the Simple Flow technique [21] is  
 421 taken as core base. This is because to its scalability to higher resolutions and  
 422 because its specialization to the concept of object flow is only natural. This  
 423 is because in the Simple Flow pipeline the smoothness localization can easily  
 424 be specified through computation masks. More specifically, the initial computation  
 425 mask is derived from the segmentation performed as prior step. The resulting  
 426 flow is then filtered only inside mask limits to enhance precision and fastening the  
 427 implementation. However, direct modifications in other optical flow methods can  
 428 be further studied. For instance, in graph-cut based minimization approaches,  
 429 the regularity constraints can be precisely targeted by disconnecting foreground  
 430 pixels from background ones.  
 431



432  
 433 **Fig. 6.** Extrapolation results from integrated flow in one frame of the Amelie Retro  
 434 sequence. From Left to Right: Annotated object, Backward object flow, Backward optical  
 435 flow, Forward object flow, Backward optical flow.  
 436  
 437  
 438  
 439

440 To evaluate the performance of the object flow in comparison with opti-  
 441 cal flow techniques, we performed a number of experiments on several video  
 442 sequences. We annotated an initial bounding box for the videos, and a segmen-  
 443 tation contour of the interest object for every frame. The experiment measures  
 444 the ability of the method to extrapolate an image from the initial frame and  
 445 the integrated flow. For every pair of frames the PSNR between the annotated  
 446

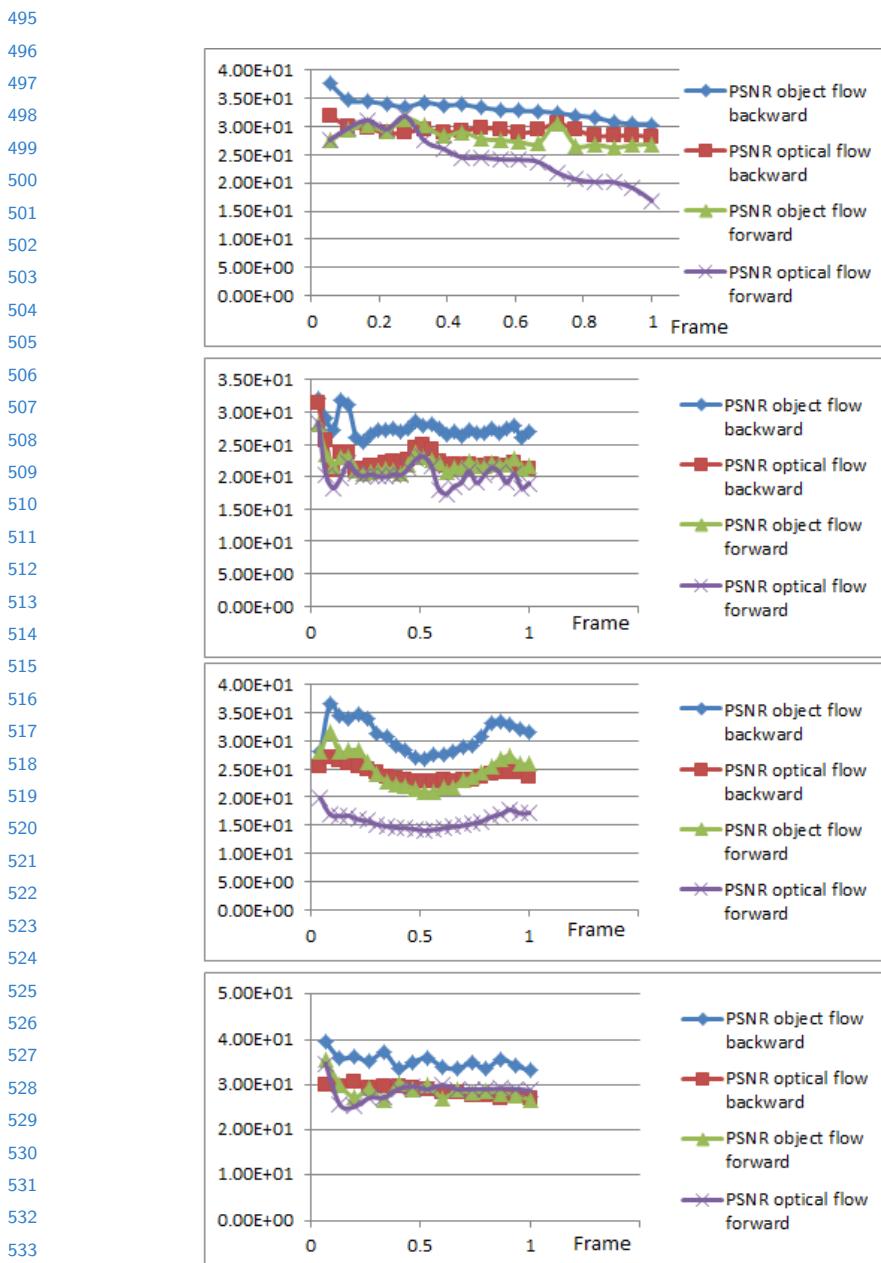
450 current state of the object and the extrapolated images is computed. The Fig.  
451 6 is a sample of the performed experiment, each column is an image generated  
452 from the given flow.

453 The Fig. 7 shows PSNR graphics for 4 different sequences. For every pair of  
454 frames an image is extrapolated, and the PSNR is computed. The measure is  
455 computed using Euler integration (Labeled as *forward* in the figs.) of the used  
456 flow (object or optical flows), and using the integration method described in [20],  
457 labeled as *backward* in the figures.

## 459 References

460

- 461 1. J. Malik and X. Ren, Learning a classification model for segmentation, *Computer*  
462 *Vision, International Conference*, 2003.
- 463 2. S. Boltz; F. Nielsen and S. Soatto, Earth mover distance on superpixels *International Conference on Image Processing*, 2010.
- 464 3. E. Boros; P. Hammer and G. Tavares, Preprocessing of unconstrained quadratic  
465 binary optimization *RUTCOR*, 2010
- 466 4. E. Boros and P. Hammer, Pseudo-boolean optimization *Discrete applied Mathematics*, 2002
- 467 5. B. Horn and B. Schunck, Determining Optical Flow *Artificial Intelligence*, 1981
- 468 6. H. Ishikawa and P. Bouthemy, Multimodal estimation of discontinuous optical flow  
469 using Markov random fields *TPAMI*, 1993
- 470 7. V. Lempitsky, S. Roth and C. Rother, Fusion Flow: Discrete-Continuos optimiza-  
471 tion for optical flow estimation *Computer Vision and Pattern Recognition*, 2008
- 472 8. M. Reso and J. Jachalsky, Temporally Consistent Superpixels *International Con-  
473 ference Computer Vision.*, 2011
- 474 9. R. Achanta; A. Shaji; K. Smith; Aurelien Lucchi; P. Fua and S. Susstrunk SLIC  
475 Superpixels compared to state of the art superpixel methods *Discrete applied Math-  
476 ematics.*, 2002
- 477 10. F. Perbet and A. Maki, Homogeneus superpixels from random walks *MVA*, 2011
- 478 11. C. Xu and J.J. Corso. Evaluation of super-voxel methods for early video proccesing.  
479 *Computer Vision and Pattern Recognition*. 2012.
- 480 12. A. Shekhovtsov, I. Kovtun and V. Hlavac. Efficient MRF deformation model for  
481 non-rigid image matching. *Computer Vision and Pattern Recognition*. 2007.
- 482 13. J. Sun, N.N Shen and H.Y. Shum. Stereo matching using propagation belief.  
483 *TPAMI*. 2003.
- 484 14. C. Rother, V. Kolmogorov and A. Blake. Grabcut: Interactive foreground extrac-  
485 tion using iterated graph cuts. *SIGGRAPH*. 2004.
- 486 15. L. Yang, Y. Guo, X. Wu and X. Wang. A new video segmentation approach:  
487 Grabcut in local window. *Soft Computing and Pattern Recognition*. 2011.
- 488 16. Y. Wu, J. Lim and M.-H. Yang. Online object tracking: A benchmark. *Computer*  
489 *Vision and Pattern Recognition*. 2013.
- 490 17. S. Baker, D. Scharstein, J.P. Lewis, S. Roth, M.J. Black and R. Szeliski. A Database  
491 and Evaluation Methodology for Optical Flow *International Journal Computer*  
492 *Vision*. 2013.
- 493 18. Y. Boykov, M-P. Jolly. Interactive Graph Cuts for Optimal Boundary & Region  
494 Segmentation of Objects in N-D images *International Conference on Computer*  
*Vision*. 2013.



**Fig. 7.** 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.

- 540 19. W. Li, D. Cosker and M. Brown. An anchor patch based optimization framework for  
541 reducing optical flow drift in long image sequences. *Asian Conference on Computer  
542 Vision*. 2012.
- 543 20. T. Crivelli, P.-H. Conze, P. Robert, M. Fradet and P. Perez. Multi-step flow fusion:  
544 towards accurate and dense correspondence in long video shots. *British Conference  
545 Machine Vision*. 2012.
- 546 21. M. Tao, J. Bai, P. Kohli, and S. Paris. SimpleFlow: A Non-iterative, Sublinear  
547 Optical Flow Algorithm. *Computer Graphics Forum, Eurographics*. 2012.

548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584