

# Self-adjustable hyper-graphs for video pose estimation based on spatial-temporal subspace construction

Jizhou MA<sup>1</sup>, Shuai LI<sup>1\*</sup>, Hong QIN<sup>2</sup>, Aimin HAO<sup>1</sup> & Qiping ZHAO<sup>1</sup>

<sup>1</sup>State Key Laboratory of Virtual Reality Technology and Systems, Beihang University, Beijing 100191, China;

<sup>2</sup>Department of Computer Science, Stony Brook University, New York 11794, USA

Received 7 October 2019/Revised 25 January 2020/Accepted 7 April 2020/Published online 20 May 2021

**Citation** Ma J Z, Li S, Qin H, et al. Self-adjustable hyper-graphs for video pose estimation based on spatial-temporal subspace construction. *Sci China Inf Sci*, 2022, 65(3): 139101, https://doi.org/10.1007/s11432-019-2869-x

Dear editor,

In recent years many supervised video pose estimation methods have achieved growing successes based on well-labeled training datasets. Nonetheless, when facing roughly-labeled training data, it still remains challenging to intrinsically encode the video contents' spatial-temporal coherency for robust video pose estimation.

Some researches aimed to directly improve and refine the existing confidence maps by combining the spatial-temporal structure models [1, 2]. Li et al. [2] suggested that fixing some reliable estimations and formulating propagation processing as a 3D trajectory completion problem. Differently, Moon et al. [1] assumed that state-of-the-art 2D human pose estimation methods have similar error distributions. Zhou and Torre [3] first learned a codebook from the motion capture dataset, and then they employed a bi-linear model to estimate poses by matching the movement mode and the dense trajectory tracing result. It enables related patterns to be expressed using sub-patterns instead of a uniform probability distribution. In [3], this flexible codebook based framework still requires a lot of extra annotated data to ensure the models' dataset-specific applicability.

To overcome this drawback, we advocate a new hierarchical hyper-graph approach based on intrinsic spatial-temporal subspace exploration and propagation. Those "mis-matched" hyper-graph subspaces, which result from imperfect data, could be adaptively improved by taking advantage of visual contents' intrinsic continuities. At the theoretic level, the key idea for subspace exploration is to design a maximum matching subspace (MMS) operator, which help propagate highly correlated action information from local video frames to all video sequences in spatial-temporal subspaces. The hyper-graph is solely built based on our MMS metric, and it could synchronously encode cross-video action similarity, inner-video temporal coherency, and synergetic relationship of different body joints.

In contrast to normal "explicit hyper-graph", we construct an "implicit hyper-graph" by hierarchically repre-

senting different-level relationships. We conceptually split "explicit hyper-graph" into a series of sub-graphs (structure), which are formulated as optimized maximum matching subspaces. Then, these subspaces ("sub-graphs") will re-contact with each other via a global MMS operator based affinity matrix (metric).

Given a set of videos belonging to the same action category, each video is divided into a group of overlapping short video segments. The initial pose extractor ResNet50 [4], and such segments are represented as  $N_P$  pose sequences  $\mathbf{P} = \{\mathbf{P}_1, \mathbf{P}_2, \mathbf{P}_3, \dots, \mathbf{P}_{N_P}\}$ . For a pose sequence  $\mathbf{P}_i \in \mathbb{R}^{n_k \times n_f}$ , it covers  $n_k$  body joints and  $n_f$  consecutive frames. To align two pose sequences, similar to [3], we apply Procrustes analysis to get a spatial transition matrix  $\mathbf{Q}$ , which conducts an affine transformation, including translating, rotating and uniformly scaling. We mark it as  $\mathbf{Q}_i(\mathbf{P}_j) : \mathbf{P}_j \xrightarrow{\mathbf{Q}_i} \mathbf{P}_j$ , which means the pose sequence  $\mathbf{P}_j$  is aligned to  $\mathbf{P}_i$  by the transformation  $\mathbf{Q}_i$ , and  $\mathbf{Q}_i(\mathbf{P}_j)$  represents the transformed matrix.

$$\alpha_{i,j} = \arg \min_{\alpha_{i,j}} \|\mathbf{P}_i - \alpha_{i,j} \mathbf{P}_j\|_F. \quad (1)$$

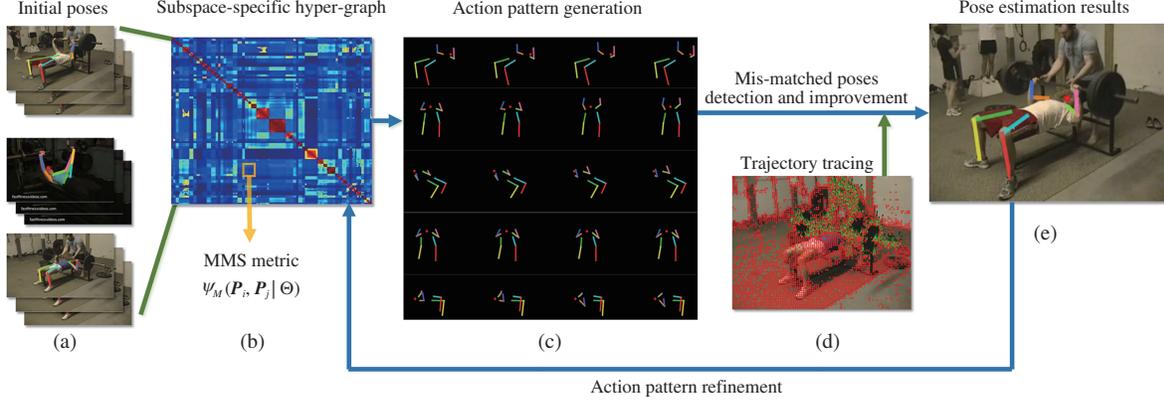
Here  $\alpha_{i,j}$  is the scaling coefficient, and it is the only parameter that needs to be optimized.  $\alpha_{i,j}$  may not equal  $\alpha_{j,i}$ , and thus  $\|\mathbf{P}_i - \alpha_{i,j} \mathbf{P}_j\|_F$  is not equal to  $\|\mathbf{P}_j - \alpha_{j,i} \mathbf{P}_i\|_F$  as well. The transformed affinities of the two pose sequences are asymmetric.

**MMS.** Most methods measure pose sequence affinities with a simple Frobenius-norm  $\|\mathbf{P}_i - \mathbf{P}_j\|_F$ , which misses a synergetic relationship analysis over joints and frames. It means that, two adjacent pose sequences in a video may be greatly different in the pose space because of temporal topology dislocation.

$$\psi_M(\mathbf{P}_i, \mathbf{P}_j | \Theta) = \|\mathbf{P}_i(K, F_i) - \mathbf{Q}_i(\mathbf{P}_j)(K, F_j)\|_F. \quad (2)$$

Here  $\Theta = \{K, F_i, F_j\}, (K, F_i)$ , and  $(K, F_j)$  are the subspace parameters, which actually is the sub-matrix index and is determined by the specific sequence pair.  $\{K, F_i, F_j\}$  respectively denotes the set (indexes) of joints, frame in  $\mathbf{P}_i$  and  $\mathbf{P}_j$ .

\* Corresponding author (email: lishuai@buaa.edu.cn)



**Figure 1** (Color online) The pipeline of our framework. (a) Input video sequences within the common action category; (b) the MMS metric governed subspace-specific hyper-graph construction; (c) generating the action pattern which consists of a series of clustered action sub-patterns; (d) the dense trajectory tracing for local details capturing; (e) detecting and improving existing mis-matched (inaccurate) poses by incorporating the semantic guidance from the action patterns and local feature tracing from the dense trajectories.

Both of the two sequences have the same number of frames ( $|F_i| = |F_j|$ ). As a result,  $\psi_M$  partitions the difference of  $P_i$  and  $P_j$  into aligned matching sub-matrix  $\Theta$  and unmatched sub-matrix. In this study,  $f/F$  indicates frames,  $k/K$  indicates body joints, and  $\Theta$  represents both of them. The  $\Theta$  can be efficiently computed via a simple path searching algorithm. The searching space encodes spatial-temporal affinity in a 3-dimension matrix  $D_{i,j} \in \mathbb{R}^{n_f \times n_f \times n_k}$ . Its element is defined as

$$D(k, f_i, f_j) = D_{k, f_i, f_j} = \|P_i(k, f_i) - Q_i(P_j)(k, f_j)\|_2. \quad (3)$$

Then, the path is partitioned at the steepest point with the maximum gradient. The first half of the path indexes the maximum matching subspaces. On the other hand, the rest of the subspaces denotes those dissimilar parts of pose sequences.

*Action pattern generation.* The action expression encoded in our hyper-graph is comprehensive but losing generality. Thus, we explore some mid-level semantic sub-patterns (sub-actions) to describe a complete action. These sub-patterns can be considered as certain words in “action codebook” to describe the pose sequences.

To begin with, we build an affinity matrix  $A \in \mathbb{R}^{N_P \times N_P}$  for all the input pose sequences, with

$$\begin{aligned} A(i, j) &= A(j, i) \\ &= \frac{1}{2} \psi_M(P_i, Q_i(P_j) | \Theta) + \frac{1}{2} \psi_M(P_j, Q_j(P_i) | \Theta). \end{aligned} \quad (4)$$

Based on  $A$ , we cluster all the pose sequences into  $N_B$  groups via spectral clustering. For the  $c$ -th group, we conduct hierarchical clustering and orderly merge the pose sequences one by one to construct a sub-pattern  $B_c$ . Yet, in the following step, we also need to trade off the contribution of those densely sampled trajectories [5] for the local detail tracing. Therefore, we adopt the 2D Gaussian function to define motional joint models.

As for “mean” pattern (pose sequence)  $\mu \in \mathbb{R}^{n_k \times n_f}$ , it has the same formulation as a pose sequence  $P_i$ .  $B_c$  is built by aggregation processing over  $n_c$  pose sequences  $P_n$  included in  $B_c$ , as

$$\begin{aligned} B_c^{(n+1)}(\Theta) &= B_c^{(n)}(\Theta) + \lambda^{(n)} \cdot P_n(\Theta), \quad n \in [1, n_c] \\ \text{s.t. } B_c^{(0)} &= \mathbf{0}^{n_k \times n_f} \text{ and } \Theta \in \psi_M(B_c^{(n)}, P_n | \Theta). \end{aligned} \quad (5)$$

Here  $\lambda^{(n)}$  denotes the hierarchical clustering distance from sub-pattern  $B_c$  to pose sequence  $P_n$ .  $\Theta$  is computed by MMS operator, which changes dynamically corresponding to current  $B_c^{(n)}$  and  $P_n$ . After constructing all sub-patterns  $B_c$ , we can obtain a complete action pattern as  $B = \{B_1, B_2, \dots, B_{N_B}\}$ .

We have defined our MMS operator  $\psi_M(P_i, P_j | \Theta)$  for two pose sequences in Eq. (2), and we also need to use it on action sub-patterns as  $\psi_M(P_i, B_c | \Theta)$ . Therefore, we redefine the distance matrix in (3) for adapting MMS to the probabilistic sub-patterns  $B_c = \{P_\mu, \rho\}$ , with

$$\begin{aligned} d &= \|P_i(k, f_i) - Q_i(B_c)(k, f_c)\|_1, \\ D_{\text{prob}}(d | B_c) &= D_{k, i, c} = \frac{1}{2\pi|\rho|^{-\frac{1}{2}}} \exp\left(-\frac{1}{2}(d^T \rho^{-1} d)\right), \end{aligned} \quad (6)$$

where  $d \in \mathbb{R}^{2 \times 1}$  denotes the distance of joint  $P_i(k, f_i)$  to the mean joint position of sub-patterns  $B_c(k, f_c)$ .  $D_{i, c}$  shows the 2D Gaussian function based affinity. In Eq. (3),  $D$  is a distance (dissimilarity) matrix, but in (6),  $D$  shows the probabilistic affinity (similarity).

Considering the confused spatial-temporal continuity over initial poses, frame-wisely refining the pose is more robust than global sequence-wise processing. The  $f$ -th frame error  $E(p_f)$  is gotten by averaging the reconstruction error overall pose sequences covering the  $f$ -th frame ( $p_f \in P_i$ ). The error  $E$  of  $f$ -th frame’s pose  $p_f$  is defined as

$$E(p_f) = E_{\text{tra}}(p_f | \dot{B}_c) + \beta \cdot E_{\text{act}}(p_f | \dot{B}_c), \quad (7)$$

where  $E_{\text{tra}}$  denotes the trajectory tracing error and  $E_{\text{act}}$  means the action pattern governed reconstruction error. The matched sub-pattern  $\dot{B}_c$  guides the computation procedure of  $E_{\text{tra}}$  and  $E_{\text{act}}$ . To compute  $E_{\text{tra}}$ ,  $\dot{B}_c$  determine the weight of each trajectory around the joints. As for  $E_{\text{act}}$ ,  $\dot{B}_c$  measures the action fitting error. We find the best matched sub-pattern for each pose sequence, with

$$(P_i, \dot{B}_c) = \arg \min_{P_i, B_c} \psi_M(P_i, B_c | \Theta), \quad (8)$$

where  $B_c \in B$  and all the related pose sequences  $P_i$  should cover the  $f$ -th frame,  $P_i(\Theta) \ni p_f$ .  $\Theta$  indexes a sub-matrix of  $P_i$ , and  $P_i$  may cover  $p_f$  but  $P_i(\Theta)$  not.

$\mathbb{T}$  is a local feature tracing function, and it gives the trajectory tracing results  $\mathbf{t} \in \mathbb{R}^{n_k \times 1}$  related to  $\mathbf{p}_f$  from the last  $(f-1)$ -th frame  $\mathbf{t}^{(0)}$  to the current  $f$ -th frame  $\mathbf{t}^{(1)}$ .  $E_{\text{tra}}$  is formulated as

$$E_{\text{tra}}(\mathbf{p}_f) = \frac{1}{|\mathbb{T}(\mathbf{p}_f)|} \cdot \sum_{\mathbf{t} \in \mathbb{T}(\mathbf{p}_f)} \|D_{\text{prob}}(\mathbf{t}^{(1)} | \dot{\mathbf{B}}_c) - D_{\text{prob}}(\mathbf{t}^{(0)} | \dot{\mathbf{B}}_c)\|_2, \quad (9)$$

where  $E_{\text{tra}}$  is able to detect the errors caused by drastic pose location changes. Moreover, we define the action pattern based reconstruction error as

$$E_{\text{act}}(\mathbf{p}_f) = \frac{1}{n_m} \sum_{\mathbf{P}_i \ni \mathbf{p}_f} \psi_M(\mathbf{P}_i, \dot{\mathbf{B}}_c | \Theta(f)). \quad (10)$$

Here  $\Theta(f)$  denotes the frame-specific parameter, which only indexes the matching joints set  $K$  at the  $f$ -th frame.  $n_m$  is the size of the candidate matching set. In practice, we do not use all eligible  $\mathbf{P}_i$  and only pick 5 ( $n_m = 5$ ) pose sequences with the minimum matching errors (Eq. (8)).

Finally, as shown in Figure 1, we formulate an iterative pose estimation framework, which alternately conducts pose improvement and action pattern refinement. The intermediately improved poses facilitate pattern amending. We fix the top 5% frames (sorted by their errors  $E(\mathbf{p}_f)$ ) for pose improvement in each loop.

*Conclusion.* We have detailed a spatial-temporal subspace involved hyper-graph method for human pose estimation in video analysis. The method can indeed improve the

existing estimated results even without the need for labeled ground-truth poses. All the exhibited advantages of the new method result from the MMS operator, which enables intrinsic encoding of the action similarity between videos, the intra-video temporal coherency, and the collaborative relevance over different body joints.

**Acknowledgements** This work was supported in part by National Key R&D Program of China (Grant No. 2018YFB-1700603), National Natural Science Foundation of China (Grant Nos. 61672077, 61532002), and Beijing Natural Science Foundation — Haidian Primitive Innovation Joint Fund (Grant No. L182016).

## References

- 1 Moon G, Chang J Y, Lee K M. Posefix: model-agnostic general human pose refinement network. In: Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019
- 2 Li Z, Wang X, Wang F, et al. On boosting single-frame 3D human pose estimation via monocular videos. In: Proceedings of IEEE/CVF International Conference on Computer Vision, 2019
- 3 Zhou F, de la Torre F. Spatio-temporal matching for human pose estimation in video. IEEE Trans Pattern Anal Mach Intell, 2016, 38: 1492–1504
- 4 Xiao B, Wu H P, Wei Y C. Simple baselines for human pose estimation and tracking. In: Proceedings of European Conference on Computer Vision, 2018
- 5 Wang H, Schmid C. Action recognition with improved trajectories. In: Proceedings of IEEE International Conference on Computer Vision, 2013