

# Temporal Visual Semantics-Induced Human Motion Understanding with Large Language Models

Zheng Xing and Weibing Zhao

**Abstract**—Unsupervised human motion segmentation (HMS) can be effectively achieved using subspace clustering techniques. However, traditional methods overlook the role of temporal semantic exploration in HMS. This paper explores the use of temporal vision semantics (TVS) derived from human motion sequences, leveraging the image-to-text capabilities of a large language model (LLM) to enhance subspace clustering performance. The core idea is to extract textual motion information from consecutive frames via LLM and incorporate this learned information into the subspace clustering framework. The primary challenge lies in learning TVS from human motion sequences using LLM and integrating this information into subspace clustering. To address this, we determine whether consecutive frames depict the same motion by querying the LLM and subsequently learn temporal neighboring information based on its response. We then develop a TVS-integrated subspace clustering approach, incorporating subspace embedding with a temporal regularizer that induces each frame to share similar subspace embeddings with its temporal neighbors. Additionally, segmentation is performed based on subspace embedding with a temporal constraint that induces the grouping of each frame with its temporal neighbors. We also introduce a feedback-enabled framework that continuously optimizes subspace embedding based on the segmentation output. Experimental results demonstrate that the proposed method outperforms existing state-of-the-art approaches on four benchmark human motion datasets.

**Index Terms**—Human motion segmentation, temporal vision semantics, subspace embedding, temporal neighbors.

## I. INTRODUCTION

**H**UMAN motion segmentation (HMS) has attracted significant attention in both industry and academic research due to its wide-ranging applications in video retrieval, virtual reality, and intelligent surveillance, particularly in human motion analysis [1]–[3]. The primary goal of unsupervised HMS is to partition frame sequences depicting human actions into non-overlapping, internally consistent groups without the need for training, serving as a preprocessing step for motion-related analysis tasks [4]. However, unsupervised HMS faces the challenge of motion primitive ambiguity due to temporal variability across different actions [5]–[9].

Subspace clustering is a well-established strategy for HMS, aiming to partition a human motion sequence into distinct groups based on the assumption that the frames originate from multiple subspaces, with frames depicting the same motion belonging to the same subspace [10]–[15]. In recent years, a

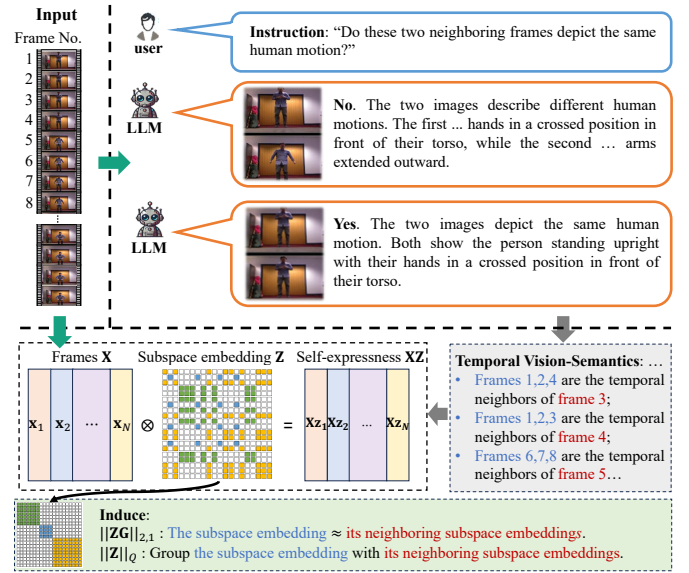


Fig. 1: Framework of the proposed method.

popular and effective approach is to first embed the human motion frames into multiple subspaces to learn the subspace structure of the data, and then apply traditional clustering algorithms to the subspace embeddings [16]–[24].

Human motion sequences inherently contain temporal information, which is crucial for HMS. As illustrated in Figure 1, a human motion sequence typically consists of multiple motion segments. For example, a person may first perform a motion with both hands clasped for a period, then extend their hands for another period, followed by squatting. However, due to the ambiguity of different motions and the complexity of temporal correlations, extracting temporal information from human motion sequences remains a significant challenge [20].

Various subspace clustering algorithms have been proposed to achieve HMS by exploring the temporal information in the data. For instance, Wang et al. [25] eliminate redundant connections between adjacent motions in the subspace embedding to extract informative instance data and capture the compact structure of human motion videos. Bai et al. [9] extract informative features during subspace embedding to capture local temporal consistency in human motions. Zhou et al. [8] employ a multi-mutual consistency learning strategy to factorize source and target data into distinct multi-layer feature spaces, thereby learning the temporal information from the source domain. Despite efforts to explore and utilize temporal information, its inaccuracy and ambiguity may lead to incor-

Weibing Zhao is with Guangdong Laboratory of Machine Perception and Intelligent Computing, Shenzhen MSU-BIT University, China. (E-mail: weibingzhao@smbu.edu.cn)

Zheng Xing is with the College of Computer Science and Software Engineering, Shenzhen University, Shenzhen 518060, China.

rect segmentation results. While neural network algorithms have shown great promise in supervised tasks, they often fall short in unsupervised tasks due to the difficulty in exploring the underlying data structures.

This paper aims to learn temporal vision semantics (TVS) from human motion sequences by leveraging the image-to-text capabilities of a pre-trained large language model (LLM) to enhance HMS performance. The key idea is to extract textual motion information from consecutive frames using the LLM and integrate this learned information into the subspace clustering framework. By incorporating TVS into the subspace clustering, we aim to ensure that the segmentation output more effectively captures the temporal dynamics inherent in the human motion sequence.

Thus, we face the following **challenges**:

- *How to learn TVM using LLM?* Although LLMs are widely used in image-to-text tasks, no research has explored how to leverage LLMs to assist in unsupervised HMS. The challenge lies in learning the textual temporal information that can be converted into a mathematical form, which can be used to induce the HMS.
- *How to integrate TVM with HMS?* We aim to learn the subspace embedding and perform segmentation based on this embedding, all induced by TVM. However, integrating TVM into both the subspace embedding and segmentation presents significant challenges.

In this study, we determine whether consecutive frames represent the same motion by querying the LLM. Based on its response, we subsequently learn the temporal relationships between frames, as illustrated in Figure 1. Building upon this, we propose a subspace clustering approach integrated with TVS, which combines subspace embedding with a temporal regularizer. This regularizer ensures that each frame shares similar subspace embeddings with its temporal neighbors. Segmentation is then performed using these subspace embeddings, with a temporal constraint that encourages the grouping of each frame with its temporal neighbors. Furthermore, we introduce a feedback-enabled framework that iteratively optimizes the subspace embeddings based on the segmentation output, ensuring continuous refinement of the model.

In summary, the main **contributions** are as follows:

- *Exploring TVS via LLMs and Integrating TVS with HMS:* This paper introduces an approach that uses LLMs to learn TVS in human motion sequences. We develop a method that applies TVS to both subspace embedding and segmentation, ensuring neighboring consistency in both processes.
- *Feedback-enabled Subspace Embedding:* We propose a feedback-enabled strategy that allows the segmentation output to inform the subspace embedding. This is not merely a combination of two methods; rather, it enables the use of HMS output to induce subspace embedding, better capturing the underlying subspace structure in the human motion sequence.

We conduct extensive experiments on four benchmark datasets for HMS. The experimental results consistently demonstrate

that our method outperforms existing state-of-the-art techniques, highlighting its superiority in HMS.

The remainder of this paper is structured as follows. Section II briefly introduces the related work including HMS and subspace clustering. Section III presents the details of the proposed method. Section IV provides the experimental settings, performance comparisons, and ablation study. Finally, Section V concludes the paper.

## II. RELATED WORKS

### A. Human Motion Segmentation

HMS is essential for accurately capturing human motion data, forming a basis for structural analysis, understanding, and practical applications [26]–[30]. Significant research efforts have led to notable achievements in this area. For instance, Zhong et al. [31] proposed a bipartite graph co-clustering framework to segment unusual activities in videos. Jenkins et al. [32] utilized zero-velocity crossing frames of angular velocity to partition motion data streams into different sequences. Barbic et al. [33] employed probabilistic principal component analysis to decompose human motion into distinct motions. Additionally, Beaudoin et al. [34] introduced a framework for distilling a motion-motif graph from motion data collections. Spatio-temporal-based Convolutional Neural Networks (CNNs) [35] and clustering-based approaches [36] have been proposed for segmenting streams of human motion into multiple activities. Despite the capability of deep learning-based motion recognition models to complete HMS tasks by training with large datasets, the unsupervised model offers significant advantages in terms of interpretability and computational efficiency. Therefore, achieving HMS tasks through an unsupervised approach is highly beneficial. [37]–[44]

However, these approaches may not fully exploit the temporal dynamics and semantic continuity inherent in human motion sequences, potentially limiting the accuracy and effectiveness of the segmentation results.

### B. Subspace Clustering

Subspace clustering discerns and segregates distinct motion types into their respective subspaces, thereby addressing human motion segmentation tasks [45]–[48]. Its capacity to manage intricate, high-dimensional motion data, combined with robustness to noise and data variability, ensures reliability in practical applications. Furthermore, by exploiting the inherent low-dimensional structures within complex motion datasets, subspace clustering facilitates the further analysis and utilization of human motion data [49]–[60].

Subspace clustering serves to ascertain the low-dimensional embedding of a high-dimensional manifold. Specifically, assuming that vectorized frames corresponding to identical actions reside within the same subspace, the subspace embedding property inherent to the data can be harnessed to derive a representative subspace embedding [5], [8], [9], [25], [61]–[63]. Consider, in particular, a matrix composed of column-wise frames  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ . Each vectorized frame may be expressed as a linear combination of all frames, that is,  $\mathbf{x}_i = \mathbf{X}\mathbf{z}_i$ , where  $\mathbf{z}_i \in$

$\mathbb{R}^N$  is the subspace embedding of  $\mathbf{x}_i$ . Hence, we may write  $\mathbf{X} = \mathbf{XZ}$ , where  $\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N] \in \mathbb{R}^{N \times N}$  constitutes the coefficient matrix. To accommodate noise, this relationship is extended to  $\mathbf{X} = \mathbf{XZ} + \mathbf{Y}$ , where  $\mathbf{Y}$  represents the noise component. The Frobenius norm of  $\mathbf{Y}$  is employed as the loss function to penalize discrepancies. Consequently, the subspace clustering problem is formulated as the minimization of  $\|\mathbf{X} - \mathbf{XZ}\|_F^2$  with respect to  $\mathbf{Z}$ . Mathematically, as delineated in [64], the formulation is expressed as minimizing  $\|\mathbf{X} - \mathbf{XZ}\|_F^2$  subject to the constraints  $\text{diag}(\mathbf{Z}) = 0$  and  $\mathbf{Z} \geq 0$ , where  $\|\cdot\|_F^2$  denotes the Frobenius norm. Herein,  $\mathbf{Z}$  embodies the sought subspace embedding.

To eliminate the effect of the noise on subspace embedding, a quadratic term  $\|\mathbf{Z}^T \mathbf{Z}\|_1$  is employed [65]. This term promotes the expression of each datum as a linear combination of other data points, whilst a regularization term enforces the nullification of reconstruction coefficients between vectors originating from distinct subspaces. Specifically, the regularization term  $\|\mathbf{Z}^T \mathbf{Z}\|_1$  compels  $\mathbf{z}_{ij} = 0$  whenever  $\mathbf{x}_i$  and  $\mathbf{x}_j$  reside in different subspaces. Given that  $\mathbf{Z}$  is nonnegative and that  $\|\mathbf{Z}^T \mathbf{Z}\|_1 = \mathbf{e}^T \mathbf{Z}^T \mathbf{Z} \mathbf{e}$ , where  $\mathbf{e}$  denotes the all-one vector, it follows that  $\|\mathbf{Z}^T \mathbf{Z}\|_1 = \sum_{i,j} \mathbf{z}_i^T \mathbf{z}_j$ . The minimization of  $\|\mathbf{Z}^T \mathbf{Z}\|_1$  thereby encourages sparsity in both  $\mathbf{z}_i$  and  $\mathbf{z}_j$ , leading the inner product  $\mathbf{z}_i^T \mathbf{z}_j$  to approach zero. Consequently, the subspace embedding problem is thus formulated as [64]:

$$\begin{aligned} & \underset{\mathbf{Z}}{\text{minimize}} \quad \|\mathbf{X} - \mathbf{XZ}\|_F^2 + \|\mathbf{Z}^T \mathbf{Z}\|_1 \\ & \text{subject to} \quad \text{diag}(\mathbf{Z}) = 0, \mathbf{Z} \geq 0 \end{aligned} \quad (1)$$

Traditional clustering methods are then applied to the subspace embedding matrix  $\mathbf{Z}$ , and the resulting clusters represent the HMS, which can be matched to specific motions using the Hungarian algorithm.

Subspace clustering has recently gained significant attention due to its effectiveness in uncovering complex data structures and improving clustering performance in high-dimensional spaces [66]–[68]. For instance, the *SIBMSC* method [16] extends the information bottleneck principle to learn view-common representations, removing redundant information and leveraging mutual information for view-specific clustering. Similarly, the *BTMSC* method [17] constructs a third-order tensor to capture high-order correlations, using the Bi-Nuclear Quasi-Norm for efficient tensor factorization. To improve robustness, the *FSMSC* method [18] integrates view-shared anchor learning with a self-guided discriminative feature selection approach, addressing noisy views and cross-view diversity. The *ARLRR* method [19] introduces affine and non-negative constraints in low-rank self-representation learning to manage affine subspaces and errors. The *DCTMSC* method [20] employs a two-step discrete cosine transform approach to simplify tensor nuclear norm calculations and enhance local structural representation. The *DCMVC* method [69] incorporates dynamic cluster diffusion and reliable neighbor-guided positive alignment to improve inter-cluster separation and within-cluster compactness. Subspace clustering methods incorporating temporal priors have proven effective in HMS tasks. For instance, the *OSC* method [5] applies a one-neighbor consistency constraint for closer representations of temporal

data, while the *TSC* method [6] uses non-negative dictionary learning and temporal Laplacian regularization. The *LTS* method [7] captures temporal correlations in both source and target data with a graph regularizer and introduces a weighted low-rank constraint to reveal clustering structures. The *CDMS* approach [8] leverages transfer subspace learning to capture multi-level information in videos. These methods, often formulated as unsupervised learning frameworks, typically adopt a self-representation strategy for motion segmentation. The *DSAE* method [70] enhances representation learning by considering temporal correlations, while the *VSDA* method [9] employs a multi-neighbor auto-encoder to extract temporal features and a long-short distance embedding/deembedding strategy to maintain representation consistency, further enhanced by a velocity-sensitive guidance mechanism.

However, existing subspace clustering approaches typically divide the process into two independent stages and often overlook the potential of incorporating temporal semantics to simultaneously enhance both subspace embedding and clustering. This oversight limits the alignment of the HMS output with the true sequential dynamics of human motion.

### C. Temporal Vision Semantics from Large Language Models

LLMs have progressed from text-only reasoning engines to unified multimodal systems capable of jointly understanding visual and linguistic information. Recent architectures such as GPT-4o, Gemini, Claude, DeepSeek, and Qwen3 integrate visual encoders with transformer-based text reasoning via large-scale contrastive pretraining, enabling them to perform complex cross-modal reasoning and semantic alignment between images and language. Unlike conventional convolutional or transformer-based visual encoders that rely on geometric or pixel-level similarity, multimodal LLMs exhibit *semantic reasoning capability*. They can compare two visual scenes and judge whether they convey the same conceptual meaning based on high-level world knowledge and contextual understanding.

Extensive research in vision–language modeling has demonstrated the strong semantic reasoning capabilities of LLMs when integrated with visual inputs. Early studies validated these capabilities in tasks such as zero-shot visual question answering and high-fidelity caption generation, where LLMs interpret visual entities, relationships, and contextual meanings directly from raw imagery [71], [72]. Building upon these foundations, subsequent works extended vision semantics to three-dimensional and dynamic scenes, leveraging language-guided scene understanding and position-aware video representations for 3D perception [73], [74]. Beyond visual applications, LLMs have also exhibited robust zero-shot reasoning and representational alignment abilities that enable general semantic understanding across modalities [75], [76].

Building upon this progress, a growing body of research has explored the integration of LLM-based visual semantics into temporal vision understanding. Recent work has examined whether video-oriented LLMs truly capture temporal reasoning or merely rely on knowledge and spatial perception [77], while subsequent studies have demonstrated that LLMs can effectively learn temporal dependencies and causal

relations across video frames [78]. Further developments employ language-guided attention mechanisms to align visual dynamics with textual motion cues, thereby enhancing spatial-temporal object understanding and fine-grained temporal reasoning [79], [80]. Beyond short-term video grounding, recent efforts extend this semantic alignment to long-term sequence modeling, revealing that LLMs can encode cross-frame dependencies with human-level temporal abstraction [81], [82].

In the domain of human motion analysis, LLMs have been increasingly adopted to reason about human activities and their semantic transitions [83]–[86]. These studies demonstrate that LLMs can interpret and describe motion in natural language, distinguish subtle phase changes, and correlate sensor or visual signals with linguistic motion descriptions. Recent studies have advanced from graph-based relational modeling to LLM-driven semantic reasoning in motion understanding. For example, some works employ LLMs to anticipate long-term actions by treating video frames as language-like tokens and enhancing vision-language interaction through cross-modal reasoning [87], while others utilize graph attention mechanisms to capture individual-group interaction dynamics in collective activities [88]. Such findings inspire the present work, where we employ LLM-based reasoning to construct temporal semantics, serving as a high-level inductive prior for unsupervised human motion segmentation.

### III. METHODOLOGY

In this section, we first introduce an LLM-based inference to identify the TVS. We then propose a feedback-enabled subspace embedding approach that incorporates TVS to efficiently determine the HMS with limited iterations.

#### A. LLM-driven Temporal Semantics Inference

Human behavior unfolds as a continuous visual process, where adjacent frames in a motion sequence often exhibit high semantic correlation. To identify segments representing the same human action, we introduce the concept of TVS, which delineates the temporal neighborhood of each frame according to semantic consistency. Specifically, for a given frame  $\mathbf{x}_i$ , we aim to discover its left and right temporal neighbor bounds  $(l_i, r_i)$  that enclose all frames depicting the same motion as  $\mathbf{x}_i$ .

Unsupervised discovery of such temporal neighborhoods is challenging using traditional machine learning methods, as the semantic boundary between motions is difficult to define purely through pixel comparisons. To address this challenge, we harness the visual reasoning capability of a LLM as a zero-shot semantic comparator. Instead of training an additional network, the LLM is instructed with a natural-language prompt to assess whether two consecutive frames represent the same human motion:

*“Do these two neighboring frames depict the same human motion? Answer Yes or No.”*

Given two adjacent frames  $(\mathbf{x}_t, \mathbf{x}_{t+1})$ , the LLM produces a binary response Yes/No, which is recorded as a Boolean variable  $\text{eq}_t \in \{0, 1\}$  indicating whether the two frames belong

---

#### Algorithm 1: Learning TVS via LLM

---

```

1 Input: Raw RGB frames sequence  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ .
2 Output: TVS matrix  $\mathbf{G} \in \mathbb{R}^{N \times N}$  and the set  $\{\mathcal{N}_i\}_{i=1}^N$ 
   1: Initialize  $\text{eq} \leftarrow \emptyset$ ;
   2: for  $i = 1$  to  $N - 1$  do
   3:   Encode images  $\mathbf{x}_i, \mathbf{x}_{i+1}$  to base64 formats separately;
   4:   Query LLM with prompt: “Do these two
     neighbouring frames depict the same human motion?
     Answer Yes or No.”;
   5:   Parse response to  $\text{eq}_i \in \{0, 1\}$  (Yes  $\rightarrow 1$ , No  $\rightarrow 0$ ).
   6:   if response ambiguous then re-query with stricter
     instruction (“Answer strictly with a single token: YES
     or NO.”) and re-parse.
   7:   Append  $\text{eq}_i$  to list  $\text{eq}$ .
   8: end for
   9: for  $i = 1$  to  $N$  do
  10:   Initialize  $l_i \leftarrow i$ ;  $r_i \leftarrow i$ .
  11:   while  $l_i > 1$  and  $\text{eq}_{l_i-1} = 1$  do  $l_i \leftarrow l_i - 1$ .
  12:   while  $r_i < N$  and  $\text{eq}_{r_i} = 1$  do  $r_i \leftarrow r_i + 1$ .
  13: end for
  14: Initialize  $\mathbf{G} \leftarrow \mathbf{0}_{N \times N}$ .
  15: for  $i = 1$  to  $N$  do
  16:    $\mathcal{N}_i \leftarrow \{j \mid j \in [l_i, r_i], j \neq i\}$ ;
  17:    $G_{ii} \leftarrow -|\mathcal{N}_i|$ ;
  18:   for each  $j \in \mathcal{N}_i$  do
  19:      $G_{ij} \leftarrow 1$ .
  20: end for
```

---

to the same motion segment. The sequence  $\{\text{eq}_1, \dots, \text{eq}_{N-1}\}$  forms the adjacency pattern of temporal consistency across the sequence. In our implementation, this is achieved through an API call to a multimodal LLM (e.g., GPT-4o, Gemini-2.0, or Claude-4.5), where both frames are provided in base64-encoded format, allowing the model to reason directly over image content.

Using the response  $\{\text{eq}_1, \dots, \text{eq}_{N-1}\}$ , we define the left and right temporal neighbor bounds for each frame  $\mathbf{x}_i$  as follows:

$$l_i = \min\{j \mid j \leq i, \mathbf{x}_j, \mathbf{x}_{j+1}, \dots, \mathbf{x}_i \text{ describe the same motion}\},$$

$$r_i = \max\{j \mid j \geq i, \mathbf{x}_i, \mathbf{x}_{i+1}, \dots, \mathbf{x}_j \text{ describe the same motion}\}.$$

Accordingly, the temporal neighborhood set  $\mathcal{N}_i$  of  $\mathbf{x}_i$  is given by

$$\mathcal{N}_i = \{j \mid j \in \{l_i, l_i + 1, \dots, r_i\}, j \neq i, j \in \mathbb{Z}^+\}.$$

Each frame  $\mathbf{x}_i$  is thus associated with a temporally and semantically coherent segment.

This structure serves as a foundation for constructing a TVS matrix  $\mathbf{G} \in \mathbb{R}^{N \times N}$ , whose entries encode the neighborhood connectivity as

$$G_{i,j} = \begin{cases} -|\mathcal{N}_i|, & \text{if } i = j, \\ 1, & \text{if } j \in \mathcal{N}_i, \\ 0, & \text{otherwise.} \end{cases}$$

In practice, the TVS matrix is implemented as a Laplacian-like structure, where the diagonal term penalizes the number



of semantic neighbors, and the off-diagonal entries reflect temporal affinity.

Algorithm 1<sup>1</sup> summarizes the complete computational procedure. For each pair of adjacent frames  $(i, i + 1)$ , the LLM is queried once and the response recorded. Subsequently, left and right neighbor bounds  $(l_i, r_i)$  are determined through a recursive traversal of the Boolean adjacency list. Finally, the TVS matrix  $\mathbf{G}$  is constructed and saved for downstream processing.

**Remark 1.** While the TVS introduces human-like temporal reasoning into motion segmentation, it only captures pairwise relationships between temporal consecutive frames, lacking global temporal dependencies. Therefore, a subsequent grouping stage is required to achieve globally consistent motion segmentation.

### B. Subspace embedding and Clustering Incorporating Vision Temporal Semantics

By leveraging matrix multiplication, we observe that the product  $\mathbf{ZG}$  captures the similarity error between the representation of a given sequential point and its neighbors. Specifically, the term

$$\mathbf{ZG} = \left[ \sum_{l \in \mathcal{N}_1} (\mathbf{z}_1 - \mathbf{z}_l), \sum_{l \in \mathcal{N}_2} (\mathbf{z}_2 - \mathbf{z}_l), \dots, \sum_{l \in \mathcal{N}_N} (\mathbf{z}_N - \mathbf{z}_l) \right].$$

measures the similarity of the  $i$ th data and its neighbors defined by  $\mathcal{N}_i$ . To encourage the subspace embedding of the subspace embedding and the embedding of its neighbors to be as similar as possible, we introduce a structural regularization term  $\|\mathbf{ZG}\|_{2,1}$ , where  $\|\cdot\|_{2,1}$  denotes the  $l_1$  norm of the vector formed by the  $l_2$  norms of each column of the matrix. This norm encourages the columns of  $\mathbf{ZG}$  to exhibit consistent behavior across neighboring points, promoting smoothness in the subspace representation. Mathematically, we express it as

$$\|\mathbf{ZG}\|_{2,1} = \sum_{i=1}^N \left\| \sum_{l \in \mathcal{N}_i} (\mathbf{z}_i - \mathbf{z}_l) \right\|_2 = \sum_{i=1}^N \sum_{l \in \mathcal{N}_i} \|\mathbf{z}_i - \mathbf{z}_l\|_2.$$

We aim to minimize the subspace embedding error  $\|\mathbf{ZG}\|_{2,1}$  to enhance the temporal consistency of the subspace embedding by promoting coherence between a subspace embedding and its neighboring subspace embeddings, which is crucial for capturing the dynamic temporal structure of the data.

**Theorem 1** (Interpretation of the TVS Regularizer). *The regularizer  $\|\mathbf{ZG}\|_{2,1}$  represents an isotropic graph total variation over the temporal graph defined by  $\{\mathcal{N}_i\}$ . Minimizing it enforces local smoothness within temporal neighborhoods while preserving discontinuities at motion boundaries, thus producing piecewise-constant embeddings consistent with human motion transitions.*

*Proof.* See Appendix A.  $\square$

**Theorem 2** (Consistency under Noisy LLM Adjacency). *If each LLM adjacency label is independently flipped with probability  $p < \frac{1}{2}$  and each motion segment has length at least  $L_{\min}$ , then the expected number of erroneous TVS boundaries scales as  $O(pN)$ . Minimizing  $\|\mathbf{ZG}\|_{2,1}$  yields piecewise-constant embeddings that smooth out isolated errors, ensuring segment-level consistency in expectation when  $p$  is small and  $L_{\min}$  is sufficiently large.*

*Proof.* See Appendix B.  $\square$

Theorem 2 implies that the proposed framework is robust to occasional LLM misjudgments: although local adjacency errors may occur, the TVS-induced regularizer preserves overall temporal coherence by enforcing smooth embeddings within segments. Thus, the method maintains consistent human motion segmentation under moderate annotation noise.

We also propose a TVS-integrated segmentation on the subspace embedding. Specifically, suppose the number of clusters is  $K$ . We introduce a cluster assignment indicator vector  $\mathbf{q}_i \in \mathbb{R}^K$  for the  $i$ -th frame, where the  $k$ -th element is set to 1 if the  $i$ -th frame is assigned to the  $k$ -th cluster, and all other elements are set to zero. We then define a clustering regularizer based on the subspace embedding  $\mathbf{Z}$  and the indicator matrix  $\mathbf{Q} = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_N] \in \mathbb{R}^{N \times K}$ :

$$\begin{aligned} \|\mathbf{Z}\|_{\mathbf{Q}} &= \frac{1}{2} \sum_{i,j} \frac{|Z_{i,j}| + |Z_{j,i}|}{2} \|\mathbf{q}_i - \mathbf{q}_j\|_2^2 \\ &= \frac{1}{2} \sum_{i,j} |Z_{i,j}| \cdot \|\mathbf{q}_i - \mathbf{q}_j\|_2^2 = \sum_{i,j} |Z_{i,j}| \cdot \frac{\|\mathbf{q}_i - \mathbf{q}_j\|_2^2}{2} \\ &= \sum_{i,j} |Z_{i,j}| \cdot \Theta_{i,j} = \sum_{i,j} |(\Theta \odot \mathbf{Z})_{ij}| = \|\Theta \odot \mathbf{Z}\|_1 \end{aligned}$$

where  $Z_{i,j}$  is the  $(i, j)$ -th element of  $\mathbf{Z}$ , and  $\Theta_{i,j} = \frac{\|\mathbf{q}_i - \mathbf{q}_j\|_2^2}{2}$ . The first equation ensures symmetry by combining both  $|Z_{i,j}|$  and  $|Z_{j,i}|$ , accounting for the interactions between off-diagonal terms, as  $Z_{i,j} \neq Z_{j,i}$  does not necessarily hold, thus incorporating these contributions into the final regularizer. The term  $\Theta_{i,j}$  measures the squared Euclidean distance between the cluster indicators  $\mathbf{q}_i$  and  $\mathbf{q}_j$ , normalized by a factor of  $\frac{1}{2}$ , capturing the dissimilarity between frames based on their clustering assignments and enforcing smoothness within clusters. The final expression  $\|\Theta \odot \mathbf{Z}\|_1$  represents the  $l_1$ -norm of the element-wise product between  $\Theta$  and  $\mathbf{Z}$ , which encourages a sparse representation of the subspace embedding while optimizing the clustering assignments, ensuring that frames assigned to the same cluster exhibit more similar representations. The term  $\|\mathbf{Z}\|_{\mathbf{Q}}$  will be minimized to optimize the clustering assignment.

To ensure that  $\mathbf{q}_i$  functions effectively as a cluster assignment indicator, we impose the condition that  $\mathbf{Q}$  is a subset of

$$\mathcal{Q} = \{\mathbf{Q} \in \{0, 1\}^{N \times K} : \mathbf{Q}\mathbf{1}_{K \times 1} = \mathbf{1}_{N \times 1}, \mathbf{q}_i = \mathbf{q}_j \forall j \in \mathcal{N}_i\}.$$

$\square$  This constraint ensures that the clustering assignment for the  $i$ -th frame and its neighbors are identical, i.e.,  $\mathbf{q}_i = \mathbf{q}_j \forall j \in \mathcal{N}_i$ , which enforces temporal consistency within cluster assignments.

<sup>1</sup><https://github.com/y66y/TVSH>

Building on the traditional subspace clustering formulation in (1), we develop a feedback-enabled framework that integrates the proposed subspace embedding, which incorporates temporal vision semantics, with the proposed clustering method. The optimization problem is formulated as follows:

$$\begin{aligned} & \underset{\mathbf{Z}, \mathbf{Q}}{\text{minimize}} \quad \|\mathbf{X} - \mathbf{XZ}\|_F^2 + \|\mathbf{Z}^T \mathbf{Z}\|_1 + \|\mathbf{ZG}\|_{2,1} + \|\mathbf{Z}\|_Q \\ & \text{subject to} \quad \text{diag}(\mathbf{Z}) = 0, \mathbf{Z} \geq 0, \mathbf{Q} \in \mathcal{Q}. \end{aligned} \quad (2)$$

**Proposition 1.** *Let  $\mathbf{X} \in \mathbb{R}^{D \times N}$  be a matrix whose columns are drawn from a union of  $K$  distinct subspaces, with the subspace assignment indicated by  $\mathbf{Q}^*$ . The optimal solution to the problem in (2) is given by  $\mathbf{Q}^*$  and  $\mathbf{Z}^*$ , where  $\mathbf{Z}^*$  is block-diagonal after permuted according to  $\mathbf{Q}^*$ .*

*Proof.* See Appendix C.  $\square$

Proposition 1 demonstrates that, under the assumption that frames are distributed across distinct subspaces, the optimal solution to problem (2) will align with the true segmentation. However, due to the influence of noise, the human motion data may not lie perfectly within the subspaces. Thus, there are inherent trade-offs in (2) due to practical dataset challenges such as image noise and subtle motions, which may cause frames not to align precisely with  $K$  subspaces. Specifically, the term  $\|\mathbf{X} - \mathbf{XZ}\|_F^2 + \|\mathbf{Z}^T \mathbf{Z}\|_1$  ensures sparsity and accurate data subspace embedding in the outputs  $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N$ , promoting distinctiveness among them. In contrast, the term  $\|\mathbf{ZG}\|_{2,1}$  requires these outputs to align with their neighbors' coefficients  $\{\mathbf{z}_l\}_{l \in \mathcal{N}_i}$ . This necessitates a balance between representing data across  $K$  clusters and maintaining temporal, aiming to segment the data sequence into smaller segments where, for instance, in a segment  $[\mathbf{x}_i, \mathbf{x}_{i+1}, \dots, \mathbf{x}_j]$  with  $j > i$ , the subspace embedding are identical, i.e.,  $\mathbf{z}_i = \mathbf{z}_{i+1} = \dots = \mathbf{z}_j$ . Additionally, the term  $\|\mathbf{Z}\|_Q$  promotes effective grouping based on  $\mathbf{Z}$  while adhering to the constraint  $\mathbf{Q} \in \mathcal{Q}$ , which stipulates that the cluster assignment of the  $i$ th frame must match that of its neighbors, introducing a further trade-off between dependent clustering and TVS considerations.

**Theorem 3** (Impact of TVS on Segmentation). *Assume each motion segment generates data lying in one of  $K$  linear subspaces with within-segment variance  $\sigma^2$  and between-subspace separation  $\Delta_{\text{sub}}^2 > 0$ . With independent LLM adjacency errors of rate  $p < \frac{1}{2}$ , the expected segmentation error satisfies  $\mathbb{E}[\text{Err}_{\text{HMS}}] \leq C_1 \frac{p}{L_{\min}} + C_2 \frac{\sigma^2}{\Delta_{\text{sub}}^2}$ . When TVS boundaries align with true actions, the optimal solution  $\mathbf{Z}^*$  becomes block-diagonal, achieving exact segmentation.*

*Proof.* See Appendix D.  $\square$

Theorem 3 establishes that TVS improves segmentation robustness by suppressing random adjacency errors and stabilizing intra-segment embeddings. The first term  $C_1 \frac{p}{L_{\min}}$  reflects the resilience to LLM-induced boundary noise, which diminishes as segment length increases, while the second term  $C_2 \frac{\sigma^2}{\Delta_{\text{sub}}^2}$  captures the dependence on subspace separability. Perfectly aligned TVS boundaries yield theoretically exact segmentation, confirming the effectiveness of LLM-guided

temporal reasoning in enhancing motion boundary localization.

### C. A Feedback-Enabled Optimization Algorithm

We employ the ADMM method [89] to solve the optimization problem formulated in (2). In order to separate the third term in (2) from the other three terms, we introduce an additional variable  $\mathbf{H} = \mathbf{ZG}$ . By incorporating an augmented Lagrangian multiplier to handle the introduced linear constraint, we can reformulate (2) as the following problem:

$$\begin{aligned} & \underset{\mathbf{Z}, \mathbf{H}, \mathbf{Q}}{\text{minimize}} \quad \|\mathbf{X} - \mathbf{XZ}\|_F^2 + \|\mathbf{Z}^T \mathbf{Z}\|_1 + \|\mathbf{H}\|_{2,1} \\ & \quad + \langle \mathbf{F}, \mathbf{H} - \mathbf{ZG} \rangle + \frac{\gamma}{2} \|\mathbf{H} - \mathbf{ZG}\|_F^2 + \|\mathbf{Z}\|_Q \\ & \text{subject to} \quad \text{diag}(\mathbf{Z}) = 0, \mathbf{Z} \geq 0, \mathbf{Q} \in \mathcal{Q} \end{aligned} \quad (3)$$

where  $\mathbf{F} \in \mathbb{R}^{N \times N}$  is the Lagrangian multiplier and  $\gamma$  is an adaptive weight parameter for enforcing the condition  $\mathbf{H} = \mathbf{ZG}$ . To solve (3), we adopt a feedback-enabled optimization strategy, where we iteratively solve three sub-problems for  $\mathbf{Z}$ ,  $\mathbf{H}$ , and  $\mathbf{Q}$  while keeping the other fixed, respectively.

1) **Z-solution:** Fixing  $\mathbf{H}$  and  $\mathbf{Q}$ , solve for  $\mathbf{Z}$  by

$$\begin{aligned} & \underset{\mathbf{Z}}{\text{minimize}} \quad \|\mathbf{X} - \mathbf{XZ}\|_F^2 + \|\mathbf{Z}^T \mathbf{Z}\|_1 + \langle \mathbf{F}, \mathbf{H} - \mathbf{ZG} \rangle \\ & \quad + \frac{\gamma}{2} \|\mathbf{H} - \mathbf{ZG}\|_F^2 + \|\mathbf{Z}\|_Q \\ & \text{subject to} \quad \text{diag}(\mathbf{Z}) = 0, \mathbf{Z} \geq 0 \end{aligned} \quad (4)$$

Since  $\mathbf{Z}$  consists of non-negative elements, we can rewrite the objective function in (4) as a function:

$$\begin{aligned} \mathcal{J}(\mathbf{Z}) = & \|\mathbf{X} - \mathbf{XZ}\|_F^2 + \mathbf{e}^T \mathbf{Z}^T \mathbf{Z} \mathbf{e} + \langle \mathbf{F}, \mathbf{H} - \mathbf{ZG} \rangle \\ & + \frac{\gamma}{2} \|\mathbf{H} - \mathbf{ZG}\|_F^2 + \|\mathbf{\Theta} \odot \mathbf{Z}\|_1 \end{aligned} \quad (5)$$

The sub-problem defined in (4) can be formulated as a convex quadratic programming problem with specific constraints for the variable  $\mathbf{Z}$ , involving the function  $\mathcal{J}(\mathbf{Z})$  from (5). In this problem, we aim to minimize  $\mathcal{J}(\mathbf{Z})$  while satisfying the given constraints. To tackle this, we employ the projected gradient method, which is a well-established approach known for its simplicity and effectiveness in solving such problems. This method is chosen as our preferred solution due to its suitability for our problem's requirements.

Consider the partial derivative of  $\|\mathbf{\Theta} \odot \mathbf{Z}\|_1$  with respect to each element  $Z_{ij}$ :

$$\frac{\partial}{\partial Z_{ij}} \left( \sum_{k,l} |\Theta_{kl} Z_{kl}| \right) = \frac{\partial}{\partial Z_{ij}} |\Theta_{ij} Z_{ij}|$$

Using the properties of the absolute value function, we get:

$$\frac{\partial}{\partial Z_{ij}} |\Theta_{ij} Z_{ij}| = \Theta_{ij} \cdot \text{sign}(\Theta_{ij} Z_{ij})$$

where  $\text{sign}(x)$  is the sign function, defined as:

$$\text{sign}(x) = \begin{cases} 1, & \text{if } x > 0 \\ 0, & \text{if } x = 0 \\ -1, & \text{if } x < 0 \end{cases}$$

**Algorithm 2:** TVSH method.

---

**1 Input:**  $\mathbf{X}$ .  
**2 Output:**  $\{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_K\}$

- 1: Initialize  $\mathbf{G}, \mathbf{F} = \mathbf{1}, \rho = 1.1, \gamma = 0.1$ .  $\mathbf{H} = \mathbf{Z}\mathbf{G}$  where  $\mathbf{Z}$  is the similarity matrix given by cosine measurement.  $\mathbf{Q}$  is initialized by K-means [11].
- 2: **repeat**
- 3: Find  $\mathbf{Z}$  by solving (6).
- 4: Calculate the projection  $\mathbf{Z} \leftarrow \prod_{\mathcal{Z}}(\mathbf{Z})$  by solving (7).
- 5: Find  $\mathbf{H}$  by solving (8);
- 6: Update  $\mathbf{F} \leftarrow \mathbf{F} + \gamma(\mathbf{H} - \mathbf{Z}\mathbf{G}), \gamma \leftarrow \rho\gamma$ .
- 7: Update  $\mathbf{Q}$  by the following steps:
- 8: Calculate  $\mathbf{L} = \mathbf{I} - \mathbf{D}^{-1/2}[(|\mathbf{Z}| + |\mathbf{Z}^T|)/2]\mathbf{D}^{-1/2}$ .
- 9: Compute the smallest  $K$  eigenvectors of  $\mathbf{L}$  denoted by  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_K]$ . Denote the row of  $\mathbf{V}$  as  $\{\mathbf{u}_i\}_{i=1}^N$ .
- 10: **repeat**
- 11: **for**  $k = 1$  **to**  $K$  **do**
- 12: Form weight matrix  $\mathbf{W}_k$  and calculate  $\mu_k$ .
- 13: **end for**
- 14: **for**  $k = 1$  **to**  $K$  **do**
- 15:  $\mathcal{C}_k \leftarrow \{i \in \{1, 2, \dots, N\} : k = \arg \min_{k \in \{1, 2, \dots, K\}} \|\mathbf{u}_i - \mu_k\|_2^2 w_{k,i}\}$
- 16: **end for**
- 17: **until**  $\mathbf{Q}$  can not be changed.
- 18: **until** The objective function value of (3) can not be decreased.

---

Therefore, the partial derivative for each element  $Z_{ij}$  is:

$$\frac{\partial \|\Theta \odot \mathbf{Z}\|_1}{\partial Z_{ij}} = \Theta_{ij} \cdot \text{sign}(\Theta_{ij} Z_{ij})$$

Combining all the partial derivatives into matrix form, the gradient of  $\|\Theta \odot \mathbf{Z}\|_1$  with respect to  $\mathbf{Z}$  is:

$$\frac{\partial \|\Theta \odot \mathbf{Z}\|_1}{\partial \mathbf{Z}} = \Theta \odot \text{sign}(\Theta \odot \mathbf{Z})$$

where  $\text{sign}(\Theta \odot \mathbf{Z})$  is the matrix obtained by applying the sign function element-wise to  $\Theta \odot \mathbf{Z}$ .

The derivative of  $\mathcal{J}(\mathbf{Z})$  with respect to  $\mathbf{Z}$  can be expressed as:  $\partial \mathcal{J}(\mathbf{Z}) = -2\mathbf{X}^T(\mathbf{X} - \mathbf{X}\mathbf{Z}) + 2\mathbf{Z}\mathbf{E} - \mathbf{F}\mathbf{G}^T - \gamma(\mathbf{H} - \mathbf{Z}\mathbf{G})\mathbf{G}^T + \Theta \odot \text{sign}(\Theta \odot \mathbf{Z})$  where  $\mathbf{E} \in \mathbb{R}^{N \times N}$  is an all-one matrix.

Setting the derivative to zero gives

$$2\mathbf{X}^T\mathbf{X}\mathbf{Z} + \mathbf{Z}(2\mathbf{E} + \gamma\mathbf{G}\mathbf{G}^T) = \mathbf{F}\mathbf{G}^T + \gamma\mathbf{H}\mathbf{G}^T + 2\mathbf{X}^T\mathbf{X} - \Theta \odot \text{sign}(\Theta \odot \mathbf{Z}). \quad (6)$$

The equation presented is a well-known Sylvester equation in the form  $\mathbf{A}\mathbf{Z} + \mathbf{Z}\mathbf{B} = \mathbf{C}$ , where  $\mathbf{A} = 2\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{B} = 2\mathbf{E} + \gamma\mathbf{G}\mathbf{G}^T$ , and  $\mathbf{C} = \mathbf{F}\mathbf{G}^T + \gamma\mathbf{H}\mathbf{G}^T + 2\mathbf{X}^T\mathbf{X} - \Theta \odot \text{sign}(\Theta \odot \mathbf{Z})$ .

We adopt Bartels-Stewart algorithm [90] to solve  $\mathbf{Z}$ . Specifically, we first perform Schur decomposition on  $\mathbf{A}$  and  $\mathbf{B}$ . The Schur decomposition of  $\mathbf{A}$  and  $\mathbf{B}$  is given by  $\mathbf{A} = \mathbf{Q}_A \mathbf{T}_A \mathbf{Q}_A^H$  and  $\mathbf{B} = \mathbf{Q}_B \mathbf{T}_B \mathbf{Q}_B^H$ , where  $\mathbf{Q}_A$  and  $\mathbf{Q}_B$  are unitary matrices and  $\mathbf{T}_A$  and  $\mathbf{T}_B$  are upper triangular matrices. Then, we use the unitary matrices from the Schur decomposition to

transform  $\mathbf{C}$  to  $\mathbf{C}' = \mathbf{Q}_A^H \mathbf{C} \mathbf{Q}_B$ . Next, we solve the simplified equation  $\mathbf{T}_A \mathbf{Z}' + \mathbf{Z}' \mathbf{T}_B = \mathbf{C}'$ . This can be done using a back-substitution method since  $\mathbf{T}_A$  and  $\mathbf{T}_B$  are upper triangular matrices. Finally, we transform  $\mathbf{Z}'$  back to  $\mathbf{Z}$  by  $\mathbf{Z} = \mathbf{Q}_A \mathbf{Z}' \mathbf{Q}_B^H$ .

However, it's worth noting that the critical frame  $\mathbf{Z}$  of the objective function may not necessarily lie within the feasible set defined in (4). To address this, we can employ a projection operator to find a feasible frame starting from the critical frame  $\mathbf{Z}$ .

$$\prod_{\mathcal{Z}}(\mathbf{Z}) = \arg \min_{\tilde{\mathbf{Z}} \in \mathcal{Z}} \|\tilde{\mathbf{Z}} - \mathbf{Z}\|_F^2 \quad (7)$$

where  $\mathcal{Z} = \{\mathbf{Z} | \text{diag}(\mathbf{Z}) = 0, \mathbf{Z} \geq 0\}$ . For a simple and quick solution to (7), we implement the projection operator  $\prod_{\mathcal{Z}}(\mathbf{Z})$  as follows:

$$z_{ij}^* = \begin{cases} z_{ij} & \text{if } z_{ij} \geq 0 \text{ and } i \neq j \\ 0 & \text{if } z_{ij} < 0 \text{ or } i = j \end{cases}$$

where  $z_{ij}$  and  $z_{ij}^*$  are the elements of  $\mathbf{Z}$  and its projection  $\prod_{\mathcal{Z}}(\mathbf{Z})$ , respectively.

2) **H-solution:** Fixing  $\mathbf{Z}$  and  $\mathbf{Q}$ , solve for  $\mathbf{H}$  by

$$\begin{aligned} \underset{\mathbf{H}}{\text{minimize}} \quad & \|\mathbf{H}\|_{2,1} + \langle \mathbf{F}, \mathbf{H} - \mathbf{Z}\mathbf{G} \rangle \\ & + \frac{\mathbf{H} - \mathbf{Z}\mathbf{G}}{2} \|\mathbf{H} - \mathbf{Z}\mathbf{G}\|_F^2 \end{aligned} \quad (8)$$

which is equivalent to minimizing  $\|\mathbf{H}\|_{2,1} + \frac{\gamma}{2} \|\mathbf{H} - (\mathbf{Z}\mathbf{G} - (1/\gamma)\mathbf{F})\|_F^2$  with respect to  $\mathbf{H}$ . Denote  $\mathbf{P} = \mathbf{Z}\mathbf{G} - (1/\gamma)\mathbf{F}$ . Then the closed-form solution to (8) will be given as follows [91]:

$$\mathbf{h}_i = \begin{cases} \frac{\|\mathbf{p}_i\| - (1/\gamma)}{\|\mathbf{p}_i\|} \mathbf{p}_i & \text{if } \|\mathbf{p}_i\| > 1/\gamma \\ 0 & \text{otherwise} \end{cases}$$

where  $\mathbf{h}_i, \mathbf{p}_i$  are the  $i$ th column of  $\mathbf{H}, \mathbf{P}$ , respectively.

**Proposition 2** (Optimality). *For fixed  $(\mathbf{Z}, \mathbf{Q}, \mathbf{F}, \gamma)$  the subproblem (8) is the proximal operator of the  $\ell_{2,1}$  norm and admits the closed-form group-shrinkage solution. Hence the  $\mathbf{H}$ -update attains the global minimizer of (8) at every iteration.*

*Proof.* See Appendix E.  $\square$

3) **Q-solution:** Fixing  $\mathbf{Z}$  and  $\mathbf{H}$ , solve for  $\mathbf{Q}$  by

$$\min_{\mathbf{Q}} \|\mathbf{Z}\|_{\mathbf{Q}}, \quad \text{subject to } \mathbf{Q} \in \mathcal{Q} \quad (9)$$

**Proposition 3.** *We have the following equivalent problem*

$$\min_{\mathbf{Q}} \|\mathbf{Z}\|_{\mathbf{Q}} \iff \min_{\mathbf{Q}} \text{Trace}(\mathbf{Q}^T (\mathbf{D} - (|\mathbf{Z}| + |\mathbf{Z}^T|)/2) \mathbf{Q})$$

where the matrix  $\mathbf{D}$  is known as the degree matrix. The degree matrix  $\mathbf{D}$  is defined as:  $D_{ii} = \sum_{j=1}^N [(|\mathbf{Z}| + |\mathbf{Z}^T|)/2]_{ij}$ . For all off-diagonal elements  $i \neq j$ ,  $D_{ij} = 0$ .

*Proof.* See Appendix F.  $\square$

The objective function in (9) is the traditional normalized cut clustering problem [92] with a TVS constraint. The Laplacian matrix  $\mathbf{L}$  can be computed using the formula  $\mathbf{L} = \mathbf{I} - \mathbf{D}^{-1/2}[(|\mathbf{Z}| + |\mathbf{Z}^T|)/2]\mathbf{D}^{-1/2}$ , where  $\mathbf{I}$  is an identity matrix. Consequently, the eigenvectors  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_K$  corresponding to the first  $K$  smallest eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_K$  of  $\mathbf{L}$  are

computed, satisfying  $\mathbf{L}\mathbf{v}_k = \tilde{\lambda}_k \mathbf{v}_k$ . These eigenvectors are arranged as columns in a matrix  $\mathbf{V} \in \mathbb{R}^{N \times K}$ .

Let  $\mathbf{u}_i \in \mathbb{R}^K$  represent the vector of the  $i$ th row of  $\mathbf{V}$ , where  $i = 1, \dots, N$ . The problem (9) can be relaxed to the following form:

$$\underset{\{\mathcal{C}_k, \boldsymbol{\mu}_k\}_{k=1}^K}{\text{minimize}} \sum_{k=1}^K \sum_{i \in \mathcal{C}_k} \|\mathbf{u}_i - \boldsymbol{\mu}_k\|_2^2, \quad \text{subject to } \mathbf{Q} \in \mathcal{Q} \quad (10)$$

where  $i \in \mathcal{C}_k$  if  $q_{i,k} = 1$ . However, the requirement  $\mathbf{q}_i = \mathbf{q}_j \forall j \in \mathcal{N}_i$  in the constraint  $\mathbf{Q} \in \mathcal{Q}$  makes solving problem (10) highly challenging. Since the constraint mandates that the clustering assignments of the  $i$ th frame and its neighbors remain consistent, we relax the constraint to that the clustering center corresponding to the  $i$ th frame should coincide with the center of its neighbors. This leads us to the formulation of the following problem:

$$\underset{\{\mathcal{C}_k, \boldsymbol{\mu}_k\}_{k=1}^K}{\text{minimize}} \sum_{k=1}^K \sum_{i \in \mathcal{C}_k} \left( \|\mathbf{u}_i - \boldsymbol{\mu}_k\|_2^2 + \eta \sum_{j \in \mathcal{N}_i} \|\mathbf{u}_j - \boldsymbol{\mu}_k\|_2^2 \right) \quad (11)$$

where the penalty coefficient  $\eta$  is set to  $1/\mathcal{N}_i$  for weight balance.

The term  $\sum_{k=1}^K \sum_{i \in \mathcal{C}_k} \|\mathbf{u}_i - \boldsymbol{\mu}_k\|_2^2$  aims to independently fit all the data with the center  $\{\boldsymbol{\mu}_k\}$ . However, the term  $\eta \sum_{k=1}^K \sum_{i \in \mathcal{C}_k} \sum_{j \in \mathcal{N}_i} \|\mathbf{u}_j - \boldsymbol{\mu}_k\|_2^2$  desires the  $\boldsymbol{\mu}_k$  to be identical temporally, i.e., the center of  $\mathbf{u}_i$  is the same as the center of  $\mathbf{u}_j$  for any  $j \in \mathcal{N}_i$ . Consequently, minimizing these two terms simultaneously leads to a trade-off between fitting data to  $K$  centers and maintaining temporal of the center assignment, where the desired outcome is to divide the data sequence into multiple small segments.

It is still challenging to solve problem (11) directly due to its NP-hard nature. We first propose the following proposition, which will be utilized to adapt problem (11) into a new form.

**Proposition 4.** *The term  $\sum_{i \in \mathcal{C}_k} (\|\mathbf{u}_i - \boldsymbol{\mu}_k\|_2^2 + \eta \sum_{j \in \mathcal{N}_i} \|\mathbf{u}_j - \boldsymbol{\mu}_k\|_2^2)$  in (11) is equivalent to  $\sum_{i=1}^N \|\mathbf{u}_i - \boldsymbol{\mu}_k\|_2^2 (\mathbb{1}(i \in \mathcal{C}_k) + \eta n_k(i))$  where  $n_k(i)$  is the number of times the frame  $\mathbf{u}_i$  appears as a sequential neighbor of a frame in the  $k$ -th cluster, i.e.,  $n_k(i) = \sum_{j \in \mathcal{C}_k} \mathbb{1}(i \in \mathcal{N}_j)$  and the indicator function  $\mathbb{1}(s) = 1$  if  $s$  is true and zero otherwise.*

*Proof.* See Appendix G.  $\square$

According to proposition 4, problem (11) can be rewritten as the following new weighted problem:  $\underset{\{\mathcal{C}_k, \boldsymbol{\mu}_k\}_{k=1}^K}{\text{minimize}} \sum_{k=1}^K \sum_{i=1}^N \|\mathbf{u}_i - \boldsymbol{\mu}_k\|_2^2 w_{k,i}$ , where  $w_{k,i} = \mathbb{1}(i \in \mathcal{C}_k) + \eta n_k(i)$ . Observing that the new weighted problem can be solved by addressing two sub-problems for  $\mathcal{C}_k$  and  $\boldsymbol{\mu}_k$  in an alternating manner when one is fixed, respectively, we first focus on solving the new problem with the given cluster assignment  $\{\mathcal{C}_k\}_{k=1}^K$ . Denote the objective function of the new weighted problem as  $\mathcal{J}_1(\{\boldsymbol{\mu}_k\}_{k=1}^K)$ . If we take the derivative of  $\mathcal{J}_1(\{\boldsymbol{\mu}_k\}_{k=1}^K)$  with respect to  $\boldsymbol{\mu}_k$  and set it to zero, i.e.,  $\frac{\partial \mathcal{J}_1(\{\boldsymbol{\mu}_k\}_{k=1}^K)}{\partial \boldsymbol{\mu}_k} = 0$ , we obtain  $\boldsymbol{\mu}_k = \frac{1}{\sum_{i=1}^N w_{k,i}} \sum_{i=1}^N w_{k,i} \mathbf{u}_i$ . We then solve the cluster assignment with the given cluster center. This is done by evaluating the weighted combination of the residual

from the frame to a given center, as well as the residuals of its sequential neighbors, so that the estimated cluster label for the frame  $\mathbf{u}_i$  is assigned to the  $l$ -th cluster, where  $l = \arg \min_{k \in \{1, 2, \dots, K\}} \|\mathbf{u}_i - \boldsymbol{\mu}_k\|_2^2 w_{k,i}$ .

The algorithm alternates between center update and cluster assignment steps until convergence. In the center update step, the resulting center represents the global optimum given a cluster assignment. This step learns the center that minimizes the distance to all frames in the cluster, including their sequential neighbors. Therefore, the center update step cannot increase the overall objective function. Similarly, in the cluster assignment step, each frame is assigned to the cluster that minimizes the distance to itself and its sequential neighbors, which also cannot increase the overall objective function. Since there is a finite number of ways the frames can be assigned, and the objective function in the new weighted problem is bounded below by zero, the proposed alternating algorithm must terminate at a locally optimal clustering result. To determine the number of clusters  $K$ , we use the silhouette score, which measures the similarity of a sample point to its own cluster in comparison to the nearest cluster. By calculating the silhouette score for different values of  $K$ , the optimal number of clusters is chosen as the value of  $K$  that maximizes the silhouette score.

By iteratively solving (4), (8), and (9), we can obtain a solution to (3). During this process, we group the subspace embeddings by solving (9) and update the embeddings based on feedback from the HMS solution of (9). The convergence of sub-problem (4), the closed-form solution of sub-problem (8), and the convergence of solving (9) ensure the overall convergence of the algorithm for (3). Algorithm 2 presents the pseudocode for our clustering method.

**Theorem 4** (Convergence). *Under bounded and lower-semicontinuous augmented Lagrangian, nondecreasing  $\gamma_t \rightarrow \gamma_\infty \in (0, \infty)$ , and bounded  $\rho > 1$ , the proposed ADMM-based alternating scheme ensures monotonic decrease of the objective and convergence of  $(\mathbf{Z}^{(t)}, \mathbf{H}^{(t)}, \mathbf{Q}^{(t)})$  to a first-order stationary point. If each  $\mathbf{Q}$ -update reaches its relaxed global optimum, every accumulation point satisfies the KKT conditions.*

*Proof.* See Appendix H.  $\square$

This theorem confirms that the alternating optimization is theoretically stable and convergent: the objective value decreases monotonically, the iterates approach a stationary solution, and, with exact subproblem updates, the algorithm attains KKT-level optimality, guaranteeing reliable convergence behaviour in practice.

#### IV. EXPERIMENTAL RESULTS AND ANALYSIS

In this section, we first introduce the human motion datasets used in our experiments (Section IV-A). We then present a comparison of our method with state-of-the-art techniques (Section IV-B). Finally, we show the effective analysis of the LLM-based TVS (Section IV-C).

TABLE I: Clustering performance of compared methods in terms of **Acc** and **NMI** on four human motion videosets. The best result is highlighted in **bold**. The improvement relative to the second-best method is depicted by  $\uparrow$ . (M) denotes the need for labeled MAD dataset assistance, and (K) indicates the requirement for labeled Keck dataset assistance.

(a) Results on Keck dataset			(b) Results on MAD dataset			(c) Results on UT dataset			(d) Results on Weiz dataset		
Method	Acc $\uparrow$	NMI $\uparrow$	Method	Acc $\uparrow$	NMI $\uparrow$	Method	Acc $\uparrow$	NMI $\uparrow$	Method	Acc $\uparrow$	NMI $\uparrow$
SSC [93]	0.3137	0.3858	SSC [93]	0.3817	0.4758	SSC [93]	0.4389	0.4998	SSC [93]	0.4576	0.6009
OSC [5]	0.4393	0.5931	OSC [5]	0.4327	0.5589	OSC [5]	0.5846	0.6877	OSC [5]	0.5216	0.7047
TSC(M) [6]	0.4653	0.6935	TSC(K) [6]	0.5473	0.7691	TSC(K) [6]	0.5213	0.7216	TSC(K) [6]	0.5931	0.7971
LTS [7]	0.4924	0.6213	LTS [7]	0.5466	0.6547	LTS [7]	0.6724	0.7435	LTS [7]	0.5674	0.6959
DSAE [70]	0.5136	0.5100	DSAE [70]	0.5898	0.6309	DSAE [70]	0.7323	0.6717	DSAE [70]	0.6120	0.6627
VSDA [9]	0.5804	0.7397	VSDA [9]	0.5606	0.7770	VSDA [9]	0.6203	0.8226	VSDA [9]	0.6287	0.7992
CDMS(M) [8]	0.6044	0.7891	CDMS(K) [8]	0.6536	0.8251	CDMS(K) [8]	0.6547	0.8267	CDMS(K) [8]	0.6465	0.8601
SIBMSC [16]	0.3886	0.4744	SIBMSC [16]	0.3639	0.4309	SIBMSC [16]	0.4477	0.4894	SIBMSC [16]	0.4127	0.5435
FSMSC [18]	0.4702	0.3970	FSMSC [18]	0.3914	0.3226	FSMSC [18]	0.4787	0.4213	FSMSC [18]	0.3914	0.3226
BTMSC [17]	0.4297	0.4862	BTMSC [17]	0.2397	0.2249	BTMSC [17]	0.4162	0.4051	BTMSC [17]	0.3638	0.4382
ARLRR [19]	0.5010	0.5270	ARLRR [19]	0.5125	0.5099	ARLRR [19]	0.5148	0.5121	ARLRR [19]	0.5436	0.5371
DCTMSC [20]	0.4723	0.4866	DCTMSC [20]	0.4885	0.5372	DCTMSC [20]	0.5569	0.5293	DCTMSC [20]	0.5592	0.5906
DCMVC [69]	0.5395	0.8049	DCMVC [69]	0.5792	0.8286	DCMVC [69]	0.5371	0.7746	DCMVC [69]	0.6030	0.8326
TVSH	<b>0.8048</b>	<b>0.8690</b>	TVSH	<b>0.8372</b>	<b>0.8438</b>	TVSH	<b>0.8723</b>	<b>0.8488</b>	TVSH	<b>0.8745</b>	<b>0.9316</b>

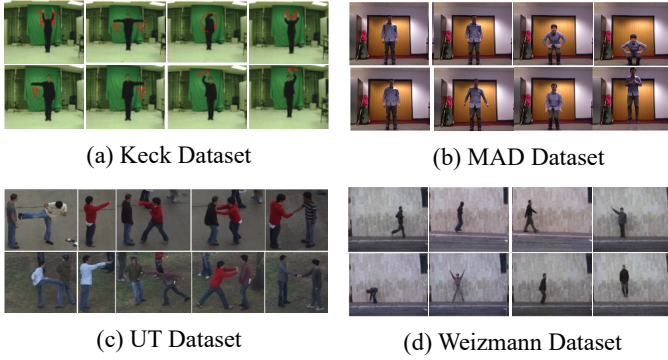


Fig. 2: Sampling frames from four human motion benchmark datasets, i.e., (a) Keck [94], (b) MAD [95], (c) UT [96], and (d) Weiz [97].

#### A. Human Motion Datasets and Experimental Setup

To provide a comprehensive evaluation of the proposed model, we perform experiments on four well-established benchmark human motion datasets. Some example frames from these datasets can be seen in Figure 2.

- *Keck Gesture Dataset (Keck)* [94] consists of 14 different motions from military signals, in which each subject is carried out 14 motions and gestures. Besides, the videos in this dataset were obtained by a fixed camera when these subjects stand out in a static background.

- *Multi-Modal Action Detection Dataset (MAD)* [95] consists of motions captured from various modalities using a Microsoft Kinect V2 system, which includes RGB images, depth cues, and skeleton formats. Specifically, the RGB images and 3D depth cues are of a size of  $240 \times 320$ . Moreover, each subject performs 35 different motions within two indoor scenes.

- *UT-Intermotion Dataset (UT)* [96] is composed of 20 videos, each of which includes six different motion types of human-human intermotions (such as punching, pushing, pointing, hugging, kicking, and handshaking).

- *Weizmann Dataset (Weiz)* [97] is composed of 90 video sequences with 10 motions (running, walking, skipping, bending, etc.) captured by nine subjects in an outdoor environment. All videos have a size of  $180 \times 144$  with 50 fps.

We evaluate clustering performance using four metrics: accuracy (Acc), normalized mutual information (NMI), precision (Pr), and adjusted rand index (ARI). These metrics assess the consistency between learned and true labels, with higher values indicating better performance. Let  $\mathcal{L} = \{l_1, l_2, \dots, l_N\}$  and  $\hat{\mathcal{L}} = \{\hat{l}_1, \hat{l}_2, \dots, \hat{l}_N\}$  represent the ground-truth and predicted labels, respectively, where  $l_i$  and  $\hat{l}_i$  denote the true and predicted labels for the  $i$ th sample. Acc is defined as the proportion of correctly clustered samples:  $\text{Acc} = \frac{1}{N} \sum_{i=1}^N \delta(l_i, \text{map}(\hat{l}_i))$ , where  $\delta(a, b)$  is the indicator function ( $\delta(a, b) = 1$  if  $a = b$ , and 0 otherwise), and  $\text{map}(\cdot)$  maps predicted labels to the best matching true labels using the Hungarian algorithm [98]. NMI quantifies the coherence between two sets. Let  $H(\mathcal{L})$  and  $H(\hat{\mathcal{L}})$  represent the entropies of the sets  $\mathcal{L}$  and  $\hat{\mathcal{L}}$ , respectively. NMI is defined as:  $\text{NMI}(\mathcal{L}, \hat{\mathcal{L}}) = \text{MI}(\mathcal{L}, \hat{\mathcal{L}}) / \sqrt{H(\mathcal{L})H(\hat{\mathcal{L}})}$ , where  $\text{MI}(\mathcal{L}, \hat{\mathcal{L}})$  measures the mutual information between the sets. Higher mutual information and lower uncertainty result in a higher NMI. If the sets are randomly distributed, NMI equals 0. Pr calculates the percentage of correctly clustered pairs among all pairs with the same clustering label. True positive (TP), false positive (FP), and false negative (FN) represent the numbers of correctly labeled samples in the positive class, misclassified samples in the positive cluster, and misclassified samples in the negative cluster, respectively. Precision is defined as:  $\text{Pr} = \frac{\text{TP}}{\text{TP} + \text{FP}}$ . ARI [99] quantifies the similarity between two sets:  $\text{ARI} = (\sum_{i,j=1}^K C_{n_{ij}}^2 - \mathbb{E}[\text{RI}]) / (C_0 - \mathbb{E}[\text{RI}])$ , where  $C_0 = \frac{1}{2} (\sum_{i=1}^K C_{|\mathcal{L}^{(i)}|}^2 + \sum_{i=1}^K C_{|\hat{\mathcal{L}}^{(i)}|}^2)$  and  $\mathbb{E}[\text{RI}] = \sum_{i=1}^K C_{|\mathcal{L}^{(i)}|}^2 \sum_{j=1}^K C_{|\hat{\mathcal{L}}^{(j)}|}^2 / C_N^2$ . Here,  $|\mathcal{L}^{(i)}|$  and  $|\hat{\mathcal{L}}^{(i)}|$  represent the number of samples in the  $i$ th cluster of the ground-truth and predicted labels, respectively. The value  $n_{ij}$  denotes the number of samples in the  $i$ th true cluster grouped into the  $j$ th predicted cluster. The notation  $C_n^m$  represents the number

TABLE II: Clustering performance of compared methods in terms of **Pr** and **ARI** on four human motion videosets.

(a) Results on Keck dataset			(b) Results on MAD dataset			(c) Results on UT dataset			(d) Results on Weiz dataset		
Method	Pr $\uparrow$	ARI $\uparrow$	Method	Pr $\uparrow$	ARI $\uparrow$	Method	Pr $\uparrow$	ARI $\uparrow$	Method	Pr $\uparrow$	ARI $\uparrow$
SSC [93]	0.3511	0.2446	SSC [93]	0.3151	0.1994	SSC [93]	0.5426	0.3772	SSC [93]	0.4469	0.3620
OSC [5]	0.3767	0.2743	OSC [5]	0.4024	0.2403	OSC [5]	0.5426	0.3966	OSC [5]	0.5126	0.4422
TSC(M) [6]	0.4214	0.3457	TSC(K) [6]	0.5116	0.3724	TSC(K) [6]	0.5864	0.4217	TSC(K) [6]	0.5667	0.5324
LTS [7]	0.4457	0.3052	LTS [7]	0.4673	0.3426	LTS [7]	0.5774	0.4457	LTS [7]	0.5991	0.5724
DSAE [70]	0.4195	0.3418	DSAE [70]	0.5492	0.3891	DSAE [70]	0.6189	0.4895	DSAE [70]	0.6233	0.5406
VSDA [9]	0.4311	0.3529	VSDA [9]	0.5667	0.3780	VSDA [9]	0.6334	0.5202	VSDA [9]	0.6180	0.5378
CDMS(M) [8]	0.5828	0.5174	CDMS(K) [8]	0.5761	0.4128	CDMS(K) [8]	0.6466	0.5539	CDMS(K) [8]	0.6316	0.5561
SIBMSC [16]	0.3773	0.2292	SIBMSC [16]	0.3227	0.1478	SIBMSC [16]	0.4972	0.3193	SIBMSC [16]	0.3643	0.2909
FSMSC [18]	0.4702	0.3970	FSMSC [18]	0.3914	0.3226	FSMSC [18]	0.5242	0.5047	FSMSC [18]	0.3914	0.3226
BTMSC [17]	0.3592	0.2418	BTMSC [17]	0.2882	0.1809	BTMSC [17]	0.5360	0.3667	BTMSC [17]	0.4397	0.3491
ARLRR [19]	0.5772	0.5046	ARLRR [19]	0.5426	0.4072	ARLRR [19]	0.6054	0.5213	ARLRR [19]	0.6211	0.5146
DCTMSC [20]	0.3753	0.2741	DCTMSC [20]	0.4508	0.2917	DCTMSC [20]	0.5635	0.4383	DCTMSC [20]	0.5502	0.4914
DCMVC [69]	0.3908	0.3136	DCMVC [69]	0.5487	0.3623	DCMVC [69]	0.6131	0.4688	DCMVC [69]	0.6045	0.5280
TVSH	<b>0.7559</b>	<b>0.7214</b>	TVSH	<b>0.7043</b>	<b>0.6973</b>	TVSH	<b>0.7477</b>	<b>0.7153</b>	TVSH	<b>0.9012</b>	<b>0.8867</b>

TABLE III: The comparison of run-time (minute) on the Keck dataset.

Method	SSC	OSC	TSC	LTS	DSAE	VSDA	CDMS
Time	1.1	2.5	1.9	3.3	4.5	4.4	5.1
Method	SIBMSC	FSMSC	BTMSC	ARLRR	DCTMSC	DCMVC	TVSH
Time	5.8	4.3	4.4	5.3	4.7	4.8	4.2

of ways to choose  $m$  items from  $n$ .

We evaluate the performance of our method through a comparative analysis with thirteen approaches, as outlined in Section II. Each method was independently tested ten times, and the average results were reported. For the proposed scheme, the TVS learning was performed using the following LLMs: GPT-o1, DeepSeek-v3-2-exp, Claude-Sonnet-4-5-20250929, Gemini-2.0-Flash-exp, Grok-4, and Qwen3-235B-a22b.

### B. HMS Performance Comparison

Tables I–II summarize results on four benchmarks (Keck, MAD, UT, Weiz) using Acc, NMI, Pr, and **ARI**. TVSH attains the best performance across all datasets and metrics.

On the Keck dataset, TVSH improves accuracy from 0.6044 (CDMS) to 0.8048 and NMI from 0.8049 (DCMVC) to 0.8690, representing a significant improvement over the best baseline. On MAD, the accuracy increases from 0.6536 to 0.8372, while NMI rises from 0.8286 to 0.8438. On UT, TVSH achieves 0.8723 accuracy and 0.8488 NMI, both higher than those of existing methods. On the more diverse Weiz dataset, TVSH reaches 0.8745 accuracy and 0.9316 NMI, improving by about 0.23 and 0.07, respectively, compared with the best previous method. Precision and ARI exhibit consistent improvement trends, confirming the robustness of TVSH in maintaining temporal coherence and enhancing motion discriminability.

Table III reports wall-clock time on Keck, which shows that TVSH maintains reasonable computational efficiency. Although it requires slightly more time than lightweight baselines such as SSC or OSC, the additional cost is modest and justified by its substantial performance gains.

#### 1) Superiority of the Proposed TVSH Method:

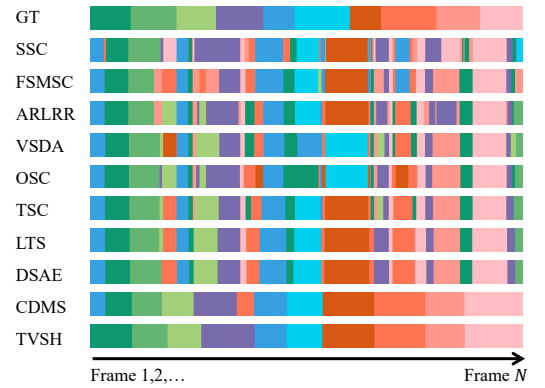


Fig. 3: Visualization of motion segmentation results of the proposed method and comparisons on Keck dataset. The different colors denote different motions. GT depicts the ground truth.

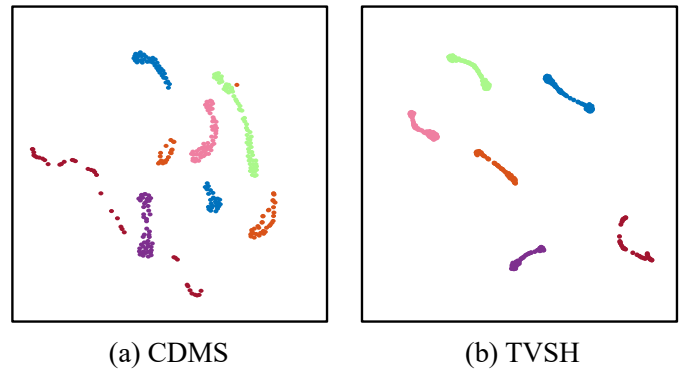


Fig. 4: Visualization of the two-dimensional t-SNE of the extracted features from six motions in the Weiz dataset. Points in different colors depict frames of different motions.

*a) Superiority 1: Temporal Modeling:* Conventional methods often fail to explicitly capture temporal dependencies in human motion sequences, making it difficult to identify gradually transitioning motions. Specifically, when a single



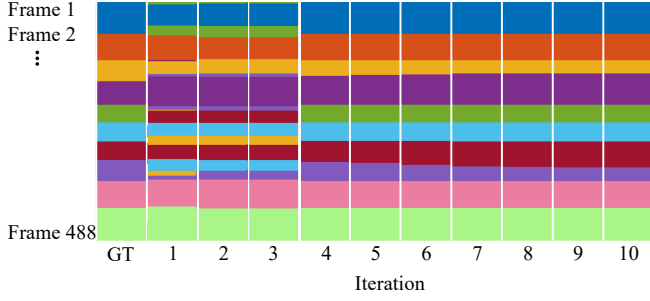


Fig. 5: Visualizations of the motion segmentation results in the different iteration of the proposed TVSH on the Weiz dataset. Different colors represent distinct motion assignments for the frames. ‘GT’ refers to the ‘ground truth’ motion segmentation results.

motion contains multiple stages with significant amplitude variations, it is often misinterpreted as multiple separate motions. Exploring temporal modeling offers a solution to this challenge. The proposed method first utilizes an LLM to obtain a TVS matrix, which explicitly encodes temporal relationships between frames. Then, the TVS matrix-based regularization is introduced to enforce temporal continuity in both the embedding space and the segmentation result, thereby reducing ambiguity in temporal motion transitions.

Figure 3 presents the motion segmentation results of the proposed method, along with comparisons on the Keck dataset. The proposed TVSH generates temporally coherent motion segmentation with well-aligned motion boundaries. Methods that do not explicitly model temporal dependencies (e.g., SSC and its variants such as FSMSC and ARLRR) achieve only limited temporal coherence. Although approaches incorporating temporal cues (e.g., OSC, TSC, LTS, DSAE, VSDA, and CDMS) perform better in terms of temporal coherence, they still fail to capture the full temporal coherence of motion sequences and cannot guarantee accurate temporal continuity in the motion segmentation results. In contrast, the proposed method leverages LLM to learn more precise temporal coherence and uses it to guide motion segmentation, producing temporally semantically accurate motion segments.

Figure 4 further visualizes the two-dimensional t-SNE embeddings of the extracted features  $u_i$  from six motions in the Weiz dataset. The proposed method produces more compact clusters than the strong baseline CDMS, highlighting the effectiveness of temporal regularization in the embedding component of the proposed TVSH. Based on the embedded features shown in Figure 4(b), achieving better motion segmentation performance becomes easier. This also explains why the proposed method can achieve more temporally semantically accurate motion segments, as shown in Figure 3.

*b) Superiority 2: Joint Optimization:* Traditional clustering-based human motion segmentation methods typically perform feature extraction and clustering of embedded features in separate stages, resulting in an embedding process that does not receive feedback from

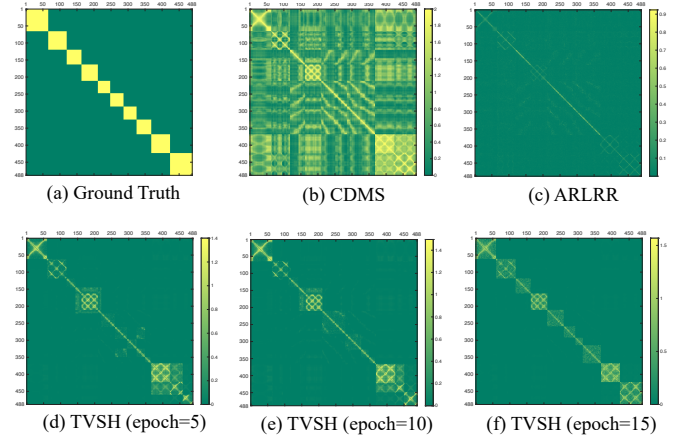


Fig. 6: Visualizations of the similarity matrix. (a) Ground-truth similarity matrix of the Ido human motion sequence in the Weiz dataset. Yellow regions indicate high similarity between frames of the same motion, while green regions denote zero similarity between frames of the same motion. (b)–(c) show the similarity matrices provided by the baselines CDMS and ARLRR. (d)–(f) show the similarity matrices generated by the proposed TVSH at different iterations.

the clustering outcome. However, a large clustering loss indicates that the embedded features are difficult to cluster, and it is meaningful to adjust the embedding to generate new features that minimize the clustering loss as much as possible. Therefore, our method adopts a feedback-enabled joint optimization framework, where the segmentation results iteratively refine the learned embeddings to achieve the smallest possible clustering loss.

Figure 5 shows the motion segmentation results of the proposed TVSH on a video from the Weiz dataset. The TVSH converges within ten iterations. During the first three iterations, the proposed TVSH focuses on identifying optimal segmentation boundaries and merging fragmented regions temporally. In the subsequent fourth to tenth iterations, it fine-tunes the boundaries, yielding stable and temporally coherent segmentation results.

Figures 6(b-c) presents the similarity matrices for different methods, while Figures 6(d)–(f) show the similarity matrices across different iterations from the joint optimization of the proposed TVSH. The similarity matrix of the proposed TVSH framework progressively exhibits a clearer block-diagonal structure, demonstrating strong alignment with the ground truth shown in Figure 6(a). In contrast, the similarity matrices obtained by baseline methods, such as CDMS and ARLRR (Figure 6(b) and (c)), are less structured and more diffuse. Obviously, a similarity matrix that is more consistent with the ground truth in Figure 6(a) facilitates more accurate motion segmentation.

*c) Ablation Study:* To assess the contribution of temporal modeling and joint optimization in TVSH, we conduct ablation experiments by selectively removing the temporal model and the joint optimization module. The variant TVSH (w/o

TABLE IV: Ablation study of the effects of temporal modeling and joint optimization.

	Keck		MAD		UT	
	Acc	NMI	Acc	NMI	Acc	NMI
TVSH (w/o joint optimization)	0.7423	0.8429	0.7848	0.8322	0.8254	0.8371
TVSH (w/o temporal model)	0.7152	0.7428	0.7211	0.7546	0.6714	0.6211
TVSH	0.8048	0.8690	0.8372	0.8438	0.8723	0.8488

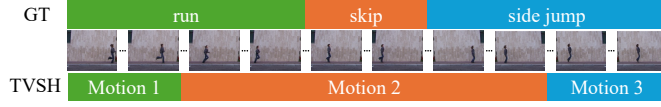


Fig. 7: Illustration of a failure case on the Weiz dataset (subject “ido”) showing gradual transition boundaries between actions.

temporal model) in Table IV removes the temporal prior by setting  $G = I$ , thereby disabling temporal regularization.

When the LLM-guided TVS is incorporated as the temporal model, performance consistently improves across all datasets. For instance, accuracy increases from 0.7152 to 0.8048 on the Keck dataset, from 0.7211 to 0.8372 on MAD, and from 0.6714 to 0.8723 on UT. These gains confirm that the TVS-based temporal modeling enhances human motion segmentation performance.

The variant TVSH (w/o joint optimization) in Table IV runs the proposed TVSH for only one iteration. This version produces weaker segmentation quality, as reflected by a 6–8% drop in accuracy across the datasets. In contrast, the full TVSH model benefits from iterative feedback between clustering and embedding, progressively refining segment boundaries and aligning the learned representation with semantic motion transitions. These results validate that joint optimization is essential for achieving temporally semantically consistent motion segmentation.

2) *Performance Bottleneck*: Although the proposed method achieves consistent improvements across all benchmarks, two primary bottlenecks preventing TVSH from reaching perfect segmentation accuracy were identified and analyzed.

a) *Bottleneck I: Gradual Transition Boundaries*: The first bottleneck arises in sequences where actions evolve smoothly without clear-cut temporal boundaries. As shown in Figure 7, when a motion transitions gradually from one motion to another (e.g., run  $\rightarrow$  skip  $\rightarrow$  side-jump), both visual and kinematic cues change continuously. These intermediate frames are semantically ambiguous, leading to minor drifts in segmentation boundaries or partial merging of adjacent segments.

b) *Bottleneck II: Look-Alike Motions*: The second bottleneck arises when motions share highly similar morphological characteristics. As depicted in Figure 8, for visually related motions such as raising one hand and raising both hands, the motion segmentation algorithm may incorrectly classify these two motions as the same “raising hand” action. This occurs because the visual cues and kinematic features between these motions overlap significantly, making it difficult for the model to distinguish between them.

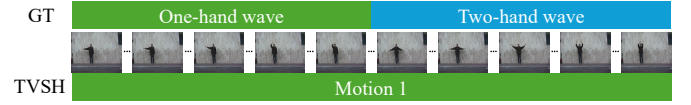


Fig. 8: Illustration of a failure case on the Weiz dataset (subject “ido”) showing visually similar or look-alike actions.

3) *Limitations and Future Directions*: Despite the effectiveness of the proposed method, several limitations remain. The first limitation arises in sequences where actions evolve smoothly without clear-cut temporal boundaries, causing minor drifts or partial merging of adjacent segments. The second limitation occurs when motions share highly similar morphological characteristics. This is due to the significant overlap in visual and kinematic features, making it difficult for the model to distinguish between them.

To address the limitations outlined above, future work will focus on enhancing the model’s ability to handle gradual motion transitions and look-alike motions. For gradual transition boundaries, we plan to incorporate uncertainty-aware temporal modeling, which can adaptively capture smooth variations and probabilistic transition boundaries. This will help the model distinguish between genuine motion transitions and intra-action fluctuations. Additionally, we aim to integrate multimodal features, such as skeletal joint trajectories, optical flow, and motion energy maps, to provide richer dynamic and geometric context. These modalities will enable the model to better differentiate between visually similar but semantically distinct actions, improving segmentation accuracy and temporal coherence in motion sequences.

### C. Effective Analysis of the LLM-Based TVS

This section analyzes the effectiveness, interpretability, and generalization of the proposed LLM-based TVS framework. Section IV-C1 visualizes the generated TVS matrices to assess their ability to capture temporal adjacency and motion coherence. Section IV-C2 compares different LLMs to evaluate how model architectures affect segmentation accuracy. Section IV-C3 examines the impact of prompt design on temporal reasoning. Section IV-C4 investigates performance across diverse motion types. Finally, Section IV-C5 summarizes the limitations of LLM-based TVS inference.

1) *Visualization of TVS from LLM*: We employ *GPT-o1* to generate the TVS matrix  $\mathbf{G}$ . Figure 9(a) presents an example TVS matrix from the *Weiz* dataset for the person “ido”. The TVS generated by the LLM exhibits minor inconsistencies. For instance, in the sixth motion, frames 271–323 correspond to the same “side jump” motion in the ground truth, yet the LLM identifies frames 271–295 as depicting the same motion and 306–323 as depicting the same motion separately. As shown in Figure 9(b), such partial inconsistencies occur in five out of ten motion segments. Nevertheless, the LLM does not introduce false distinctions, as it never explicitly labels segments 271–395 and 306–323 as different motions. Consequently, these minor inconsistencies do not mislead the subsequent TVSH algorithm. For the motions 1, 3, 5, 8, and 10, the LLM produces nearly perfectly consistent segmentations, demonstrating

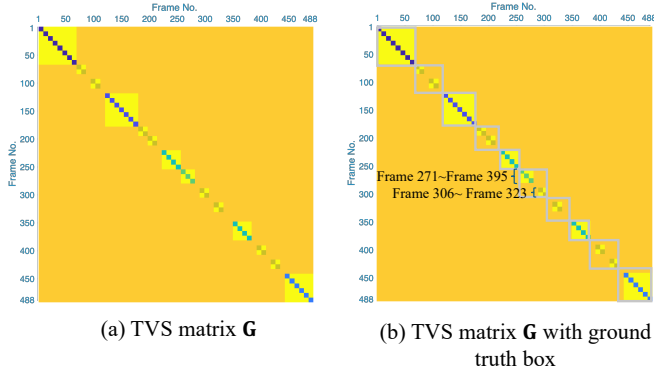


Fig. 9: (a) The TVS matrix  $G$  on Weiz dataset (person 'ido') generated by GPT-o1. (b) The TVS matrix  $G$  with ground truth label. The gray boxes indicate groups of frames that describe the same motion.

its effectiveness in capturing semantic motion coherence and temporal adjacency.

2) *Comparison with Different LLM Models*: We further evaluate six widely used multimodal LLMs for TVS learning, including *GPT-o1*, *DeepSeek-v3-2-exp*, *Claude-Sonnet-4-5-20250929*, *Gemini-2.0-Flash-exp*, *Grok-4*, and *Qwen3-235B-a22b*. Table V summarizes the parameter scales, runtime, and segmentation accuracy of the proposed TVSH, evaluated on the Weiz “ido” video (488 frames) using these LLM models under identical prompt and API configurations. All experiments were conducted on a MacBook Air equipped with an M4 chip and 32 GB of memory. Overall, *Gemini-2.0-Flash-exp* achieves the fastest runtime, followed by *Claude-Sonnet-4-5-20250929* and *DeepSeek-v3-2-exp*, whereas *Grok-4* is the slowest due to its extremely large parameter size. *GPT-o1* and *Qwen3-235B-a22b* fall in the mid range, balancing accuracy and computational cost.

In terms of segmentation accuracy, *GPT-o1* attains the best overall performance (88.27%), closely followed by *Qwen3-235B-a22b* and *Gemini-2.0-Flash-exp*, which achieve comparable results. *Claude-Sonnet-4-5-20250929* and *DeepSeek-v3-2-exp* exhibit slightly lower accuracy, while *Grok-4*, despite its largest parameter scale, yields the lowest performance. These results indicate that a larger model size does not necessarily guarantee better temporal reasoning or motion understanding; rather, architectural design and multimodal alignment play a more decisive role. Furthermore, model runtime generally increases with parameter size, reflecting the trade-off between computational complexity and inference precision. Nevertheless, all evaluated LLMs enhance TVS quality over non-LLM baselines (76.51%), confirming that integrating multimodal reasoning effectively strengthens temporal semantics and motion segmentation performance.

3) *Comparison with Different Prompts*: To examine the effect of prompt design on temporal semantic inference, we compared several variants of the prompt used to instruct the LLM. While the baseline prompt offers simplicity and generalization, additional prompt designs can enhance precision, robustness, and interpretability in identifying temporal

TABLE V: Parameter scale, runtime, and accuracy of different LLMs for TVS learning on the Weiz dataset (subject “ido”).

Method	Parameters (B)	Runtime/question (s)	Acc (%)
GPT-o1	175	1.84	88.27
DeepSeek-v3-2-exp	67	0.95	85.71
Claude-Sonnet-4-5-20250929	70	0.72	86.88
Gemini-2.0-Flash-exp	120	0.36	87.47
Grok-4	314	3.69	84.17
Qwen3-235B-a22b	235	2.20	87.52

TABLE VI: Accuracy of different prompts for TVS Learning on Weiz dataset.

Prompt	a)	b)	c)	d)	e)	f)
Acc	88.27	89.14	88.89	89.58	89.77	89.94

consistency between frames.

(a) *Baseline prompt*. This concise binary question directly queries the LLM’s perception of motion similarity:

“Do these two neighboring frames depict the same human motion? Answer Yes or No.”

Although simple and generalizable, this form provides limited guidance on how the model should evaluate visual similarity. Therefore, we explored more detailed formulations that explicitly direct attention toward motion dynamics and semantic continuity.

(b) *Attribute-focused prompt*. The LLM was asked to compare explicit aspects of human motion, such as global body posture, limb configuration, contact state (e.g., feet or hand support), and motion direction, while ignoring irrelevant visual variations like background and illumination:

“Carefully compare the two human figures. Focus on body posture, limb angles, contact with the ground, and movement direction. Ignore lighting, clothing, and background. Decide if they represent the same stage of an action. Answer Yes or No.”

(c) *Confidence-based prompt*. To quantify uncertainty in LLM judgment, we introduced a structured response that requests a confidence score:

“Do these two frames depict the same human motion? Provide your answer (Yes/No) and a confidence score between 0 and 1.”

This allows adaptive thresholding during TVS construction and enables selective re-querying of low-confidence pairs.

(d) *Step-aware prompt*. For temporally distant frames, the model was instructed to reason about motion continuity across a temporal gap:

“Compare frame  $i$  and frame  $i + \Delta t$ . Decide whether they correspond to the same stage of motion despite intermediate movement. Ignore viewpoint and background differences.”

(e) *Phase-aware prompt*. Incorporating explicit temporal reasoning, the LLM is asked to determine whether the two frames belong to the same *action phase* (e.g., preparation, execution, or completion):

“Identify whether these two frames occur in the same phase of an action (preparation, execution,



or completion). Focus on body posture and motion trajectory. Answer Yes or No.”

This helps capture fine-grained transitions within a continuous action.

(f) *Causal-motion prompt*. To leverage the model’s reasoning ability, we designed a prompt that emphasizes causal understanding of movement progression:

“Analyze how the motion evolves between these two frames. Determine if the second frame naturally follows from the first as part of the same continuous action. Answer Yes or No.”

This causal formulation improves temporal coherence by aligning LLM reasoning with the physical progression of motion.

Table VI reports the accuracy of six prompt variants used to guide the LLM in HMS. The results exhibit a steady improvement from the baseline formulation to the more context-aware and causality-driven designs, indicating that richer semantic cues lead to better temporal reasoning. The *causal-motion prompt* (f) achieves the highest accuracy (89.94%), confirming that prompting the LLM to reason about physical motion progression enhances its ability to capture temporal continuity and human-intuitive semantics. The *phase-aware prompt* (e) yields a comparable result (89.77%), showing that explicitly considering action phases (such as preparation, execution, and completion) helps distinguish fine-grained temporal transitions within continuous motions.

Both the *step-aware* (d) and *attribute-focused* (b) prompts demonstrate stable performance, as emphasizing motion continuity or detailed physical attributes effectively reduces ambiguity and reinforces local consistency in semantic comparison. The *confidence-based prompt* (c) provides moderate improvement by quantifying uncertainty in LLM judgments, which benefits reliability but contributes less to deeper semantic reasoning. In contrast, the simple *baseline prompt* (a) performs the weakest, as it lacks guidance on how the LLM should interpret motion similarity, relying only on its implicit visual understanding.

Overall, the results reveal a clear semantic progression: as the prompts evolve from generic and perception-based instructions to structured and reasoning-oriented formulations, the inferred temporal semantics become increasingly coherent. This progression (from the baseline to attribute-focused, confidence-based, step-aware, phase-aware, and finally causal-motion prompts) reflects the shift from surface-level perceptual matching toward a deeper, causality-driven understanding of human motion dynamics.

4) *Performance on Diverse Motion Types*: Representative motion examples in the Weiz dataset are shown in Figure 10. We observe that motions characterized by slower and more structured movements, such as *bend*, *side walk*, *walk*, *one hand wave*, and *two hands wave*, produce more precise and complete TVS representations. In contrast, faster and more dynamic motions, including *jumping jack*, *side jump*, *jump*, *run*, and *skip*, often result in incomplete TVS coverage due to rapid posture transitions and higher motion variability. As previously discussed for the TVS in Figure 9(a), such motions may



Fig. 10: Frame examples of the ten motions in the Weiz dataset (person ‘ido’).

introduce partial inconsistencies in temporal segmentation. For example, frames 271–323 correspond to the same “side jump” motion in the ground truth, yet the LLM identifies frames 271–295 as one continuous motion and frames 306–323 as one continuous motion.

5) *Limitation*: A notable limitation of the LLM-driven TVS is its tendency to misidentify a single motion as multiple distinct motions. This issue is evident in Figure 9(a), where the LLM splits a continuous motion into several segments. While using appropriate prompts can help mitigate this problem, it cannot be completely avoided. As a result, LLM-driven TVS cannot be directly applied to human motion segmentation in a straightforward manner. Additional steps, such as the proposed TVSH, are required. This issue becomes particularly pronounced when a motion involves multiple stages and occurs at high speed, but the camera’s frame rate is low. Furthermore, different LLM models may lead to TVS matrices with varying accuracy, and calling the LLM API is typically time-consuming, with each query requiring between 0.36 and 3.69 seconds.

## V. CONCLUSION

In this paper, we introduced a novel feedback-enabled subspace embedding approach for HMS, leveraging TVS embedded in human motion videos. We formulated the subspace embedding problem by integrating a temporal regularizer to capture the underlying temporal structure. Furthermore, we incorporated clustering with a temporal constraint to ensure that the clustering assignments reflect temporal characteristics. Finally, we developed a feedback-enabled framework to optimize the subspace embedding based on the segmentation results. Experimental results on benchmark datasets for HMS consistently demonstrated the superior performance of our approach compared to existing state-of-the-art techniques.

## APPENDIX

### A. Proof of Theorem 1

Let  $\{\mathcal{N}_i\}_{i=1}^N$  be the temporal neighborhoods and define the (possibly directed) adjacency  $A \in \{0, 1\}^{N \times N}$  by  $A_{i\ell} = 1$  iff  $\ell \in \mathcal{N}_i$ , otherwise  $A_{i\ell} = 0$ . Let  $\deg(i) = \sum_{\ell} A_{i\ell}$  and let  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N] \in \mathbb{R}^{d \times N}$  denote the column-wise embeddings. Recall the TVS matrix  $\mathbf{G} \in \mathbb{R}^{N \times N}$  is defined nodewise by

$$G_{ii} = -|\mathcal{N}_i|, G_{i\ell} = 1 \text{ if } \ell \in \mathcal{N}_i, G_{i\ell} = 0 \text{ otherwise.} \quad (12)$$

Hence each row  $i$  of  $\mathbf{G}$  encodes a *star Laplacian* centered at  $i$ .

**Lemma 1** (Discrete divergence-of-differences identity). *For every  $i \in \{1, \dots, N\}$ ,*

$$(\mathbf{ZG})_{:i} = \sum_{\ell \in \mathcal{N}_i} (\mathbf{z}_i - \mathbf{z}_\ell). \quad (13)$$

*Proof.* By (12), the  $i$ -th column of  $\mathbf{G}$  is  $G_{ii} = -\deg(i)$  and  $G_{\ell i} = 1$  for  $\ell \in \mathcal{N}_i$ , zero otherwise. Therefore

$$\begin{aligned} (\mathbf{ZG})_{:i} &= \sum_{j=1}^N \mathbf{z}_j G_{ji} = \mathbf{z}_i G_{ii} + \sum_{\ell \in \mathcal{N}_i} \mathbf{z}_\ell G_{\ell i} \\ &= -\deg(i) \mathbf{z}_i + \sum_{\ell \in \mathcal{N}_i} \mathbf{z}_\ell = \sum_{\ell \in \mathcal{N}_i} (\mathbf{z}_\ell - \mathbf{z}_i). \end{aligned}$$

Changing sign inside the sum yields (13).  $\square$

Define, for each node  $i$ , the *edge-difference stack*

$$\mathbf{v}_i = [(\mathbf{z}_i - \mathbf{z}_\ell) : \ell \in \mathcal{N}_i] \in \mathbb{R}^{d \deg(i)},$$

i.e., concatenate all incident differences at node  $i$ . The *node-wise isotropic graph total variation* is then

$$\text{GTV}_{\text{iso}}(\mathbf{Z}; A) = \sum_{i=1}^N \|\mathbf{v}_i\|_2 = \sum_{i=1}^N \left( \sum_{\ell \in \mathcal{N}_i} \|\mathbf{z}_i - \mathbf{z}_\ell\|_2^2 \right)^{1/2}. \quad (14)$$

Using Lemma 1, we can write

$$\|\mathbf{ZG}\|_{2,1} = \sum_{i=1}^N \|(\mathbf{ZG})_{:i}\|_2 = \sum_{i=1}^N \left\| \sum_{\ell \in \mathcal{N}_i} (\mathbf{z}_i - \mathbf{z}_\ell) \right\|_2. \quad (15)$$

Observe that (14) aggregates edge differences at node  $i$  by a *stack-then- $\ell_2$*  operation, while (15) aggregates them by a *sum-then- $\ell_2$*  operation. These two node-wise aggregations are *equivalent up to degree-dependent constants*: by the triangle inequality and Cauchy–Schwarz, for any collection  $\{\mathbf{a}_\ell\}_{\ell=1}^m$ ,

$$\left\| \sum_{\ell=1}^m \mathbf{a}_\ell \right\|_2 \leq \sum_{\ell=1}^m \|\mathbf{a}_\ell\|_2 \leq \sqrt{m} \left( \sum_{\ell=1}^m \|\mathbf{a}_\ell\|_2^2 \right)^{1/2}.$$

Applying this to  $\mathbf{a}_\ell = \mathbf{z}_i - \mathbf{z}_\ell$  with  $m = \deg(i)$  gives, nodewise,

$$\begin{aligned} \left\| \sum_{\ell \in \mathcal{N}_i} (\mathbf{z}_i - \mathbf{z}_\ell) \right\|_2 &\leq \sum_{\ell \in \mathcal{N}_i} \|\mathbf{z}_i - \mathbf{z}_\ell\|_2 \\ &\leq \sqrt{\deg(i)} \left( \sum_{\ell \in \mathcal{N}_i} \|\mathbf{z}_i - \mathbf{z}_\ell\|_2^2 \right)^{1/2}. \end{aligned} \quad (16)$$

Summing (16) over  $i$  yields the sandwich bound

$$\begin{aligned} \|\mathbf{ZG}\|_{2,1} &\leq \sum_{i=1}^N \sum_{\ell \in \mathcal{N}_i} \|\mathbf{z}_i - \mathbf{z}_\ell\|_2 \\ &\leq \sqrt{\deg_{\max}} \text{GTV}_{\text{iso}}(\mathbf{Z}; A), \deg_{\max} = \max_i \deg(i). \end{aligned} \quad (17)$$

Conversely, Jensen’s inequality implies  $(\sum_{\ell} \|\mathbf{z}_i - \mathbf{z}_\ell\|_2^2)^{1/2} \leq \sum_{\ell} \|\mathbf{z}_i - \mathbf{z}_\ell\|_2$ , and combining with the triangle inequality in the other direction gives constants  $c_1, c_2 > 0$  (depending only on  $\deg_{\max}$ ) such that

$$c_1 \text{GTV}_{\text{iso}}(\mathbf{Z}; A) \leq \|\mathbf{ZG}\|_{2,1} \leq c_2 \text{GTV}_{\text{iso}}(\mathbf{Z}; A). \quad (18)$$

Therefore,  $\|\mathbf{ZG}\|_{2,1}$  is *equivalent* (up to fixed multiplicative constants on degree-bounded temporal graphs) to the isotropic

Graph-TV (14). This establishes the first statement in Theorem 1.

Two key properties follow immediately: (i) *Zero iff node-wise constancy*. By (13),  $\|\mathbf{ZG}\|_{2,1} = 0$  iff for every  $i$ ,  $\sum_{\ell \in \mathcal{N}_i} (\mathbf{z}_i - \mathbf{z}_\ell) = \mathbf{0}$ . Since all summands are nonnegative in norm and the graph is degree-bounded, this holds iff  $\mathbf{z}_i = \mathbf{z}_\ell$  for all  $\ell \in \mathcal{N}_i$ . Thus  $\mathbf{z}$  is constant on every connected component induced by  $A$  (i.e., within each TVS-consistent segment). (ii) *Local smoothness with boundary preservation*. Minimizing (14) (and hence  $\|\mathbf{ZG}\|_{2,1}$  by (18)) penalizes *intra-segment* variations  $\|\mathbf{z}_i - \mathbf{z}_\ell\|_2$  along edges while not penalizing differences across *absent* edges. On temporal data, edges are present within action-consistent neighborhoods and absent across motion transitions. Therefore, the minimum-energy configurations are *piecewise-constant* on segments separated by motion boundaries (where edges vanish), exactly capturing the “smooth-inside / sharp-across” behavior typical of Graph-TV.

The identity (13) and the norm-equivalence (18) show that  $\|\mathbf{ZG}\|_{2,1}$  is an isotropic Graph-TV (up to degree-dependent constants) on the temporal graph defined by  $\{\mathcal{N}_i\}$ . Consequently, minimizing  $\|\mathbf{ZG}\|_{2,1}$  enforces local smoothness within temporal neighborhoods while allowing discontinuities at motion boundaries, yielding piecewise-constant embeddings aligned with human motion transitions, thus proving Theorem 1.

### B. Proof of Theorem 2

Consider a sequence of  $N$  frames with true binary adjacency labels  $\text{eq}_k^* \in \{0, 1\}$  on consecutive pairs  $(k, k+1)$ ,  $k = 1, \dots, N-1$ , where  $\text{eq}_k^* = 1$  indicates same motion and  $\text{eq}_k^* = 0$  indicates a true motion boundary. Let the observed labels be  $\text{eq}_k$ , obtained by flipping each  $\text{eq}_k^*$  independently with probability  $p < \frac{1}{2}$ . Let the minimum true segment length be  $L_{\min} \geq 1$ . The TVS neighborhoods  $\{\mathcal{N}_i\}$  and matrix  $\mathbf{G}$  are built from  $\{\text{eq}_k\}$  via the rule described in the method section, and embeddings  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N]$  are obtained by minimizing the TVS regularizer  $\|\mathbf{ZG}\|_{2,1}$  together with other convex terms of the objective.

Define the set of indices of flipped adjacencies  $\mathcal{F} = \{k \in \{1, \dots, N-1\} : \text{eq}_k \neq \text{eq}_k^*\}$  and let  $F = |\mathcal{F}|$ . By independence and linearity of expectation,

$$\mathbb{E}[F] = \sum_{k=1}^{N-1} \mathbb{P}(\text{eq}_k \neq \text{eq}_k^*) = p(N-1) = O(pN).$$

Each flip can at most create one *spurious* boundary (when a true within-segment edge is flipped from 1 to 0) or remove one *true* boundary (when a boundary edge is flipped from 0 to 1). Hence, if  $B_{\text{err}}$  denotes the number of erroneous TVS boundaries inferred from  $\{\text{eq}_k\}$ , then deterministically  $B_{\text{err}} \leq F + F = 2F$ , so

$$\mathbb{E}[B_{\text{err}}] \leq 2\mathbb{E}[F] = 2p(N-1) = O(pN).$$

This proves the first claim.

A single flipped within-segment edge (changing a run of 1’s to 1, 0, 1 locally) splits a true segment into two pieces whose lengths sum to the original length and differ by at

most 1. Since segment lengths are at least  $L_{\min}$ , any isolated flip creates at worst a one-edge “notch” inside a long run. Similarly, flipping a boundary edge can merge two adjacent true segments but only across a single location.

By Lemma (Graph-TV identity) used in Theorem 1, the TVS penalty can be written nodewise as

$$\|\mathbf{Z}\mathbf{G}\|_{2,1} = \sum_{i=1}^N \left\| \sum_{\ell \in \mathcal{N}_i} (\mathbf{z}_i - \mathbf{z}_\ell) \right\|_2,$$

which is equivalent (up to degree-dependent constants on a degree-bounded temporal graph) to an isotropic graph total variation that penalizes within-neighborhood differences  $\|\mathbf{z}_i - \mathbf{z}_\ell\|_2$ . Consider a true segment  $S = \{s, \dots, t\}$  of length  $|S| \geq L_{\min}$ . If all adjacencies inside  $S$  are correct ( $\text{eq}_k = 1$  for  $k \in \{s, \dots, t-1\}$ ), the minimal TVS energy within  $S$  is attained at *constant*  $\mathbf{z}_i$  over  $S$  (zero variation). If there is a single flip at index  $k \in \{s, \dots, t-1\}$  producing a spurious cut, the TVS graph inside  $S$  loses one local edge near  $k$ , but all remaining adjacent edges still connect the two sides across many nodes. Any nonconstant jump that honors the spurious cut introduces at least one additional nonzero difference along the many surviving edges across the two sides; keeping the embedding constant on  $S$  sets all those differences to zero. Thus, for an isolated spurious cut, the constant solution on  $S$  weakly dominates any split solution in TVS energy.

Formally, let  $S$  be partitioned into  $S_1 = \{s, \dots, k\}$  and  $S_2 = \{k+1, \dots, t\}$  by a single flipped edge at  $k$ . Let  $\mathbf{z}_{S_1}$  and  $\mathbf{z}_{S_2}$  denote the respective constants of a piecewise-constant candidate on  $S_1, S_2$ . For every surviving edge  $(i, \ell)$  with  $i \in S_1, \ell \in S_2$  that remains in  $\mathcal{N}_i$  or  $\mathcal{N}_\ell$  (these are the edges not removed by the single flip and they are  $\Omega(|S|)$  many when  $L_{\min}$  is large), the TVS contribution adds  $\|\mathbf{z}_{S_1} - \mathbf{z}_{S_2}\|_2$ . Hence the total TVS cost increases by at least  $c \|\mathbf{z}_{S_1} - \mathbf{z}_{S_2}\|_2$  for some  $c = \Omega(L_{\min})$ . Setting  $\mathbf{z}_{S_1} = \mathbf{z}_{S_2}$  brings this increase to zero, so the minimum is achieved by the constant solution unless opposed by a substantially large data term. In our stated theorem, the conclusion is in *expectation* and for *small*  $p$  with *large*  $L_{\min}$ ; the probability that multiple adjacent flips accumulate to remove most cross-side edges within  $S$  decays geometrically in the number of required flips, hence their contribution is negligible for small  $p$ .

Combining Steps 2–3 with the bound  $\mathbb{E}[B_{\text{err}}] = O(pN)$ , the flips appear sparsely along the chain for small  $p$ . With high probability, the flips are isolated at  $\Theta(1/p)$  spacing, while segment lengths are at least  $L_{\min}$ . When  $L_{\min}$  is sufficiently large relative to the typical spacing and the TVS weight is positive, the TVS minimization favors constant embeddings over each true segment and suppresses the local notches caused by isolated flips. Therefore the recovered embeddings are piecewise-constant on the true segments in expectation, yielding *segment-level consistency* for small  $p$  and large  $L_{\min}$ .

The expected number of erroneous TVS boundaries scales as  $O(pN)$  by linearity of expectation. Moreover, because  $\|\mathbf{Z}\mathbf{G}\|_{2,1}$  acts as an isotropic graph total variation on the temporal graph, isolated boundary errors are smoothed out by the regularizer on sufficiently long segments, so that the final embeddings remain piecewise constant per true segment in

expectation when  $p$  is small and  $L_{\min}$  is large. This completes the proof of Theorem 2.

### C. Proof of Proposition 1

We first consider the first two term of (2), which is denoted as

$$f(\mathbf{Z}) = \|\mathbf{X}\mathbf{Z} - \mathbf{X}\|_{\text{F}}^2 + \|\mathbf{Z}^T \mathbf{Z}\|_1 = \|\mathbf{X}\mathbf{Z} - \mathbf{X}\|_{\text{F}}^2 + \mathbf{e}^T \mathbf{Z}^T \mathbf{Z} \mathbf{e}$$

The columns of  $\mathbf{X}$  are in general position:  $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_K]$ , where all the columns of submatrix  $\mathbf{X}_\alpha$  lie in the same subspace  $\mathcal{S}_\alpha$ .

Assume  $\mathbf{Z}^*$  minimizes the function  $f(\mathbf{Z})$ , and we decompose  $\mathbf{Z}^*$  to be the sum of two matrices

$$\mathbf{Z}^* = \mathbf{Z}^D + \mathbf{Z}^C$$

$$= \begin{bmatrix} \mathbf{Z}_{11}^* & & & \mathbf{0} \\ & \mathbf{Z}_{22}^* & & \\ & & \ddots & \\ \mathbf{0} & & & \mathbf{Z}_{KK}^* \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{Z}_{12}^* & \cdots & \mathbf{Z}_{1K}^* \\ \mathbf{Z}_{21}^* & \mathbf{0} & \cdots & \mathbf{Z}_{2K}^* \\ \vdots & & \ddots & \vdots \\ \mathbf{Z}_{K1}^* & \mathbf{Z}_{K2}^* & \cdots & \mathbf{0} \end{bmatrix}$$

where  $\mathbf{Z}_{ij}^* \in \mathbb{R}^{\mathcal{N}_i \times \mathcal{N}_j}$ . Note that both  $\mathbf{Z}^D$  and  $\mathbf{Z}^C$  are non-negative.

According to the decomposition of  $\mathbf{Z}^*$ , any column of  $\mathbf{Z}^*$  can be written as  $\mathbf{z}_i^* = \mathbf{z}_i^D + \mathbf{z}_i^C$ , with  $\mathbf{z}_i^D$  and  $\mathbf{z}_i^C$  supported on disjoint subset of indices. We can write  $\|\mathbf{X}\mathbf{Z}^* - \mathbf{X}\|_{\text{F}}^2$  as

$$\begin{aligned} \|\mathbf{X}\mathbf{Z}^* - \mathbf{X}\|_{\text{F}}^2 &= \sum_{i=1}^N \|\mathbf{X}\mathbf{z}_i^* - \mathbf{x}_i\|_2^2 = \sum_{i=1}^N \|\mathbf{X}\mathbf{z}_i^D + \mathbf{X}\mathbf{z}_i^C - \mathbf{x}_i\|_2^2 \\ &= \sum_{i=1}^N \|\mathbf{X}\mathbf{z}_i^D - \mathbf{x}_i\|_2^2 + \sum_{i=1}^N \|\mathbf{X}\mathbf{z}_i^C\|_2^2 + 2 \sum_{i=1}^N \cos \theta_i \|\mathbf{X}\mathbf{z}_i^D - \mathbf{x}_i\|_2 \|\mathbf{X}\mathbf{z}_i^C\|_2 \end{aligned}$$

where  $\theta_i$  is the angle between vector  $\mathbf{X}\mathbf{z}_i^D - \mathbf{x}_i$  and  $\mathbf{X}\mathbf{z}_i^C$ .

Since the matrix  $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_K]$  is well arranged, any column  $\mathbf{x}_i \in \mathbf{X}_\alpha$  and  $\mathbf{x}_j \in \mathbf{X}_\beta$  lie in different subspaces if  $\alpha \neq \beta$ . Let  $\mathbf{x}_i \in \mathcal{S}_\alpha$ , according to the definition of  $\mathbf{z}_i^D$  and  $\mathbf{z}_i^C$ , we have  $\mathbf{X}\mathbf{z}_i^D \in \mathcal{S}_\alpha$  and  $\mathbf{X}\mathbf{z}_i^C \notin \mathcal{S}_\alpha$ . Based on the orthogonal subspace assumption, we have  $(\mathbf{X}\mathbf{z}_i^D - \mathbf{x}_i) \perp \mathbf{X}\mathbf{z}_i^C$  and  $\theta_i = \pi/2$ , thus

$$\begin{aligned} \|\mathbf{X}\mathbf{Z}^* - \mathbf{X}\|_{\text{F}}^2 &= \|\mathbf{X}\mathbf{Z}^D - \mathbf{X}\|_{\text{F}}^2 + \|\mathbf{X}\mathbf{Z}^C\|_{\text{F}}^2 \\ &\geq \|\mathbf{X}\mathbf{Z}^D - \mathbf{X}\|_{\text{F}}^2 \end{aligned} \quad (19)$$

Based on the nonnegativity of  $\mathbf{Z}^*$ ,  $\mathbf{Z}^C$ , and  $\mathbf{Z}^D$ , we have

$$\begin{aligned} \|\mathbf{Z}^*\|_1 &= \sum_{i,j} |(\mathbf{z}_i^*)^T \mathbf{z}_j^*| = \sum_{i,j} (\mathbf{z}_i^*)^T \mathbf{z}_j^* = \sum_{i,j} (\mathbf{z}_i^D + \mathbf{z}_i^C)^T (\mathbf{z}_j^D + \mathbf{z}_j^C) \\ &\geq \sum_{i,j} (\mathbf{z}_i^D)^T \mathbf{z}_j^D + \sum_{i,j} (\mathbf{z}_i^C)^T \mathbf{z}_j^C = \|(\mathbf{z}^D)^T \mathbf{z}^D\|_1 + \|(\mathbf{z}^C)^T \mathbf{z}^C\|_1 \\ &\geq \|(\mathbf{z}^D)^T \mathbf{z}^D\|_1 \end{aligned} \quad (20)$$

From inequalities (19) and (20) we have  $f(\mathbf{Z}^*) \geq f(\mathbf{Z}^D)$ . Because  $\mathbf{Z}_{ij}^* \in \mathbb{R}^{\mathcal{N}_i \times \mathcal{N}_j}$ , we have  $f(\mathbf{Z}^*) = f(\mathbf{Z}^D)$  and  $\mathbf{Z}^C = \mathbf{0}$ , thus  $\mathbf{Z}^* = \mathbf{Z}^D$ .

We then consider the third term in (2), namely,  $g(\mathbf{Z}) = \lambda_2 \|\mathbf{Z}\mathbf{G}\|_{2,1} = \sum_{i=1}^N \sum_{\ell \in \mathcal{N}_i} \|\mathbf{z}_i - \mathbf{z}_\ell\|_2$ . Given the subspace assumption for the samples, the construction of the neighbor set  $\mathcal{N}_i$  for the  $i$ th sample based on cosine measurements captures all the temporal neighbors of the  $i$ th sample. It is therefore straightforward to demonstrate that  $g(\mathbf{Z}) \geq \sum_{i=1}^N \sum_{\ell \in \mathcal{N}_i} \|\mathbf{z}_i^* - \mathbf{z}_\ell^*\|_2$ , confirming that  $\mathbf{Z}^*$  also minimizes  $g(\mathbf{Z})$ .



Finally, we address the last term in (2). With  $\mathbf{Z}$  being block-diagonal, the  $(i, j)$ th element of  $\mathbf{Z}$  is nonzero if and only if the  $i$ th and  $j$ th samples are situated in the same subspace. Consequently,  $\|\mathbf{Z}^*\|_{\mathbf{Q}^*} = 0$ .

Thus, Proposition 1 is upheld.

#### D. Proof of Theorem 3

Let the sequence be partitioned into ground-truth motion segments  $\{\mathcal{S}_g\}_{g=1}^K$ , where each segment  $\mathcal{S}_g$  has length  $|\mathcal{S}_g| \geq L_{\min}$  and generates observations from a linear subspace  $\mathcal{U}_g \subset \mathbb{R}^D$  with within-segment variance  $\sigma^2$ . The between-subspace separation is  $\Delta_{\text{sub}}^2 := \min_{g \neq h} \text{dist}^2(\mathcal{U}_g, \mathcal{U}_h) > 0$ . Let  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N]$  denote the learned subspace embeddings and  $\mathbf{Q}$  the cluster indicator matrix obtained by solving problem (2) with the TVS matrix  $\mathbf{G}$ . Assume LLM adjacency on consecutive pairs is independently flipped with probability  $p < \frac{1}{2}$  when constructing  $\mathbf{G}$ . We evaluate segmentation error  $\text{Err}_{\text{HMS}}$  as the normalized mis-segmentation rate (e.g., the Hamming error of predicted boundaries or the misclustering fraction, up to label permutation).

The objective (2) couples a data-fitting term and two regularizers: (i) the *subspace fidelity* driven by  $\|\mathbf{X} - \mathbf{XZ}\|_F^2 + \|\mathbf{Z}^\top \mathbf{Z}\|_1$  and the clustering regularizer  $\|\mathbf{Z}\|_{\mathbf{Q}}$ ; and (ii) the *temporal smoothness* driven by the TVS term  $\lambda_G \|\mathbf{ZG}\|_{2,1}$ . Accordingly, we bound

$$\mathbb{E}[\text{Err}_{\text{HMS}}] \leq \underbrace{\mathbb{E}[\text{Err}_{\text{TVS}}]}_{\text{temporal adjacency noise}} + \underbrace{\mathbb{E}[\text{Err}_{\text{sub}}]}_{\text{subspace separability/noise}},$$

and control each component in turn.

Let  $\text{eq}_k^* \in \{0, 1\}$  be the true adjacency on  $(k, k+1)$  and  $\text{eq}_k$  the observed (noisy) label used to build  $\mathbf{G}$ . The number of flipped adjacencies  $F = \sum_{k=1}^{N-1} \mathbb{I}\{\text{eq}_k \neq \text{eq}_k^*\}$  satisfies  $\mathbb{E}[F] = p(N-1)$  by independence. Each flip can create at most one spurious cut or remove one true cut, hence the number of erroneous TVS boundaries  $B_{\text{err}} \leq 2F$  and  $\mathbb{E}[B_{\text{err}}] = O(pN)$ .

However, the TVS penalty is an isotropic graph total-variation on the temporal graph (Theorem 1):

$$\|\mathbf{ZG}\|_{2,1} = \sum_{i=1}^N \left\| \sum_{\ell \in \mathcal{N}_i} (\mathbf{z}_i - \mathbf{z}_\ell) \right\|_2,$$

which favors constant embeddings on long runs and penalizes isolated notches. Inside a true segment  $\mathcal{S}_g$  with  $|\mathcal{S}_g| \geq L_{\min}$ , a single flipped edge breaks one local link but leaves  $\Omega(L_{\min})$  many cross-links among neighbors intact; any nonconstant split of  $\mathbf{Z}$  across that notch incurs at least  $c \|\Delta\|_2$  extra TVS cost with  $c = \Omega(L_{\min})$ . Thus, for sufficiently large  $\lambda_G$  (bounded away from zero and not exceeding the data term scale), the optimizer prefers to *heal* isolated flips and keep  $\mathbf{z}_i$  constant within  $\mathcal{S}_g$ . Since the flips are sparse in expectation and segments are long, the fraction of frames affected at the *segment level* scales as the number of flips divided by the segment length, yielding

$$\mathbb{E}[\text{Err}_{\text{TVS}}] \leq C_1 \frac{p}{L_{\min}},$$

for a constant  $C_1$  depending on neighborhood width and the TVS weight.

Within each true segment, the data lie near a subspace  $\mathcal{U}_g$  with variance  $\sigma^2$ , and different segments correspond to subspaces separated by  $\Delta_{\text{sub}}^2$ . Standard perturbation arguments for subspace clustering (and nearest-subspace assignment) imply a misassignment probability bounded by  $C_2 \sigma^2 / \Delta_{\text{sub}}^2$  when the separation dominates the noise (large-margin regime). In our formulation, the terms  $\|\mathbf{X} - \mathbf{XZ}\|_F^2 + \|\mathbf{Z}^\top \mathbf{Z}\|_1$  and  $\|\mathbf{Z}\|_{\mathbf{Q}}$  promote embeddings that are (approximately) block-sparse/diagonal across subspaces and cluster-coherent; hence the induced clustering error satisfies

$$\mathbb{E}[\text{Err}_{\text{sub}}] \leq C_2 \frac{\sigma^2}{\Delta_{\text{sub}}^2},$$

with  $C_2$  depending on the regularization weights and the conditioning of  $\{\mathcal{U}_g\}$ .

By the decomposition above,

$$\mathbb{E}[\text{Err}_{\text{HMS}}] \leq C_1 \frac{p}{L_{\min}} + C_2 \frac{\sigma^2}{\Delta_{\text{sub}}^2},$$

which proves the claimed bound.

If the TVS neighborhoods  $(l_i, r_i)$  coincide with true segments, then the TVS graph has no spurious cross-segment edges. Minimizing  $\|\mathbf{ZG}\|_{2,1}$  forces  $\mathbf{z}_i$  to be constant within each segment, and the data-fitting plus sparsity terms make inter-segment connections in  $\mathbf{Z}$  suboptimal when  $\Delta_{\text{sub}}^2 > 0$ . Therefore, at any optimum  $(\mathbf{Z}^*, \mathbf{Q}^*)$ , the matrix  $\mathbf{Z}^*$  is (after permutation) block-diagonal with blocks aligned to  $\{\mathcal{S}_g\}$ , which yields exact segmentation up to label permutation.

The constants  $C_1, C_2$  absorb factors due to neighborhood width, TVS weight  $\lambda_G$ , and spectral clustering relaxation tightness. The bound holds for any fixed choice of regularization weights in a compact interval  $[\underline{\lambda}_G, \bar{\lambda}_G]$  that preserves the healing effect of TVS while not overwhelming the data-fitting terms.

#### E. Proof of Proposition 2

For fixed  $(\mathbf{Z}, \mathbf{Q}, \mathbf{F}, \gamma)$ , the  $\mathbf{H}$ -subproblem

$$\min_{\mathbf{H}} \|\mathbf{H}\|_{2,1} + \langle \mathbf{F}, \mathbf{H} - \mathbf{ZG} \rangle + \frac{\gamma}{2} \|\mathbf{H} - \mathbf{ZG}\|_F^2 \quad (21)$$

is the proximal operator of the  $\ell_{2,1}$  norm and admits the closed-form group-shrinkage solution. Hence the  $\mathbf{H}$ -update attains the global minimizer of (21) at every iteration.

Define

$$\mathbf{P} := \mathbf{ZG} - \frac{1}{\gamma} \mathbf{F}.$$

Expanding the quadratic and completing the square yields

$$\begin{aligned} & \frac{\gamma}{2} \|\mathbf{H} - \mathbf{ZG}\|_F^2 + \langle \mathbf{F}, \mathbf{H} - \mathbf{ZG} \rangle \\ &= \frac{\gamma}{2} \left\| \mathbf{H} - \left( \mathbf{ZG} - \frac{1}{\gamma} \mathbf{F} \right) \right\|_F^2 - \frac{1}{2\gamma} \|\mathbf{F}\|_F^2 \\ &= \frac{\gamma}{2} \|\mathbf{H} - \mathbf{P}\|_F^2 + \text{const.} \end{aligned}$$

Since the additive constant does not affect the minimizer, problem (21) is equivalent to

$$\min_{\mathbf{H}} \|\mathbf{H}\|_{2,1} + \frac{\gamma}{2} \|\mathbf{H} - \mathbf{P}\|_F^2, \quad (22)$$

which is the proximal mapping of the  $\ell_{2,1}$  norm at point  $\mathbf{P}$  with parameter  $1/\gamma$ :

$$\mathbf{H}^* = \text{prox}_{(1/\gamma)\|\cdot\|_{2,1}}(\mathbf{P}) = \arg \min_{\mathbf{H}} \frac{\gamma}{2} \|\mathbf{H} - \mathbf{P}\|_F^2 + \|\mathbf{H}\|_{2,1}.$$

Let  $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_N]$  and  $\mathbf{P} = [\mathbf{p}_1, \dots, \mathbf{p}_N]$  denote the column partitions. Since

$$\|\mathbf{H}\|_{2,1} = \sum_{i=1}^N \|\mathbf{h}_i\|_2, \quad \|\mathbf{H} - \mathbf{P}\|_F^2 = \sum_{i=1}^N \|\mathbf{h}_i - \mathbf{p}_i\|_2^2,$$

the objective in (22) decomposes into  $N$  independent vector problems:

$$\mathbf{h}_i^* = \arg \min_{\mathbf{h} \in \mathbb{R}^D} \|\mathbf{h}\|_2 + \frac{\gamma}{2} \|\mathbf{h} - \mathbf{p}_i\|_2^2 \quad (i = 1, \dots, N). \quad (23)$$

Thus, it suffices to solve (23) for a single column and apply the solution to each  $i$ .

Consider the convex function  $\phi(\mathbf{h}) = \|\mathbf{h}\|_2 + \frac{\gamma}{2} \|\mathbf{h} - \mathbf{p}\|_2^2$  with  $\mathbf{p} \in \mathbb{R}^D$  fixed. A vector  $\mathbf{h}^*$  is optimal iff

$$\mathbf{0} \in \partial \|\mathbf{h}^*\|_2 + \gamma(\mathbf{h}^* - \mathbf{p}).$$

The subdifferential of  $\|\cdot\|_2$  is

$$\partial \|\mathbf{h}\|_2 = \begin{cases} \left\{ \frac{\mathbf{h}}{\|\mathbf{h}\|_2} \right\}, & \mathbf{h} \neq \mathbf{0}, \\ \{\mathbf{u} \in \mathbb{R}^D : \|\mathbf{u}\|_2 \leq 1\}, & \mathbf{h} = \mathbf{0}. \end{cases}$$

*Case A* ( $\mathbf{h}^* \neq \mathbf{0}$ ). Then there exists  $\alpha > 0$  such that  $\mathbf{h}^* = \alpha \mathbf{p}$  (the solution must align with  $\mathbf{p}$  by symmetry). Plugging into the optimality condition:

$$\begin{aligned} \mathbf{0} &= \frac{\mathbf{h}^*}{\|\mathbf{h}^*\|_2} + \gamma(\mathbf{h}^* - \mathbf{p}) = \frac{\alpha \mathbf{p}}{\alpha \|\mathbf{p}\|_2} + \gamma(\alpha - 1)\mathbf{p} \\ &= \left( \frac{1}{\|\mathbf{p}\|_2} + \gamma(\alpha - 1) \right) \mathbf{p}, \end{aligned}$$

which yields

$$\alpha = 1 - \frac{1}{\gamma \|\mathbf{p}\|_2}.$$

Feasibility requires  $\alpha > 0$ , i.e.,  $\|\mathbf{p}\|_2 > 1/\gamma$ .

*Case B* ( $\mathbf{h}^* = \mathbf{0}$ ). The optimality condition becomes  $\mathbf{0} \in \partial \|\mathbf{0}\|_2 - \gamma \mathbf{p}$ , i.e., there exists  $\mathbf{u}$  with  $\|\mathbf{u}\|_2 \leq 1$  such that  $\mathbf{u} = \gamma \mathbf{p}$ , which is possible iff  $\|\mathbf{p}\|_2 \leq 1/\gamma$ .

Combining the two cases gives the *block (group) soft-thresholding* operator

$$\mathbf{h}^* = \begin{cases} \left(1 - \frac{1}{\gamma \|\mathbf{p}\|_2}\right) \mathbf{p}, & \text{if } \|\mathbf{p}\|_2 > \frac{1}{\gamma}, \\ \mathbf{0}, & \text{otherwise.} \end{cases}$$

Applying this column-wise with  $\mathbf{p}_i = \mathbf{P}_{:,i}$  yields

$$\mathbf{h}_i^* = \max \left( 1 - \frac{1}{\gamma \|\mathbf{p}_i\|_2}, 0 \right) \mathbf{p}_i, \quad i = 1, \dots, N. \quad (24)$$

The objective in (22) is strictly convex in each column due to the strongly convex quadratic term. Therefore, the solution in (24) is the unique global minimizer for each subproblem (23); stacking the columns gives the unique global minimizer of (22), and hence of (21). This establishes that the  $\mathbf{H}$ -update is a *global* proximal step at every iteration, proving Proposition 2.

## F. Proof of Proposition 3

To prove that the given problem is equivalent to the spectral clustering problem, that is, solving  $\min_{\mathbf{Q}} \text{Tr}(\mathbf{Q}^\top (\mathbf{D} - \mathbf{A}) \mathbf{Q})$ , where  $\mathbf{D}$  is a diagonal matrix with elements  $\mathbf{D}_{j,j} = \sum_i A_{i,j}$ , we begin by expanding the objective function. Recall that our problem is  $\min_{\mathbf{Q}} \frac{1}{2} \sum_{i,j} A_{i,j} \|\mathbf{q}_i - \mathbf{q}_j\|_2^2$ , which can be expanded as  $\frac{1}{2} \sum_{i,j} A_{i,j} \|\mathbf{q}_i - \mathbf{q}_j\|_2^2 = \frac{1}{2} \sum_{i,j} A_{i,j} (\mathbf{q}_i^\top \mathbf{q}_i - 2\mathbf{q}_i^\top \mathbf{q}_j + \mathbf{q}_j^\top \mathbf{q}_j)$ . This expression can be broken into three parts:  $\frac{1}{2} \sum_{i,j} A_{i,j} \mathbf{q}_i^\top \mathbf{q}_i - \sum_{i,j} A_{i,j} \mathbf{q}_i^\top \mathbf{q}_j + \frac{1}{2} \sum_{i,j} A_{i,j} \mathbf{q}_j^\top \mathbf{q}_j$ . Now, observing each part, we can express it in matrix form. First,  $\sum_{i,j} A_{i,j} \mathbf{q}_i^\top \mathbf{q}_i = \sum_i \mathbf{q}_i^\top \mathbf{q}_i \sum_j A_{i,j} = \sum_i \mathbf{q}_i^\top \mathbf{q}_i \mathbf{D}_{i,i} = \text{Tr}(\mathbf{Q}^\top \mathbf{D} \mathbf{Q})$ , and similarly,  $\sum_{i,j} A_{i,j} \mathbf{q}_j^\top \mathbf{q}_j = \sum_j \mathbf{q}_j^\top \mathbf{q}_j \sum_i A_{i,j} = \sum_j \mathbf{q}_j^\top \mathbf{q}_j \mathbf{D}_{j,j} = \text{Tr}(\mathbf{Q}^\top \mathbf{D} \mathbf{Q})$ . Finally, we have  $-\sum_{i,j} A_{i,j} \mathbf{q}_i^\top \mathbf{q}_j = -\text{Tr}(\mathbf{Q}^\top \mathbf{A} \mathbf{Q})$ . Thus, combining these parts, we arrive at the following:  $\frac{1}{2} \sum_{i,j} A_{i,j} \|\mathbf{q}_i - \mathbf{q}_j\|_2^2 = \frac{1}{2} \text{Tr}(\mathbf{Q}^\top \mathbf{D} \mathbf{Q}) + \frac{1}{2} \text{Tr}(\mathbf{Q}^\top \mathbf{D} \mathbf{Q}) - \text{Tr}(\mathbf{Q}^\top \mathbf{A} \mathbf{Q}) = \text{Tr}(\mathbf{Q}^\top \mathbf{D} \mathbf{Q}) - \text{Tr}(\mathbf{Q}^\top \mathbf{A} \mathbf{Q}) = \text{Tr}(\mathbf{Q}^\top (\mathbf{D} - \mathbf{A}) \mathbf{Q})$ . This establishes the equivalence between the given problem and the spectral clustering problem.

## G. Proof of Proposition 4

The cost associated with a sample located in the  $k$ -th cluster, whose center is denoted by  $\mu_k$ , can be expressed as the sum of the squared Euclidean distances between the sample and the cluster center, augmented by a term accounting for the influence of neighboring samples. Specifically, the cost is given by  $\sum_{i \in \mathcal{C}_k} \left( \|\mathbf{u}_i - \mu_k\|_2^2 + \eta \sum_{j \in \mathcal{N}_i} \|\mathbf{u}_j - \mu_k\|_2^2 \right)$ , which can be expanded as  $\sum_{i=1}^N \mathbb{1}(i \in \mathcal{C}_k) \|\mathbf{u}_i - \mu_k\|_2^2 + \sum_{i=1}^N \mathbb{1}(i \in \mathcal{C}_k) \sum_{j=1}^N \mathbb{1}(j \in \mathcal{N}_i) \|\mathbf{u}_j - \mu_k\|_2^2$ . To simplify, we define  $n_k(i)$ , the number of times the  $i$ -th frame is considered a neighbor of samples in the  $k$ -th cluster, as  $n_k(i) = \sum_{j \in \mathcal{C}_k} \mathbb{1}(i \in \mathcal{N}_j)$ . Substituting this into the previous expression, the cost becomes  $\sum_{i=1}^N \mathbb{1}(i \in \mathcal{C}_k) \|\mathbf{u}_i - \mu_k\|_2^2 + \sum_{j=1}^N \eta n_k(j) \|\mathbf{u}_j - \mu_k\|_2^2$ . This formulation can be further compacted into the following expression  $\sum_{i=1}^N \|\mathbf{u}_i - \mu_k\|_2^2 (\mathbb{1}(i \in \mathcal{C}_k) + \eta n_k(i))$ . This equation reveals the total cost, which is the sum of the direct distance between each sample and the cluster center, and the weighted influence of its neighbors, with the weight determined by  $\eta$ . The term  $n_k(i)$  quantifies how many times sample  $i$  is considered a neighbor within the  $k$ -th cluster.

## H. Proof of Theorem 4

Recall the augmented formulation

$$\begin{aligned} \min_{\mathbf{Z}, \mathbf{H}, \mathbf{Q}} \quad & \Phi(\mathbf{Z}, \mathbf{H}, \mathbf{Q}) := \underbrace{\|\mathbf{X} - \mathbf{X}\mathbf{Z}\|_F^2}_{\triangleq f(\mathbf{Z})} + \underbrace{\|\mathbf{Z}\|_{\mathbf{Q}}}_{\triangleq g(\mathbf{H})} + \underbrace{\iota_{\mathcal{Q}}(\mathbf{Q})}_{\triangleq h(\mathbf{Q})} \\ \text{s.t.} \quad & \mathbf{H} = \mathbf{Z}\mathbf{G}, \quad \text{diag}(\mathbf{Z}) = \mathbf{0}, \quad \mathbf{Z} \geq \mathbf{0}, \end{aligned} \quad (25)$$

where  $\iota_{\mathcal{Q}}$  is the indicator of the feasible set  $\mathcal{Q}$ . The (scaled) augmented Lagrangian is

$$\begin{aligned} \mathcal{L}_\gamma(\mathbf{Z}, \mathbf{H}, \mathbf{Q}; \mathbf{F}) &:= f(\mathbf{Z}) + g(\mathbf{H}) + h(\mathbf{Q}) + \langle \mathbf{F}, \mathbf{H} - \mathbf{Z}\mathbf{G} \rangle \\ &\quad + \frac{\gamma}{2} \|\mathbf{H} - \mathbf{Z}\mathbf{G}\|_F^2, \end{aligned} \quad (26)$$

with  $\gamma > 0$  and multiplier  $\mathbf{F}$ . One outer iteration performs:

- (i) **Z-update**: minimize  $\mathcal{L}_{\gamma_t}(\cdot, \mathbf{H}^t, \mathbf{Q}^t; \mathbf{F}^t)$  over  $\mathbf{Z} \in \mathcal{Z} := \{\text{diag}(\mathbf{Z}) = 0, \mathbf{Z} \geq 0\}$  (via a Sylvester step followed by the projection  $\Pi_{\mathcal{Z}}$ ).
- (ii) **H-update**: minimize  $\mathcal{L}_{\gamma_t}(\mathbf{Z}^{t+1}, \cdot, \mathbf{Q}^t; \mathbf{F}^t)$  over  $\mathbf{H}$ , i.e., a proximal  $\ell_{2,1}$  step with closed form.
- (iii) **Q-update**: minimize  $h(\mathbf{Q})$  for fixed  $\mathbf{Z}^{t+1}$  (normalized-cut relaxation with TVS).
- (iv) **Dual update**:  $\mathbf{F}^{t+1} = \mathbf{F}^t + \gamma_t(\mathbf{H}^{t+1} - \mathbf{Z}^{t+1}\mathbf{G})$ ; update  $\gamma_{t+1} \geq \gamma_t$ ,  $\gamma_t \rightarrow \gamma_\infty \in (0, \infty)$ ; keep  $\rho > 1$  bounded.

We assume:

- A1  $\mathcal{L}_\gamma$  is proper, lower-semicontinuous, and bounded below on the feasible set.
- A2 Each block subproblem admits a minimizer; the **Z**-step is solved exactly for the quadratic subproblem followed by the exact projection onto  $\mathcal{Z}$ ; the **H**-step is the exact proximal minimizer; the **Q**-step attains a (relaxed) global minimizer of  $h(\cdot)$  or at least a value not exceeding  $h(\mathbf{Q}^t)$ .
- A3  $\{\gamma_t\}$  is nondecreasing with  $\gamma_t \rightarrow \gamma_\infty \in (0, \infty)$ , and the penalty growth factor  $\rho > 1$  is bounded.

**Lemma 2** (Blockwise descent). *For any  $t$ , the updates (i)–(iii) satisfy*

$$\begin{aligned} \mathcal{L}_{\gamma_t}(\mathbf{Z}^{t+1}, \mathbf{H}^t, \mathbf{Q}^t; \mathbf{F}^t) &\leq \mathcal{L}_{\gamma_t}(\mathbf{Z}^t, \mathbf{H}^t, \mathbf{Q}^t; \mathbf{F}^t), \\ \mathcal{L}_{\gamma_t}(\mathbf{Z}^{t+1}, \mathbf{H}^{t+1}, \mathbf{Q}^t; \mathbf{F}^t) &\leq \mathcal{L}_{\gamma_t}(\mathbf{Z}^{t+1}, \mathbf{H}^t, \mathbf{Q}^t; \mathbf{F}^t), \\ \mathcal{L}_{\gamma_t}(\mathbf{Z}^{t+1}, \mathbf{H}^{t+1}, \mathbf{Q}^{t+1}; \mathbf{F}^t) &\leq \mathcal{L}_{\gamma_t}(\mathbf{Z}^{t+1}, \mathbf{H}^{t+1}, \mathbf{Q}^t; \mathbf{F}^t). \end{aligned}$$

*Proof.* Each inequality holds because each block is minimized exactly (or to a value no worse than the current one) while other blocks are fixed.  $\square$

**Lemma 3** (Dual ascent keeps augmented value nonincreasing). *With the dual update  $\mathbf{F}^{t+1} = \mathbf{F}^t + \gamma_t(\mathbf{H}^{t+1} - \mathbf{Z}^{t+1}\mathbf{G})$  and nondecreasing  $\gamma_t$ , there exists  $c > 0$  (independent of  $t$ ) such that*

$$\begin{aligned} \mathcal{L}_{\gamma_{t+1}}(\mathbf{Z}^{t+1}, \mathbf{H}^{t+1}, \mathbf{Q}^{t+1}; \mathbf{F}^{t+1}) \\ \leq \mathcal{L}_{\gamma_t}(\mathbf{Z}^{t+1}, \mathbf{H}^{t+1}, \mathbf{Q}^{t+1}; \mathbf{F}^t) - c \|\mathbf{H}^{t+1} - \mathbf{Z}^{t+1}\mathbf{G}\|_F^2. \end{aligned}$$

*Proof.* This is the standard ADMM identity obtained by expanding (26) at the two pairs  $(\gamma_t, \mathbf{F}^t)$  and  $(\gamma_{t+1}, \mathbf{F}^{t+1})$ , using the dual update, and the nondecreasing penalty. The quadratic penalty dominates the linear coupling, yielding the negative quadratic term.  $\square$

Combining Lemmas 2 and 3,

$$\begin{aligned} \mathcal{L}_{\gamma_{t+1}}(\mathbf{Z}^{t+1}, \mathbf{H}^{t+1}, \mathbf{Q}^{t+1}; \mathbf{F}^{t+1}) \\ \leq \mathcal{L}_{\gamma_t}(\mathbf{Z}^t, \mathbf{H}^t, \mathbf{Q}^t; \mathbf{F}^t) - c \|\mathbf{H}^{t+1} - \mathbf{Z}^{t+1}\mathbf{G}\|_F^2, \end{aligned}$$

so the augmented Lagrangian value is monotonically nonincreasing along the iterates. By A1, it is bounded below; hence it converges to a finite limit, and

$$\|\mathbf{H}^{t+1} - \mathbf{Z}^{t+1}\mathbf{G}\|_F \rightarrow 0. \quad (27)$$

**Lemma 4** (Boundedness of iterates). *The sequence  $\{(\mathbf{Z}^t, \mathbf{H}^t, \mathbf{Q}^t, \mathbf{F}^t)\}$  is bounded.*

*Proof.* Since  $\mathcal{L}_{\gamma_t}$  decreases and is coercive in each block due to the quadratic penalty and the constraints (nonnegativity, zero

diagonal) restricting  $\mathbf{Z}$ , the proximal term controlling  $\mathbf{H}$ , and the indicator  $\iota_{\mathcal{Q}}$  restricting  $\mathbf{Q}$ , each block remains bounded. The dual sequence is bounded because  $\mathbf{F}^{t+1} - \mathbf{F}^t = \gamma_t(\mathbf{H}^{t+1} - \mathbf{Z}^{t+1}\mathbf{G})$  and (27).  $\square$

Thus there exists a convergent subsequence (not relabeled) with

$$(\mathbf{Z}^t, \mathbf{H}^t, \mathbf{Q}^t, \mathbf{F}^t) \rightarrow (\mathbf{Z}^*, \mathbf{H}^*, \mathbf{Q}^*, \mathbf{F}^*), \quad \gamma_t \rightarrow \gamma_\infty.$$

Moreover, (27) implies  $\mathbf{H}^* = \mathbf{Z}^*\mathbf{G}$ .

Consider the first-order optimality (variational inequality) of each exact block update at iteration  $t$ :

$$0 \in \partial_{\mathbf{Z}} \left( f(\mathbf{Z}) - \langle \mathbf{F}^t, \mathbf{Z}\mathbf{G} \rangle + \frac{\gamma_t}{2} \|\mathbf{H}^t - \mathbf{Z}\mathbf{G}\|_F^2 + \iota_{\mathcal{Z}}(\mathbf{Z}) \right) \Big|_{\mathbf{Z}^{t+1}}, \quad (28)$$

$$0 \in \partial_{\mathbf{H}} \left( g(\mathbf{H}) + \langle \mathbf{F}^t, \mathbf{H} \rangle + \frac{\gamma_t}{2} \|\mathbf{H} - \mathbf{Z}^{t+1}\mathbf{G}\|_F^2 \right) \Big|_{\mathbf{H}^{t+1}}, \quad (29)$$

$$0 \in \partial_{\mathbf{Q}} h(\mathbf{Q}) \Big|_{\mathbf{Q}^{t+1}}. \quad (30)$$

Passing to the limit along the convergent subsequence, using: (i) outer semicontinuity of subdifferentials for proper l.s.c. functions, (ii)  $\gamma_t \rightarrow \gamma_\infty$ , (iii)  $\mathbf{H}^{t+1} - \mathbf{Z}^{t+1}\mathbf{G} \rightarrow \mathbf{0}$ , and (iv)  $\mathbf{F}^{t+1} - \mathbf{F}^t \rightarrow \mathbf{0}$ , we obtain the KKT-type stationary conditions at  $(\mathbf{Z}^*, \mathbf{H}^*, \mathbf{Q}^*, \mathbf{F}^*)$ :

$$\begin{aligned} 0 &\in \partial_{\mathbf{Z}} \left( f(\mathbf{Z}) - \langle \mathbf{F}^*, \mathbf{Z}\mathbf{G} \rangle + \iota_{\mathcal{Z}}(\mathbf{Z}) \right) \Big|_{\mathbf{Z}^*}, \\ 0 &\in \partial_{\mathbf{H}} \left( g(\mathbf{H}) + \langle \mathbf{F}^*, \mathbf{H} \rangle \right) \Big|_{\mathbf{H}^*}, \\ 0 &\in \partial_{\mathbf{Q}} h(\mathbf{Q}) \Big|_{\mathbf{Q}^*}, \quad \mathbf{H}^* = \mathbf{Z}^*\mathbf{G}. \end{aligned}$$

These are precisely the first-order (primal-dual) stationary conditions for (25) with the linear constraint  $\mathbf{H} = \mathbf{Z}\mathbf{G}$ . Therefore, any limit point is a first-order stationary point.

If in addition each **Q**-update solves its relaxed subproblem globally (Assumption A2 strengthened), then at  $(\mathbf{Z}^*, \mathbf{H}^*, \mathbf{Q}^*)$ ,  $\mathbf{Q}^*$  satisfies the global optimality condition of the convex relaxation. In this case, the limit point satisfies the KKT conditions of the relaxed problem; hence every accumulation point is a KKT point.

We have shown (i) monotone decrease and convergence of the augmented Lagrangian values, (ii) boundedness of the iterates and vanishing primal residual, and (iii) that every limit point is a first-order stationary point; with globally optimal **Q**-updates for the relaxed subproblem, every accumulation point is a KKT point. This completes the proof of Theorem 4.

## REFERENCES

- [1] F. Zhou, F. De la Torre, and J. K. Hodgins, “Hierarchical aligned cluster analysis for temporal clustering of human motion,” *IEEE Trans. Pattern Anal. Mach.*, vol. 35, no. 3, pp. 582–596, 2012.
- [2] M. Keuper, S. Tang, B. Andres, T. Brox, and B. Schiele, “Motion segmentation & multiple object tracking by correlation co-clustering,” *IEEE Trans. Pattern Anal. Mach.*, vol. 42, no. 1, pp. 140–153, 2018.
- [3] R. Poppe, “Vision-based human motion analysis: An overview,” *Comput. Vis. Image Underst.*, vol. 108, no. 1–2, pp. 4–18, 2007.
- [4] J. F.-S. Lin, M. Karg, and D. Kulić, “Movement primitive segmentation for human motion modeling: A framework for analysis,” *IEEE Trans. Human-Mach. Syst.*, vol. 46, no. 3, pp. 325–339, 2016.
- [5] S. Tierney, J. Gao, and Y. Guo, “Subspace clustering for sequential data,” in *Proceedings of IEEE conference on computer vision and pattern recognition*, 2014, pp. 1019–1026.

- [6] S. Li, K. Li, and Y. Fu, "Temporal subspace clustering for human motion segmentation," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 4453–4461.
- [7] L. Wang, Z. Ding, and Y. Fu, "Low-rank transfer human motion segmentation," *IEEE Trans. Image Process.*, vol. 28, no. 2, pp. 1023–1034, 2018.
- [8] T. Zhou, H. Fu, C. Gong, L. Shao, F. Porikli, H. Ling, and J. Shen, "Consistency and diversity induced human motion segmentation," *IEEE Trans. Pattern Anal. Mach.*, vol. 45, no. 1, pp. 197–210, 2022.
- [9] Y. Bai, L. Wang, Y. Liu, Y. Yin, H. Di, and Y. Fu, "Human motion segmentation via velocity-sensitive dual-side auto-encoder," *IEEE Trans. Image Process.*, vol. 32, pp. 524 – 536, 2022.
- [10] X. Cao, C. Zhang, C. Zhou, H. Fu, and H. Foroosh, "Constrained multi-view video face clustering," *IEEE Trans. Image Process.*, vol. 24, no. 11, pp. 4381–4393, 2015.
- [11] A. Ng, M. Jordan, and Y. Weiss, "On spectral clustering: Analysis and an algorithm," *Advances in neural information processing systems*, vol. 14, 2001.
- [12] M. Rahmani and G. Atia, "Innovation pursuit: A new approach to the subspace clustering problem," in *International conference on machine learning*, 2017, pp. 2874–2882.
- [13] Z. Xing, W. Liu, B. Li, J. Tian, M. Chu, and J. Chen, "HMM-based CSI embedding for trajectory recovery via feature engineering on MIMO-OFDM channels in LOS/NLOS regions," in *Proceedings of IEEE/CIC International Conference on Communications in China (ICCC)*, 2025, pp. 1–6.
- [14] Z. Xing and W. Zhao, "Clustering structure identification with ordering graph," 2023. [Online]. Available: <https://openreview.net/forum?id=HG0SwOmlaEo>
- [15] Z. Xing and J. Chen, "Constructing indoor region-based radio map without location labels," *IEEE Trans. Signal Process.*, vol. 72, p. 2512–2526, 2024.
- [16] S. Wang, C. Li, Y. Li, Y. Yuan, and G. Wang, "Self-supervised information bottleneck for deep multi-view subspace clustering," *IEEE Trans. Image Process.*, vol. 32, pp. 1555–1567, 2023.
- [17] S. Wang, Z. Lin, Q. Cao, Y. Cen, and Y. Chen, "Bi-nuclear tensor Schatten-p norm minimization for multi-view subspace clustering," *IEEE Trans. Image Process.*, vol. 32, pp. 4059–4072, 2023.
- [18] Z. Chen, X.-J. Wu, T. Xu, and J. Kittler, "Fast self-guided multi-view subspace clustering," *IEEE Trans. Image Process.*, vol. 32, pp. 6514–6525, 2023.
- [19] Y. Tang, Y. Xie, and W. Zhang, "Affine subspace robust low-rank self-representation: from matrix to tensor," *IEEE Trans. Pattern Anal. Mach.*, vol. 45, no. 8, pp. 9357–9373, 2023.
- [20] Y. Chen, S. Wang, Y.-P. Zhao, and C. P. Chen, "Double discrete cosine transform-oriented multi-view subspace clustering," *IEEE Trans. Image Process.*, vol. 33, pp. 2491–2501, 2024.
- [21] Z. Xing and W. Zhao, "Trajectory map-matching in urban road networks based on RSS measurements," *IEEE Trans. Intell. Transp. Syst.*, vol. 26, no. 4, pp. 4647–4660, 2025.
- [22] Z. Xing, J. Chen, and Y. Tang, "Integrated segmentation and subspace clustering for RSS-based localization under blind calibration," in *Proc. IEEE Global Commun. Conf. (GlobeCom)*, 2022, pp. 5360–5365.
- [23] Z. Xing and W. Zhao, "Block-diagonal structure learning for subspace clustering," *Expert Systems with Applications*, vol. 285, no. 0957-4174, pp. 127 767–127 767, 2025.
- [24] Z. Xing, H. Li, W. Liu, Z. Ren, J. Chen, J. Xu, and C. Qin, "Spectrum efficiency prediction for real-world 5g networks based on drive testing data," in *Proc. IEEE Wireless Commun. Netw. Conf. (WCNC)*, 2022, pp. 2136–2141.
- [25] X. Wang, D. Guo, and P. Cheng, "Support structure representation learning for sequential data clustering," *Pattern Recognition*, vol. 122, p. 108326, 2022.
- [26] J. K. Aggarwal and Q. Cai, "Human motion analysis: A review," *Computer vision and image understanding*, vol. 73, no. 3, pp. 428–440, 1999.
- [27] R. Lan and H. Sun, "Automated human motion segmentation via motion regularities," *The Visual Computer*, vol. 31, no. 1, pp. 35–53, 2015.
- [28] L. Wang, W. Hu, and T. Tan, "Recent developments in human motion analysis," *Pattern recognition*, vol. 36, no. 3, pp. 585–601, 2003.
- [29] S. Schulz and A. Woerner, "Automatic motion segmentation for human motion synthesis," in *International Conference on Articulated Motion and Deformable Objects*. Springer, 2010, pp. 182–191.
- [30] R. Li, Z. Liu, and J. Tan, "Human motion segmentation using collaborative representations of 3d skeletal sequences," *IET Computer Vision*, vol. 12, no. 4, pp. 434–442, 2018.
- [31] H. Zhong, J. Shi, and M. Visontai, "Detecting unusual activity in video," in *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2, 2004, pp. II–II.
- [32] A. Fod, M. J. Matarić, and O. C. Jenkins, "Automated derivation of primitives for movement classification," *Autonomous robots*, vol. 12, pp. 39–54, 2002.
- [33] J. Barbič, A. Safonova, J.-Y. Pan, C. Faloutsos, J. K. Hodgins, and N. S. Pollard, "Segmenting motion capture data into distinct behaviors," in *Proceedings of Graphics Interface*, 2004, pp. 185–194.
- [34] P. Beaudoin, S. Coros, M. Van de Panne, and P. Poulin, "Motion-motif graphs," in *Proceedings of ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, 2008, pp. 117–126.
- [35] C. Lea, A. Reiter, R. Vidal, and G. D. Hager, "Segmental spatiotemporal cnns for fine-grained action segmentation," in *European Conference-Computer Vision*, 2016, pp. 36–52.
- [36] F. De la Torre, J. Campoy, Z. Ambadar, and J. F. Cohn, "Temporal segmentation of facial behavior," in *IEEE International Conference on Computer Vision*, 2007, pp. 1–8.
- [37] H. Gao, F. Guo, J. Zhu, Z. Kan, and X. Zhang, "Human motion segmentation based on structure constraint matrix factorization," *Science China Information Sciences*, vol. 65, no. 1, p. 119103, 2022.
- [38] Q. Jiang, M. Liu, X. Wang, M. Ge, and L. Lin, "Human motion segmentation and recognition using machine vision for mechanical assembly operation," *SpringerPlus*, vol. 5, no. 1, p. 1629, 2016.
- [39] Y. Liu, L. Feng, S. Liu, and M. Sun, "Sensor network oriented human motion segmentation with motion change measurement," *IEEE Access*, vol. 6, pp. 9281–9291, 2017.
- [40] J. F.-S. Lin, V. Joukov, and D. Kulic, "Human motion segmentation by data point classification," in *2014 36th Annual International Conference of the IEEE Engineering in Medicine and Biology Society*. IEEE, 2014, pp. 9–13.
- [41] Z. Xing and W. Zhao, "Segmentation and completion of human motion sequence via temporal learning of subspace variety model," *IEEE Trans. Image Process.*, vol. 33, pp. 5783–5797, 2024.
- [42] Z. Xing and J. Chen, "HMM-based CSI embedding for trajectory recovery from RSS measurements of non-cooperative devices," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2024, pp. 7060–7064.
- [43] Z. Xing and W. Zhao, "Calibration-free indoor positioning via regional channel tracing," *IEEE Internet Things J.*, vol. 12, no. 5, pp. 5449–5461, 2025.
- [44] W. Zhao, X. Yan, J. Gao, R. Zhang, J. Zhang, Z. Li, S. Wu, and S. Cui, "Pointlie: Locally invertible embedding for point cloud sampling and recovery," in *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence (IJCAI)*, 2021, pp. 1345–1351.
- [45] J. F.-S. Lin and D. Kulić, "Online segmentation of human motion for automated rehabilitation exercise analysis," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 22, no. 1, pp. 168–180, 2013.
- [46] M. Hoai, Z.-Z. Lan, and F. De la Torre, "Joint segmentation and classification of human actions in video," in *CVPR 2011*. IEEE, 2011, pp. 3265–3272.
- [47] Y. Guo, G. Xu, and S. Tsuji, "Understanding human motion patterns," in *Proceedings of the 12th IAPR International Conference on Pattern Recognition, Vol. 3-Conference C: Signal Processing (Cat. No. 94CH3440-5)*, vol. 2. IEEE, 1994, pp. 325–329.
- [48] D. Gehrig, T. Stein, A. Fischer, H. Schwameder, and T. Schultz, "Towards semantic segmentation of human motion sequences," in *Annual Conference on Artificial Intelligence*. Springer, 2010, pp. 436–443.
- [49] G. R. Bradschi and J. W. Davis, "Motion segmentation and pose recognition with motion history gradients," *Machine Vision and Applications*, vol. 13, no. 3, pp. 174–184, 2002.
- [50] D. Gong, G. Medioni, and X. Zhao, "Structured time series analysis for human action segmentation and recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 7, pp. 1414–1427, 2013.
- [51] M. Dimiccoli, L. Garrido, G. Rodriguez-Corominas, and H. Wendt, "Graph constrained data representation learning for human motion segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1460–1469.
- [52] L. Shao, L. Ji, Y. Liu, and J. Zhang, "Human action segmentation and recognition via motion and shape analysis," *Pattern Recognition Letters*, vol. 33, no. 4, pp. 438–445, 2012.
- [53] K. Kahol, P. Tripathi, S. Panchanathan, and T. Rikakis, "Gesture segmentation in complex motion sequences," in *Proceedings 2003 International Conference on Image Processing (Cat. No. 03CH37429)*, vol. 2. IEEE, 2003, pp. II–105.

- [54] J. F.-S. Lin, V. Bonnet, A. M. Panchea, N. Ramdani, G. Venture, and D. Kulić, "Human motion segmentation using cost weights recovered from inverse optimal control," in *2016 IEEE-RAS 16th International Conference on Humanoid Robots (Humanoids)*. IEEE, 2016, pp. 1107–1113.
- [55] C. Lu and N. J. Ferrier, "Repetitive motion analysis: Segmentation and event classification," *IEEE transactions on pattern analysis and machine intelligence*, vol. 26, no. 2, pp. 258–263, 2004.
- [56] Z. Xing and J. Chen, "Blind construction of angular power maps in massive MIMO networks," *IEEE Trans. Signal Process.*, vol. 00, no. 00, pp. 00–00, 2025.
- [57] W. Zhao, H. Zhang, C. Zheng, X. Yan, S. Cui, and Z. Li, "Cpu: Codebook lookup transformer with knowledge distillation for point cloud upsampling," in *Proceedings of the ACM International Conference on Multimedia*, 2023, p. 3917–3925.
- [58] Z. Xing and J. Chen, "Constructing angular power maps in massive MIMO networks using measurements without location labels," in *Proc. IEEE Int. Conf. Commun. (ICC)*, vol. 0, 2025, pp. 0–0.
- [59] Z. Xing and W. Zhao, "Block-diagonal guided DBSCAN clustering," *IEEE Trans. Knowl. Data Eng.*, vol. 36, no. 11, pp. 5709–5722, 2024.
- [60] Z. Xing and J. Chen, "Unsupervised radio map construction in mixed los/nlos indoor environments," in *Proc. IEEE Global Commun. Conf. (GlobeCom)*, 2025.
- [61] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma, "Robust recovery of subspace structures by low-rank representation," *IEEE Trans. Pattern Anal. Mach.*, vol. 35, no. 1, pp. 171–184, 2012.
- [62] Y. Qin, X. Zhang, L. Shen, and G. Feng, "Maximum block energy guided robust subspace clustering," *IEEE Trans. Pattern Anal. Mach.*, vol. 45, no. 2, 2022.
- [63] L. Wang, Z. Ding, and Y. Fu, "Learning transferable subspace for human motion segmentation," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, no. 1, 2018.
- [64] J. Xu, K. Xu, K. Chen, and J. Ruan, "Reweighted sparse subspace clustering," *Comput. Vis. Image Underst.*, vol. 138, pp. 25–37, 2015.
- [65] S. Wang, X. Yuan, T. Yao, S. Yan, and J. Shen, "Efficient subspace segmentation via quadratic programming," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 25, no. 1, 2011, pp. 519–524.
- [66] Z. Xing and W. Zhao, "K-means clustering: A review of the past 70 years," *Available at SSRN 5842722*, 2025.
- [67] Z. Xing and J. Chen, "Blind radio mapping via spatially regularized bayesian trajectory inference," *arXiv preprint arXiv:2512.13701*, 2025.
- [68] Z. Xing and W. Zhao, "K-means clustering: A review of the past 70 years," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024, pp. 6270–6278.
- [69] J. Cui, Y. Li, H. Huang, and J. Wen, "Dual contrast-driven deep multi-view clustering," *IEEE Trans. Image Process.*, vol. 33, pp. 4753–4764, 2024.
- [70] Y. Bai, L. Wang, Y. Liu, Y. Yin, and Y. Fu, "Dual-side auto-encoder for high-dimensional time series segmentation," in *IEEE International Conference on Data Mining*, 2020, pp. 918–923.
- [71] J. Guo, J. Li, D. Li, A. M. H. Tiong, B. Li, D. Tao, and S. Hoi, "From images to textual prompts: Zero-shot visual question answering with frozen large language models," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 10 867–10 877.
- [72] Y. Ge, X. Zeng, J. S. Huffman, T.-Y. Lin, M.-Y. Liu, and Y. Cui, "Visual fact checker: Enabling high-fidelity detailed caption generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 14 033–14 042.
- [73] H. Zhi, P. Chen, J. Li, S. Ma, X. Sun, T. Xiang, Y. Lei, M. Tan, and C. Gan, "Lscenellm: Enhancing large 3d scene understanding using adaptive visual preferences," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 3761–3771.
- [74] D. Zheng, S. Huang, and L. Wang, "Video-3d llm: Learning position-aware video representation for 3d scene understanding," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 8995–9006.
- [75] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, "Large language models are zero-shot reasoners," *Advances in neural information processing systems*, vol. 35, pp. 22 199–22 213, 2022.
- [76] D. Wang, Y. Zuo, F. Li, and J. Wu, "LLMs as zero-shot graph learners: Alignment of GNN representations with LLM token embeddings," in *Advances in Neural Information Processing Systems*, vol. 37. Curran Associates, Inc., 2024, pp. 5950–5973.
- [77] B. Feng, Z. Lai, S. Li, Z. Wang, S. Wang, P. Huang, and M. Cao, "Breaking down video LLM benchmarks: Knowledge, spatial perception, or true temporal understanding?" *arXiv preprint arXiv:2505.14321*, 2025.
- [78] R. Liu, C. Li, H. Tang, Y. Ge, Y. Shan, and G. Li, "St-llm: Large language models are effective temporal learners," in *European Conference on Computer Vision*. Springer, 2024, pp. 1–18.
- [79] Y. Yuan, H. Zhang, W. Li, Z. Cheng, B. Zhang, L. Li, X. Li, D. Zhao, W. Zhang, Y. Zhuang *et al.*, "Videorefer suite: Advancing spatial-temporal object understanding with video llm," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 18 970–18 980.
- [80] M. Nie, D. Ding, C. Wang, Y. Guo, J. Han, H. Xu, and L. Zhang, "Slowfocus: Enhancing fine-grained temporal understanding in video llm," *Advances in Neural Information Processing Systems*, vol. 37, pp. 81 808–81 835, 2024.
- [81] X. Ding and L. Wang, "Do language models understand time?" in *Companion Proceedings of the ACM on Web Conference 2025*, 2025, pp. 1855–1868.
- [82] A. Deng, Z. Gao, A. Choudhuri, B. Planche, M. Zheng, B. Wang, T. Chen, C. Chen, and Z. Wu, "Seq2time: Sequential knowledge transfer for video LLM temporal grounding," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 13 766–13 775.
- [83] L.-H. Chen, S. Lu, A. Zeng, H. Zhang, B. Wang, R. Zhang, and L. Zhang, "Motionllm: Understanding human behaviors from human motions and videos," *arXiv preprint arXiv:2405.20340*, 2024.
- [84] L. Li, S. Jia, J. Wang, Z. Jiang, F. Zhou, J. Dai, T. Zhang, Z. Wu, and J.-N. Hwang, "Human motion instruction tuning," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 17 582–17 591.
- [85] Z. Li, S. Deldari, L. Chen, H. Xue, and F. D. Salim, "Sensorllm: Aligning large language models with motion sensors for human activity recognition," *arXiv preprint arXiv:2410.10624*, 2024.
- [86] Y. Wang, S. Zheng, B. Cao, Q. Wei, W. Zeng, Q. Jin, and Z. Lu, "Scaling large motion models with million-level human motions," *arXiv preprint arXiv:2410.03311*, 2024.
- [87] B. Wang, Y. Tian, S. Wang, and L. Yang, "Multimodal large models are effective action anticipators," *IEEE Transactions on Multimedia*, vol. 27, pp. 2949–2960, 2025.
- [88] L. Lu, Y. Lu, R. Yu, H. Di, L. Zhang, and S. Wang, "GAIM: Graph attention interaction model for collective activity recognition," *IEEE Transactions on Multimedia*, vol. 22, no. 2, pp. 524–539, 2019.
- [89] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein *et al.*, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Found. Trends Mach. Learn.*, vol. 3, no. 1, pp. 1–122, 2011.
- [90] R. H. Bartels and G. W. Stewart, "Solution of the matrix equation  $ax+xb=c$  [4]," *Communications of the ACM*, vol. 15, no. 9, pp. 820–826, 1972.
- [91] K. Haynes, P. Fearnhead, and I. A. Eckley, "A computationally efficient nonparametric approach for changepoint detection," *Statistics and computing*, vol. 27, pp. 1293–1305, 2017.
- [92] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE Trans. Pattern Anal. Mach.*, vol. 22, no. 8, pp. 888–905, 2000.
- [93] E. Elhamifar and R. Vidal, "Sparse subspace clustering: Algorithm, theory, and applications," *IEEE Trans. Pattern Anal. Mach.*, vol. 35, no. 11, pp. 2765–2781, 2013.
- [94] Z. Jiang, Z. Lin, and L. Davis, "Recognizing human actions by learning and matching shape-motion prototype trees," *IEEE Trans. Pattern Anal. Mach.*, vol. 34, no. 3, pp. 533–547, 2012.
- [95] D. Huang, S. Yao, Y. Wang, and F. De La Torre, "Sequential max-margin event detectors," in *European Conference Computer Vision*, 2014, pp. 410–424.
- [96] M. S. Ryoo and J. K. Aggarwal, "Spatio-temporal relationship match: Video structure comparison for recognition of complex human activities," in *IEEE international conference on computer vision*, 2009, pp. 1593–1600.
- [97] L. Gorelick, M. Blank, E. Shechtman, M. Irani, and R. Basri, "Actions as space-time shapes," *IEEE Trans. Pattern Anal. Mach.*, vol. 29, no. 12, pp. 2247–2253, 2007.
- [98] H. W. Kuhn, "The hungarian method for the assignment problem," *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.
- [99] K. Zhan, F. Nie, J. Wang, and Y. Yang, "Multiview consensus graph clustering," *IEEE Trans. Image Process.*, vol. 28, no. 3, pp. 1261–1270, 2018.