

Supervoxel-based segmentation of 3D imagery with optical flow integration for spatiotemporal processing

Xiaohui Huang, Chengliang Yang, Sanjay Ranka and Anand Rangarajan

Dept. of CISE, University of Florida, Gainesville, FL, USA

Abstract

The past twenty years has seen a progressive evolution of computer vision algorithms for unsupervised 2D image segmentation. While earlier efforts relied on Markov random fields and efficient optimization (graph cuts etc.), the next wave of methods beginning in the early part of this century were, in the main, stovepiped. Of these 2D segmentation efforts, one of the most popular and indeed, one that comes close to being a state of the art method, is the ultrametric contour map (UCM). The pipelined methodology consists of (i) computing local, oriented responses, (ii) graph creation, (iii) eigenvector computation (globalization), (iv) integration of local and global information, (v) contour extraction and (vi) superpixel hierarchy construction. UCM performs well on a range of 2D tasks. Consequently, it is somewhat surprising that no 3D version of UCM exists at the present time. To address that lack, we present a novel 3D supervoxel segmentation method, dubbed 3D UCM, which closely follows its 2D counterpart while adding 3D relevant features. The methodology, driven by supervoxel extraction, combines local and global gradient-based features together to first produce a low level supervoxel graph. Subsequently, an agglomerative approach is used to group supervoxel structures into a segmentation hierarchy with explicitly imposed containment of lower level supervoxels in higher level supervoxels. Comparisons are conducted against state of the art 3D segmentation algorithms. The considered applications are 3D spatial and 2D spatiotemporal segmentation scenarios. For the latter comparisons, we present results of 3D UCM with and without optical flow video pre-processing. As expected, when motion correction beyond a certain range is required, we demonstrate that 3D UCM in conjunction with optical flow is a very useful addition to the pantheon of video segmentation methods.

1 Introduction

The ready availability of 3D datasets and video has opened up the need for more 3D computational tools. Further, the big data era with ever improving computational resources allows us to envisage complex 3D processing methodologies with the anticipation of increased practicality. The focus in the present work is on 3D segmentation. There is a clear and pressing need to extract coherent, volumetric and spatiotemporal structures from high dimensional datasets—specifically lightfields, RGBD videos and regular video data, particle-laden turbulent flow data comprising dust storms and snow avalanches and finally 3D medical images like MRIs. It is natural to represent coherent 3D structures in the form of locally homogeneous segments or irregularly shaped regions with the expectation that such structures correspond to objects of interest.

A particular focus of the present work is the parsing of volumetric and spatiotemporal datasets into supervoxels. While there exists previous work on this topic, the supervoxel literature is considerably sparser than the corresponding superpixel literature. Supervoxels perform a similar function to superpixels: the codification of locally coherent, homogeneous regions. Superpixels and supervoxels have a conceptual connection to clustering since the latter is a popular way to extract coherent structures from a set of patterns (using some imposed metric). However, the principal difference is that supervoxels can benefit from the gridded structure of the data (which is not the case for general patterns). Consequently, successful supervoxel estimation techniques and the grouping algorithms at their center can be expected to leverage the “dense” nature of gridded 3D patterns.

The popular ultrametric contour map (UCM) framework [1] has established itself as the state of the art in 2D superpixel segmentation. However, a direct extension of this framework has not been seen in the case of 3D supervoxels. The current popular frameworks in 3D segmentation start by constructing regions based on some kind of agglomerative clustering (graph based or otherwise) of single pixel data. However, the 2D counterparts of these methods on which they are founded are not as accurate as UCM. This is because the UCM algorithm does not begin by grouping pixels into regions. Instead, it combines local feature information and global grouping in a two step process which is responsible for its success. Therefore, it follows that a natural extension of this two step process in the case of 3D should be a harbinger for success.

UCM derives its power from a combination of local and global cues which complement each other in order to detect highly accurate boundaries. However, all the other methods, with the exception of normalized cuts [2] (which directly obtains regions) are inherently local and do not incorporate global image information. In sharp contrast, UCM includes global image information by estimating eigenfunction scalar fields of graph Laplacians formed from local image features. Subsequently, the two pieces of information (local and global) are integrated to produce a good segmentation. Despite this inherent advantage, the high computational cost of UCM (mainly in its globalization step) is a bottleneck in the development of its 3D analog. Recent work in [3] has addressed this issue in 2D by providing an efficient GPU-based implementation but the huge number of voxels involved in the case of spatiotemporal volumes remains an issue. The work in [4] provides an efficient CPU-only implementation by leveraging the structure of the underlying problem and providing a reduced order normalized cuts approach to address the original computationally inefficient eigenvector problem. Consequently, this opens up a plausible route to 3D as a reduced order approach effectively solving a closely related globalization problem but at a fraction of the cost. It should be noted that we do not expect the current rapid improvements in hardware and distributed processing to significantly improve the tractability of the original eigenvector problem in 3D, hence our foray into an approximate, reduced order, normalized cuts approach.

The main work of this paper has a tripartite structure. First, we design filters for 3D volumes (either space-time video or volumetric data) which provide local cues at multiple scales and at different orientations—akin to the approach in 2D UCM. Second, we introduce a new method for graph affinity matrix construction—henceforth called the *oriented intervening contour cue*—which extends the idea of the intervening contour cue in [5]. Third, we solve a reduced order eigenvector problem by leveraging ideas from the approach in [4]. After this globalization step, the local and global fields are merged to obtain a new 3D scalar field. Subsequent application of a watershed transform on this 3D field yields supervoxel tessellations (represented as relatively uniform polyhedra that tessellate the region). The next step in this paper is to build a hierarchy of supervoxels. While 2D UCM merges regions based on their boundary strengths using an oriented watershed transform approach, we did not find this approach to work well in 3D. Instead, to obtain the supervoxel hierarchy, we follow the approach of [6] which performs a graph-based merging of regions using the method of internal variation [7].

The contributions of this work are summarized below:

- We present the first 3D UCM method for the segmentation problem. Our approach extends the popular and highly accurate *gPb*-UCM framework (a technical term which will be elaborated below) to 3D resulting in a highly scalable framework based on reduced order normalized cuts. To the best of our knowledge, there does not exist a surface detection method in 3D which uses a graph Laplacian-based globalization approach for estimating supervoxels.
- It is a fully unsupervised method and it achieved state of the art performance on supervoxel benchmarks, exhibiting lesser fragmentation, better region quality and better compactness. Further, our results indicate that owing to the deployment of a high quality boundary detector as the initial step, our method maintains a distinction between the foreground and background regions by not causing unnecessary oversegmentation of the background at the lower levels of the hierarchy. This is especially clear in the video segmentation results presented in the paper.
- We expect applications of 3D UCM in a wide range of vision tasks, including video semantic understanding, video object tracking and labeling in high-resolution medical imaging, etc.

The organization of the paper is as follows. Section 2 describes related work on 2D and 3D segmentation. Since this work is almost impossible to adequately summarize, we paint a broad brush while focusing on

work that had a direct impact on the present work. Section 3 describes the extension of the 2D gPb -UCM framework to 3D while clearly pointing out the compromises and design choices. This is followed in Section 4 by a brief foray into video segmentation and integration with optical flow. Section 5 evaluates the proposed approach and other state of the art supervoxel methods on two different types of 3D volumetric datasets both qualitatively and quantitatively. Section 6 concludes by summarizing our contributions and also discusses the scope for future work. Throughout the paper, we use the term pixel and voxel interchangeably when referring to the basic “atom” of 2D/3D images.

2 Related Work

We first address previous work based on the normalized cuts framework. Subsequently, we attempt to summarize previous work on volumetric image segmentation with the caveat that this literature is almost impossible to summarize.

Normalized Cuts and gPb -UCM

Normalized cuts [2] gained immense popularity in 2D image segmentation by treating the image as a graph and computing hard partitions. The gPb -UCM framework [1] leveraged this approach in a soft manner to introduce globalization into their contour detection process. This led to a drastic reduction in the over-segmentation arising out of gradual texture or brightness changes in the previous approaches. Being a computationally expensive method, the last decade saw the emergence of several techniques to speed up the underlying spectral decomposition process in [1]. These techniques range from various optimization techniques like multilevel solvers [8, 9], to systems implementations exploiting GPU parallelism [3], and to approximate methods like reduced order normalized cuts [4, 10].

3D Volumetric Image Segmentation

Segmentation techniques applied to 3D volumetric images can be found in the literature of medical imaging (usually MRI) and 2D+time video sequence segmentation. In 3D medical image segmentation, unsupervised techniques like region growing [11] have been well studied. In [12], normalized cuts were applied to MRIs but gained little attention. Recent literature mostly focuses on supervised techniques [13, 14]. In video segmentation, there are two streams of works. On the one hand, [15, 16, 17, 18] are frameworks that rely on segmenting each frame into superpixels in the first place. Therefore they are not applicable to general 3D volumes. On the other hand, a variety of other methods treat video sequences as spatiotemporal volumes and try to segment them into supervoxels. These methods are mostly extensions of popular 2D image segmentation approaches, including graph cuts [19, 20], SLIC [21], mean shift [22], graph-based methods [6] and normalized cuts [23, 9]. Besides these, temporal superpixels [24, 25, 26] are another set of effective approaches that extract high quality spatiotemporal volumes. A comprehensive review of supervoxel methods can be found in [27, 28]. All of this development in the area of 3D supervoxels has led to a general consensus in the video segmentation community that supervoxels have favorable properties which can be leveraged later in the pipeline. These characteristics are summarized as follows: (i) supervoxel boundaries should stick to meaningful image boundaries; (ii) regions within one supervoxel should be homogeneous while inter-supervoxel differences should be substantially larger; (iii) it is important that supervoxels have regular topologies; (iv) a hierarchy of supervoxels is favored as different applications have different supervoxel granularity preferences.

Deep Neural Networks

Recently, deep learning, a re-branding of neural networks (NNs), has seen successful application in numerous areas, especially in computer vision and natural language processing. Typical approaches, such as Support Vector Machines (SVMs), Gaussian Mixture Models, Hidden Markov Models, and Conditional Random Fields, mostly have “shallow” architectures that have only one to two processing layers between the features and the output. Deep neural networks, e.g. convolutional neural networks (CNNs) and recurrent neural

networks (RNNs), usually have many layers which allow for a rich variety of highly complex, nonlinear and hierarchical features to be learned from the data. Video segmentation problems involve several subproblems which can be classified as semantic segmentation [29] and general object segmentation [30, 31]. For video object segmentation problems, convolutional neural networks (CNNs) based approaches have made great advances in both unsupervised learning [30, 32, 33, 34, 35] and semi-supervised learning [36, 37, 38]. Convolutional Neural Networks (CNNs) have recently been integrated with gPb -UCM. Convolutional Oriented Boundaries (COB) [39, 40] is a generic CNN architecture that allows end-to-end learning of multiscale oriented contours and provides state of the art contours and region hierarchies. It uses a pipeline similar to gPb -UCM but replaces the contour probability maps by CNN maps from VGGNet [41] and ResNet [42]. We introduced transfer learning to 3D UCM without the intermediate supervised CNN layers as in COB. Our very preliminary investigation into using transfer learning (re-application of discriminative deep learning features) in 3D UCM was unsuccessful in outperforming unsupervised 3D UCM. However, due to rapid developments that are ongoing in this space, a definitive conclusion cannot be reached at this time. Since the basic approach of adding deep network features to our 3D UCM pipeline is extremely straightforward, we plan to revisit this in the future to see if significant improvements can be achieved over and beyond our current unsupervised learning UCM pipeline.

3 The supervoxel extraction framework

As mentioned in the Introduction, our work is a natural extension of UCM, the state of the art 2D superpixel extraction approach. We now briefly describe the steps in 2D UCM since these are closely followed (for the most part) in our 3D effort.

UCM begins by introducing the Pb gradient-based detector which assigns a “probability” value $Pb(x, y, \theta)$ at every location (x, y) and orientation θ . This is obtained from an oriented gradient detector $G(x, y, \theta)$ applied to an intensity image. A multiscale extension $mPb(x, y, \theta)$ to the Pb detector is deployed by executing the gradient operator at multiple scales followed by cue combination over color, texture and brightness channels. Thus far, only local information (relative to scale) has been used. In the next globalization step, a weighted graph is constructed with graph edge weights being set proportional to the evidence for a strong contour connecting the nodes. This is performed for all pairs of pixels resulting in an $N^2 \times N^2$ graph given N^2 pixels. The top eigenvectors of the graph are extracted, placed in image coordinates followed by gradient computation. This results in the $sPb(x, y, \theta)$ detector which carries global information as it is derived from eigenvector “images.” Finally, the globalized probability detector $gPb(x, y, \theta)$ is computed via a weighted linear combination of mPb and sPb . While this completes the pipeline in terms of information accrued for segmentation, UCM then proceeds to obtain a set of closed regions using gPb as the input via the application of the oriented watershed transform (OWT). Watershed based flood filling is performed at the lowest level of a hierarchy leading to an oversegmentation. A graph-based region merging algorithm (with nodes, edges and weights corresponding to regions, separating arcs and measures of dissimilarity respectively) is deployed resulting in an entire hierarchy of contained segmentations (respecting an ultrametric). The UCM pipeline can be broadly divided into (i) the gPb detector, (ii) the oriented watershed transform and (iii) graph-based agglomeration.

The pipeline in our 3D UCM framework closely follows that of 2D gPb -UCM with one major exception which will be clarified and elaborated upon below. The greater voxel cardinality forces us to revamp the globalization step above wherein a graph is constructed from every pair of voxels: if the volumetric dataset is $N \times N \times N$, the graph is $N^3 \times N^3$ which is too large for eigenvector computation. Therefore, computational considerations force us to adopt reduced order eigensystem solvers. A second exception concerns the agglomeration step. 2D UCM merges regions together by only considering contour edge pixels at the base level and not all pixels. This approach leads to the creation of fragmented surfaces in 3D. To overcome this problem, we perform graph-based agglomeration using all voxels following recent work. With these changes to the pipeline, the 3D UCM framework is broadly subdivided into (i) local, volume gradient detection, (ii) globalization using reduced order eigensolvers and (iii) graph-based agglomeration to reflect the emphasis on the changed subsystems. The upside is that 3D UCM becomes scalable to handle sizable datasets.

Local gradient feature extraction

The UCM framework begins with gradient-based edge detection to quantify the presence of boundaries. Most gradient-based edge detectors in 2D [43, 44] can be extended to 3D for this purpose. The 3D gradient operator used in this work is based on the *mPb* detector proposed in [1, 45] which has been empirically shown to have superior performance in 2D.

The building block of the 3D *mPb* gradient detector is an oriented gradient operator $G(x, y, z, \theta, \varphi, r)$ described in detail in Fig. 1. To be more specific, in a 3D volumetric or spatiotemporal intensity field, we place a sphere centered at each voxel to denote its neighborhood. An equatorial plane specified by its normal vector $\vec{t}(\theta, \varphi)$ splits the sphere into two half spheres. We compute the intensity histograms for both half spheres, denoted as \mathbf{g} and \mathbf{h} . Then we define the gradient magnitude in the direction $\vec{t}(\theta, \varphi)$ as the χ^2 distance between \mathbf{g} and \mathbf{h} :

$$\chi^2(\mathbf{g}, \mathbf{h}) = \frac{1}{2} \sum_i \frac{(g(i) - h(i))^2}{g(i) + h(i)}. \quad (1)$$

In order to capture gradient information at multiple scales, this gradient detector is executed for different radius values r of the neighborhood sphere. Gradients obtained from different scales are then linearly combined together using

$$G_s(x, y, z, \theta, \varphi) \equiv \sum_r \alpha_r G(x, y, z, \theta, \varphi, r) \quad (2)$$

where α_r weighs the gradient contribution at different scales. For multi-channel 3D images like video sequences, $G_s(x, y, z, \theta, \varphi)$ is separately calculated from each channel and summed up using equal weights. Finally, the measure of boundary strength at (x, y, z) is computed as the maximum response over all directions $\vec{t}(\theta, \varphi)$:

$$mPb(x, y, z) \equiv \max_{\theta, \varphi} G_s(x, y, z, \theta, \varphi). \quad (3)$$

In our experiments, θ and φ take values in $\{0, \frac{\pi}{4}, \frac{\pi}{2}, \frac{3\pi}{4}\}$ and $\{-\frac{\pi}{4}, 0, \frac{\pi}{4}\}$ respectively and in one special case, $\varphi = \frac{\pi}{2}$. Therefore we compute local gradients in 13 different directions. Neighborhood values of 2, 4 and 6 voxels were used for r . Equal weights α_r were used to combine gradients across different scales. Also, as is standard, we always apply an isotropic Gaussian smoothing filter with $\sigma = 3$ voxels before any gradient operation.

Globalization using a reduced order eigensystem

The globalization core of *gPb*-UCM is driven by nonlinear dimensionality reduction (closely related to spectral clustering). The local cues obtained from the gradient feature detection phase are globalized (and therefore emphasize the most salient boundaries in the image) by computing generalized eigenvectors of a graph Laplacian (obtained from the normalized cuts principle) [2]. However, this approach depends on solving a sparse eigensystem at the scale of the number of pixels in the image. Thus as the size of the image grows large, the globalization step becomes the computational bottleneck of the entire process. This problem is even more severe in the 3D setting because the voxel cardinality far exceeds the pixel cardinality of our 2D counterparts. An efficient approach was proposed in [4] to reduce the size of the eigensystem while maintaining the quality of the eigenvectors used in globalization. We generalize this method to 3D so that our approach becomes scalable to handle sizable datasets.

In the following, we describe the globalization steps: (i) graph construction and oriented intervening contour cue, (ii) reduced order normalized cuts and eigenvector computation, (iii) scale-space gradient computation on the eigenvector image and (iv) the combination of local and global gradient information. For the most part, this pipeline mirrors the transition from *mPb* to *gPb* with the crucial difference being the adoption of reduced order normalized cuts.

Graph construction and the oriented intervening contour cue

In 2D, the normalized cuts approach begins with sparse graph construction obtained by connecting pixels that are spatially close to each other. *gPb*-UCM [1] specifies a sparse symmetric affinity matrix W using the

intervening contour cue [5] which is the maximal value of mPb along a line connecting the two pixels i, j at the ends of relation W_{ij} . However, this approach does not utilize all of the useful information obtained from the previous gradient feature detection step. Fig. 2 describes a potential problem and our resolution. To improve the accuracy of the affinity matrix, we take the direction vector of the maximum gradient magnitude into consideration when calculating the pixel-wise affinity value. This new variant is termed the *oriented intervening contour cue*. For any spatially close voxels i and j , we use $\vec{i\bar{j}}$ to denote the line segment connecting i and j . \vec{d} is defined as the unit direction vector of $\vec{i\bar{j}}$. Assume P is a set of voxels that lie close to $\vec{i\bar{j}}$. For any $p \in P$, \vec{n} is the unit direction vector associated with its mPb value. We define the affinity value W_{ij} between i and j as follows:

$$W_{ij} = \exp(-\max_{p \in P} \{mPb(p) |\langle \vec{d}, \vec{n} \rangle| \} / \rho) \quad (4)$$

where $\langle \cdot, \cdot \rangle$ is the inner product operator of the vector space and ρ is a scaling constant. In our experiments, the set P contains the voxels that are at most 1 voxel away from $\vec{i\bar{j}}$. ρ is set to 0.1. In the affinity matrix W , each voxel is connected to voxels that fall in the $5 \times 5 \times 5$ cube centered at that voxel. Thus the graph defined by W is very sparse.

Reduced order normalized cuts and eigenvector computation

At this point, standard 2D gPb -UCM solves for the generalized eigenvectors of the sparse eigensystem

$$(D - W)\vec{v} = \lambda D\vec{v} \quad (5)$$

where D is a diagonal matrix defined by $D_{ii} = \sum_j W_{ij}$. However, this eigenvector problem is computationally very intensive. It becomes the bottleneck, both in time and memory efficiency, of the normalized cuts based segmentation algorithms. To overcome this problem, an efficient and highly parallel GPU implementation was provided in [3]. However, this approach requires us to use GPU-based hardware and software suites—an unnecessary restriction at this stage of development. A clever alternative in [4, 10] builds the graph on superpixels instead of pixels to reduce the size of the eigensystem. We chose to generalize [4]’s approach to 3D as (i) the superpixel solution is more scalable than the GPU solution in terms of memory requirements, (ii) specialized GPU co-processors are not commonly available in many computing platforms like smart phones and wearable devices, and (iii) the approach in [10] is specifically designed for superpixels in each frame in video segmentation, thus not easily generalizable. Finally, the approach in [4] constructs a reduced order normalized cuts system which is easier to solve. We denote m as the number of supervoxels, and n as the number of voxels. The reduced order eigensystem is denoted by

$$(L^T(D - W)L)\vec{x} = \lambda' L^T D L \vec{x} \quad (6)$$

where $L \in \mathbb{R}^{m \times n}$, $\vec{x} \in \mathbb{R}^m$ and $L\vec{x} = \vec{v}$. The purpose of L is to assign each pixel to a superpixel/supervoxel. In our approach, the supervoxels are generated by a watershed transform on the mPb image obtained from the volumetric gradient feature detection step. Obviously the number of supervoxels m is much smaller than the number of voxels n in the 3D volumetric/spatio-temporal image. In practice, there are usually two to three orders reduction in the size of the eigensystem (from millions of voxels to few thousands of supervoxels). Therefore it is much more efficient to solve Equation (6) than Equation (5).

Scale space gradient computation on the eigenvector image

We solve for the generalized eigenvectors $\{\vec{x}_0, \vec{x}_1, \dots, \vec{x}_n\}$ of the system in (6) corresponding to the smallest eigenvalues $\{\lambda'_0, \lambda'_1, \dots, \lambda'_n\}$. As stated in [4], λ_i in (5) will be equal to λ'_i and $L\vec{x}_i$ will match \vec{v}_i modulo an irrelevant scale factor, where \vec{v}_i are the eigenvectors of the original eigensystem (5). Similar to the 2D scenario [1], eigenvectors \vec{v}_i carry surface information. Fig. 3 shows several example eigenvectors obtained from two types of 3D volumetric datasets. In both cases, the eigenvectors distinguish salient aspects of the original image. Based on this observation, we apply the gradient operator mPb defined in (3) to the

eigenvector images. The outcome of this procedure is denoted as sPb because it represents the *spectral* component of the boundary detector, following the convention established in [1]:

$$sPb(x, y, z) = \sum_{i=1}^K \frac{1}{\sqrt{\lambda_i}} mPb_{\vec{v}_i}(x, y, z). \quad (7)$$

Note that this weighted summation starts from $i = 1$ because λ_0 always equals 0 and \vec{v}_0 is a vanilla image. The weighting by $1/\sqrt{\lambda_i}$ is inspired by the mass-spring system in mechanics [46, 1]. In our experiments, we use 16 eigenvectors, i.e. $K = 16$.

The combination of local and global gradient information

The last step is to combine local cues mPb and global cues sPb . mPb tries to capture local variations while sPb aims to obtain salient boundary surfaces. By linearly combining them together, we get a *globalized* boundary detector gPb :

$$gPb(x, y, z) = \omega mPb(x, y, z) + (1 - \omega) sPb(x, y, z). \quad (8)$$

In practice, we use equal weights for mPb and sPb . After obtaining the gPb values, we apply a post-processing step of non-maximum suppression [43] to get thinned boundary surfaces when the resulting edges from mPb are too thick. Fig. 4 shows some examples of mPb , sPb and gPb .

Supervoxel agglomeration

At this point, 2D gPb -UCM proceeds with the oriented watershed transform (OWT) [47, 48, 1] to create a hierarchical segmentation of the image resulting in the ultrametric contour map. However, we find that the same strategy does not work well in 3D. The reasons are two-fold. First, because of the presence of irregular topologies, it is more difficult to approximate boundary surfaces with square or triangular meshes in 3D than to approximate boundary curves with line segments in 2D. Second, following OWT, during superpixel merging, only the information of the pixels on boundaries are used during greedy boundary removal. This is not a robust design especially when we take into account the fragmentation of boundary surfaces in 3D.

For the above reasons, we turn to the popular graph based image and video segmentation methods [7, 6] to create the segmentation hierarchy. We first apply a watershed transform to the gPb field obtained from the previous step to get an oversegmentation. Next we iteratively merge the adjacent segments beginning with this oversegmentation. The output of this procedure is a segmentation hierarchy represented by a tree-structure whose lower level segments are always contained in higher level segments. As in [6], the merge rules run on a graph. The nodes of the graph are regions and the edges represent region to region relations. First, for any two adjacent regions R_i and R_j , we assign an edge e_{ij} to connect them on the graph. The weight of e_{ij} is set to the χ^2 distance between (Lab space) intensity value histograms of R_i and R_j with 20 bins used. Also, for any region R , a quantity named the relaxed internal variation $RInt(R)$ is defined as follows:

$$RInt(R) \equiv Int(R) + \frac{\tau}{|R|} \quad (9)$$

where $Int(R)$ is defined as the maximum edge weight of its minimum spanning tree (MST). For the lowest level regions, i.e. the regions of oversegmentation obtained from the watershed transform, $Int(R)$ is set to 0. $|R|$ is the voxel cardinality of region R . τ is a parameter which triggers the merging process and controls the granularity of the regions. In each iteration of merging, all the edges are traversed in ascending order. For any edge e_{ij} , we merge incident regions R_i and R_j if the weight of e_{ij} is less than the minimum of the relaxed internal variation of the two regions. Thus the merging condition is written as

$$weight(e_{ij}) < \min\{RInt(R_i), RInt(R_j)\}. \quad (10)$$

In practice, we increase the granularity parameter τ by a factor of 1.1 in each iteration. This agglomeration process iteratively progresses until no edge meets the merging criterion. The advantage of graph based methods is that they make use of the information in *all* voxels in the merged regions. Furthermore, as shown

in the experiments below, we see that it overcomes the weakness of fragmented supervoxels of previous graph based methods. This is because traditional graph based methods are built on voxel-level graphs.

Finally, we obtain a supervoxel hierarchy represented by a bottom-up tree structure. This is the final output of the 3D UCM algorithm. The granularity of the segmentation is a user driven choice guided by the application.

4 Integration with optical flow

Integration of 3D UCM with optical flow is quite natural in the spatiotemporal setting. We would like to establish spatiotemporal coherence and stability of structures across time and one way to achieve this is via the application of optical flow. Estimation of optical flow across multiple frames is a difficult problem: motion models and optimization strategies are required to resolve the ambiguities and discontinuities prevalent. However, when multiple frames are simultaneously considered, we are afforded with the possibility of coupling supervoxel estimation with multiframe optical flow. While this is straightforward to visualize in video segmentation, we argue that the registration of multiple frames across time can be of benefit in far flung applications such as particle flows in computational fluid dynamics and change detection in remote sensing.

Our work is an extension of 3D UCM supervoxel segmentation to the spatiotemporal case. The pipeline in the 3D UCM framework is to sequentially perform gradient feature extraction, globalization and supervoxel agglomeration on 3D volumetric images. In the spatiotemporal case, we precede 3D UCM by first bringing the multiple frames into register using optical flow. This integration is dubbed 3D OF UCM to denote the use of optical flow.

Optical flow computation

Optical flow has a complex relationship to the actual 3D motion field since it represents apparent motion caused by intensity change [49, 50]. Optical flow has been a central topic in computer vision since the 80s and has considerably matured since then [51]. The basic approach consists of estimating a vector field (velocity vectors at each grid point) which is subsequently related to object motion.

In general, optical flow can be divided into two types, sparse optical flow and dense optical flow [51, 52]. Briefly, the sparse optical flow method only computes flow vectors at some “interesting” pixels in the image, while the dense optical flow method estimates the flow field at every pixel and is more accurate [53, 54]. We use dense optical flow here since we wish to integrate with scene segmentation.

Supervoxel segmentation

In recent work on video segmentation problems, motion estimation has emerged as a key ingredient for state of the art approaches [7, 55, 56] for both deep learning methods and typical methods. We introduced optical flow estimation to the pipeline of 3D OF UCM. On top of 3D UCM, we first compute optical flow and then consider to employ it as the source of an additional cue to guide the segmentation.

Given a video sequence, we compute the optical flow, and combine RGB channels with the magnitude of the optical flow as input to 3D UCM. We use the state of the art Classic+NL-fast method [57] to estimate the optical flow. Fig. 5 shows the examples of optical flow computed for video sequences from [58, 31]. This is repeated across multiple frames followed by interpolation to bring the set of frames into register [59]. We used color coding in [60] for optical flow in order to aid visualization. Subsequently, 3D UCM is applied to obtain supervoxel segmentation of spatiotemporal imagery.

Since our input data has additional optical flow information (a motion channel), we usually achieve better segmentation results especially for some small regions. As optical flow provides complementary information to RGB images, 3D OF UCM improves the overall performance. Datasets, evaluation metrics and comparison results are presented in Section 5.

5 Evaluation

We perform quantitative and qualitative comparisons between 3D UCM, 3D OF UCM and state of the art supervoxel methods on different types of 3D volumetric and spatiotemporal datasets.

Experimental Setup

Datasets

We use three video datasets for quantitative and qualitative comparisons. The most typical use cases of 3D segmentation are medical images like MRIs and video sequences. We use the publicly available Internet Brain Segmentation Repository (IBSR) [61] for our medical imaging application. It contains 18 volumetric brain MRIs with their white matter, gray matter and cerebro-spinal fluid labeled by human experts. These represent cortical and subcortical structures of interest in neuroanatomy.

The second dataset is BuffaloXiph [58]. It is a subset of the well known xiph.org videos and the 8 video sequences in the dataset have 69-85 frames with dense semantic pixel labels. The third dataset is the DAVIS dataset (Densely Annotated Video Segmentation) [31], a benchmark dataset designed for the task of video object segmentation. DAVIS contains 55 high quality and full HD video sequences and densely annotated, pixel-accurate ground truth masks, spanning multiple occurrences of common video object segmentation challenges such as occlusions, motion blur and appearance changes. For quantitative comparisons, we use the BuffaloXiph dataset and carry out experiments with a thorough quantitative analysis on a well accepted benchmark, LIBSVX [28].

Methods

A comprehensive comparison of current supervoxel and video segmentation methods is available in [27, 28]. Building on their approach, in our experiments, we compare 3D UCM and 3D OF UCM to two state of the art supervoxel methods: (i) hierarchical graph based (GBH) [6] and (ii) segmentation by weighted aggregation (SWA) [8, 9, 62]. GBH is the standard graph based method for video segmentation and SWA is a multilevel normalized cuts solver that also generates a hierarchy of segmentations. Some quantitative superpixel evaluation metrics have been recently used [63] for frame-based 2D images. For 3D spacetime validation, a comprehensive comparison of current supervoxel and video segmentation methods is available in [27, 28].

Qualitative Comparisons

We present some example segmentation results for three datasets: IBSR, BuffaloXiph and DAVIS. To compare GBH, SWA and 3D UCM, we present some slices from both the IBSR and BuffaloXiph datasets in Fig. 6. It shows that, through different levels, there are significant differences in the segmentation results of all three methods. However, it is difficult to say which method is the best compared with the ground truth. It should be noted that GBH has a fragmentation problem in the IBSR and BuffaloXiph datasets. On the other hand, SWA has a clean segmentation in IBSR but suffers from fragmentation in video sequences. In contrast, 3D UCM has the most regular segmentation.

And in Fig. 7 and Fig. 8, from top to bottom, we show fine level to coarse level supervoxel segmentation comparison of these methods on IBSR and BuffaloXiph datasets respectively. The results show that all four methods differ markedly in the kind of segmentations they obtain. For brain frames in the IBSR dataset, the visual comparison shows that 3D UCM performs well especially on the globalization aspect, such as having no oversegmentation in the background. As can be seen, GBH and SWA both have a fragmentation problem in the BuffaloXiph datasets through different levels. In contrast, 3D UCM has the most regular segmentation and 3D OF UCM has the finest segmentation even in small detailed regions. In general, 3D UCM generates meaningful and compact supervoxels at all levels of the hierarchy while GBH and SWA have a fragmentation problem at the lower levels. To demonstrate the improved supervoxel agglomeration of 3D OF UCM compared to 3D UCM, examine Fig. 9 which shows the original images, optical flow images, segmentation from 3D UCM and 3D OF UCM. The optical flow images generated prior to segmentation extract the boundaries of different objects and improve the supervoxel segmentation. For instance, the two

legs of the person in the center of the image in Fig. 9 are grouped as the same supervoxel object by 3D OF UCM but segmented as different objects by 3D UCM. Also for the snowboard frames, 3D OF UCM generates finer supervoxels for both the moving object and the background.

And in the DAVIS dataset, we selected some examples of optical flow and segmentation results of 3D OF UCM in Fig. 10. It shows that 3D OF UCM generates fine supervoxels and avoids unnecessary over-segmentation of the background for the moving object because of information from the optical flow motion channel. The result can be regarded as early video processing and be employed to other vision tasks, e.g. semantic segmentation and object segmentation. To sum up, 3D UCM leverages the strength of a high quality boundary detector and based on that, 3D OF UCM leverages additional flow field information to perform better. Hence, from qualitative comparisons on two benchmark datasets, it suffices to say that 3D OF UCM generates more meaningful and compact supervoxels at all levels of the hierarchy than three other methods.

Quantitative Measures

As both the IBSR and BuffaloXiph datasets are densely labeled, we are able to compare the supervoxel methods on a variety of characteristic measures. Besides, we re-examined our 3D OF UCM experiments to quantify improvement by using optical flow as complementary cue. In the previous experiment integrating optical flow with 3D UCM, we only compute optical flow for the first 15 frames and added the flow vectors to the original RGB channels. In the next experiment, we process optical flow estimation for the complete video in the BuffaloXiph dataset, and the results show more improvements. We posit the improvements mostly benefit from consistency of motion flow in video. For performance comparisons between 3D OF UCM and 3D UCM, the main improvements are demonstrated on quantitative measures. From the figures of quantitative measures, we notice that, expect at very small granularity, 3D OF UCM perform better than 3D UCM on boundary quality, region quality and supervoxel compactness.

Boundary Quality

The precision-recall curve is the most recommended measure and has found widespread use in comparing image segmentation methods [45, 1]. This was introduced into the world of video segmentation in [16]. We use it as the boundary quality measure in our benchmarks. It measures how well the machine generated segments stick to ground truth boundaries. More importantly, it shows the tradeoff between the positive predictive value (precision) and true positive rate (recall). For a set of machine generated boundary pixels S_b and human labeled boundary pixels G_b , precision and recall are defined as follows:

$$\text{precision} \equiv \frac{|S_b \cap G_b|}{|S_b|}, \quad \text{recall} \equiv \frac{|S_b \cap G_b|}{|G_b|}. \quad (11)$$

We show the precision-recall curves on IBSR and BuffaloXiph datasets in Fig. 11a and 11d respectively. 3D UCM performs the best on IBSR and is the second best on BuffaloXiph while 3D OF UCM perform the best. GBH does not perform well on BuffaloXiph while SWA is worse on IBSR. One limitation of 3D OF UCM and 3D UCM is that there is an upper limit of its boundary recall because it is based on supervoxels. GBH and SWA can have arbitrarily fine segmentations. Thus they can achieve a recall rate arbitrarily close to 1, although the precision is usually low in these situations. Since our 3D OF UCM adds a motion channel to 3D UCM, its boundary quality is better than 3D UCM.

Region Quality

Measures based on overlaps of regions such as Dice’s coefficients are widely used in evaluating region covering performances of voxelwise segmentation approaches [64, 65]. We use the 3D segmentation accuracy introduced in [27, 28] to measure the average fraction of ground truth segments that is correctly covered by the machine generated supervoxels. Given that $G_v = \{g_1, g_2, \dots, g_m\}$ are ground truth volumes, $S_v = \{s_1, s_2, \dots, s_n\}$ are supervoxels generated by the algorithms and V represents the whole volume, the

3D segmentation accuracy is defined as

$$\text{3D segmentation accuracy} \equiv \frac{1}{m} \sum_{i=1}^m \frac{\sum_{j=1}^n |s_j \cap g_i| \times \mathbf{1}(|s_j \cap g_i| \geq |s_j \cap \bar{g}_i|)}{|g_i|} \quad (12)$$

where $\bar{g}_i = V \setminus g_i$. We plot the 3D segmentation accuracy against the number of supervoxels in Figs. 11b and 11e. GBH and SWA again perform differently in IBSR and BuffaloXiph datasets. But 3D UCM consistently performs well and 3D OF UCM shows the best performance on the BuffaloXiph dataset, especially when the number of supervoxels is low.

Supervoxel Compactness

Compact supervoxels of regular shapes are always favored because they benefit further higher level tasks in computer vision. The compactness of superpixels generated by a variety of image segmentation algorithms in 2D was investigated in [66]. It uses a measure inspired by the isoperimetric quotient to measure compactness. We use another quantity defined in a similar manner to specific surface area in material science and biology. In essence, specific surface area and isoperimetric quotient both try to quantify the total surface area per unit mass or volume. Formally, given that $S_v = \{s_1, s_2, \dots, s_n\}$ are supervoxels generated by algorithms, the specific surface area used to measure supervoxel compactness is defined as

$$\text{specific surface area} = \frac{1}{n} \sum_{i=1}^n \frac{\text{Surface}(s_i)}{\text{Volume}(s_i)} \quad (13)$$

where Surface() and Volume() count voxels on the surfaces and inside the supervoxels respectively. Lower values of specific surface area imply more compact supervoxels. The compactness comparisons on IBSR and BuffaloXiph are shown in Figs. 11c and 11f. We see that 3D UCM always generates the most compact supervoxels except at small supervoxel granularities on IBSR. GBH does not perform well in this measure because of its fragmentation problem, which is consistent with our qualitative observations.

Our quantitative measures cover boundary quality, region quality and supervoxel compactness and we feel these are the most important aspects of supervoxel quality. 3D OF UCM always performs better than 3D UCM and achieves best or the second best in all measures on the BuffaloXiph dataset. In contrast, GBH and SWA fail on some measures. In conclusion, we have demonstrated that 3D OF UCM has improved supervoxel characteristics (relative to 3D UCM) according to various measures on spatiotemporal datasets while 3D UCM is a very competitive supervoxel method on 3D volumetric datasets.

6 Discussion

In this paper, we presented the 3D UCM supervoxel framework, an extension of the most successful 2D image segmentation technique, gPb -UCM, to 3D. Experimental results show that our approach outperforms the current state of the art in most benchmark measures on two different types of 3D volumetric datasets. When combined with optical flow, 3D OF UCM performs very well on spatiotemporal datasets. The principal difference between 3D UCM and its 2D counterpart is a reduced order normalized cuts technique to overcome a computational bottleneck of the traditional gPb -UCM framework when directly extended to 3D. This immediately allows our method to scale up to larger datasets. When we jointly consider supervoxel quality and computational efficiency, we believe that 3D UCM can become the standard bearer for massive 3D volumetric datasets. We expect applications of 3D UCM in a wide range of vision tasks, including video semantic understanding, video object tracking and labeling in high-resolution medical imaging, etc.

Considering recent advances in deep learning-based methods on vision tasks, we investigated state of the art video segmentation methods using deep neural networks and also performed additional experiments on introducing transfer learning to the 3D UCM framework. For the BuffaloXiph dataset, we extracted its low-level features (first 1-4 layers from VGG-19 [41]) from the state of the art neural network, added to its original RGB channels, and used the new data as input to our 3D UCM framework. The results did not show much improvement. In general, the transfer learning based 3D UCM has very similar performance as 3D

UCM. In some cases, the transfer learning based 3D UCM has lower segmentation accuracy than 3D UCM since the CNN provided additional features are effectively noise w.r.t. supervoxel extraction. So we did not include this method in the paper. And this remained the case even after we did extra experiments with manually selected features. We posit that it is hard to get hierarchical segmentation from end-to-end deep learning networks because of the lack of multilevel dense annotation. We also did experiments similar to the pipeline of [67] using 3D UCM for video object segmentation on the Densely Annotated Video Segmentation (DAVIS) dataset. Using supervoxels generated by our proposed work, with optical flow and saliency map, our video object segmentation results demonstrate its effectiveness and top rank performance in comparison with other state of the art unsupervised methods.

Since this is a new and fresh approach to 3D segmentation, 3D UCM still has several limitations from our perspective. First, because it is a general purpose technique, the parameters of 3D UCM have not been tuned using supervised learning for specific applications. In immediate future work, we plan to follow the metric learning framework as in [1] to deliver higher performance. Second, since it is derived from the modular framework of *gPb*-UCM, 3D UCM has numerous alternative algorithmic paths that can be considered. In image and video segmentation, a better boundary detector was proposed in [18], with different graph structures deployed in [68, 69, 17] followed by graph partitioning alternatives in [70, 71]. A careful study of these alternative options may result in improved versions of 3D UCM. A GPU implementation will almost certainly permit us to replace the reduced order normalized cuts eigensystem with the original eigensystem leading to better accuracy. Finally, 3D UCM at the moment does not incorporate prior knowledge for segmentation [72, 73]. As transfer learning methods develop in 3D leading to the publication of useful 3D filters, these can be readily incorporated into the local feature extraction step of 3D UCM. These avenues represent relatively straightforward extension paths going forward.

Availability of data and material

IBSR: <https://www.nitrc.org/projects/ibsr/>

BuffaloXiph: <https://www.cse.buffalo.edu/~jcorso/r/labelprop.html>

DAVIS: <http://davischallenge.org/>

List of abbreviations

- UCM: ultrametric contour map
- 3D OF UCM: three dimensional optical flow UCM
- *gPb*-UCM: globalized probability map UCM
- hierarchical graph based (GBH)
- segmentation by weighted aggregation (SWA)
- IBSR: Internet Brain Segmentation Repository
- DAVIS: Densely Annotated Video Segmentation
- CNNs: Convolutional Neural Networks

Acknowledgements

We acknowledge Manu Sethi for assisting us with the optical flow integration. This work is partially supported by NSF IIS 1743050 to A.R. and S.R.

References

- [1] Arbelaez, P., Maire, M., Fowlkes, C., Malik, J.: Contour detection and hierarchical image segmentation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **33**(5), 898–916 (2011)
- [2] Shi, J., Malik, J.: Normalized cuts and image segmentation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **22**(8), 888–905 (2000)
- [3] Catanzaro, B., Su, B.-Y., Sundaram, N., Lee, Y., Murphy, M., Keutzer, K.: Efficient, high-quality image contour detection. In: *Computer Vision, 2009 IEEE 12th International Conference on*, pp. 2381–2388 (2009). IEEE
- [4] Taylor, C.: Towards fast and accurate segmentation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1916–1922 (2013)
- [5] Fowlkes, C., Martin, D., Malik, J.: Learning affinity functions for image segmentation: Combining patch-based and gradient-based approaches. In: *Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer Society Conference on*, vol. 2, pp. 54–61 (2003). IEEE
- [6] Grundmann, M., Kwatra, V., Han, M., Essa, I.: Efficient hierarchical graph-based video segmentation. In: *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pp. 2141–2148 (2010). IEEE
- [7] Felzenszwalb, P.F., Huttenlocher, D.P.: Efficient graph-based image segmentation. *International Journal of Computer Vision* **59**(2), 167–181 (2004)
- [8] Sharon, E., Brandt, A., Basri, R.: Fast multiscale image segmentation. In: *Computer Vision and Pattern Recognition, 2000. Proceedings. IEEE Conference on*, vol. 1, pp. 70–77 (2000). IEEE
- [9] Sharon, E., Galun, M., Sharon, D., Basri, R., Brandt, A.: Hierarchy and adaptivity in segmenting visual scenes. *Nature* **442**(7104), 810–813 (2006)
- [10] Galasso, F., Keuper, M., Brox, T., Schiele, B.: Spectral graph reduction for efficient image and streaming video segmentation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 49–56 (2014)
- [11] Pohle, R., Toennies, K.D.: Segmentation of medical images using adaptive region growing. In: *SPIE Medical Imaging 2001*, pp. 1337–1346 (2001). International Society for Optics and Photonics
- [12] Carballido-Gamio, J., Belongie, S.J., Majumdar, S.: Normalized cuts in 3-D for spinal MRI segmentation. *Medical Imaging, IEEE Transactions on* **23**(1), 36–44 (2004)
- [13] Fischl, B., Salat, D.H., Busa, E., Albert, M., Dieterich, M., Haselgrove, C., Van Der Kouwe, A., Killiany, R., Kennedy, D., Klaveness, S., *et al.*: Whole brain segmentation: automated labeling of neuroanatomical structures in the human brain. *Neuron* **33**(3), 341–355 (2002)
- [14] Deng, Y., Rangarajan, A., Vemuri, B.C.: Supervised learning for brain MR segmentation via fusion of partially labeled multiple atlases. In: *IEEE International Symposium on Biomedical Imaging (ISBI)*, pp. 633–637 (2016)
- [15] Galasso, F., Cipolla, R., Schiele, B.: Video segmentation with superpixels. In: *Computer Vision—ACCV 2012*, pp. 760–774. Springer, Berlin, Heidelberg. (2012)
- [16] Galasso, F., Nagaraja, N.S., Cardenas, T.J., Brox, T., Schiele, B.: A unified video segmentation benchmark: Annotation, metrics and analysis. In: *IEEE International Conference on Computer Vision* (2013)
- [17] Khoreva, A., Galasso, F., Hein, M., Schiele, B.: Classifier based graph construction for video segmentation. In: *Computer Vision and Pattern Recognition (CVPR), 2015 IEEE Conference on*, pp. 951–960 (2015). IEEE

- [18] Khoreva, A., Benenson, R., Galasso, F., Hein, M., Schiele, B.: Improved image boundaries for better video segmentation. arXiv preprint arXiv:1605.03718 (2016)
- [19] Boykov, Y., Veksler, O., Zabih, R.: Fast approximate energy minimization via graph cuts. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **23**(11), 1222–1239 (2001)
- [20] Boykov, Y.Y., Jolly, M.-P.: Interactive graph cuts for optimal boundary & region segmentation of objects in N-D images. In: *Computer Vision, 2001. ICCV 2001. Proceedings. Eighth IEEE International Conference on*, vol. 1, pp. 105–112 (2001). IEEE
- [21] Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P., Susstrunk, S.: SLIC superpixels compared to state-of-the-art superpixel methods. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **34**(11), 2274–2282 (2012)
- [22] Paris, S., Durand, F.: A topological approach to hierarchical segmentation using mean shift. In: *Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on*, pp. 1–8 (2007). IEEE
- [23] Fowlkes, C., Belongie, S., Malik, J.: Efficient spatiotemporal grouping using the Nystrom method. In: *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*, vol. 1, pp. 231–238 (2001). IEEE
- [24] Chang, J., Wei, D., Fisher, J.: A video representation using temporal superpixels. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2051–2058 (2013)
- [25] Reso, M., Jachalsky, J., Rosenhahn, B., Ostermann, J.: Temporally consistent superpixels. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 385–392 (2013)
- [26] Bergh, M., Roig, G., Boix, X., Manen, S., Gool, L.: Online video seeds for temporal window objectness. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 377–384 (2013)
- [27] Xu, C., Corso, J.J.: Evaluation of super-voxel methods for early video processing. In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pp. 1202–1209 (2012). IEEE
- [28] Xu, C., Corso, J.J.: LIBSVX: A supervoxel library and benchmark for early video processing. *International Journal of Computer Vision*, 119(3): 272–290. (2016)
- [29] Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3213–3223. (2016)
- [30] Fragkiadaki, K., Arbelaez, P., Felsen, P., Malik, J.: Learning to segment moving objects in videos. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4083–4090 (2015)
- [31] Perazzi, F., Pont-Tuset, J., McWilliams, B., Van Gool, L., Gross, M., Sorkine-Hornung, A.: A benchmark dataset and evaluation methodology for video object segmentation. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 724–732. (2016)
- [32] Fragkiadaki, K., Zhang, G., Shi, J.: Video segmentation by tracing discontinuities in a trajectory embedding. In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pp. 1846–1853 (2012). IEEE
- [33] Koh, Y.J., Kim, C.-S.: Primary object segmentation in videos based on region augmentation and reduction. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7417–7425 (2017)
- [34] Tokmakov, P., Alahari, K., Schmid, C.: Learning video object segmentation with visual memory. arXiv preprint arXiv:1704.05737 (2017)
- [35] Jain, S.D., Xiong, B., Grauman, K.: Fusionseg: Learning to combine motion and appearance for fully automatic segmentation of generic objects in videos. arXiv preprint arXiv:1701.05384 **2**(3), 6 (2017)

- [36] Perazzi, F., Khoreva, A., Benenson, R., Schiele, B., Sorkine-Hornung, A.: Learning video object segmentation from static images. In: IEEE Conference on Computer Vision and Pattern Recognition, pp. 3491–3500. (2017)
- [37] Caelles, S., Maninis, K.-K., Pont-Tuset, J., Leal-Taixé, L., Cremers, D., Van Gool, L.: One-shot video object segmentation. In: IEEE Conference on Computer Vision and Pattern Recognition, pp. 5320–5329, (2017). IEEE
- [38] Cheng, J., Tsai, Y.-H., Wang, S., Yang, M.-H.: Segflow: Joint learning for video object segmentation and optical flow. In: 2017 IEEE International Conference on Computer Vision (ICCV), pp. 686–695 (2017). IEEE
- [39] Maninis, K.K., Pont-Tuset, J., Arbeláez, P., Gool, L.V.: Convolutional oriented boundaries. In: European Conference on Computer Vision (ECCV), Lecture Notes in Computer Science, vol. 9905, Springer, Cham, pp. 580–596. (2016)
- [40] Maninis, K.K., Pont-Tuset, J., Arbeláez, P., Gool, L.V.: Convolutional oriented boundaries: From image segmentation to high-level tasks. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), 40(4): 819–833. (2017)
- [41] Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
- [42] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778 (2016)
- [43] Canny, J.: A computational approach to edge detection. Pattern Analysis and Machine Intelligence, IEEE Transactions on (6), 679–698 (1986)
- [44] Meer, P., Georgescu, B.: Edge detection with embedded confidence. Pattern Analysis and Machine Intelligence, IEEE Transactions on **23**(12), 1351–1365 (2001)
- [45] Martin, D.R., Fowlkes, C.C., Malik, J.: Learning to detect natural image boundaries using local brightness, color, and texture cues. Pattern Analysis and Machine Intelligence, IEEE Transactions on **26**(5), 530–549 (2004)
- [46] Belongie, S., Malik, J.: Finding boundaries in natural images: A new method using point descriptors and area completion. In: Computer Vision-ECCV’98, pp. 751–766. Springer, (1998)
- [47] Arbelaez, P., Maire, M., Fowlkes, C., Malik, J.: From contours to regions: An empirical evaluation. In: Computer Vision and Pattern Recognition, IEEE Conference on, pp. 2294–2301 (2009). IEEE
- [48] Najman, L., Schmitt, M.: Geodesic saliency of watershed contours and hierarchical segmentation. Pattern Analysis and Machine Intelligence, IEEE Transactions on **18**(12), 1163–1173 (1996)
- [49] Horn, B.K., Schunck, B.G.: Determining optical flow. Artificial intelligence **17**(1-3), 185–203 (1981)
- [50] Fleet, D., Weiss, Y.: Optical flow estimation. In: Handbook of Mathematical Models in Computer Vision, pp. 237–257. Springer, Boston, MA, (2006)
- [51] Black, M.J., Anandan, P.: The robust estimation of multiple motions: Parametric and piecewise-smooth flow fields. Computer Vision and Image Understanding **63**(1), 75–104 (1996)
- [52] Brox, T., Bregler, C., Malik, J.: Large displacement optical flow. In: Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on, pp. 41–48 (2009). IEEE
- [53] Xu, L., Jia, J., Matsushita, Y.: Motion detail preserving optical flow estimation. IEEE Transactions on Pattern Analysis and Machine Intelligence **34**(9), 1744–1757 (2012)
- [54] Yilmaz, A., Javed, O., Shah, M.: Object tracking: A survey. ACM computing surveys (CSUR) **38**(4), 45 pages, (2006)

- [55] Tsai, Y.-H., Yang, M.-H., Black, M.J.: Video segmentation via object flow. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3899–3908 (2016)
- [56] Ramakanth, S.A., Babu, R.V.: Seamseg: Video object segmentation using patch seams. In: 2014 IEEE Conference on Computer Vision and Pattern Recognition, pp. 376–383 (2014).
- [57] Sun, D., Roth, S., Black, M.J.: Secrets of optical flow estimation and their principles. In: Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on, pp. 2432–2439 (2010). IEEE
- [58] Chen, A.Y., Corso, J.J.: Propagating multi-class pixel labels throughout video frames. In: Image Processing Workshop (WNYIPW), 2010 Western New York, pp. 14–17 (2010). IEEE
- [59] Thirion, J.-P.: Image matching as a diffusion process: an analogy with Maxwell’s demons. *Medical image analysis* **2**(3), 243–260 (1998)
- [60] Baker, S., Scharstein, D., Lewis, J., Roth, S., Black, M.J., Szeliski, R.: A database and evaluation methodology for optical flow. *International Journal of Computer Vision* **92**(1), 1–31 (2011)
- [61] Worth, A.: The Internet Brain Segmentation Repository (IBSR). <http://www.nitrc.org/projects/ibsr> (2016)
- [62] Corso, J.J., Sharon, E., Dube, S., El-Saden, S., Sinha, U., Yuille, A.: Efficient multilevel brain tumor segmentation with integrated Bayesian model classification. *Medical Imaging, IEEE Transactions on* **27**(5), 629–640 (2008)
- [63] Levinshtein, A., Stere, A., Kutulakos, K.N., Fleet, D.J., Dickinson, S.J., Siddiqi, K.: Turbopixels: Fast superpixels using geometric flows. *IEEE transactions on Pattern Analysis and Machine Intelligence* **31**(12), 2290–2297 (2009)
- [64] Menze, B.H., Jakab, A., Bauer, S., Kalpathy-Cramer, J., Farahani, K., Kirby, J., Burren, Y., Porz, N., Slotboom, J., Wiest, R., *et al.*: The multimodal brain tumor image segmentation benchmark (BRATS). *Medical Imaging, IEEE Transactions on* **34**(10), 1993–2024 (2015)
- [65] Malisiewicz, T., Efros, A.A.: Improving spatial support for objects via multiple segmentations. In: Proceedings of the British Machine Vision Conference (BMVC), pp. 1–10 (2007)
- [66] Schick, A., Fischer, M., Stiefelhagen, R.: Measuring and evaluating the compactness of superpixels. In: Pattern Recognition (ICPR), 2012 21st International Conference on, pp. 930–934 (2012). IEEE
- [67] Griffin, B.A., Corso, J.J.: Video object segmentation using supervoxel-based gerrymandering. arXiv preprint arXiv:1704.05165 (2017)
- [68] Ren, X., Malik, J.: Learning a classification model for segmentation. In: Computer Vision, 2003. Proceedings. Ninth IEEE International Conference on, pp. 10–17 (2003). IEEE
- [69] Briggman, K., Denk, W., Seung, S., Helmstaedter, M.N., Turaga, S.C.: Maximin affinity learning of image segmentation. In: Advances in Neural Information Processing Systems, pp. 1865–1873 (2009)
- [70] Brox, T., Malik, J.: Object segmentation by long term analysis of point trajectories. In: Computer Vision–ECCV 2010, pp. 282–295. Springer, (2010)
- [71] Palou, G., Salembier, P.: Hierarchical video representation with trajectory binary partition tree. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2099–2106 (2013)
- [72] Vu, N., Manjunath, B.: Shape prior segmentation of multiple objects with graph cuts. In: Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on, pp. 1–8 (2008). IEEE
- [73] Chen, F., Yu, H., Hu, R., Zeng, X.: Deep learning shape priors for object segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1870–1877 (2013)

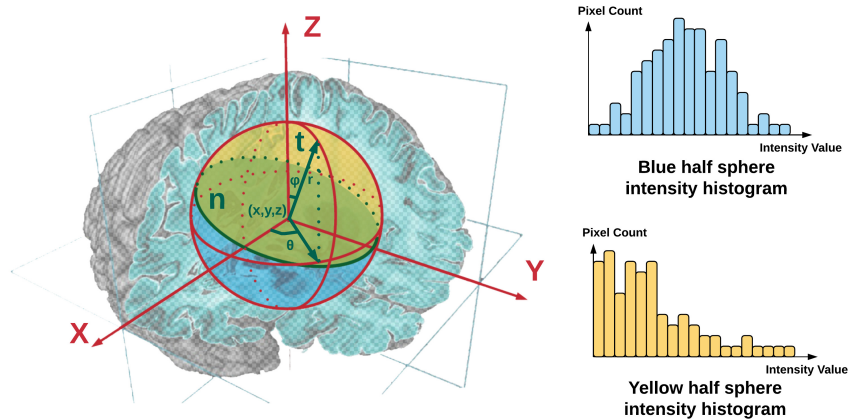


Figure 1: The **oriented gradient operator** $G(x, y, z, \theta, \varphi, r)$: At location (x, y, z) , the local neighborhood is defined by a sphere with radius r . An equatorial plane \vec{n} (shaded green) along with its normal vector \vec{t} splits the sphere into two half spheres. The one above \vec{n} is shaded yellow and the one below is shaded blue. We histogram the intensity values of voxels that fall in the yellow and blue half spheres respectively. Finally, the local gradient magnitude is set to the χ^2 distance between the yellow and blue histograms in the direction \vec{t} at scale r .

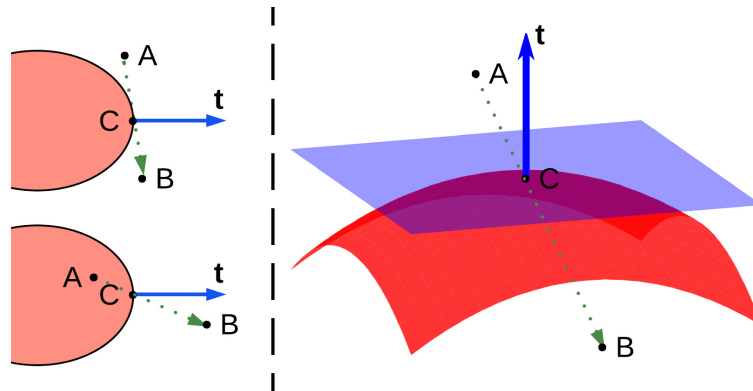


Figure 2: **Left:** Suppose we want to calculate the affinity value between voxels A and B . C is the voxel with maximal mPb value that lies on the line segment \bar{AB} . In the upper left case, A and B belong to one region. In the lower left case, A and B are in two different regions. But the intervening contour cue of UCM [?] gives the same affinity value in both cases, which is not very satisfactory. Obviously it would be better if we consider the direction \vec{t} of C 's mPb . **Right:** In our *oriented intervening contour cue* approach, when calculating affinity values, we always take the product of $mPb(C)$ with the absolute value of the inner product $\langle \vec{t}, \vec{n} \rangle$, where \vec{n} is the unit direction vector of line segment \bar{AB} . If A and B are on different sides of a boundary surface, $|\langle \vec{t}, \vec{n} \rangle|$ will be large, leading to small affinity value and vice-versa.

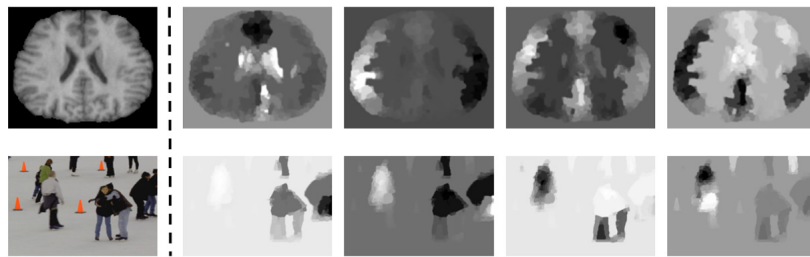


Figure 3: **Upper Left:** One slice of a brain MRI from the IBSR dataset [?].
Lower Left: One frame of a video sequence from the BuffaloXiph dataset [?].
Right: The corresponding slices/frames of the first 4 eigenvectors.

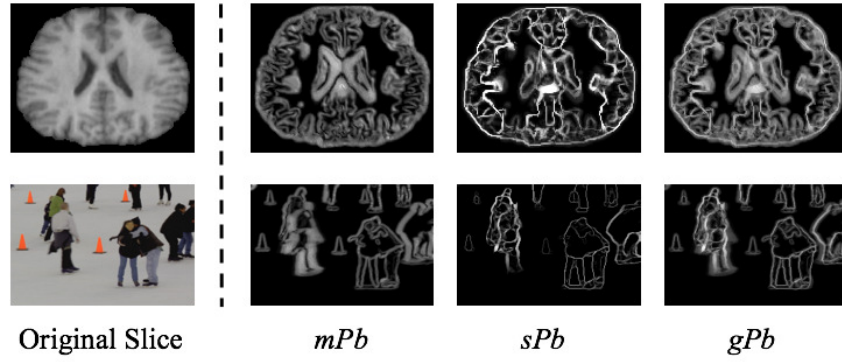


Figure 4: *sPb* augments the strength of the most salient boundaries in *gPb*.



Figure 5: Examples of optical flow images. Top: RGB images. Bottom: corresponding motion magnitude estimates.

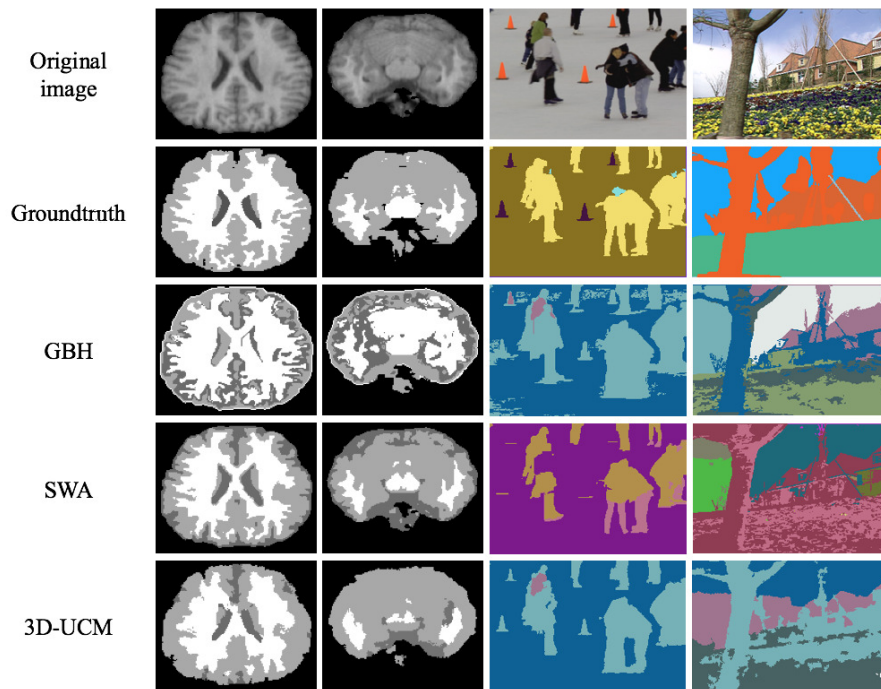


Figure 6: The level of the segmentation hierarchy is chosen to be similar to the ground truth granularity. **Left:** IBSR dataset results. The white, gray and dark gray regions are white matter, gray matter and cerebro-spinal fluid (CSF) respectively. **Right:** BuffaloXiph dataset results.

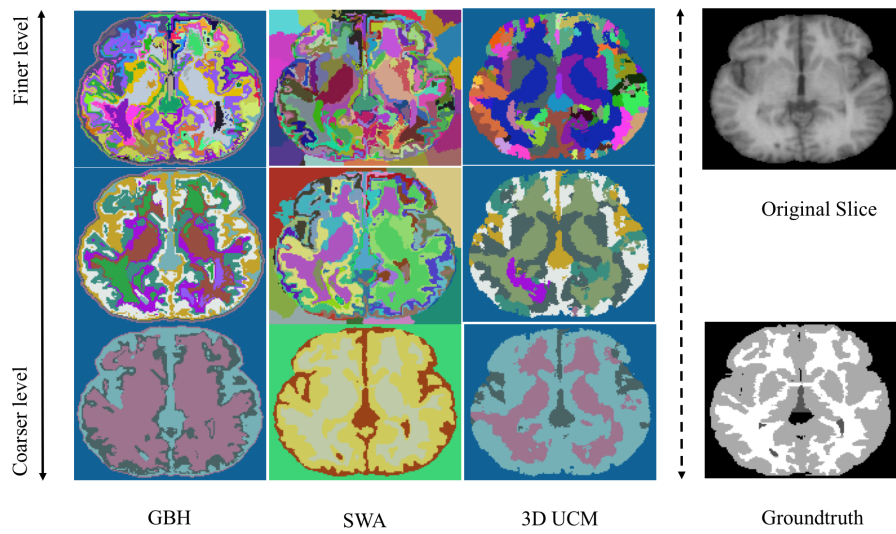


Figure 7: The hierarchy of segmentations for the IBSR dataset. **Left:** GBH. **Middle Left:** SWA. **Middle Right:** 3D UCM. **Top Right:** Original Slice. **Bottom Right:** Ground-truth.

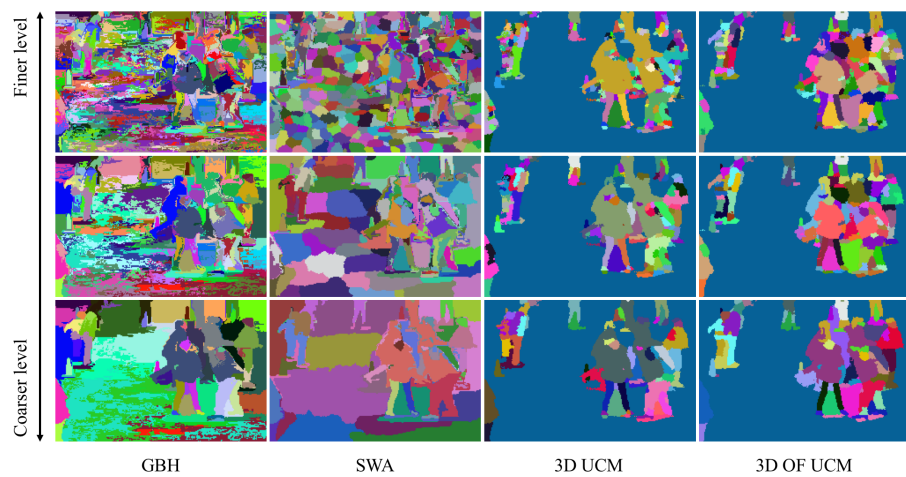


Figure 8: The hierarchy of segmentations for the BuffaloXiph dataset. **Left:** GBH. **Middle Left:** SWA. **Middle Right:** 3D UCM. **Right:** 3D OF UCM.

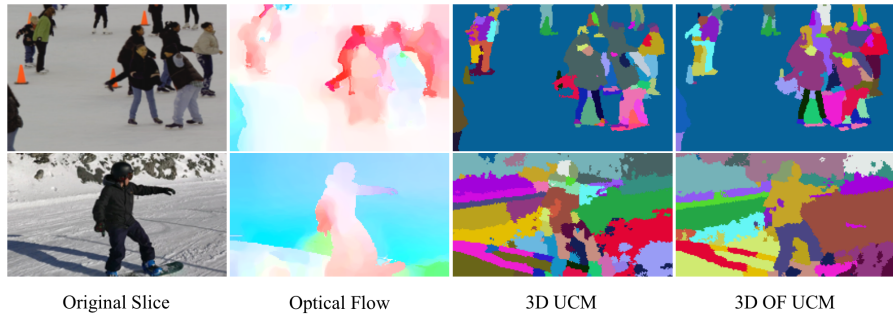


Figure 9: Comparison of 3D UCM and 3D OF UCM on the BuffaloXiph (top) and DAVIS (bottom) dataset. **Left:** Original frames. **Middle Left:** Optical flow frames. **Middle Right:** 3D UCM segmentation for selected level. **Right:** 3D OF UCM segmentation for selected level.

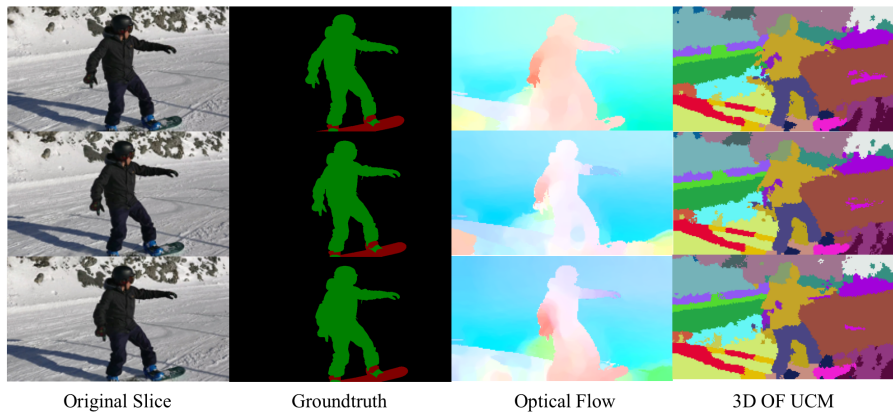


Figure 10: Optical flow estimation and supervoxel segmentation for multiple temporal frames on DAVIS dataset.

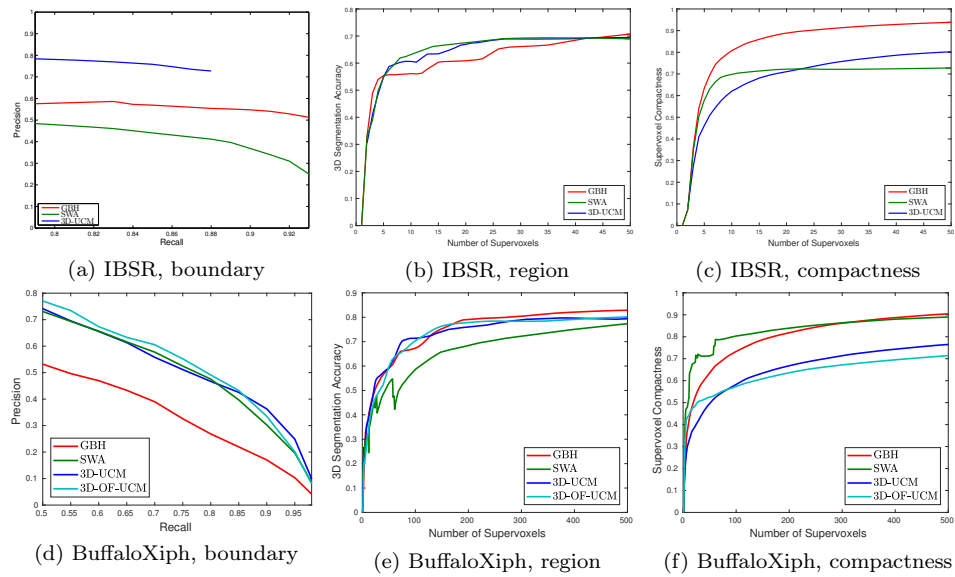


Figure 11: Quantitative measures.