

|              |   |
|--------------|---|
| Title        | 人の動作認識のための可分な意味特徴の自動抽出に関する研究  |
| Author(s)    | Tran, Thang Thanh   |
| Citation     |   |
| Issue Date   | 2014-09   |
| Type         | Thesis or Dissertation  |
| Text version | ETD   |
| URL          | <a href="http://hdl.handle.net/10119/12289">http://hdl.handle.net/10119/12289</a> |
| Rights       |   |
| Description  | Supervisor:小谷 一孔, 情報科学研究科, 博士   |

## **Thesis title**

Automatic Extraction of Discriminative Semantic Features for Action Recognition

## **Abstract**

In our everyday lives, we are constantly confronted with human motion: we watch other people as they move, and we perceive our own movements. Motion, in terms of physics, is a change in the position of a body with respect to time. Since the human body is not simply a rigid block but rather a complex aggregation of flexibly connected limbs and body parts, human motion can have a very complex spatial-temporal structure. Deeper understanding on human actions is required in many applications, e.g., action recognition (security), animation (sport, 3D cartoon movies and virtual world), etc.

With the development of the technology like 3D specialized markers, we could capture the moving signals from marker joints and create a huge set of 3D action motion capture (MOCAP) data which is capable of accurately digitizing a motion's spatial-temporal structure for further processing on a computer. Recently, motion capture data have become publicly available on a larger scale, e.g. CMU, HDM05. The task of automatic extraction of semantic action features is gaining in importance. The underlying questions are how to measure similarity between motions or how to compare motions in a meaningful way. The main problem is that the granularity of MOCAP data is too fine for our purpose: Human actions typically exhibit global and local temporal deformation, i.e. different movement speed and timing difference. The similar types of motions may exhibit significant spatial and temporal variations. The irrelevant details (like noise) as well as spatial pose deformations may interfere with the actual semantics that we are trying to capture. The problems require the identification and extraction of logically related motions scattered within the data set. This leads us to the field of motion analysis for identifying the significant spatial and temporal features of an action.

To automatically extract the action features from 3D MOCAP data, we proposed two approaches dealing at features levels: 1) Extract of Discriminate Patterns from Skeleton Sequences approach provides a foundation in lower dimensional representation for the movement sequence analysis, retrieval, identification and synthesis; and 2) Automatic Extraction of Semantic Action Features approach which focuses on solving the high-dimensional computational problems arising from the human motion sequences. They support the follow-up stages of processing the human movement on a natural language level. As one common underlying concept, the proposed approaches contain a retrieval component for extracting the above-mentioned features.

Firstly, we extract the discriminative patterns as local features and the utilization of a statistical approach in text classification to recognize actions. In text classification, documents are presented as vectors where each component is associated to a particular word from the code book. Traditional weighting methods like Term Frequency Inverse Document Frequency (TF-IDF) are used to estimate the importance of each word in the document. In this approach, we use the beyond TF-IDF weighting method to extract discrimination patterns which obtain a set of characteristics that remain relatively constant to separate different categories. This weighting defines the importance of words in representing specific categories of documents. It not only reduces the number of feature dimension compared to the original 3D sequence of skeletons, but also reduces the viewing time of browsing, bandwidth, and computational requirement of retrieval.

Secondly, we propose the semantic annotation approach of the human motion capture data and use the Relational Feature concept to automatically extract a set of action features. For each action class, we propose a statistical method to extract the common sets. The features extracted is used to recognize the action in real-time. We extract the set of action features automatically based on the velocity feature of body joints. We consider this set as action spatial information. We combine both spatial and temporal

processes to extract the action features and use them for action recognition. In our experiments, we use the 3D motion capture database HDM05 for performance validation. With few training samples, our experiment shows that the features extracted by this method achieve high accuracy in recognizing actions on testing data. Our proposed method gets high accuracy comparing to others state-of-art approaches.

### **Keywords**

Discriminate Semantic Features; Automatic Extraction Features; Semantic Action Features; Action Recognition; Joint Velocity.