

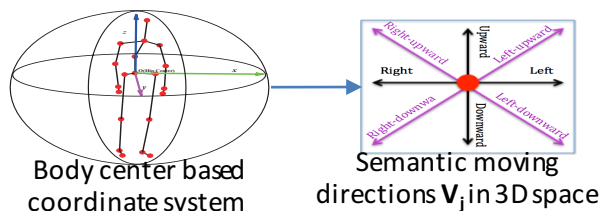
1. Abstract

Human-human interaction recognition has attracted increasing attention in recent years due to its wide applications in computer vision fields. However, it still remains a challenge due to mutual occlusion, various subject appearance or body size and complex context. In this paper, a novel feature descriptor based on spatial relationship and semantic motion trend similarity between body parts is proposed for human-human interaction recognition.

2. Feature Extraction

The moving direction of each joint is firstly quantified into several semantic words and then the motion trend similarity is captured by histogram intersection.

2.1 Motion trend feature

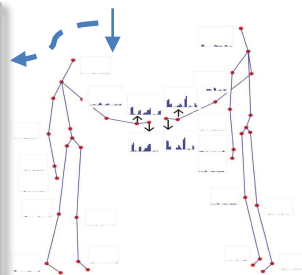


- Frame-direction:

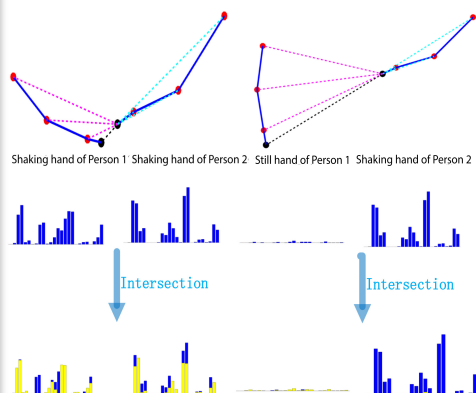
$$v_i^k = \{x_{p_i^k} - x_{p_{i-1}^k}, y_{p_i^k} - y_{p_{i-1}^k}, z_{p_i^k} - z_{p_{i-1}^k}\}$$
- Soft direction encoding depending on cosine similarity and displacement:

$$\cos\theta_j^k(t) = \frac{v_j \cdot v_i^k}{\|v_i^k\| \|v_j\|}, j \in [1, m]$$

$$\begin{cases} bin_{first} = bin_{first} + Dis^k(t) \times \cos\theta_{first}^k(t) \\ bin_{second} = bin_{second} + Dis^k(t) \times \cos\theta_{second}^k(t) \end{cases}$$



2.2 Motion trend similarity



$$InterSectionJ(b_i^k, b_j^k) = \min(|b_i^m|, |b_j^m|)$$

$$H = \left(\underbrace{1, \dots, 1, 0, \dots, 0}_{N-H_1}, \underbrace{1, \dots, 1, 0, \dots, 0}_{N-H_2}, \dots, \underbrace{1, \dots, 1, 0, \dots, 0}_{N-H_m} \right)$$

$$K_{type} = \sum_{p=1}^8 \sum_{q=1}^8 \sum_{k=1}^{26} InterSectionJ(b_p^k, b_q^k)$$

$$= \sum_{p=1}^8 \sum_{q=1}^8 K(H(p), H(q))$$

Where,

$$K(H(i), H(j)) = \sum_{k=1}^m InterSectionJ(b_i^k, b_j^k) = H(i) \cdot H(j)$$

3. Experiment Results

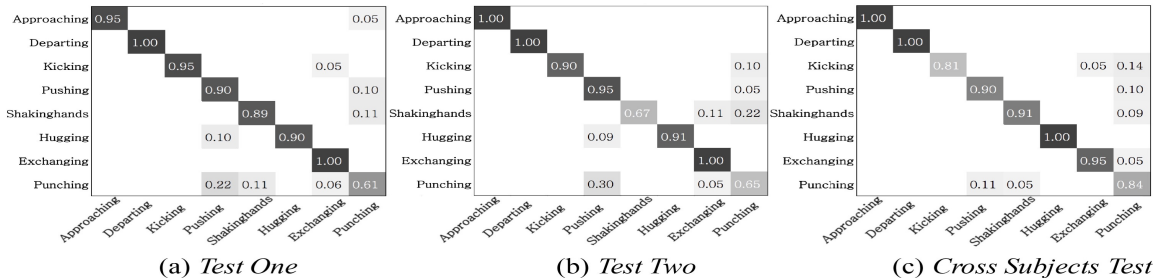


Fig. 1 Confusion matrices

Table 1: Recognition Accuracy (%) on SBU dataset.

State-of-the-art	Joint features [13] CFDM [9] [18]	80.3 89.4 86.9
Proposed Method	Test One	90.58
	Test Two	90.28
	Cross Subjects Test Average	92.50 91.12

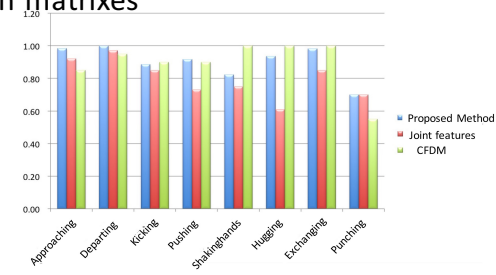


Fig. 2 Comparison between categories