

## خلاصه سازی مقیاس پذیر ویدئویی با استفاده از انتخاب واژه نامه تنک

پوریا اعتضادی فر<sup>۱</sup>، حسن فرسی<sup>۲\*</sup>

۱- دانشجوی دکترا، ۲- استاد، دانشگاه بیرجند

(دریافت: ۹۴/۱۲/۰۸، پذیرش: ۹۵/۰۶/۲۸)

## چکیده

یکی از حوزه های مهم در پدافند غیرعامل، شناسایی تهدیدات و اعلام هشدار است. یکی از روش های پرکاربرد در حوزه شناسایی بررسی داده های ویدئویی به منظور شناسایی اهداف ناشناس و اعلام هشدار است. به منظور بررسی سریع، با دقت بالا روش های خلاصه سازی ویدئو ارائه شده است. همچنین در طول سال های گذشته، ایجاد ویدئو دیجیتالی منجر به رشد نمایی محتوای ویدئویی شده است. به منظور افزایش قابلیت استفاده از این حجم بالای ویدئو، تحقیقات بسیاری به انجام رسیده و خلاصه سازی ویدئو به جهت مرور سریع این مجموعه ویدئویی بزرگ و برای کمک به فهم سریع محتوای داده های ویدئویی پیشنهاد شده است. در خلاصه سازی ویدئو، تصاویری به عنوان نماینده از هر صحنه انتخاب می شود تا مروری تصویری از تمام فیلم به دست آید. اخیراً روش هایی با استفاده از فرمول بندی تنک برای خلاصه سازی ویدئو، داده های ویدئویی را به میزان زیادی نسبت به دیگر روش ها خلاصه نموده اند. در این مقاله به خلاصه سازی ویدئو، به عنوان یک مسئله انتخاب واژه نامه تنک پرداخته می شود. بدین منظور، با استفاده از روشی جدید بر پایه کدینگ تنک، می توان به میزان زیادی خلاصه سازی داده های ویدئویی را نسبت به دیگر روش های خلاصه سازی ویدئو که با روش تنک و یا روش های دیگر پیشنهاد شده اند بهبود بخشید. اساس این روش بر پایه حل معادله بهینه سازی با استفاده از آستانه گذاری نرم است که پیچیدگی کمتری نسبت به روش های پیشنهاد شده اخیر است. این امر را می توان با بررسی میزان پیچیدگی روش پیشنهادی با روش های متداول اخیر متوجه شد. در انتها نتایج آزمایش برای مجموعه داده های معیاری زمین حقیقت و با روش های State of the art، ادعای ما در بهبود میزان خلاصه سازی روش پیشنهادی را نشان می دهد.

**کلیدواژه ها:** خلاصه سازی ویدئو، انتخاب واژه نامه، فریم کلیدی، آنالیز ویدئو، کدینگ تنک

## Scalable Videos Summarization Using Sparse Dictionary Selection

P. Etezadifar, H. Farsi\*

University of Birjand

(Received: 27/02/2016; Accepted: 18/09/2016)

## Abstract

One of the important topics of passive defense is threats detection and warning alarm. One of the most widely used methods in detection field is video data investigation in order to identify unknown targets and warning alarm. In order to evaluate a fast and high-precision technique, video summarization is presented. Also, during the past years, creation of digital videos has caused exponential growth of video content. To increase the high volume of video usability, a lot of researches have been done and video summarization has been proposed to quick view of large video collection and quick understanding of the content of video data. In the video summarization, pictures are selected as a representative of each scene to obtain a visual overview of whole video. Recently, new methods using sparse formulation are suggested for video summarization being more effective in video data summarization than other methods. In this paper, video summarization is presented as a sparse dictionary selection problem. For this purpose, using a new method based on sparse coding, have been able to improve video data summarization compared to other video summarization methods based on sparse or other coding. Finally, the results for the ground truth data collection and State of the art methods, shows improvement our claim in the video summary on proposed method.

**Key words:** Video Summarization, Dictionary Selection, Key Frame, Video Analysis, Sparse Coding.

\*Corresponding Author E-mail: hfarsi@birjand.ac.ir

## ۱. مقدمه

همان طور که می‌دانید، یکی از حوزه‌های مهم در پدافند غیرعامل، شناسایی تهدیدات و اعلام خبر است. به منظور پیشگیری از حملات دشمنان می‌بایست هر لحظه اقدامات احتمالی مورد بررسی و جستجو قرار گیرد. یکی از این حوزه‌ها رصد هوایی با استفاده از سامانه‌های هوشمند پردازشی است. با توجه به احتمال تجاوز در هر لحظه از زمان به وسیله هواپیماهای رادار گریز و موشک‌های کروز که اصولاً به دلیل ارتفاع پایین پروازی قادر به ردیابی توسط سامانه‌های راداری نمی‌باشند؛ می‌بایست از روش‌های تصویربرداری برخط (ویدیو) استفاده کرد. به دلیل حجم بالای داده‌های به دست آمده توسط دوربین‌های ویدئویی نمی‌توان اطلاعات آن‌ها را به صورت برخط پردازش کرده و سامانه قادر به اطلاع‌رسانی در زمان خبر نیست. یکی از روش‌هایی که اخیراً مطرح شده است، خلاصه‌سازی ویدیو است که با استفاده از آن بتوان نسبت داده به اطلاعات را افزایش داده و خروجی به دست آمده را بتوان به عنوان داده‌ای با میزان اطلاعات بالا به سامانه تصمیم‌گیر وارد کرده و با استفاده از آن بتوان به وجود و یا عدم وجود اشیای متخاصم هوایی پی برد. همچنین از روش خلاصه‌سازی ویدیو در حملات آفندی به منظور عدم آشکارسازی به وسیله سنسورهای راداری استفاده می‌شود [۳-۱].

به عنوان مثال در موشک‌های کروز و هواپیماهای بدون سرنشین به منظور مسیریابی از سامانه‌های ویدئویی استفاده می‌شود که اساس کار آن‌ها با استفاده از روش‌های خلاصه‌سازی ویدیو است. بنابراین، یکی از روش‌ها به منظور مقابله با یک سامانه آفندی روش‌های مورد استفاده به منظور تعیین هدف در فاز نهایی موشک کروز استفاده از ویژگی‌های بصری است. بدین صورت که در فاز انتهایی، موشک کروز شروع به جست‌وجو در فریم‌های ویدئویی اخذشده از سطح زمین می‌کند و پس از تطبیق منطقه مورد نظر با تصاویر موجود در حافظه‌اش، به منطقه هدف اصابت می‌کند [۴-۶]. به دلیل نرخ بالای تصاویر به دست آمده از دوربین تعبیه شد، در موشک کروز، نمی‌توان تمام فریم‌ها را مورد بررسی قرار داد. یکی از راه‌های کاهش نرخ فریم‌های جست‌وجو استفاده از روش‌های خلاصه‌سازی ویدیو است. در این مقاله به بهبود روش خلاصه‌سازی پرداخته شده است. همچنین، می‌توان در زمینه‌های صنعتی و تجاری نیز استفاده کرد. باید به این نکته اشاره نمود که خلاصه‌سازی ویدیو به یک ابزار کلیدی، برای مرور و دسترسی کارآمد و نیز برای دست‌کاری مجموعه بزرگ ویدئویی، تبدیل شده و توجه بیشمار محققین را به خود جذب کرده است. بیشتر

فعالیت‌هایی که در حوزه خلاصه‌سازی سیگنال ویدئویی متمرکز شده‌اند، بر پایه ویژگی‌های گوناگونی مثل حرکت<sup>۱</sup> [۷]، شنیداری<sup>۲</sup> [۸] و یا روش‌های ترکیبی<sup>۳</sup> [۹] است. مقاله مروری در مورد سامانه اتیک از نمونه عملیات نیز ارائه شده اند [۱۰]. در مدیریت مجموعه ویدئویی، ویدیوهای مشابه و یا مرتبط، معمولاً مطابق با محتوای آن‌ها به فصل‌هایی گروه‌بندی می‌شوند. بنابراین، کاربران این ویدیوها، به دیگر روابط بین گروه‌هایی از این فایل‌های ویدئویی علاقه‌مند می‌شوند. برای مثال، در آرشیوهای اخبار ویدئویی، تمام کلیپ‌هایی که در مورد یک موضوع از کانال‌های مختلف و جلسات خبری پخش می‌شود، هر دو با هم تکامل موضوع را نشان می‌دهند. از طرفی با پیشرفت فناوری در زمینه تصویربرداری و افزایش اطلاعات ویدئویی ساخته شده، افزایش تقاضا مردم برای ذخیره‌سازی، به اشتراک گذاشتن اطلاعات ویدئویی در اینترنت و شبکه‌های اجتماعی، بازیابی ویدیو<sup>۴</sup> [۱۱]، مراقبت داخلی<sup>۵</sup> [۱۲]، آشکارسازی ناهنجاری<sup>۶</sup> [۱۳] و ... بیش از پیش، مورد توجه قرار گرفته است. در این مورد خلاصه‌سازی ویدیو با حذف داده‌های حشو، با استفاده از فشرده‌سازی و چگال کردن آن‌ها، کمک بسیاری به افزایش سرعت برنامه‌ها، کاهش فضای ذخیره‌سازی و سادگی ارسال از طریق اینترنت می‌شود. روش‌های خلاصه‌سازی ویدئویی را می‌توان به دو صورت کلی خلاصه‌سازی استاتیکی<sup>۷</sup> و دینامیکی<sup>۸</sup> مجزا نمود [۱۴]. خلاصه‌سازی استاتیکی ویدیو بر اساس پیدا کردن فریم‌های کلیدی کار می‌کند. منظور از فریم‌های کلیدی، تصاویری ثابت از ویدیو هستند که مهم‌ترین محتوای آن را شامل می‌شوند [۱۵]. ولی خلاصه دینامیکی ویدیو، شامل دنباله‌ای از شات‌های کوچک است که پشت سر هم قرار گرفته‌اند. با این وجود، استخراج فریم کلیدی، انعطاف بیشتری برای نشان دادن محتوای ویدیو دارد. به طور معمول، استخراج فریم کلیدی به دو صورت مبتنی بر متن و محتوا انجام می‌گیرد [۱۶]. در روش مبتنی بر متن، کلمه و یا جمله‌ای دلخواه انتخاب می‌شود و فریمی کلیدی متناسب با آن متن استخراج می‌شود. این روش زمان‌بر و پرهزینه است. روش مبتنی بر محتوا، بر اساس ویژگی‌های استخراج شده از ویدیو صورت می‌پذیرد. این روش نسبت به روش مبتنی بر متن ارزان‌تر و سریع‌تر است ولی به دلیل تفاوت ادراک انسان و ماشین، زیاد قابل اطمینان نیست.

<sup>1</sup> Motion

<sup>2</sup> Audio

<sup>3</sup> Multi-Modality

<sup>4</sup> Video Retrieval

<sup>5</sup> Indoor Surveillance

<sup>6</sup> Anomaly Detection

<sup>7</sup> Static

<sup>8</sup> Dynamic

دهد. بنابراین یک صحنه نسبتاً مشابه که در زمان‌های متفاوت نمایش داده می‌شود به عنوان صحنه‌ای جدید در نظر گرفته شده و فریم‌های کلیدی تکراری تولید می‌گردد. در روشی دیگر که برای خلاصه‌سازی محتوایی ویدیو استریو پیشنهاد شده در ابتدا ویدیو را به شات‌هایی تقسیم‌بندی می‌کند. سپس فریم‌هایی که ویژگی بصری شبیه به یک دیگر دارند با یک دیگر جمع می‌شوند. سپس ناهمخوانی زمینه<sup>۲</sup>، مناطق مسدود<sup>۳</sup> و نقشه عمق<sup>۴</sup> تخمین زده می‌شود. در گام بعد، بردار ویژگی با استفاده از کلاس‌بندی فازی ویژگی‌ها ساخته شده و انتخاب شات با استفاده از کلاس‌بندی شات‌های مشابه صورت می‌گیرد. در انتها فریم‌های کلیدی با استفاده از روش بهینه‌سازی برای مکان فریم‌ها با پیدا کردن کمترین شباهت بین فریمی استخراج می‌شوند [۲۱]. این روش نیز مشکلاتی نزدیک به آنچه برای روش پن گفته شد را دارا است. ساتو و لی [۲۲] الگوریتم دسته‌بندی برای گروه‌بندی فریم‌ها با رنگ‌های مشابه ارائه دادند. در این روش، دنباله ویدئویی به قطعه‌هایی تقسیم می‌شود که این کار با گروه‌بندی فریم‌های مشابه بر اساس ممان‌های رنگ انجام می‌شود. سپس، به دسته‌بندی این قطعه‌ها می‌پردازد. پس از آن، قطعات پشت سرهم که متعلق به یک دسته‌اند با هم ترکیب می‌شوند و قطعات مجاور به منظور کاهش افزونگی، ادغام می‌شوند. در نهایت، زیرمجموعه‌ای از فریم‌ها از طولانی‌ترین قسمت‌ها برای خلاصه ویدئویی سخت، انتخاب می‌شوند. در روشی دیگر، با استفاده از انتخاب واژه‌نامه تنک به استخراج فریم‌های کلیدی پرداخته است، البته این روش به جای استفاده از نرم ۱، نرم ۲ و ۱ را برای شرط تنکی خود استفاده کرده و با استفاده از روش نسترو<sup>۵</sup> [۲۳] به انتخاب واژه‌نامه تنک برای مسئله محدب ولی غیر مسطح<sup>۶</sup> پرداخته است [۲۴]. این روش دارای پیچیدگی نسبتاً بالایی است که از معایب آن محسوب می‌شود. روشی دیگر، با استفاده از محاسبه میزان تفاوت بین هیستوگرام رنگ دو فریم متوالی فریم کلیدی استخراج شده است [۲۵]. این روش بسیار نسبت به تغییرات شدت روشنایی حساس بوده و صحنه‌های شبیه به یک دیگر که دارای تغییرات شدت روشنایی هستند را به عنوان صحنه‌های جدید انتخاب کرده و فریم‌ها کلیدی شبیه به یک دیگر انتخاب می‌کند. فو و همکارانش [۲۶]، روشی برای خلاصه‌سازی ویدیوهای چند نمایشی<sup>۷</sup> پیشنهاد کرده‌اند. در این

بنابراین، محققین برای کاهش این تفاوت، از توصیف‌گرهای مختلفی استفاده می‌کنند. در ادامه ساختار کلی مقاله به صورت خلاصه مورد بررسی قرار می‌گیرد. در اولین گام به فرموله کردن مسئله خلاصه‌سازی ویدیو با استفاده از انتخاب واژه‌نامه تنک پرداخته می‌شود. در مرحله دوم به حل مسئله محاسبه شده برای خلاصه‌سازی ویدیو، در مرحله سوم به چگونگی پیاده‌سازی حل مسئله و الگوریتم آن پرداخته می‌شود. در مرحله بعد، به بررسی ویژگی‌هایی که برای آموزش واژه‌نامه<sup>۱</sup> استفاده می‌شود پرداخته شده و در انتها نتیجه روش پیشنهادی با دیگر روش‌های پیشنهاد شده در خلاصه‌سازی ویدیو مقایسه می‌شود. در ادامه به بررسی کارهای صورت گرفته در این زمینه توسط محققین پرداخته می‌شود.

در زمینه خلاصه‌سازی ویدیو، کارهای زیادی انجام شده است [۱۷ و ۱۸] که مجال بررسی همه آن‌ها نیست. اساس روش‌های پیشنهاد شده بدین صورت است که ویژگی خاصی را برای مجموعه فریم‌های ویدئویی در نظر گرفته و فریم‌هایی را میزان تغییرات این ویژگی نسبت به فریم قبل خود بیشتر بوده به عنوان فریم کلیدی استخراج می‌شوند. در ادامه به بررسی چند مقاله معروف در این زمینه پرداخته می‌شود. کلبان [۱۹] روشی برای مقابله با مشکلات خلاصه‌سازی ویدیو معرفی کرد که این روش بر اساس به کارگیری ویژگی‌های تلفیقی سطح بالا بود. در روش پیشنهادی کلبان، در ابتدا نمونه‌برداری به منظور حذف داده‌های اضافی، انجام می‌شود و تنها زیرمجموعه‌ای از فریم‌های ویدئویی باقی می‌ماند. سپس پنج ویژگی سطح بالا از فریم‌های باقی‌مانده استخراج می‌شود. در مرحله بعد، وزن‌های بهینه‌ای برای ترکیب این ویژگی‌ها، توسط روش گرادیان نزولی استنتاج می‌شوند. در نهایت، از الگوریتم وزنی k-means برای تشخیص مهم‌ترین قسمت‌های یک خلاصه‌سازی نهایی، استفاده شده است. در این مقاله، از روش دسته‌بندی ضعیفی به منظور انتخاب فریم کلیدی استفاده شده است که ضعف این روش محسوب می‌گردد روش پیشنهادی توسط پن [۲۰] روشی برای خلاصه‌سازی ویدیو بر اساس الگوریتم آشکارسازی مرز شات است. در ابتدا، مرز شات توسط مقایسه هیستوگرام رنگ فریم‌های متوالی، به دست می‌آید. پس از آن، با استفاده از گروه‌بندی شات‌های مشابه و انتخاب یک شات به عنوان نماینده از هر گروه و حذف شات‌های تکراری عمل می‌کند. در انتها، قسمت‌های انتخاب شده بر اساس اهمیت آن‌ها، خلاصه ویدیو را می‌سازند. در این روش به دلیل استفاده از مرز شات نمی‌توان به طور کلی داده ویدئویی را مورد بررسی قرار

<sup>2</sup> Disparity Field<sup>3</sup> Occluded Areas<sup>4</sup> Depth Map<sup>5</sup> Nesterov<sup>6</sup> Non-Smooth<sup>7</sup> Multi-View Videos<sup>1</sup> Dictionary

در این مسئله، فرض می‌شود مشاهده‌ای از بردار  $x$  به صورت زیر در دسترس باشد:

$$y = Dx + e \quad \& \quad \begin{cases} x \in \mathcal{R}^n \\ y \in \mathcal{R}^m \\ D \in \mathcal{R}^{m \times n} \end{cases} \quad (1)$$

در معادله (۱) منظور از  $y$ ، مشاهدات و منظور از  $e$ ، نویز است. همچنین،  $D$  به عنوان واژه‌نامه و معلوم بوده و  $x$  بردار مجهول است. همچنین، توزیع پیشین هر یک از درایه‌های  $x$  در این مسئله یک توزیع لاپلاس با میانگین صفر است و درایه‌ها مستقل از هم<sup>۳</sup> (i.i.d) می‌باشند. دلیل استفاده از توزیع لاپلاس به این لحاظ است که این توزیع حول مقادیر صفر تیز بوده و به جواب مسئله که باید جوابی تنک باشد نزدیک‌تر است تا بقیه توزیع‌ها از قبیل گوسی و ... همچنین نویز دارای توزیع گوسی با میانگین صفر و واریانس  $\sigma^2 I$  است. در این مسئله چون توزیع پیشین متغیر مجهول  $x$  در دسترس است بهترین تخمین MAP است [۳۲]. بنابراین برای تخمین  $x$  باید معادله (۲) حل گردد.

$$\hat{x}_{map} = \arg \max_x \{Ln(P(y|x)) + Ln(P(x))\} \quad (2)$$

با توجه به فرضیات مسئله می‌دانید که توزیع  $y|x$ ، همان توزیع نویز با میانگین  $Dx$  است که برابر با  $N(Dx, \sigma^2 I)$  است. همچنین می‌دانید اگر  $e \in \mathcal{R}^m$ ،  $e \sim N(\mu, \Sigma)$ ، آنگاه دارید [۳۴]:

$$P(y|x) = \frac{1}{(2\pi)^{m/2} \sigma^m} e^{-\frac{\|y-Dx\|_2^2}{2\sigma^2}} \quad (3)$$

با استفاده از معادله‌های (۲) و (۳) و با حذف مقادیر ثابت، با توجه به اینکه مقادیر ثابت در بهینه‌سازی تأثیری ندارد، دارید:

$$\hat{x}_{map} = \arg \min_x \left\{ \|y - Dx\|_2^2 + \lambda \|x\|_1 \right\} \quad (4)$$

می‌توان معادله (۴) را با استفاده از مدل‌سازی واژه‌نامه استفاده‌شده در Group Lasso [۳۳ و ۳۴]، به صورت زیر بازنویسی کرد:

$$\hat{X} = \arg \min_x \left\{ \|Y - DX\|_F^2 + \lambda \|X\|_1 \right\} \quad (5)$$

در معادله (۵)، منظور از  $X$ ، ماتریس ضرایب تعقیب است، نرم فروبنیوس  $\|X\|_F$ ، به صورت  $\left( \sum_{i,j} X_{ij}^2 \right)^{1/2}$  تعریف می‌شود.

مقاله در گام اول گراف شات فضا-زمانی<sup>۱</sup> ساخته می‌شود و با فرموله کردن عملیات برچسب‌گذاری گراف به خلاصه‌سازی ویدیو می‌پردازد. در روشی دیگر ژانگ و همکارانش [۲۷]، ابتدا کل فضای ویدیو کلاس‌بندی کرده و پس از آن فریم‌هایی که در مرکز کلاس قرار دارند را به عنوان فریم‌های کلیدی انتخاب کرده؛ ولی اگر فریمی در مرکز خوشه قرار نداشته باشد، نزدیک‌ترین فریم به مرکز به عنوان فریم کلیدی استخراج می‌شود. در این روش، کلاس‌بندی استفاده‌شده ضعیف بوده و به خوبی قادر به مجزا نمودن صحنه‌ها از یک دیگر ندارد. در روشی دیگر، با استفاده از محاسبه آنتروپی توأم و اطلاعات مشترک بین فریم‌های متوالی به آشکارسازی محوشدگی و برش شات می‌پردازد و سپس در هر شات فریم‌های کلیدی را با استفاده از ویژگی‌های آنتروپی محاسبه می‌کند [۲۸]. این روش نیز شبیه به روش پن قادر به بررسی کلی فیلم ناست و نمی‌تواند به درستی صحنه‌های مشابه را حذف نماید.

## ۲. الگوریتم پیشنهادی

در این مقاله، هدف خلاصه‌سازی ویدیو با استفاده از انتخاب واژه‌نامه تنک است. در ادامه به صورت خلاصه فرضیات اولیه استفاده‌شده در این مقاله را توضیح داده خواهد شد و پس از آن به فرموله کردن مسئله پرداخته می‌شود. در این مقاله بردار با حرف کوچک و به صورت  $x$  و همچنین ماتریس با حرف بزرگ و به صورت  $X$  نشان داده می‌شود. همچنین فرض شده است که سایز پایگاه داده ویدئویی  $n$  است به این معنا که در این پایگاه داده  $n$  فریم ویدئویی وجود دارد ( $D \in \mathcal{R}^{d \times n}$ ). همچنین، پایگاه داده به صورت  $D = \{\underline{d}_1, \underline{d}_2, \dots, \underline{d}_n\}$  نمایش داده می‌شود. در این پایگاه داده فریم‌های دوبعدی به بردارهای ویژگی یک‌بعدی تبدیل شده‌اند. بنابراین،  $\underline{d}_i \in \mathcal{R}^d$ ، نشان‌دهنده  $i$  امین تصویر از پایگاه داده با ابعاد  $d \times 1$  است. برای این کار می‌توان از روش کیف کلمات<sup>۲</sup> Gist [۲۹]، و CENTRIST [۳۰] و [۳۱] استفاده کرد و یا اینکه پیکسل‌ها را به صورت برداری زیر هم قرار داد. هدف این مقاله استخراج فریم‌های کلیدی منظور از  $\mathfrak{F}_{key}(D) = \{d_1, d_2, \dots, d_k\}$  به طول  $k$  است. همچنین  $D_{key}$  پایگاه داده فریم‌های کلیدی استخراج‌شده به صورت  $D_k \in \mathcal{R}^{d \times k}$  است که باید تا حد امکان مقدار  $k = n$  باشد. این شرط باعث ارضای تنک بودن فریم‌های کلیدی می‌شود.

<sup>3</sup> Independent and Identically Distributed

<sup>1</sup> Spatio-Temporal Shot Graph

<sup>2</sup> Bag of Words (BOW)

بنابراین به جای حل معادله (۷) به حل معادله زیر پرداخته می‌شود:

$$\hat{X} = \arg \min_x \left\{ \begin{array}{l} \|Y - DX\|_F^2 + \lambda \|X\|_{2,1} \\ + \frac{c}{2} \|\underline{x} - \underline{x}_0\|_2^2 \\ + \frac{1}{2} \|D\underline{x} - D\underline{x}_0\|_2^2 \end{array} \right\} \quad (۹)$$

در ادامه، با محاسبه  $\nabla \tilde{f}(\underline{x}, \underline{x}_0) = 0$ ، مقدار  $\underline{x}$  را به دست می‌آید.

$$\begin{aligned} \nabla \tilde{f}(\underline{x}, \underline{x}_0) &= -D^T(Y - DX) + \\ &\lambda \operatorname{sgn}(\underline{x}) + C(\underline{x} - \underline{x}_0) - \\ D^T(D\underline{x} - D\underline{x}_0) &= 0 \\ \Rightarrow -D^T Y + D^T DX + \lambda \operatorname{sgn}(\underline{x}) &+ \\ C(\underline{x} - \underline{x}_0) - D^T D\underline{x} + D^T D\underline{x}_0 &= 0 \\ \Rightarrow \underline{x} &= \left( \frac{1}{c} D^T(Y - DX) + \underline{x}_0 \right) + \\ \frac{1}{c} \operatorname{sgn}(\underline{x}) \end{aligned} \quad (۱۰)$$

با توجه به [۳۶ و ۳۷] می‌توان مقدار  $\underline{x}$  را به صورت زیر به دست آورد:

$$\underline{x} = \operatorname{soft}_{\frac{1}{c}} \left[ \frac{1}{c} D^T(Y - DX) + \underline{x}_0 \right] = \begin{cases} \frac{1}{c} D^T(Y - DX) + \underline{x}_0 - \frac{1}{c} & ; \frac{1}{c} < \frac{1}{c} D^T(Y - DX) + \underline{x}_0 \\ 0 & ; \frac{1}{c} \geq \left| \frac{1}{c} D^T(Y - DX) + \underline{x}_0 \right| \\ \frac{1}{c} D^T(Y - DX) + \underline{x}_0 + \frac{1}{c} & ; \frac{1}{c} > \frac{1}{c} D^T(Y - DX) + \underline{x}_0 \end{cases} \quad (۱۱)$$

در معادله (۱۱)،  $C$  مقداری ثابت است که در این آزمایش برابر با  $\lambda(D^T D)$  در نظر گرفته شده است. باید به این نکته اشاره کرد که برای همگرا شدن الگوریتم آستانه‌گذاری نرم باید  $C \geq \lambda(D^T D)$  باشد [۳۶]. بنابراین مقدار هر سطر از ماتریس  $X$  با استفاده از رابطه (۱۱) محاسبه می‌شود.

در الگوریتم پیشنهادی، شرط همگرایی کوچک بودن مقدار نرم ۲ سطرهای ماتریس  $X$  است. پس از اجرای الگوریتم پیشنهادی، نرم ۲ سطرهای ماتریس  $X$  محاسبه شده و طبقه‌بندی<sup>۴</sup> می‌شوند. در مرحله بعد مقادیر بیشتر به عنوان فریم‌های کلیدی ویدیو مورد نظر استخراج می‌گردند. در این مقاله با انجام هم‌زمان آموزش و انتخاب تنک واژه‌نامه به خروجی بهتری نسبت به روش‌های دیگر خلاصه‌سازی ویدیو دست پیدا کرده است. شبه کد الگوریتم پیشنهادی در ادامه نشان داده شده است.

همچنین، نرم ۱ برای ماتریس به صورت  $\|X\|_1 = \sum_{i,j} |X_{ij}|$  محاسبه می‌شود. جواب معادله (۵) در صورت کوچک بودن مقدار  $\lambda$  به صورت حدی به ماتریس همانی  $I$  میل می‌کند. این بدان معناست که حل معادله به سمتی پیش می‌رود تا فریم‌های کلیدی برابر با کل فریم‌های ویدیوی ورودی شوند، اما با افزایش مقدار  $\lambda$ ، میزان عناصر صفر ماتریس  $X$  افزایش پیدا می‌کند و معادله (۵) به سمتی می‌رود که تمامی مقادیر  $X$  به سمت صفر میل کنند. این بدان معناست که فریم‌های کلیدی استخراج شده را به کمترین تعداد فریم‌ها میل می‌دهد. در این مقاله مقادیر واژه‌نامه  $D \in \mathfrak{R}^{d \times k}$  با  $D_{key} \in \mathfrak{R}^{d \times k}$  که واژه‌نامه آموزشی نامیده می‌شود، با  $Y \in \mathfrak{R}^{d \times n}$  که همان واژه‌نامه اولیه ساخته شده با استفاده از ویژگی‌ها و همچنین  $X \in \mathfrak{R}^{k \times n}$  ماتریس ضرایب پیگیری<sup>۱</sup> می‌باشند. قسمت دوم معادله (۵) که تنک بودن ماتریس  $X$  را با استفاده از نرم ۱ محاسبه می‌کند؛ به دنبال آن است تا سطرهای زیادی از ماتریس  $X$  را تمام صفر کند و با استفاده از آن سطرهای غیر صفر را به عنوان فریم‌های کلیدی انتخاب کند. بنابراین می‌توان به جای استفاده از نرم ۱ از نرم ۲ و ۱ استفاده کرد [۳۵] که به صورت زیر تعریف می‌شود:

$$\|X\|_{2,1} = \sum_{i=1}^k \left( \sum_{j=1}^n x_{ij}^2 \right)^{\frac{1}{2}} \quad (۶)$$

با توجه به توضیحات قبلی منظور از  $\|\underline{x}_i\|_2$  نرم ۲، سطر  $i$ ام ماتریس  $X$  است. بنابراین با جایگزینی مقدار نرم ۲ و ۱ با نرم ۱ معادله (۵) به صورت زیر بازنویسی می‌شود:

$$\hat{X} = \arg \min_x \left\{ \|Y - DX\|_F^2 + \lambda \|X\|_{2,1} \right\} \quad (۷)$$

در این مقاله با استفاده از الگوریتم آستانه‌گذاری نرم<sup>۲</sup> از خانواده الگوریتم‌های انقباضی<sup>۳</sup> به حل این موضوع پرداخته می‌شود. الگوریتم آستانه‌گذاری نرم به حل معادله (۷) می‌پردازد. می‌توان معادله (۷) را به صورت زیر به منظور حذف ضریب  $D$  از  $X$  نوشت:

$$\tilde{f}(\underline{x}, \underline{x}_0) = f(\underline{x}) + d(\underline{x}, \underline{x}_0) \quad (۸)$$

$$\text{if} \begin{cases} f(\underline{x}) = \|Y - DX\|_F^2 \\ \quad + \lambda \|X\|_{2,1} \\ d(\underline{x}, \underline{x}_0) = \frac{c}{2} \|\underline{x} - \underline{x}_0\|_2^2 \\ \quad + \frac{1}{2} \|D\underline{x} - D\underline{x}_0\|_2^2 \end{cases}$$

<sup>۱</sup> Pursuit Coefficient Matrix

<sup>۲</sup> Soft Thresholding

<sup>۳</sup> Shrinkage Algorithms

<sup>۴</sup> Sort

در این مقاله از دو سطح مکانی استفاده شده است. در یک سطح ۵ تکه<sup>۴</sup> از فریم ویدئویی و در سطح دیگر ۱ تکه فریم ویدئویی در نظر گرفته شده است. در این روش به هر تکه از تصاویر ۴۲ ویژگی نسبت داده می‌شود که ۴۰ ویژگی مربوط به مقادیر ویژه هر تکه است و ۲ ویژگی برای محاسبه میانگین و واریانس مقادیر ویژه مورد استفاده قرار می‌گیرد. تا این مرحله برای هر فریم ویدئویی ۴۲×(۱+۵) = ۲۵۲ ویژگی استخراج شد. در مرحله بعد با توجه به اینکه روش CENTRIST مستقل از ویژگی‌های رنگ است،

در برخی گزارش‌ها [۳۶] از سه ممان مرکزی<sup>۵</sup> برای محاسبه توزیع احتمال رنگ در هر فریم استفاده شده است. این سه ممان برابر با میانگین، واریانس و کجی<sup>۶</sup> هستند در این مقاله علاوه بر سه ممان مرکزی که در بالا گفته شد از ممان درجه ۴ نیز استفاده می‌شود. دلیل استفاده از این ممان آن است که سه ممان قبلی قادر به تشخیص میزان کشیدگی<sup>۷</sup> در توزیع ویژگی‌های استخراج شده ندارند. بنابراین به سه ممان گفته شده در [۴۳] ممان درجه ۴ را نیز اضافه می‌شود. می‌توان میزان کشیدگی را با استفاده از معادله (۱۲) محاسبه کرد.

لازم به ذکر است که معادله (۱۲) بر روی فضای رنگ سه‌بعدی HSI<sup>۸</sup> محاسبه می‌شود. این سامانه رنگی، اطلاعات رنگ C را از اطلاعات شدت روشنایی آن جدا می‌کند. در این سامانه رنگ، اطلاعات رنگ در قالب پرده رنگ<sup>۹</sup> و اشباع<sup>۱۰</sup> ارائه می‌شود در حالی که روشنایی<sup>۱۱</sup>، مقدار نور تصویر را بیان می‌کند. این سامانه رنگی توانایی بالایی در ارائه درک رنگ توسط انسان را دارد. زیرا سامانه بینایی انسان به راحتی قادر به تشخیص پرده‌های رنگ متفاوت، بدون توجه به مقدار روشنایی و اشباع است. فضای رنگ HSI با استفاده از روابط زیر از فضای رنگ RGB به دست می‌آید [۴۳].

برای محاسبه ویژگی، فریم ویدئویی را به (۳×۴) قسمت تقسیم می‌شود. بردار ویژگی ساخته شد با استفاده از ویژگی‌های رنگی و آماری برابر با ۱۴۴ = (۳×۴)<sup>۲</sup> می‌شود. بردار ویژگی جدیدی با استفاده از CENTRIST و ویژگی‌هایی رنگی و آماری ساخته می‌شود که این مقدار برابر با ۳۹۶ ویژگی می‌شود. در ادامه به ارزیابی روش پیشنهادی پرداخته می‌شود:

الگوریتم پیشنهادی: انتخاب واژه‌نامه تنک

Input: The Video Set  $D \in \mathfrak{R}^{d \times n}$ ,  $\lambda, c, x_0 = 0$

Initialization:  $X = X^{(0)}$

Output: The Summarized Frames, Which is

Used from  $X \in \mathfrak{R}^{k \times n}$

repeat

for  $j = 1, 2, \dots, k$  do

$$\underline{x} = \text{soft}_{\%c} \left[ \%c D^T (Y - DX) + \underline{x}_0 \right]$$

$$\text{if } \left( \%c < \%c D^T (Y - DX) + \underline{x}_0 \right)$$

$$\underline{x}_l = \%c D^T (Y - DX) + \underline{x}_0 - \%c$$

$$\text{Else if } \left( \%c > \%c D^T (Y - DX) + \underline{x}_0 \right)$$

$$\underline{x}_l = \%c D^T (Y - DX) + \underline{x}_0 + \%c$$

Else

$$\underline{x}_l = 0$$

end if

end for

$$X \leftarrow X^{(k+1)}$$

Until  $\|X\|_{2,1} < \delta$

همان طور که در شبه‌کد الگوریتم مورد نظر نشان داده شده است، به ازای مقادیر مختلف  $\delta$ ، فریم‌های کلیدی به دست آمده متفاوت می‌باشند. این نکته به عنوان مقیاس‌پذیری<sup>۱</sup> روش پیشنهادی است. هرچه مقدر  $\delta$ ، بزرگ‌تر باشد، تعداد فریم‌های کلیدی خروجی بیشتر و هرچه مقدار  $\delta$ ، کوچک‌تر باشد، تعداد فریم‌های کلیدی کمتر می‌باشند.

### ۳. معرفی ویژگی مبتنی بر هیستوگرام تبدیل Sensus<sup>۲</sup> با ترکیب ویژگی رنگ و ویژگی آماری استفاده شده در این مقاله

ویژگی که در این مقاله مورد استفاده قرار گرفته است ترکیبی از ویژگی CENTRIST [۳۸]، رنگ [۳۹] و آماری [۴۰-۴۲] است. ویژگی‌های به دست آمده با استفاده از روش CENTRIST با قسمت‌بندی تصویر به چند بلوک و اعمال تبدیل Sensus روی آن‌ها به دست می‌آید. البته در این روش از فیلترهای مشتق‌گیر نیز برای فیلتر کردن تصویر مورد نظر استفاده می‌شود. اساس این روش برای ساختن قالبی<sup>۳</sup> از روی تصاویر بدون در نظر گرفتن فضای رنگ آن‌ها است. این الگوریتم از ساختار هرم مکانی استفاده می‌کند.

<sup>4</sup> Patch

<sup>5</sup> Central Moments

<sup>6</sup> Skewness

<sup>7</sup> Kurtosis

<sup>8</sup> Hue, Saturation, Intensity (HSI)

<sup>9</sup> Hue

<sup>10</sup> Saturation

<sup>11</sup> Intensity

<sup>1</sup> Scalable

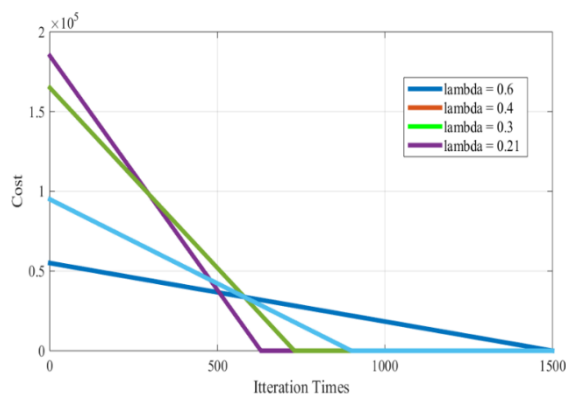
<sup>2</sup> CENsus Transform hISTogram (CENTRIST)

<sup>3</sup> Template

که بهترین مقدار  $\lambda$  برای خلاصه‌سازی ویدیو  $\lambda = 0.21$  در نظر گرفته می‌شود.

در ادامه به مقایسه روش پیشنهادی با روش‌های پیشنهادشده در زمینه خلاصه‌سازی ویدیو پرداخته می‌شود.

در اولین آزمایش، روش پیشنهادی را با روش<sup>۹</sup> DSVS [۳۱] مقایسه می‌شود. این مرجع بر روی پایگاه داده VSUMM<sup>۰</sup> برای انجام آزمایش‌های خود بهره برده است. علاوه بر روش فوق، روش پیشنهادی با چند روش دیگر مقایسه می‌شود که عبارت‌اند از روش<sup>۱۰</sup> OCFE [۴۸] است که با استفاده کلاس‌بندی برخط با استفاده از ویژگی‌های استفاده‌شده در [۳۱] ویدیو را خلاصه می‌کند. همچنین، روش UCF [۴۹] که این روش در چند بخش به خلاصه‌سازی ویدیو می‌پردازد.



شکل ۱. راندمان الگوریتم پیشنهادی با استفاده از تغییر مقدار  $\lambda$  به صورت تجربی که محور افقی آن مربوط به تعداد تکرار الگوریتم و محور عمودی آن مربوط به تابع ارزش  $\|X\|_{2,1}$  است

در بخش اول با استفاده از ویژگی محتوایی بر مبنای حرکت شات و شباهت رنگ بین شات‌ها به آشکارسازی شات و پیدا کردن شات‌های مشابه پرداخته است. در مرحله بعد استفاده از ساختن گراف شات‌های مشابه و تشخیص چهره به بازسازی و آشکارسازی صحنه به خلاصه‌سازی ویدیو منجر می‌شود. در این آزمایش به دلیل برابر بودن شرایط آزمایش از ویژگی مبتنی بر هیستوگرام تبدیل Sensus با ترکیب ویژگی رنگ و ویژگی آماری استفاده شده است. در ادامه به نحوه امتیازدهی پرداخته می‌شود. برای امتیاز دادن به فریم‌های خلاصه استخراج شده بدین گونه عمل می‌شود که اگر فریم انتخاب شده درست باشد به خروجی ۱ واحد و اگر اشتباه باشد صفر واحد اضافه می‌شود. البته برای تطبیق‌های ضعیف نیز مقدار  $0.5$  واحد به خروجی اضافه می‌شود.

<sup>۹</sup> Dictionary Selection Based Video Summarization (DSVS)

<sup>۱۰</sup> Online Clustering Key Frames Extraction (OCFE)

Kurtosis :  $k_i =$

$$\left\{ \frac{N(N+1)}{(N-1)(N-2)(N-3)} \sum_{j=1}^N \left( \frac{p_{ij} - \bar{X}_i}{\sigma_i} \right)^4 \right\} - \frac{3(N-1)^2}{(N-2)(N-3)} \quad (12)$$

#### ۴. ارزیابی روش پیشنهادی

در ابتدا به بررسی پایگاه داده‌های استفاده‌شده در این مقاله پرداخته می‌شود.

##### ۴-۱. پایگاه داده

در این مقاله از ۲ پایگاه داده ویدیویی استفاده شده است.

**پایگاه داده VSUMM:** این پایگاه داده، شامل ۵۰ ویدیو انتخاب شده از OVP [۴۴] است و شامل صحنه‌های مختلف ویدیویی که در قسمت قبل توضیح داده شد، است [۴۵]. مدت زمان این کلیپ‌ها بین ۱ تا ۴ دقیقه است. خلاصه زمین حقیقت<sup>۱</sup> این پایگاه داده در با استفاده از ۵۰ کاربر به دست آمده است [۴۶]. برای هر ویدیو از ۵ کاربر مختلف کمک گرفته شده است. بنابراین برای هر ویدیو ۵ مجموعه خلاصه وجود دارد. در کل به تعداد ۲۵۰ داده خلاصه‌شده ویدیویی برای ۵۰ کلیپ ویدیویی در دسترس است.

**پایگاه داده OVP<sup>۲</sup>:** این پایگاه داده از مرجع [۴۳] به دست آمده است که شامل ۶ گروه کلیپ ویدیویی است که شامل: مستند<sup>۳</sup>، آموزشی<sup>۴</sup>، ویدیوهای روزمره<sup>۵</sup>، تاریخی<sup>۶</sup>، کنفرانس<sup>۷</sup> و خدمات عمومی<sup>۸</sup> است. برای انجام آزمایش‌ها بر روی این پایگاه داده از ۱۰ ویدیو کلیپ که در جدول (۱) از مرجع [۴۷] معرفی شده‌اند، استفاده شده است. مجموعه فریم‌های کلیدی شبیه به مرجع [۴۸]، به صورت دستی و با استفاده از ۵ دانشجو با سابقه کار در زمینه پردازش ویدیو انتخاب شده است. در ادامه به بررسی مقدار بهینه  $\lambda$  پرداخته می‌شود.

##### ۴-۲. مقایسه روش پیشنهادی با روش‌های دیگر

در اولین گام به بررسی مقدار انتخابی برای پارامتر  $\lambda$  پرداخته می‌شود. با بررسی مقادیر مختلف  $\lambda$  بر روی مقادیر مختلف ویدیویی که در شکل (۱) نشان داده شده است به این نتیجه رسید

<sup>۱</sup> Ground Truth

<sup>۲</sup> Open Video Project (OVP)

<sup>۳</sup> Documentary

<sup>۴</sup> Educational

<sup>۵</sup> Ephemeral

<sup>۶</sup> Historical

<sup>۷</sup> Lecture

<sup>۸</sup> Public Service

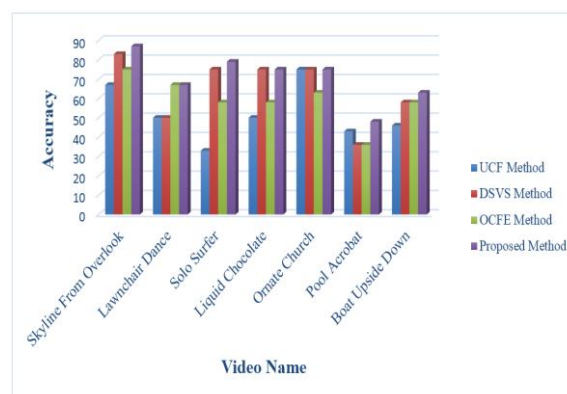
در آزمایش دوم، روش پیشنهادی با استفاده از پایگاه داده OVP و بر روی ۱۰ فایل ویدئویی صورت پذیرفته است که در جدول (۱) از مرجع [۴۹] معرفی شده‌اند. در این آزمایش به مقایسه روش پیشنهادی با روش‌های BoIVS<sup>۱</sup>، DT<sup>۱</sup> [۵۰]، STIMO<sup>۲</sup> [۵۱]، VSUMM [۴۶]، KBKS<sup>۳</sup> [۵۲] و DSVS [۳۱] می‌پردازیم. این مقایسه بر روی تمام پایگاه داده VSUUM صورت پذیرفته است و نتیجه آن در شکل (۴) نشان داده شده است.

همان‌طور که در شکل (۴) نشان داده شده است، روش پیشنهادی به میزان ۱٫۳٪ خلاصه‌های به دست آمده از پایگاه داده VSUMM را بهبود داده است. در ادامه یک نمونه از نتایج به دست آمده برای پایگاه داده VSUMM در شکل (۵) نشان داده شده است.

همان‌طور که در شکل (۵) نشان داده شده است فریم‌های خلاصه شده به وسیله روش پیشنهادی به فریم‌های زمین مرجع نزدیک‌تر می‌باشند. روش‌های STIMO و BoIVS نسبتاً دارای خطای بیشتری نسبت به روش VSUMM می‌باشند. ولی در مجموع روش پیشنهادی توانسته است روش‌های گفته شده را بهبود بخشد.

در ادامه خروجی روش پیشنهادی با ۳ روش گفته شده در شکل (۲) نشان داده شده است.

به عنوان مثال نتیجه روش پیشنهادی روی مجموعه ویدئویی Liquid Chocolate که در مرجع نیز مورد بررسی قرار گرفته [۳۱] و در شکل (۳) نشان داده شده است.



شکل ۲. مقایسه روش پیشنهادی با روش‌های UCF، DSVS و OCFE. برای ۱۸ فایل ویدئویی از پایگاه داده VSUMM

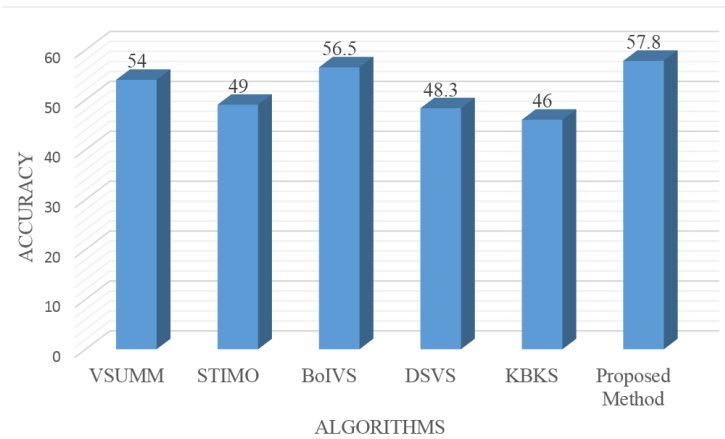


شکل ۳. نمونه‌ای از فریم‌های خلاصه شده برای نمونه ویدئویی Liquid Chocolate که سطر اول نمونه زمین مبنا، سطر دوم، سوم، چهارم و پنجم به ترتیب فریم‌های به دست آمده از OCFE، DSVS، UCF و روش پیشنهادی می‌باشند. در سطرهای دوم تا پنجم، فریم‌های استخراج شده‌ای که درست می‌باشند با رنگ سبز و آن دسته از فریم‌هایی که نزدیک به زمین مبنا است با رنگ قرمز مشخص شده است.

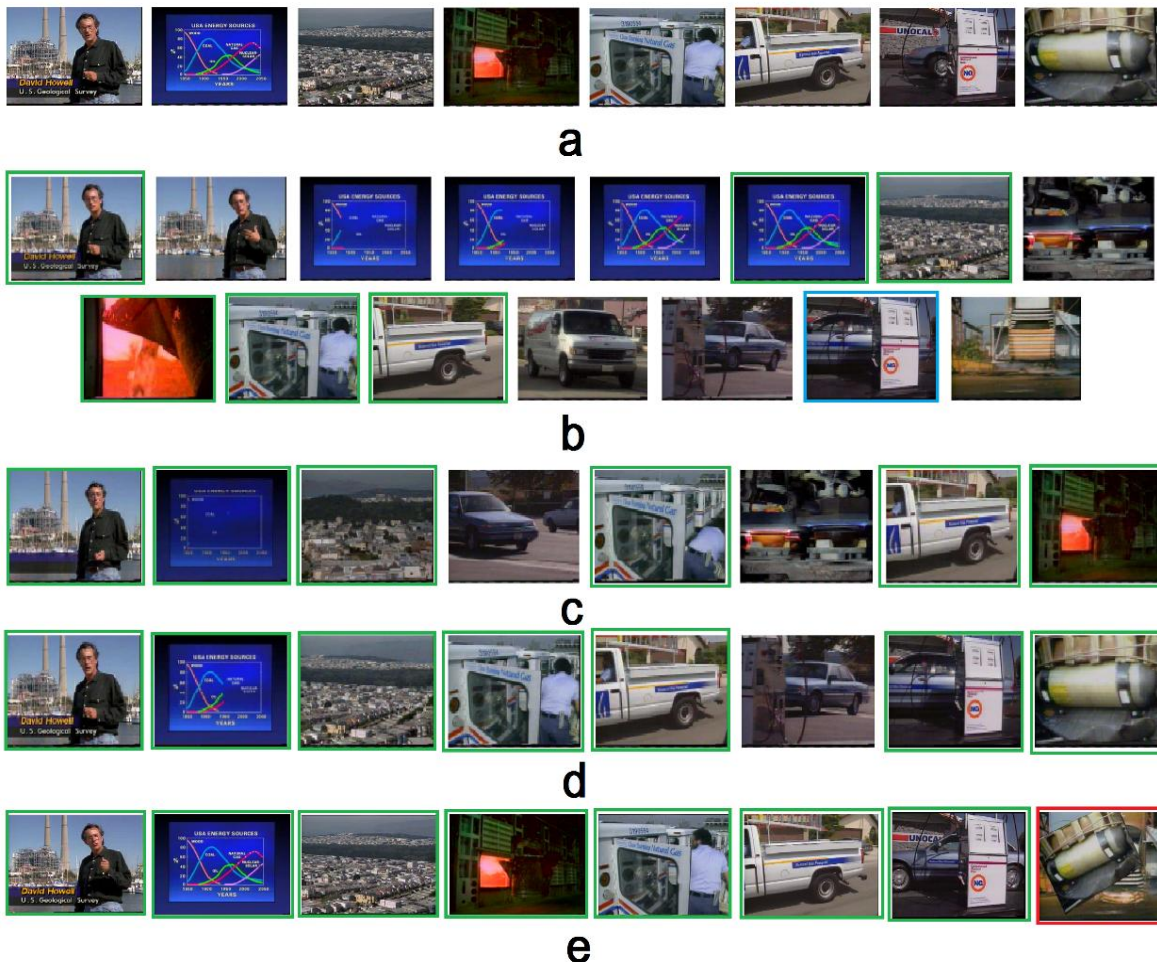
<sup>۱</sup> STII and MOving video storyboard (STIMO)

<sup>۱</sup> Keypoint Based Keyframe Selection (KBKS)





شکل ۴. مقایسه عملکرد روش پیشنهادی با روش‌های VSUMM، STIMO، BoIVS، DSVS و KBKS بر روی تمامی ۵۰ فایل ویدیویی پایگاه VSUMM



شکل ۵. نمونه‌ای از فریم‌های خلاصه‌شده برای ۲ نمونه ویدیویی از پایگاه داده VSUMM که a نمونه زمین مبنا، b، c، d و e به ترتیب فریم‌های به دست آمده از VSUMM، BoIVS، STIMO و روش پیشنهادی می‌باشند. در b تا e، فریم‌های استخراج‌شده‌ای که درست می‌باشند با رنگ سبز و آن دسته از فریم‌هایی که در نمونه زمین مبنا وجود نداشته ولی از لحاظ محتوایی به عنوان فریم خلاصه‌شده می‌توان در نظر گرفت با رنگ قرمز مشخص شده است.

<sup>1</sup>Delaunay Triangulation (DT)

## ۵. نتیجه‌گیری

با توجه به اینکه روش‌های مراقبت ویدئویی، جزو روش‌های پر استفاده در حوزه مراقبت و آشکارسازی است. در این مقاله با استفاده از روش‌های جدید در حوزه پردازش ویدئو به بهبود عملکرد این سامانه‌ها پرداخته شد. بنابراین با استفاده از انتخاب واژه‌نامه تنک با الگوریتم انقباضی آستانه نرم می‌توان میزان خلاصه‌سازی ویدئو را بهبود بخشید. بدین منظور ابتدا کدینگ تنک را برای خلاصه‌سازی ویدئو فرموله شد و پس از آن با توجه به آنچه از حل معادلات کلاسیک به دست آمد می‌توان الگوریتمی برای خلاصه‌سازی با استفاده از این روش پیشنهاد داد. برخلاف عملیات ریاضی در حل معادلات تنک؛ پس از به دست آمدن عبارت نهایی، مشاهده می‌شود که دارای پیچیدگی کمی نسبت به روش‌های دیگر بر پایه تنک است. در انتها، روش پیشنهادی به صورت یک الگوریتم نمایش داده شده است که منجر به سادگی پیاده‌سازی آن می‌شود. پس از اجرای الگوریتم بر روی ۲ پایگاه داده OVP، VSUMM و مقایسه فریم‌های کلیدی به دست آمده از روش پیشنهادی و روش‌های STIMO، VSUMM، DSVS، BoIVS و KBKS می‌توان میزان خلاصه‌سازی را بهبود بخشید که این امر ادعای فرض شده در این مقاله، مبنی بر بهبود روش پیشنهادی نسبت به دیگر روش‌ها را نشان می‌دهد.

## ۶. مراجع

- [9] Ma, Y. F.; Lu, L.; Zhang, H. J.; Li, M. "A User Attention Model for Video Summarization"; ACM Multimedia'02 2002, 533-542.
- [10] Truong, B.; Venkatesh, S. "Video Abstraction: A Systematic Review and Classification"; ACM Trans. on Multimedia Computing, Communication and Applications 2007, 3, 1-37.
- [11] Hu, W.; Xie, N.; Zeng, X.; Maybank, S. "A Survey on Visual Content-Based Video Indexing and Retrieval"; IEEE Trans. Syst., Man, and Cyber. 2011, 41, 797-819.
- [12] Tsai, D. M.; Lai, S. C. "Independent Component Analysis-Based Background Subtraction for Indoor Surveillance"; IEEE Trans. Image Process 2009, 18, 158-167.
- [13] Xiang, T.; Gong, S. "Video Behavior Profiling for Anomaly Detection"; IEEE Trans. Pattern Anal. Mach. Intell. 2008, 30, 5, 893-908.
- [14] Li, Y.; Zhang, T.; Tretter, D. "An Overview of Video Abstraction Techniques"; Technical Report HPL-2001-191, 2001.
- [15] Ciocca, G.; Schettini, R. "Innovative Algorithm for Key Frame Extraction in Video Summarization"; J. Real-Time Image Process 2006, 1, 1, 69-88.
- [16] Massimiliano, A. "Extracting and Summarizing Information from Large Data Repositories"; Ph.D. Dissertation, University of Naples Federico II, Italia, 2006.
- [17] Money, A. G.; Agius, H. "Video Summarization: A Conceptual Framework and Survey of the State of the Art"; J. Vis. Com. Image Rep. 2008, 19, 121-143.
- [18] Li, Y.; Lee, S. H.; Yeh, C. H.; Kuo, C. C. "Techniques for Movie Content Analysis and Skimming: Tutorial and Overview on Video Abstraction Techniques"; IEEE Signal Process. Magazine 2006, 23, 79-89.
- [19] Kleban, J.; Sarkar, A.; Moxley, E.; Mangiat, S.; Joshi, S.; Kuo, T.; Manjunath, B. S. "Feature Fusion and Redundancy Pruning for Rush Video Summarization"; Proc. of the ACM Int. Workshop on Video Summarization, 2007, 84-88.
- [20] Pan, C. M.; Chuang, Y. Y.; Hsu, W. H. "Fast Rushes Summarization System"; Proc. of the ACM Int. Workshop on Video Summarization, 2007, 74-78.
- [21] Doulamis, N. D.; Doulamis, A. D.; Avrithis, Y. S.; Ntalianis, K. S.; Kollias, S. D. "Efficient Summarization of Stereoscopic Video Sequences"; IEEE Trans. Circuits Syst. Video Tech. 2000, 10, 501-517.
- [22] Le, D. D.; Satoh, S. "National Institute of Informatics"; Proc. of the ACM Int. Workshop on Video Summarization, 2007, 70-73.
- [23] Nesterov, Y. "Gradient Methods for Minimizing Composite Objective Function"; Louvain-la-Neuve, Belgium, CORE, 2007.
- [24] Cong, Y.; Yuan, J.; Luo, J. "Towards Scalable Summarization of Consumer Videos via Sparse Dictionary Selection"; IEEE Trans. Multimedia 2012, 14, 66-75.
- [25] Hanjalic, A.; Langendijk, R. L.; Biemond, J. "A New Key Frame Allocation Method for Representing Stored Video Streams"; 1<sup>st</sup> Int. Workshop on Image Databases & Multi. Search, Amsterdam, 1996, 67-74.
- [26] Fu, Y.; Guo, Y.; Zhu, Y.; Liu, F.; Song, Ch.; Zhou, Zh. "Multi View Video Summarization"; IEEE Trans. Multimedia. 2010, 12, 717-729.
- [1] Weiming, H.; Tieniu, T.; Iang, W.; Steve, M. "A Survey on Visual Surveillance of Object Motion and Behaviors"; IEEE Trans. Syst. Cybernetics 2004, 3, 334-352.
- [2] Hong, L.; Ruan, Y.; Wicker, D.; Layne, J. "Energy-Based Video Tracking Using Joint Target Density Processing with an Application to Unmanned Aerial Vehicle Surveillance"; IET Image Processing 2008, 2, 1-12.
- [3] Hampapur, A.; Connell, J.; Haas, N.; Merkl, H.; Shu, C. F. "Smart Video Surveillance"; IEEE Signal Magazin 2005, 38-51.
- [4] Conte, G.; Dohetry, P. "Vision-based Unmanned Aerial Vehicle Navigation Using Geo-Referenced Information"; J. on Advances in Signal Processing 2009, 10, 1-18.
- [5] Carr, J. R.; Sobek, J. S. "Digital Scene Matching Area Correlator (DSMAC)"; Proc. Conf. Image Processing for Missile Guidance, 1980, 36.
- [6] Tsai, S. X. "Introduction to the Scene Matching Missile Guidance Technologies"; Guided Missile, 1996.
- [7] Liu, T.; Zhang, H. J.; Qi, F. "A Novel Video Key Frame Extraction Algorithm Based on Perceived Motion-Energy Model"; IEEE Trans. on Circuits and Syst. for Video Tech. 2003, 13, 1006-1013.
- [8] Taskiran, C. M.; Pizlo, Z.; Amir, A.; Ponceleon, D.; Delp, E. "Automated Video Program Summarization Using Speech Transcripts"; IEEE Trans. on Multimedia 2006, 8, 775-791.

- [39] Stricker, M.; Orengo, M. "Similarity of Color Images"; Proc. SPIE Storage and Retrieval for Image and Video Databases, 1995, 2420, 381–392.
- [40] Rayner, J. C. W.; Best, D. J.; Matthews, K. L. "Interpreting the Skewness Coefficient"; Communications in Statistics Theory and Methods 1995, 24, 593–600.
- [41] Groeneveld, R. A.; Meeden, G. "Measuring Skewness and Kurtosis"; J. Roy. Statist. Soc. 1984, 33, 391–399.
- [42] Arnold, B. C.; Groeneveld, R. A. "Measuring Skewness with Respect to the Mode"; Am. Statistician J. 1995, 49, 34–38.
- [43] Gonzalez, C.; Woods, R. E. "Digital Image Processing"; Prentice-Hall, Second Edition, 2002.
- [44] "The Open Video Project"; <http://www.open-video.org>.
- [45] "VSUMM Database"; <https://sites.google.com/site/vsummsite/results>.
- [46] Avila, S. E. F.; Lopes, A. P. B.; Daluz, A.; Araújo, A. "Vsumm: A Mechanism Designed to Produce Static Video Summaries and a Novel Evaluation Method"; Pattern Rec. Let. 2011, 32, 56–68.
- [47] Lu, S.; Wnag, Z.; Mei, T.; Guan, G.; Feng, D. D. "A Bag of Important Model with Locality Constrained Coding Based Feature Learning for Video Summarization"; IEEE Trans, Multimedia 2014, 16, 1497–1509.
- [48] Luo, J.; Papin, C.; Costello, K. "Towards Extracting Semantically Meaningful Key Frames From Personal Video Clips: From Humans to Computers"; IEEE Trans. Circuits Syst. Video Tech. 2009, 19, 289–301.
- [49] Rasheed, Z.; Shah, M. "Detection and Representation of Scenes in Videos"; IEEE Trans. Multimedia. 2005, 7, 1097–1105.
- [50] Mundur, P.; Rao, Y.; Yesha, Y. "Keyframe Based Video Summarization Using Delaunay Clustering"; Int. J. Digit. Libraries 2006, 6, 2, 219–232.
- [51] Furini, M.; Geraci, F.; Montangero, M.; Pellegrini, M. "Stimo: Still and Moving Video Storyboard for the Web Scenario"; Multimedia Tools Appl. 2010, 46, 47–69.
- [52] Guan, G.; Wang, Z.; Lu, S.; Deng, J. D.; Feng, D. "Keypoint Based Keyframe Selection"; IEEE Trans. Circuits Syst. Video Tech. 2013, 23, 729–734.
- [27] Zhuang, Y.; Rui, Y.; Huang, T.; Mehrotra, S. "Adaptive Key Frame Extraction Using Unsupervised Clustering"; Proc. Int. Conf. Image Process, 1998, 866–870.
- [28] Cernekova, Z.; Pitas, I.; Nikou, C. "Information Theory-Based Shot Cut/Fade Detection and Video Summarization"; IEEE Trans. Circuits & Sys for Video Tech. 2006, 16, 82–91.
- [29] Oliva, A.; Torralba, A. "Modeling the Shape of the Scene: A Holistic Representation of the Spatial Envelope"; Int. J. Comput. Vis. 2001, 42, 145–175.
- [30] Wu, J.; Christensen, H.; Rehg, J. "Visual Place Categorization: Problem, Dataset, and Algorithm"; Proc. IROS, 2009.
- [31] Wu, J.; Rehg, J. "Centrist: A Visual Descriptor for Scene Categorization"; IEEE Trans. Pattern Anal. Mach. Intell. 2010, 33, 1489–1501.
- [32] Gallager, R. G. "Information Theory and Reliable Communication"; John Wiley and Sons, New York, 1968.
- [33] Tibshirani, R. "Regression Shrinkage and Selection via the Lasso"; J. Roy. Statist. Soc. Series B (Methodological), 1996, 58, 267–288.
- [34] Yuan, M.; Lin, Y. "Model Selection and Estimation in Regression with Grouped Variables"; J. Roy. Statist. Soc. 2006, 68, 49–67.
- [35] Cong, Y.; Yuan, J.; Liu, J. "Sparse Reconstruction Cost for Abnormal Event Detection"; Proc. IEEE Conf. Com. Vision & Pattern Recognition (CVPR), 2011, 3449–3456.
- [36] Fornasier, M.; Rauhut, H. "Iterative Thresholding Algorithms"; Science Direct, Applied and Computational Harmonic Analysis 2008, 25, 187–208.
- [37] Selesnick, I. "A Derivation of the Soft-Thresholding Function"; Ph.D. Thesis, Polytechnic Institute of New York University, 2009.
- [38] Wu, J.; Christensen, H.; Rehg, J. "Visual Place Categorization: Problem, Dataset, and Algorithm"; Proc. Intelligent Robots and Systems, 2009, 4763–47760.