

Content-Sensitive Supervoxels via Uniform Tessellations on Video Manifolds

Ran Yi, Yong-Jin Liu
Tsinghua University, China
{yr16, liuyongjin}@tsinghua.edu.cn

Yu-Kun Lai
Cardiff University, UK
Yukun.Lai@cs.cf.ac.uk

Abstract

Supervoxels are perceptually meaningful atomic regions in videos, obtained by grouping voxels that exhibit coherence in both appearance and motion. In this paper, we propose content-sensitive supervoxels (CSS), which are regularly-shaped 3D primitive volumes that possess the following characteristic: they are typically larger and longer in content-sparse regions (i.e., with homogeneous appearance and motion), and smaller and shorter in content-dense regions (i.e., with high variation of appearance and/or motion). To compute CSS, we map a video to a 3-dimensional manifold \mathcal{M} embedded in \mathbb{R}^6 , whose volume elements give a good measure of the content density in \mathcal{M} . We propose an efficient Lloyd-like method with a splitting-merging scheme to compute a uniform tessellation on \mathcal{M} , which induces the CSS in \mathcal{M} . Theoretically our method has a good competitive ratio $O(1)$. We also present a simple extension of CSS to stream CSS for processing long videos that cannot be loaded into main memory at once. We evaluate CSS, stream CSS and seven representative supervoxel methods on four video datasets. The results show that our method outperforms existing supervoxel methods.

1. Introduction

Supervoxels are perceptually meaningful atomic regions obtained by grouping similar voxels (i.e., exhibiting coherence in both appearance and motion) in the spatiotemporal domain, which over-segment a video while well preserving its structural content. Supervoxels can greatly reduce the computational complexity and have been widely used as a preprocessing step in many vision applications, such as video segmentation [12, 17, 34], foreground object segmentation [13], action segmentation and recognition [14, 21], spatiotemporal object detection [23], spatiotemporal closure in videos [15], and many others.

Depending on the size of video data, methods to compute supervoxels can be classified into *off-line* and *stream-*

ing methods. Off-line methods require the video to be short enough such that all voxels can be loaded into the memory. On the other hand, streaming methods do not have such a limitation on the video length, i.e., video data is accessed sequentially in blocks and the memory available to streaming methods only needs to fit a block. Many methodologies, such as clustering, hierarchical, generative, statistical and graph partitioning methods, have been applied to compute supervoxels; see an excellent survey in [31]. Xu and Corso [31] further select seven representative methods, including five off-line [8, 9, 24, 6, 12] and two streaming [32, 4] methods, to represent the state of the art.

To measure the quality of supervoxels, the following principles are taken into consideration: (1) *Feature preservation*: supervoxels align well with object boundaries in a video; (2) *Spatiotemporal uniformity*: in non-feature regions, supervoxels are uniform and regular in the spatiotemporal domain; (3) *Performance*: computing supervoxels is time-and-space efficient and scales well with large video data; (4) *Easy to use*: users simply specify the desired number of supervoxels and should not be bothered by tuning other parameters; (5) *Parsimony*: the above principles should be maintained with as few supervoxels as possible.

So far, none of existing methods satisfy all these principles. In this paper, we propose a content-sensitive approach to address the *parsimony* principle, and therefore, achieve a good balance among all principles. Our method is motivated by an important observation: the scene layouts and motions of different objects in a video usually exhibit large diversity, and thus the density of video content often varies significantly in different parts of the video. Applying spatiotemporally uniform distribution of supervoxels indiscriminately to the whole video often leads to under-segmentation in content-dense regions (i.e., with high variation of appearance and/or motion), and over-segmentation in content-sparse regions (i.e., with homogeneous appearance and motion). Therefore, computing supervoxels adaptively with respect to the density of video content can achieve the best performance (see Figure 1).

To compute *content-sensitive supervoxels* (CSS), we extend the image manifold concept [19, 20] to a video man-

Corresponding author

where ρ is a density function on C_M . An RV T(S_K, M) is a restricted centroidal Voronoi tessellation (RCVT) if and only if each generator $s_i \in S_K$ is the mass centroid of $C_M(s_i)$.

Theorem 1. [19] Let M be an l -manifold embedded in R^d , $d > l - 2$, and $K \in \mathbb{Z}_+$ be a positive integer. For an arbitrary set S_K of points $\{s_i\}_{i=1}^K$ in R^d and an arbitrary tessellation $\{C_i\}_{i=1}^K$ on M , $\bigcup_{i=1}^K C_i = M$, $C_i \cap C_j = \emptyset$, $i \neq j$, define the tessellation energy functional as follows:

$$E(\{(s_i, C_i)\}_{i=1}^K) = \sum_{i=1}^K \int_{C_i} \|x - s_i\|^2 dx \quad (5)$$

Then the necessary condition for E to be minimized is that $\{(s_i, C_i)\}_{i=1}^K$ is an RCVT of M .

Theorem 1 indicates that RCVT is a uniform tessellation on M , which minimizes the energy E .

2.2. Bi-criteria approximation algorithms

The discretized counterpart of RCVT is the solution to the K -means problem on the manifold domain M . Given a fixed K , denote by $S_K^{\text{opt}} = \{s_i^{\text{opt}}\}_{i=1}^K$ and $\{C_i^{\text{opt}}\}_{i=1}^K$ the (unknown) optimal generator set and tessellation on M respectively, which minimize the energy E . Let S_K and $\{C_i\}_{i=1}^K$ be the generator set and tessellation output from an algorithm A . The competitive ratio is defined to be the worst-case ratio $r = \frac{E(\{(s_i, C_i)\}_{i=1}^K)}{E(\{(s_i^{\text{opt}}, C_i^{\text{opt}})\}_{i=1}^K)}$ and A is said to be a b -approximation algorithm if $r \leq b$. An algorithm is called (a, b) -approximation, if it outputs $\{(s_i, C_i)\}_{i=1}^K$ with aK generators and tessellation cells, such that $r = \frac{E(\{(s_i, C_i)\}_{i=1}^{aK})}{E(\{(s_i^{\text{opt}}, C_i^{\text{opt}})\}_{i=1}^K)} \leq b$, where $a > 1$ and $b > 1$.

3. Overview

Our method relies on a *video manifold* representation. Denote by \mathcal{V} an input video with N voxels. Each voxel is represented by $v(x, y, t)$, where (x, y) is the pixel location and t the frame index. Inspired by the success of image manifolds [19, 20], we represent the color in the CIELAB color space. Let $c(v) = (l(v), a(v), b(v))$ be the color at the voxel v . We define a video manifold M using a stretching map $\gamma: M \rightarrow R^6$ that maps all voxels in \mathcal{V} into a 3-manifold M embedded in the 6-dimensional space:

$$\gamma(v) = (\gamma_1 x, \gamma_1 y, \gamma_2 t, \gamma_3 l(v), \gamma_3 a(v), \gamma_3 b(v)), \quad (6)$$

where we follow [20] to set global stretching factors $\gamma_1 = 0.435$ and $\gamma_3 = 1$.

For each voxel $v(x, y, t)$, denote by \mathcal{V}_v the unit cube (i.e., of size $1 \times 1 \times 1$) centered at v . Refer to Figure 2. Let a_1, a_2, \dots, a_8 be eight corner points of \mathcal{V}_v , each of which is determined by an average of its eight

Figure 2. The stretching map $\gamma: M \rightarrow R^6$ maps the unit cube \mathcal{V}_v (red box) centered at the voxel $v(x, y, t)$ into a 3-manifold M . Each corner a_i of \mathcal{V}_v , $i = 1, 2, \dots, 8$, is the center of its eight neighboring voxels; e.g., $a_1 = \frac{1}{8} \sum_{i=-1}^0 \sum_{j=0}^1 \sum_{k=-1}^0 v(x+i, y+j, t+k)$.

neighboring voxels. We decompose \mathcal{V}_v into six non-overlapped tetrahedrons tet_i , $i = 1, 2, \dots, 6$, denoted by $\{(a_5, a_6, a_3, a_2), (a_5, a_6, a_3, a_7), (a_5, a_8, a_3, a_4), (a_5, a_8, a_3, a_7), (a_5, a_1, a_3, a_2), (a_5, a_1, a_3, a_4)\}$. We approximate the volume of the *curved* tetrahedron $(\text{tet}(a, b, c, d))$ in M as

$$V(\text{tet}(a, b, c, d)) = \frac{1}{3} A(\triangle_{abc}) \cdot \text{dist}(d, \triangle_{abc}) \quad (7)$$

where $\text{dist}(d, \triangle_{abc})$ is the distance from d to the subspace spanned by vectors a, b and c , and $A(\triangle_{abc})$ is the area of the *curved* triangle (\triangle_{abc}) , which is approximated by

$$A(\triangle_{abc}) = \frac{1}{2} \sqrt{(\overline{a(b)})^2 (\overline{a(c)})^2 \sin \theta} \quad (8)$$

is the angle between vectors $\overline{a(b)}$ and $\overline{a(c)}$, and $\|u\|_2$ is the Euclidean length of the vector u in R^6 . Using the least squares method, $\text{dist}(d, \triangle_{abc})$ can be efficiently obtained as the length of the residual vector r :

$$r = d - L(L^T L)^{-1} L^T d, \quad (9)$$

where L is a 6×3 matrix $L = (a \ b \ c)$, and a, b, c, d are column 6-vectors. Then the volume of $(\mathcal{V}_v) \cap M$ is approximated by

$$V((\mathcal{V}_v) \cap M) = \sum_{i=1}^6 V(\text{tet}_i) \quad (10)$$

For any region $\mathcal{R} \subset M$, the volume of $\mathcal{R} \cap M$ is simply the sum $\sum_{v_j \in \mathcal{R}} V((\mathcal{V}_{v_j}) \cap M)$. We assume the density is 1 everywhere in M .

Our method is based on an important characteristic in the video manifold M : the volume of a region $\mathcal{R} \subset M$ depends on both the volume of \mathcal{R} and the color variation in \mathcal{R} . The higher variation of colors in \mathcal{R} , the larger the volume of $\mathcal{R} \cap M$ and vice versa. We propose an algorithm in Section 4, which is theoretically a constant-factor bi-criteria approximation, to quickly and efficiently compute a uniform tessellation $\{(s_i, C_i)\}_{i=1}^K$ in M . Then the inverse mapping γ^{-1} will transform $\{C_i\}_{i=1}^K$ into *content-sensitive* supervoxels in \mathcal{V} . See Figure 3 for an illustration.

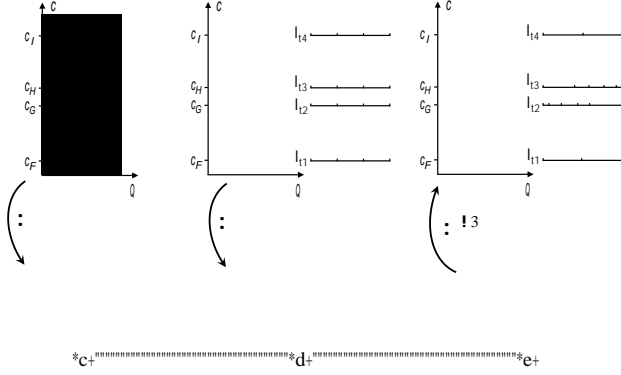


Figure 3. Overview of CSS generation on a synthetic, degenerate gray video, for easy illustration. In (a), the image frame at time t is a degenerate (1D) gray line image I_t . Supervoxels in I_t can be regarded as a tessellation of I_t . Clipping supervoxels in each line image I_t results in a set of induced superpixels in I_t . (a) We represent M as a video manifold $M = (I, t) \in \mathbb{R}^3$, whose area elements give a good measure of content density in I_t . (b) A uniform tessellation of 12 generating points on the domain of I_t leads to a non-uniform tessellation on M . The induced superpixels on line images I_{t_1} , I_{t_2} , I_{t_3} and I_{t_4} are also uniform. (c) We generate CSS by computing a uniform tessellation $\{(s_i, C_i)\}_{i=1}^{12}$ on M and the induced superpixels on line images are also content-sensitive. The real examples of induced superpixels on video frames are illustrated in Figure 1.

4. $(O(1), O(1))$ -Approximation Algorithm

To obtain a uniform tessellation $\{(s_i, C_i)\}_{i=1}^K$ in M , our algorithm consists of the following two steps:

- **Initialization (Section 4.1).** We apply a variant of the K -means++ algorithm¹ [2, 1, 30] to determine the initial positions of generating points $S_K = \{s_i\}_{i=1}^K$.
- **Lloyd refinement (Section 4.2).** Observing that the classic Lloyd method converges only to a local minimum with a large number of iterations, we propose an efficient splitting-merging scheme to compute $\text{RCVT}(S_K, M)$, which helps move the solution out of local minima and ensures a good competitive ratio.

Our algorithm is easy to implement and can obtain high-quality CSS in very few iterations. Theoretically our algorithm is $(O(1), O(1))$ -approximation. To ease reading, all proofs in this paper are presented in supplemental material.

4.1. Initialization

We apply a variant of K -means++ algorithm to obtain a provable high-quality initialization of generating points $S_K = \{s_i\}_{i=1}^K$. The pseudo-code is summarized in Algorithm 1. In each step, a point in M is picked up with

¹In literature, K -means++ algorithm includes a seeding step and a Lloyd refinement step. In our method, we only use the seeding step.

Algorithm 1 Initialization

Input: A video M of N voxels and the desired number of supervoxels K .

Output: The initial positions of K generating points $S_K = \{s_i\}_{i=1}^K$.

- 1: Compute $V(\cdot, \cdot)$ for each voxel $v \in M$ (Eq.(10)).
- 2: Choose a point v_1 from all voxels $v \in M$ with probability proportional to $V(\cdot, \cdot)$.
- 3: Set $s_1 = v_1$, $S_1 = \{s_1\}$ and $i = 1$.
- 4: **while** $i < K$ **do**
- 5: Choose a point v_{i+1} from all voxels $v \in M$ with probability proportional to the score $p_{S_i}(v)$ (Eq.(12)).
- 6: Set $s_{i+1} = v_{i+1}$, $S_{i+1} = S_i \cup \{s_{i+1}\}$ and $i = i + 1$.
- 7: **end while**

probability proportional to its current score (defined as its squared distance to the nearest generator picked so far), and added as a new generator. To compute the required probability in the manifold domain M , we consider the positions of mapped voxels $(v) \in M$, $v \in M$. With respect to an existing generator (v_i) , the score of a mapped voxel $(v_j) \in M$, $j = i$, is

$$p_{v_i}(v_j) = \frac{\int_M \|x - (v_i)\|^2 dx}{V(M) \cdot \|v_j - (v_i)\|^2} \quad (11)$$

Then the score of picking (v_j) with respect to an existing generator set S is

$$p_S(v_j) = \min_{v_i \in S} p_{v_i}(v_j) \quad (12)$$

Algorithm 1 runs in $O(NK)$ time. A simple adaption of the proofs in [2, 30] shows that using $\text{RCVT}(S_K, M)$ as the tessellation, Algorithm 1 is an expected $(\log K)$ -approximation algorithm, and furthermore, if we sample $K (> 1)$ generators, the expected approximation ratio of Algorithm 1 is bounded by

$$E(r) < 8 \cdot \left(1 + \frac{1 + \sqrt{5}}{2(\sqrt{5} - 1)}\right). \quad (13)$$

4.2. Lloyd refinement with splitting and merging

Given the initial generators $S_K = \{s_i\}_{i=1}^K$, $s_i \in M$, the classic Lloyd method computes $\text{RCVT}(S_K, M)$ iteratively by alternating the following two steps:

- **Step 1:** Fixing the generator set S_K , compute $\text{RCVT}(S_K, M)$;
- **Step 2:** For each cell C_M in $\text{RCVT}(S_K, M)$, update its generator to be the mass centroid of C_M .

Algorithm 2 CSS generation

Input: A video of N voxels, the desired number of supervoxels K , $\text{num}_{\text{random}}$ the number of repeated random sampling of generators in each iteration, and the maximum number of iterations iter_{max} .

Output: K content-sensitive supervoxels.

```

1: Initialize the generators  $S_K = \{s_i\}_{i=1}^K$  (Algorithm 1).
2: Set  $\text{iter} = 0$ .
3: while  $\text{iter} < \text{iter}_{\text{max}}$  do
4:   Compute  $\text{RV T}(S_K, M)$ .
5:   Set  $n = 0$ .
6:   while  $n < \text{num}_{\text{random}}$  do
7:     Randomly pick three generators  $s_m, s_i, s_j$  in  $S_K$ ,
       in which  $s_i$  and  $s_j$  are neighbors (Algorithm 4).
8:     Check the splitting-merging feasibility of
        $(s_m, s_i, s_j)$  (Algorithm 3) and put the return
       values in  $(\text{Flag}, s_p, s_q, s_k)$ .
9:     if  $\text{Flag} == \text{TRUE}$  then
10:      Update  $S_K$  by splitting  $s_m$  into  $(s_p, s_q)$  and
        merging  $(s_i, s_j)$  into  $s_k$ .
11:    end if
12:     $n = n + 1$ ;
13:  end while
14:  Compute  $\text{RV T}(S_K, M)$ .
15:  for each cell  $C_M(s_i)$  in  $\text{RV T}$  do
16:    Update the generator  $s_i$  to be the mass centroid of
     $C_M(s_i)$ .
17:  end for
18:   $\text{iter} = \text{iter} + 1$ ;
19: end while
20: Compute  $^{-1}(\text{RV T}(S_K, M))$  to obtain  $K$  CSS.

```

This method converges only to a local minimum with a large number of iterations [7].

To compute a high quality $\text{RCV T}(S_K, M)$ in very few iterations, we propose a splitting-merging scheme in the Lloyd iteration. The pseudo-code is summarized in Algorithm 2, in which line 4 and lines 15-16 correspond to Steps 1 and 2 in the classic Lloyd method respectively, and lines 5-14 summarize our splitting-merging scheme.

The splitting operation $S : s_m \rightarrow (s_p, s_q)$ splits a cell $C_M(s_m)$ into two new cells $C(s_p)$ and $C(s_q)$. Conversely, the merging operation $M : (s_i, s_j) \rightarrow s_k$ merges two cells $C_M(s_i)$ and $C_M(s_j)$ into a new cell $C(s_k)$. Splitting reduces the tessellation energy E and merging increases it. The number of generators does not change by applying a pair of splitting and merging operations $(S, M) : (s_m, (s_i, s_j)) \rightarrow ((s_p, s_q), s_k)$. Our goal is to design a pair (S, M) such that E does not increase.

Definition 1. The diameter d_i of a cell $C_M(s_i)$, $s_i \in S_K$, is the maximum Euclidean distance between pairs of points

Algorithm 3 Check splitting-merging feasibility

Input: Three generators (s_m, s_i, s_j) in S_K and an $\text{RV T}(S_K, M)$.

Output: A Boolean variable Flag indicating the feasibility and three new generators (s_p, s_q, s_k) .

```

1: Compute the mass centroids  $s_m, s_i$  and  $s_j$  of  $C_M(s_m)$ ,
    $C_M(s_i)$  and  $C_M(s_j)$ , respectively.
2: Compute the diameter  $d_m$  of the cell  $C_M(s_m)$  and the
   points  $p_1(d_m)$  and  $p_2(d_m)$  (see Definition 1).
3: Compute two new cells  $C(p_1(d_m))$  and  $C(p_2(d_m))$ ,
   which are the Voronoi cells of  $p_1(d_m)$  and  $p_2(d_m)$  in
   the domain  $C_M(s_m)$ .
4: Compute the mass centroids  $s_k, s_p$  and  $s_q$  of  $C_M(s_i)$ ,
    $C_M(s_j)$ ,  $C(p_1(d_m))$  and  $C(p_2(d_m))$ , respectively.
5: Compute  $m_{i,j}$  in Eq. (15).
6: if  $s_p - s_m \geq m_{i,j}$  and  $s_q - s_m \geq m_{i,j}$  then
7:   return  $\text{TRUE}$  and  $(s_p, s_q, s_k)$ .
8: else
9:   return  $\text{FALSE}$  and  $(\text{NULL}, \text{NULL}, \text{NULL})$ .
10: end if

```

in the cell, i.e.,

$$d_i = \max_{x,y \in C_M(s_i)} \|x - y\|_2 \quad (14)$$

Denote by $p_1(d_i)$ and $p_2(d_i)$ the two points in $C_M(s_i)$ satisfying $\|p_1(d_i) - p_2(d_i)\|_2 = d_i$.

Theorem 2. Let s_m, s_i, s_j be three generators in an $\text{RV T}(S_K, M)$. For the cells $C_M(s_m)$, $C_M(s_i)$ and $C_M(s_j)$, let m_m, m_i, m_j be their masses, s_m, s_i, s_j be their mass centroids, respectively. For any partitioning of $C_M(s_m)$ into two new cells $C(p_1(d_m))$ and $C(p_2(d_m))$, which satisfies $p_1(d_m) \in C(p_1(d_m))$, $p_2(d_m) \in C(p_2(d_m))$, $C(p_1(d_m)) \cap C(p_2(d_m)) = \emptyset$ and $C(p_1(d_m)) \cup C(p_2(d_m)) = C_M(s_m)$, let s_k, s_p and s_q be the mass centroids of $C_M(s_i)$, $C_M(s_j)$, $C(p_1(d_m))$ and $C(p_2(d_m))$, respectively. If $s_p - s_m \geq m_{i,j}$ and $s_q - s_m \geq m_{i,j}$, where

$$m_{i,j} = \frac{m_i m_j}{m_m(m_i + m_j)} \|s_i - s_j\|_2 \quad (15)$$

then the pair of operations $(S, M) : (s_m, (s_i, s_j)) \rightarrow ((s_p, s_q), s_k)$ does not increase the tessellation energy E .

By Theorem 2, we check the splitting-merging feasibility condition at line 8 of Algorithm 2 and this condition is summarized in Algorithm 3. Note that computing the diameter of an arbitrary region (line 2 of Algorithm 3) is time-consuming. In practice, we compute the axis-aligned bounding box B of $C_M(s_m)$. B is determined by two supporting points p_1 and p_2 in $C_M(s_m)$ and we use them as fast approximations to $p_1(d_m)$ and $p_2(d_m)$.

Algorithm 4 Randomly pick three generators

Input: An RV T(S_K, M) and an expected cell volume $E(V(C_M))$ (Eq.(18)).

Output: Three generators (s_m, s_i, s_j) in S_K .

```

1: Set  $S_{dense} = \emptyset$  and  $S_{sparse} = \emptyset$ .
2: for each cell  $C_M(s_i)$  in RV T( $S_K, M$ ) do
3:   Compute the volume  $V(C_M(s_i))$ .
4:   if  $V(C_M(s_i)) > 4E(V(C_M))$  then
5:      $S_{dense} = S_{dense} \cup \{s_i\}$ .
6:   else if  $V(C_M(s_i)) < E(V(C_M))/4$  then
7:      $S_{sparse} = S_{sparse} \cup \{s_i\}$ .
8:   end if
9: end for
10: if  $|S_{dense}| > 2$  then
11:   Randomly pick a generator  $s_m$  in  $S_{dense}$ .
12: else
13:   Randomly pick a generator  $s_m$  in  $S_K$ .
14: end if
15: Collect all neighboring pairs in  $S_{sparse}$  in the set  $N$ .
16: if  $|N| > 2$  then
17:   Randomly pick a pair  $(s_i, s_j)$  in  $N$ .
18: else
19:   Randomly pick two neighboring generators  $s_i$  and  $s_j$  in  $S_K$ .
20: end if
21: return ( $s_m, s_i, s_j$ ).
```

Corollary 1. Denote a voxel in the video as v_i and let

$$m_{\min} = \min_{v_i} \{V(C(v_i))\}, \quad (16)$$

$$d_{\max} = \max_{v_i} \{d(C(v_i))\}, \quad (17)$$

where $d(C(v_i))$ is the diameter of $C(v_i)$ in M . Let s_m, s_i, s_j be three generators in an RV T(S_K, M) and s_m be the mass centroid of $C_M(s_m)$. Let P be the hyperplane which passes through s_m and is perpendicular to the line connecting $p_1(d_m)$ and $p_2(d_m)$. P partitions $C_M(s_m)$ into $C(p_1(d_m))$ and $C(p_2(d_m))$. If $d_m = w(\frac{m_m}{m_{\min}}, m_{i,j} + d_{\max})$, where $m_m = V(C_M(s_m))$, $w = \max\{1 + \frac{1}{m_{i,j}}, \frac{1}{m_{i,j}}\}$, $= \frac{p_1(d_m) - s_m}{p_2(d_m) - s_m} \cdot 2$ and $m_{i,j}$ is defined in Eq.(15), then the pair of operations $(S, M) : (s_m, (s_i, s_j)) \rightarrow ((s_p, s_q), s_k)$ does not increase the tessellation energy E , where s_p, s_q and s_k are the mass centroids of $C(p_1(d_m))$, $C(p_2(d_m))$ and $C_M(s_i) \cup C_M(s_j)$, respectively.

Note that for a region with a fixed volume, the higher variation of colors in it, the larger the volume of it in M and vice versa. In Algorithm 2, s_i and s_j are chosen to be neighboring generators. Then by Corollary 1, Algorithm 2 has the following characteristics:

- the smaller the volumes of cells $C_M(s_i)$ and $C_M(s_j)$,

Algorithm 5 Streaming CSS generation

Input: A video of N voxels and the desired number of supervoxels K .

Output: K content-sensitive supervoxels.

```

1: Compute the discretized manifold representation  $M = \{(x_i, y_i, t_i, w_i)\}_{i=1}^N$ .
2: Initialize  $S = M$ .
3: while  $S$  cannot be loaded into main memory do
4:   Set  $S = \emptyset$ .
5:   Divide  $S$  into  $l$  disjoint pieces  $s_1, \dots, s_l$ , such that each piece can be loaded into main memory.
6:   for each piece  $s_i$  do
7:     Apply Algorithm 2 to compute 1.2K generators  $S_K(s_i)$ .
8:     Compute RV T( $S_K(s_i), s_i$ ).
9:     for each new generator  $g_j$  in  $S_K(s_i)$  do
10:      Compute the total weight of all points in the cell corresponding to  $g_j$  in RV T( $S_K(s_i), s_i$ ) and assign it to  $g_j$  as the weight  $w_j$ ;
11:    end for
12:     $S = S \cup (g_j, w_j)$ ,  $g_j \in S_K(s_i)$ .
13:  end for
14:   $S = S$ .
15: end while
16: Apply Algorithm 2 to  $S$  for obtaining  $K$  supervoxels.
```

the more likely they are merged, thus reducing the number of generators in content-sparse regions;

- the larger the volume of a cell $C_M(s_m)$, the more likely it is split, thus producing more generators in content-rich regions.

Accordingly, to increase the feasibility of the splitting-merging operation at line 8 of Algorithm 2, we estimate content-dense and content-sparse regions in RV T(S_K, M) and collect their corresponding generators into subsets S_{dense} and S_{sparse} in Algorithm 4. If S_{dense} and S_{sparse} contain sufficient generators, we randomly pick two neighboring generators in S_{sparse} to be merged and pick one generator in S_{dense} to be split; otherwise, we randomly pick three generators in S_K . To estimate the content density of cells, we compute the expected cell volume as the average of K cells over the total volume of video manifold M :

$$E(V(C_M)) = \frac{\sum_{v_i} V(C(v_i))}{K} \quad (18)$$

For each cell C_M in RV T(S_K, M), we compare its volume $V(C_M)$ with $E(V(C_M))$: (1) if $V(C_M) > 4E(V(C_M))$, we put the generator of this cell into S_{dense} , and (2) if $V(C_M) < E(V(C_M))/4$, we put the generator of this cell into S_{sparse} . Algorithm 4 summarizes the pseudo-code.

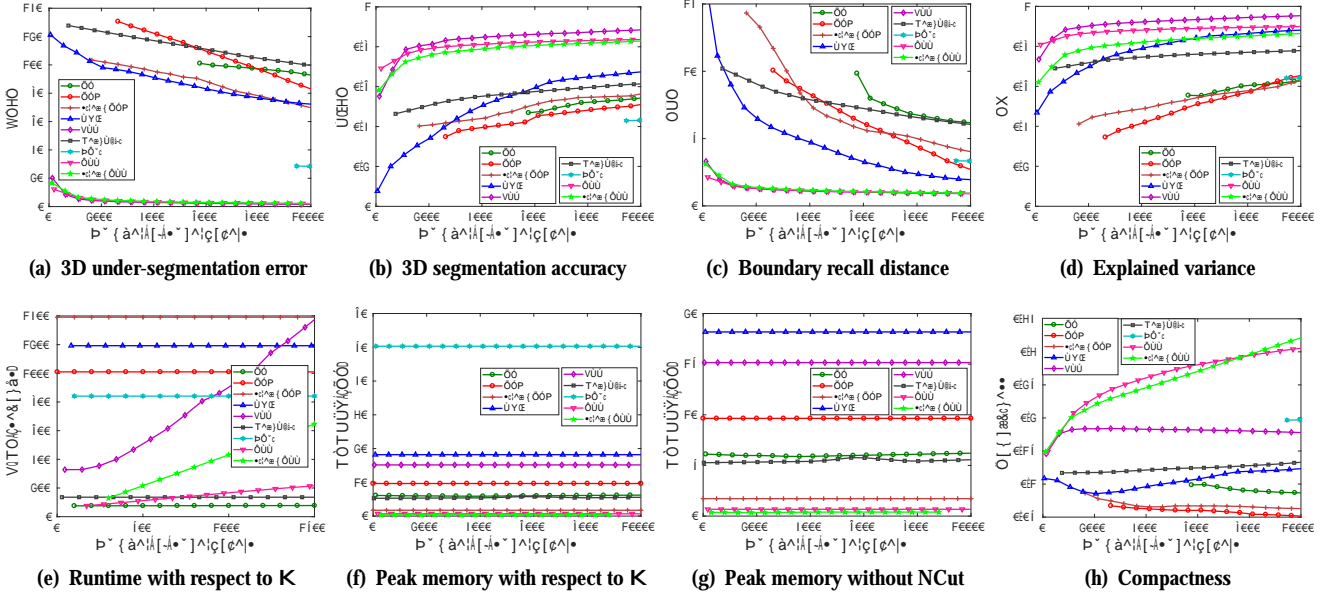


Figure 4. Evaluation of seven representative methods and our methods (CSS and streamCSS) on the BuffaloXiph dataset. Due to its high computational cost, NCut is run at a fixed frame resolution of 240×160 downsampled from original videos. CSS, streamCSS and TSP have the best performance in terms of UE3D, SA3D, BRD and EV. CSS and streamCSS is better than TSP in terms of compactness, time and space efficiency. Similar performances are observed on the other three video datasets (SegTrack v2, BVDS and CamVid), which are reported in supplemental material.

In all our experiments, we set $\text{iter}_{\max} = 20$ and $\text{num}_{\text{random}} = 20$ in Algorithm 2. We show in Section 5 that our algorithm can obtain high-quality CSS in 20 iterations. We have the following theoretical results.

Theorem 3. *By selecting $(1 + \epsilon)K$ generators, $\epsilon > 0$, Algorithm 2 is a bi-criteria $(1 + \epsilon, 1 + \frac{1+\epsilon}{2})$ -approximation algorithm in expectation.*

Theorem 4. *By selecting $(1 + \epsilon)K$ generators, $0 < \epsilon < 1$, the time and space complexities of Algorithm 2 are $O(NK)$ and $O(N + K)$, respectively.*

4.3. Streaming CSS algorithm for long videos

Thanks to the streaming K-means algorithm [1], Algorithm 2 is readily extended to a streaming version for handling long videos. The streaming CSS algorithm represents the video manifold \mathcal{M} by an ordered, discretized sequence of weighted points $\mathcal{M} = \{(x_i, y_i, t_i, w_i)\}_{i=1}^N$, where (x_i, y_i, t_i) is the position of voxel v_i in \mathcal{V} and $w_i = V((x_i, y_i, t_i))$. Pseudo-code is summarized in Algorithm 5.

Theorem 5. *If $(1 + \epsilon)K$ generators, $0 < \epsilon < 1$, are selected by Algorithm 2, Algorithm 5 is $(O(1), O(1))$ -approximation.*

5. Experiments

We implemented CSS (Algorithm 2) and streamCSS (Algorithm 5) in C++ and tested them on a PC with Intel

Core E5-2683V3 and 256GB RAM running Linux. Source code is available². We compare CSS and streamCSS with seven representative methods selected in [31], including NCut [27, 10, 9], SWA [25, 26, 6], MeanShift [24], GB [8], GBH [12], streamGBH [32] and TSP [4]. Since CSS and streamCSS adopt a random initialization, we report the average results of 20 initializations. The performances are evaluated on four video datasets, i.e., BuffaloXiph [5], SegTrack v2 [18], BVDS [29, 11] and CamVid [3], which have groundtruth labels drawn by human annotators.

In the original implementation of above seven methods collected in the LIBSVX benchmark [31], the clustering of supervoxels does not consider the connectivity of voxels. Therefore, the voxels that have the same supervoxel label (corresponding to a cluster) can have many disjoint components. In video applications such as detecting spatiotemporal closure for foreground object segmentation [15], the characteristic of multiple disjoint components in one label makes supervoxels' performance very bad. In our evaluation, for fairness we extract all the connected components in each supervoxel and relabel them. Therefore, more supervoxels are produced and accordingly the supervoxels' range in Figure 4 are adjusted and different from the ones in the LIBSVX benchmark.

Adherence to object boundaries. As perceptually

²<http://cg.cs.tsinghua.edu.cn/people/~Yongjin/Yongjin.htm>

meaningful atomic regions in videos, supervoxels should well preserve the object boundaries of ground-truth segmentation. 3D under-segmentation error (UE3D), 3D segmentation accuracy (SA3D) and boundary recall distance (BRD) are standard metrics in this aspect [4, 16, 31]. UE3D and SA3D are complementary to each other and both measure the tightness of supervoxels that overlap with ground-truth segmentation. BRD measures how well the groundtruth boundaries are successfully retrieved by the supervoxel boundaries. As shown in Figures 4(a)-(c), CSS, streamCSS and TSP have the highest SA3D, and the smallest UE3D and BRD, demonstrating their ability to adhere to object boundaries.

Explained variation (EV). EV is a standard metric that measures the color variations in supervoxels [22, 31]. A large EV means the color in each supervoxel is close to homogeneity. As shown in Figure 4(d), CSS, streamCSS and TSP have the largest EV.

Computational cost. We record runtime and peak memory of all nine methods. All methods are implemented in C or C++ except NCut (Matlab running with 8 threads) and TSP (Matlab with MEX). As shown in Figure 4(e), GB, CSS and MeanShift are three fastest methods. In particular, CSS is $5\times$ to $10\times$ faster than TSP. As shown in Figures 4(f)-(g), streamCSS and CSS are two methods that use smallest peak memory. In particular, the peak memory of streamCSS and CSS is $22\times$ smaller than TSP.

Three more metrics – mean size variation (MSV), temporal extent (TEX) and label consistency (LC) – are used in [31]. MSV and TEX measure the size variation and average temporal extent of all supervoxels in a video. Since our work advocates to adapt the size of supervoxels according to video content density, these two metrics are no longer suitable. Instead, we qualitatively compare the *content sensitivity* and propose to use the *compactness* measure below. LC is evaluated using groundtruth optical flow. However, as aforementioned, optical flow is only an optional preprocessing tool to video applications and may introduce extra error into supervoxel evaluation.

Content sensitivity. Figure 1 (last column) shows a typical result of CSS. More qualitative results are illustrated in supplemental material. By clipping supervoxels in image frames, these results clearly show that CSS well captures object boundaries in a video and the supervoxels are content sensitive, i.e., small in content-dense regions and large in content-sparse regions. The content sensitive feature is due to the characteristic that regions of high appearance and motion variance have large volumes in \mathbf{M} .

Compactness. In many real-world video applications, the solution relies on minimizing an energy function defined on a spatiotemporal supervoxel graph in a video clip. The shape regularity of supervoxels has a direct influence on the complexity of this spatiotemporal supervoxel graph,

and thus, affects the application performance. It was observed that compact supervoxels usually have better segmentation performance than non-compact ones [33]. Note that for any connected region R^3 , the isoperimetric inequality holds:

$$A(\cdot) \geq 3V(\cdot)^{\frac{2}{3}}V(B_1)^{\frac{1}{3}} \quad (19)$$

where B_1 is a unit sphere, $A(\cdot)$ and $V(\cdot)$ are bounding surface area and volume of \cdot respectively, and the equality holds when \cdot is a sphere. Therefore, for a given supervoxel over-segmentation $\tilde{S} = \{\tilde{s}_1, \tilde{s}_2, \dots, \tilde{s}_K\}$, the compactness metric C is defined as [33]

$$C(\tilde{S}) = \frac{1}{|\tilde{S}|} \sum_{\tilde{s}_i \in \tilde{S}} Q(\tilde{s}_i) \frac{|\tilde{s}_i|}{N}, \text{ where } Q(\tilde{s}_i) = \frac{6^{\frac{1}{2}} V(\tilde{s}_i)}{A(\tilde{s}_i)^{\frac{3}{2}}}, \quad (20)$$

$|\tilde{s}_i|$ is the number of voxels in \tilde{s}_i . The larger the compactness value is, the more regular the shape of supervoxels is. As shown in Figure 4(h), CSS and streamCSS have the largest compactness values.

6. Conclusion

In this paper, we propose a simple yet efficient algorithm which obtains content-sensitive supervoxels by computing RCVT – a uniform tessellation – on a video manifold \mathbf{M} . \mathbf{M} is constructed by mapping a video into a combined color and spatiotemporal space R^6 and the volume elements on \mathbf{M} reflect the density of video content. We propose a splitting-merging scheme and use it in the classic Lloyd method such that a high quality RCVT can be computed in very few iterations. Our algorithm is easily extended to a stream version for handling long videos. In addition to its easy implementation, our algorithm is theoretically an $(O(1), O(1))$ -approximation. Experimental results show that our method and TSP outperform other six representative methods (NCut, SWA, MeanShift, GB, GBH and streamGBH) in terms of metrics UE3D, SA3D, BRD and EV. Our method is better than TSP in terms of compactness and time and space efficiency.

Acknowledgment

This work was supported by the Natural Science Foundation of China (61725204, 61521002, U1736220, 61661130156), Royal Society-Newton Advanced Fellowship (NA150431) and Beijing National Research Center for Information Science and Technology.

References

- [1] N. Ailon, R. Jaiswal, and C. Monteleoni. Streaming k-means approximation. In *Proceedings of the 22nd International Conference on Neural Information Processing Systems*, NIPS '09, pages 10–18, 2009. 2, 4, 7

- [2] D. Arthur and S. Vassilvitskii. K-means++: The advantages of careful seeding. In *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '07, pages 1027–1035, 2007. 2, 4
- [3] G. J. Brostow, J. Shotton, J. Fauqueur, and R. Cipolla. Segmentation and recognition using structure from motion point clouds. In *European Conference on Computer Vision*, ECCV '08, pages 44–57. Springer, 2008. 7
- [4] J. Chang, D. Wei, and J. W. Fisher III. A video representation using temporal superpixels. In *Proceedings of the 2013 IEEE Conference on Computer Vision and Pattern Recognition*, CVPR '13, pages 2051–2058, 2013. 1, 2, 7, 8
- [5] A. Y. Chen and J. J. Corso. Propagating multi-class pixel labels throughout video frames. In *Proceedings of Western New York Image Processing Workshop*, 2010. 7
- [6] J. J. Corso, E. Sharon, S. Dube, S. El-Saden, U. Sinha, and A. L. Yuille. Efficient multilevel brain tumor segmentation with integrated Bayesian model classification. *IEEE Trans. Med. Imaging*, 27(5):629–640, 2008. 1, 2, 7
- [7] Q. Du, V. Faber, and M. Gunzburger. Centroidal Voronoi tessellations: Applications and algorithms. *SIAM Review*, 41(4):637–676, 1999. 5
- [8] P. F. Felzenszwalb and D. P. Huttenlocher. Efficient graph-based image segmentation. *International Journal of Computer Vision*, 59(2):167–181, 2004. 1, 2, 7
- [9] C. Fowlkes, S. Belongie, F. Chung, and J. Malik. Spectral grouping using the nyström method. *IEEE Trans. Pattern Anal. Mach. Intell.*, 26(2):214–225, 2004. 1, 7
- [10] C. C. Fowlkes, S. J. Belongie, and J. Malik. Efficient spatiotemporal grouping using the nyström method. In *2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, CVPR '01, pages 231–238, 2001. 7
- [11] F. Galasso, N. S. Nagaraja, T. J. Cárdenas, T. Brox, and B. Schiele. A unified video segmentation benchmark: Annotation, metrics and analysis. In *Proceedings of the 2013 IEEE International Conference on Computer Vision*, ICCV '13, pages 3527–3534, 2013. 7
- [12] M. Grundmann, V. Kwatra, M. Han, and I. A. Essa. Efficient hierarchical graph-based video segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, CVPR '10, pages 2141–2148, 2010. 1, 2, 7
- [13] S. D. Jain and K. Grauman. Supervoxel-consistent foreground propagation in video. In *13th European Conference on Computer Vision*, ECCV '14, pages 656–671, 2014. 1
- [14] Y. Ke, R. Sukthankar, and M. Hebert. Spatio-temporal shape and flow correlation for action recognition. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, CVPR'07, 2007. 1
- [15] A. Levinstein, C. Sminchisescu, and S. J. Dickinson. Optimal image and video closure by superpixel grouping. *International Journal of Computer Vision*, 100(1):99–119, 2012. 1, 7
- [16] A. Levinstein, A. Stere, K. N. Kutulakos, D. J. Fleet, S. J. Dickinson, and K. Siddiqi. Turbopixels: Fast superpixels using geometric flows. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 31(12):2290–2297, 2009. 8
- [17] C. Li, L. Lin, W. Zuo, S. Yan, and J. Tang. SOLD: sub-optimal low-rank decomposition for efficient video segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, CVPR '15, pages 5519–5527, 2015. 1
- [18] F. Li, T. Kim, A. Humayun, D. Tsai, and J. M. Rehg. Video segmentation by tracking many figure-ground segments. In *Proceedings of the IEEE International Conference on Computer Vision*, ICCV '13, pages 2192–2199, 2013. 7
- [19] Y.-J. Liu, C. Yu, M. Yu, and Y. He. Manifold SLIC: a fast method to compute content-sensitive superpixels. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2016)*, pages 651–659, 2016. 1, 2, 3
- [20] Y.-J. Liu, M. Yu, B.-J. Li, and Y. He. Intrinsic manifold SLIC: A simple and efficient method for computing content-sensitive superpixels. *IEEE Trans. Pattern Anal. Mach. Intell.*, 40(3):653–666, 2018. 1, 3
- [21] J. Lu, R. Xu, and J. J. Corso. Human action segmentation with hierarchical supervoxel consistency. In *IEEE Conference on Computer Vision and Pattern Recognition*, CVPR '15, pages 3762–3771, 2015. 1
- [22] A. P. Moore, S. Prince, J. Warrell, U. Mohammed, and G. Jones. Superpixel lattices. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, CVPR '08, 2008. 8
- [23] D. Oneata, J. Revaud, J. Verbeek, and C. Schmid. Spatio-temporal object detection proposals. In D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, editors, *13th European Conference on Computer Vision*, ECCV '14, pages 737–752, 2014. 1
- [24] S. Paris and F. Durand. A topological approach to hierarchical segmentation using mean shift. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, CVPR '07, 2007. 1, 2, 7
- [25] E. Sharon, A. Brandt, and R. Basri. Fast multiscale image segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, CVPR '00, pages 1070–1077, 2000. 2, 7
- [26] E. Sharon, M. Galun, D. Sharon, R. Basri, and A. Brandt. Hierarchy and adaptivity in segmenting visual scenes. *Nature*, 442(7104):810–813, 2006. 2, 7
- [27] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 22(8):888–905, 2000. 7
- [28] D. Sun, S. Roth, and M. J. Black. A quantitative analysis of current practices in optical flow estimation and the principles behind them. *International Journal of Computer Vision*, 106(2):115–137, 2014. 2
- [29] P. Sundberg, T. Brox, M. Maire, P. Arbeláez, and J. Malik. Occlusion boundary detection and figure/ground assignment from optical flow. In *IEEE Conference on Computer Vision and Pattern Recognition*, CVPR '11, pages 2233–2240. IEEE, 2011. 7
- [30] D. Wei. A constant-factor bi-criteria approximation guarantee for k-means++. In *Annual Conference on Neural Information Processing Systems (NIPS) 2016*, pages 604–612, 2016. 2, 4

- [31] C. Xu and J. J. Corso. Libsvx: A supervoxel library and benchmark for early video processing. *International Journal of Computer Vision*, 119(3):272–290, 2016. **1, 7, 8**
- [32] C. Xu, C. Xiong, and J. J. Corso. Streaming hierarchical video segmentation. In *Proceedings of the 12th European Conference on Computer Vision - Volume Part VI*, ECCV '12, pages 626–639, 2012. **1, 2, 7**
- [33] R. Yi, Y.-J. Liu, and Y.-K. Lai. Evaluation on the compactness of supervoxels. *Submitted for publication*, 2018. **8**
- [34] C.-P. Yu, H. Le, G. Zelinsky, and D. Samaras. Efficient video segmentation using parametric graph partitioning. In *IEEE International Conference on Computer Vision (ICCV 2015)*, pages 3155–3163, 2015. **1**