

Automated Modularization of Human Motion into Actions and Behaviors

Odest Chadwicke Jenkins
Robotics Research Laboratory
Department of Computer Science
University of Southern California
Los Angeles, CA 90089-0781
cjenkins@cs.usc.edu

Maja J Matarić
Robotics Research Laboratory
Department of Computer Science
University of Southern California
Los Angeles, CA 90089-0781
mataric@cs.usc.edu

Abstract

In this paper we address the problem of automatically deriving vocabularies of motion modules from human motion data, taking advantage of the underlying structure in motion. We approach this problem with a data-driven methodology for modularizing a motion stream (or time-series of human motion) into a vocabulary of parameterized actions and a set of high-level behaviors for sequencing actions. Central to this methodology is the discovery of spatio-temporal structure in a motion stream. We estimate this structure by using a spatio-temporal dimension reduction method based on extended Isomap. The utility of the derived vocabularies is validated through their use in synthesizing new humanoid motion.

1 Introduction

Human motion data are used broadly in fields such as character animation and humanoid robotics. However, working with such data can be a tedious process requiring extensive skilled manual effort. One of the major challenges is modification of existing, captured motion to create a new, desired resulting motion. Several research efforts have been devoted to developing methods for addressing this problem, such as [3] for parameterized robot trajectory planning and [4] for constraint-based motion editing. Typically, motion reuse approaches are restricted to a single type of motion and/or require significant user intervention.

In this paper, we present an approach (summarized in Figure 1) to automatically deriving a vocabulary of human motion modules that allow new motion to be synthesized. Our approach is applicable to multiple types of human motion and greatly reduces the cost of manual effort involved in synthesizing new motion. We assume that the motion of an individual human performer has structure and our aim is to provide a feasible estimate of this structure from which a vocabulary of motion modules can be built. Each module represents a set of motion with a common theme or meaning (e.g., punch, jab, reach) and is also capable of realizing new motions that are variations on the same theme. Once the vocabulary is established, it can

be used to synthesize new motion and sequences of movements.

We address several issues involved in building motion vocabularies. First, we address the derivation of structure from human motion using dimension reduction, similar to [2]. Most existing dimension reduction techniques assume spatial data which have no temporal order. However, we seek to take advantage of both spatial and temporal dependencies in movement data. Tenenbaum et al. [9] alluded to use of temporal order within the context of Isomap. Our approach extends Isomap to extract spatio-temporal structure. With additional clustering and interpolation mechanisms, we derive parameterized *action primitives*. These primitives are similar to “verbs” in manually-derived *Verbs and Adverbs* vocabularies [5]. Drawing an analogy to linguistic grammars, action primitives could be considered terminals. Further dimension reduction iterations allow for the derivation of *parent behavior* modules that represent sequences of action primitives. Parent behaviors are similar to “verb graphs” in Verb and Adverbs and could be considered non-terminals in a grammar.

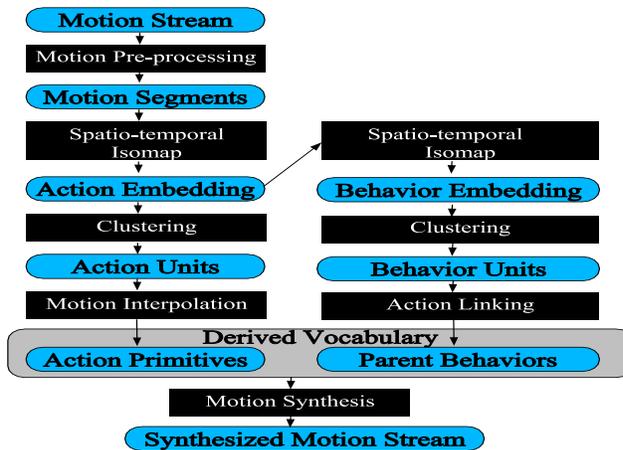


Figure 1: *Flowchart of approach. The input to the system is a motion stream, which is segmented (using one of several approaches). Dimension reduction, clustering, and interpolation are applied to the segments to derive action primitives. Using the initial embedding, another iteration of dimension reduction and clustering are applied to find behavior units. Parent behaviors are formed by determining component actions from a behavior unit and linking those with derived transition probabilities. The resulting action and behavior vocabulary is used to synthesize novel motion streams.*

2 Spatio-temporal dimension reduction

Central to our motion derivation methodology is the ability to transform a set of motion data so that its underlying structure can be estimated. Our approach is to use dimension reduction to extract structure-indicative features. We use joint angle data as input, and consider a motion segment as a data point in a space of joint angle trajectories with dimensionality equal to the product of the number of frames in the segment and the number of degrees of freedom in the kinematics of the humanoid data source. The dimension reduction mechanism transforms the motion data into a space in which the extraction of features can be performed.

Recently, [2] used *Principal Components Analysis (PCA)* to derive movement primitives from human arm motion. While the reduction of dimension using PCA is useful, the extraction of features from the linear PCA-embedded subspace is unintuitive. The derived principal components (PCs) are not convenient features, as they do not have a meaningful interpretation at least in part because the input data are non-linear. Alternatively, features could be constructed by clustering the PCA dimension-reduced data. However, without an a priori specification for the number of clusters, this also produced unintuitive features. Furthermore, the best result from clustering will provide a discretization of the space of joint angle trajectories. This limitation restricts our ability to extract features that vary across a wide volume in the space of joint angle trajectories, potentially overlapping with other features. Several nonlinear dimension reduction techniques, including [9], [6], and [7] address the linear limitation of PCA. However, these approaches perform a *spatial* nonlinear dimension reduction. Consequently, our attempts to apply these approaches to the same arm motion data resulted in the same feature extraction problems as described for PCA, except with fewer and more appropriate extracted features.

Human motion data have a meaningful temporal ordering that can be utilized for feature extraction. In our approach, we use long streams of motion as input, which we then segment (automatically or manually), thus retaining the natural temporal sequence in which the segments occur. This temporal structure needs to now be retained and utilized by Isomap. Figure 2 illustrates the differences between PCA, Isomap, and our spatio-temporal extension of Isomap for sequentially-ordered S-Curve data describing the motion of a point in 3D.

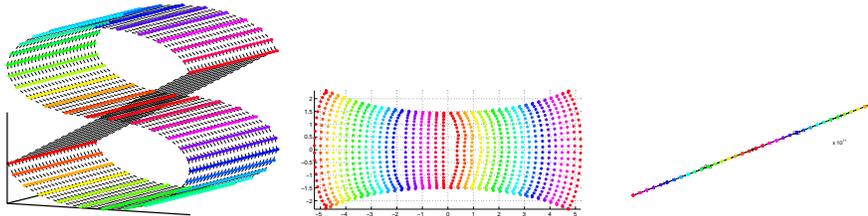


Figure 2: “S”-curve trajectory example motivated by [6]. (a) plot of 3D trajectory of a point moving along an S-Curve with its temporal order specified by the dashed line. The point traverses the S and returns to its initial position along the “S”, shifted slightly off of the previous S. (b) The embedding produced by spatial Isomap removes the spatial nonlinearity of the S. (c) The embedding produced by spatio-temporal Isomap removes the spatial nonlinearity and collapses the multiple traversals of the S to a line.

In the remainder of this section, we describe our extension of the Isomap algorithm for spatio-temporal dimension reduction to allow for the extraction of meaningful features. The procedure for *spatio-temporal Isomap* is as follows:

1. Determine a local neighborhood for at each point in the data set.
 - (a) Set the distance $D(i, j)$ between point i and a neighboring point j based on a chosen distance metric
2. Adjust D matrix for temporal dependencies
3. Compute all-pairs shortest paths, with D containing pairwise distances
4. Construct d -dimensional embedding by an eigenvalue decomposition of D , given d

This algorithm is the same as Isomap, except for step 2 (underlined) above. We provide spatio-temporal embedding by the adjustment of the similarity matrix D , through *Weighted Temporal Neighbors (WTN)*, or a different distance metric, through a *phase-space distance metric*. The phase-space distance metric was implemented as the Euclidean distance between two data points with concatenated spatial and velocity information. A comparison of the two methods showed that WTN provided more meaningful embeddings for deriving motion vocabularies.

2.1 Weighted Temporal Neighbors

In spatial Isomap, neighborhoods local to each point x_i are formed by the spatially closest points to x_i . In WTN, these spatial neighbors and adjacent temporal neighbors, points x_{i-1} and x_{i+1} , form local neighborhoods. By including adjacent temporal neighbors, a single connected component can be realized in the D matrix and, thus, include all of the data points in the embedding. We use the constant c_{atn} to regulate the distance in the D matrix between a point and its adjacent temporal neighbors.

WTN also modifies the D matrix based on Common Temporal Neighbors. We define two data points, t_x and t_y , as *common temporal neighbors (ctn)* if $t_y \in nbhd(t_x)$ and $t_{ya} \in nbhd(t_{xa})$, where $nbhd(t_x)$ is the spatial neighborhood of t_x , t_{x+1} and t_{y+1} are data points temporally adjacent to t_x and t_y , respectively. Common temporal neighbors are used to identify points in the local spatial neighborhood that are more likely to be grouped in the same feature. We use a constant c_{ctn} to specify how much to reduce the distance in D between two common temporal neighbors. By providing a significant distance reduction between two common temporal neighbors, we ensure that these two points will be *proximal* in the resulting embedding. Two points that are not common temporal neighbors, but are linked by common temporal neighbors, will also be proximal in the embedding. We define a set of points in which all pairs in the set are connected by a path of common temporal neighbors as a *similar connected component*. Points belonging to a single similar connected component will be proximal in the embedding and thus distinguishable through clustering.

3 Deriving actions and behaviors

We iteratively apply spatio-temporal Isomap to human motion data to produce embedding from which actions and behaviors can be extracted. The first iteration of embedding yields clusterable groups of motion, called *action units*. An interpolation mechanism is combined with each action unit to form an *action primitive* capable of producing new instances of the action. Spatio-temporal Isomap is applied to the first embedding to yield clusterable groups of actions *behavior units*. From a behavior units, a parent behavior is formed to encapsulate its component action primitives and linking them with transition probabilities.

3.1 Motion pre-processing

The first step in the derivation of actions and behaviors is the segmentation of a single motion stream into a set of motion segments. The motion streams consisted of human upper-body motion with 27 degrees of freedom (DOFs). The streams were segmented manually, for ground truth, and also by using *Kinematic Centroid Segmentation (KCS)*. KCS segments the motion of a kinematic substructure (e.g., an arm) based on the motion of a centroid feature that is the average of a set of Cartesian features along the arm. KCS determines motion segment boundaries in a greedy fashion using the following procedure:

1. Set current segment to the first frame
2. Compute distance between centroid at current segment boundary and centroid at every subsequent frame
3. Find first local maximum in centroid distance function
 - (a) Traverse frames with a moving window until current frame exceeds a distance threshold and is the maximum value in the window
4. Place new segment at the found local maximum, goto step 2

KCS was applied to each kinematic substructure (left arm, right arm, and torso) independently and proximal segment boundaries were merged. For both segmentation methods, the resulting segments were normalized to a fixed number of frames using cubic spline interpolation.

3.2 Extracting action units and action primitives

Spatio-temporal Isomap was performed on the set of motion segments using WTN for temporal adjustment of the pairwise distance matrix D . For extracting actions, we aim to collapse motion segments which are in the same similar connected component into proximity in the resulting embedding. This can be achieved with a significantly large value assigned to c_{ctn} . However, c_{atn} cannot be set to a single constant because the adjacent temporal neighbor of a motion segment may or may not be included in the same action unit. We set c_{atn} to a negative value as a flag for spatio-temporal Isomap to set the distance between adjacent temporal neighbors as their spatial distance.

The resulting embedding produces linearly separable groups that can be extracted automatically by clustering. The clustering method is implemented based on the one-dimensional “sweep-and-prune” technique introduced in [1] for detecting overlapping axis-aligned bounding boxes. This clustering method does not require the number of clusters to be specified *a priori*, but rather a separating distance for distinguishing cluster intervals along each dimension. Once clustering is applied, each cluster is considered an action unit.

Next, each action unit is generalized to an action primitive using interpolation. An action unit contains a set of motion segments and corresponding coordinates in a reduced space. Thus, we can select new coordinates in that space and interpolate within the joint angle trajectory space to produce new motion variations of the action. A variety of interpolation mechanisms could be used. We chose Shepard’s interpolation [8] because of the simplicity of its implementation. The locations from the embedding are used for interpolation coordinates, although these locations could be repositioned manually or automatically (potentially through spatial Isomap) to improve interpolation results.

3.3 Extracting behavior units and parent behaviors

Using the first embedding, we perform another iteration of spatio-temporal Isomap to extract behavior units. From the first embedding, motion segments from the same action were proximal. In the second embedding, we aim to collapse actions typically performed sequentially into behaviors. Consequently, we must collapse motion segments from actions in this sequence into the same behavior unit. For this to take place, we set c_{ctn} to a large constant and c_{atn} to a constant small enough to collapse action units but large enough to not split our existing action units. We then perform bounding-box clustering on this embedding.

We generalized behavior units into parent behaviors by determining component actions, providing transition probabilities, and providing the ability to determine valid transitions between component behaviors. Behavior units only specify motion segments belonging to a particular behavior. By associating action units and behavior units with motion segments in common, we can determine what actions are components of a certain behavior. We also used the segments in the behavior units to determine the transition probabilities between component actions. A transition from one action to another is valid if no large “hops” in joint space occur in the resulting motion. In order to determine valid transitions, we densely sampled each action primitive for interpolated motions. When contemplating a transition to a certain action, the parent behavior can examine the interpolated motions for the action to determine if a valid transition to the action exists.

4 Motion synthesis

As an example of the usefulness of vocabularies we automatically derived, we implemented a mechanism to synthesize a stream of human motion from a user-selected parent behavior. The input to this mechanism is the selection of a parent behavior. An initial action is selected, by the user or randomly, from the component actions of the selected behavior. An interpolated motion segment, sampled from the current action primitive, is the initial piece in a synthesized motion stream. Given the currently selected action and interpolated motion, we collect a set of *append candidates* from the interpolated motion segments of the component actions. Append candidates are motion segments that would provide a valid transition if appended to the end of the current synthesized motion stream. Transition validity is enforced with a threshold on the joint-space distance between the last frame of the current synthesized motion stream and the first frame of the interpolated motion segment. The append candidates are weighted by the transition probabilities from the current action and the action represented by the candidate segment. The current synthesized motion stream is updated by appending a randomly selected append candidate. Using the current action and synthesized stream, the process of selecting append candidates and updating the stream is repeated until no append candidates are available or a stopping condition is reached.

5 Results and discussion

Using the implementation of our method, we derived action and behavior vocabularies for two different streams of motion. The first stream contained motion of a human performing various activities, including several types of punching, dancing, arm waving, semaphores, and circular hand movements. The second stream contained only two-arm reaching motions to various goal positions. Each stream used 27 DOFs to describe the upper body of the performer and contained 22,549 and 9,145 frames, respectively. The streams were segmented using KCS and segments were time-normalized to 100 frames. We derived 56 actions and 14 behaviors for the first stream and 2 actions and 2 behaviors for the second stream. The results for the first stream are shown in Figure 3. The vocabulary derived from the second stream was as expected; the two actions, “reach out” and “return to idle”, formed one reaching behavior. The second behavior was irrelevant because it contained only a single transient segment. The derived vocabulary from the first stream also produced a desirable results, including distinct behaviors for the circular, punching, and some of the dancing behaviors and no distinct behavior for the semaphores, which had no meaningful motion pattern. However, the arm waving behavior merged with several dancing behaviors due to their similarity. Each derived action was sampled

for 200 new interpolated instances. Our motion synthesis mechanism was used to sequence new motion streams for each derived behavior. Each stream synthesizes plausible motion for the structure of each behavior. Excerpts from these streams for a few behaviors are annotated in Figure 3.

There is a distinct trade-off in exchanging convenience of our automated approach for the elegance of a manual approach. The common sense and skill of a human animator or programmer allow for an intuitive semantic meaning and expressiveness to be applied to motion. A significant cost, however, is incurred in terms of time, effort, and training. By using our automated approach, we reduce the cost of manual intervention in exchange for useful motion modules with an observable, but not explicitly stated, meaning. Even used with unsophisticated techniques for clustering, interpolation, and synthesis, our derived vocabularies were able to produce plausible motion requiring manual effort for observing the types of motion produced by each derived behavior.

6 Conclusion

We have proposed an approach for deriving vocabularies consisting of action and behavior modules from streams of human motion data. We were able to derive these vocabularies based on embeddings produced by our extension of Isomap for spatio-temporal data. Using these derived motion vocabularies, we demonstrated the usefulness of our approach with respect to synthesizing new human motion. Our vocabulary derivation and motion synthesis procedures required little manual effort and intervention in producing useful results.

References

- [1] Jonathan D. Cohen, Ming C. Lin, Dinesh Manocha, and Madhav K. Ponamgi. I-COLLIDE: An interactive and exact collision detection system for large-scale environments. In *Proceedings of the 1995 symposium on Interactive 3D graphics*, pages 189–196, 218. ACM Press, 1995.
- [2] A. Fod, M. Mataric, and O. Jenkins. Automated derivation of primitives for movement classification. *Autonomous Robots*, 12(1):39–54, January 2002.
- [3] A. J. Ijspeert, J. Nakanishi, and S. Schaal. Trajectory formation for imitation with nonlinear dynamical systems. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS2001)*, pages 752–757, 2001.
- [4] Zoran Popović and Andrew P. Witkin. Physically based motion transformation. In *Proceedings of SIGGRAPH 99*, Computer Graphics Proceedings, Annual Conference Series, pages 11–20, Los Angeles, California, August 1999. ACM SIGGRAPH / Addison Wesley Longman. ISBN 0-20148-560-5.
- [5] Charles Rose, Michael F. Cohen, and Bobby Bodenheimer. Verbs and adverbs: Multi-dimensional motion interpolation. *IEEE Computer Graphics & Applications*, 18(5):32–40, September - October 1998. ISSN 0272-1716.
- [6] S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000.
- [7] Bernhard Scholkopf, Alex J. Smola, and Klaus-Robert Muller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10(5):1299–1319, 1998.
- [8] D. Shepard. A two-dimensional interpolation function for irregularly-spaced data. In *Proceedings of the ACM national conference*, pages 517–524. ACM Press, 1968.
- [9] J. B. Tenenbaum, V. de Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.

