

# Deriving Action and Behavior Primitives from Human Motion Data

Odest Chadwicke Jenkins and Maja J Matarić

*Robotics Research Labs, University of Southern California*

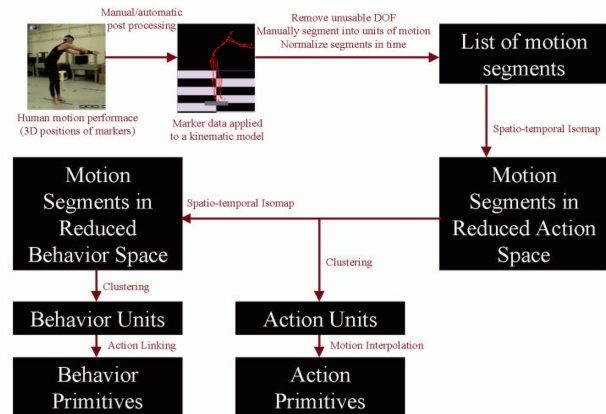
## Abstract

We address the problem of creating basis behaviors for modularizing humanoid robot control and representing human activity. These behaviors, called *perceptual-motor primitives*, serve as a substrate for linking a system's perception of human activities and the ability to perform those activities. We present a data-driven method for deriving perceptual-motor action and behavior primitives from human motion capture data. In order to find these primitives, we employ a spatio-temporal non-linear dimension reduction technique on a set of motion segments. From this transformation, motions representing the same action can be clustered and generalized. Further dimension reduction iterations are applied to derive extended-duration behaviors.

## 1 Introduction

Controlling a humanoid robot and perceiving human activity are both complicated problems. In order for imitation, the process of learning new movement patterns and skills by observation, to be possible, such mechanisms for both perception and control must be present. Although algorithms used for the separate problems differ vastly, both problems at their core assume that some structure exists in the motion being generated or analyzed. If we assume that the underlying motion structure for perception is the same motion structure used for control, we can address both problems in a synergistic fashion.

In our approach, we assume that motion is structured for both perception and control by a vocabulary of parameterized motion modules, called *perceptual-motor primitives*. These primitives generally refer to a set of behaviors that are capable of the recognition and generation of complex movement through sequencing and/or superposition of the behaviors [10]. Through the use of such structures, we can address the problems of robot control and activity recognition in the context of our model for imitation [4] [9]. For perception, a system can encode a teacher's previously unseen demonstration into a set of primi-



**Figure 1:** Flowchart for primitive derivation methodology

tives (activity recognition). From this encoding, the system can synthesize an imitation of the demonstration from the controllers of the classified behaviors (motor control).

While primitive behaviors provide a valuable structuring of perception and control, a major question remains as to what behaviors should be used as the set of primitives. Primitive behaviors could be chosen by hand and crafted manually [11] [9]. However, motion data could contain structure that is unanticipated by human decision making. As stated in [4], determining a basis set of primitives is a difficult problem, but could be addressed by deriving a primitive set in a data-driven fashion through dimension reduction. In this work, we extend a technique for intrinsic (non-linear) dimension reduction [18] to emphasize temporal as well as non-linear spatial structure.

Using dimension reduction (DR), we describe a method for deriving sets of perceptual-motor primitives as atomic-level action primitives and meta-level behavior primitives. We apply spatio-temporal intrinsic DR to manually segmented human motion

capture data. Segments representing instances of same action are proximate in the reduced space and cluster into *action units*. Action units are then generalized into parameterized primitives through motion interpolation [12]. By iteratively applying DR, atomic-level action units form extended duration meta-level *behavior units*. Behavior units can be generalized into behavior primitives by linking action units. The flowchart for our approach is shown in Figure 1.

## 2 Related Work

The work presented in this paper closely follows the framework of [4] and is based on the same model for imitation. In [4], however, Principal Components Analysis (PCA) was used to learn a set principal components (PCs) to serve as perceptual-motor primitives. The main problem with using PCA in such a manner is that linear PCs cannot represent the non-linear nature of human motion. Even though the PCs can span the space of possible motion, the principal component primitives do not capture the underlying structure of the motion. Thus, creating parameterized motor controllers for these primitives becomes non-trivial. In addition, motion segments in [4] are implicitly assumed to have no temporal dependencies, which is not an accurate assumption. These problems are handled in this work through non-linear spatial DR and accounting for temporal dependencies.

Other approaches to imitation have used sensory-motor mapping representations such as dynamical systems [8], programmable pattern generators [14], forward models [3], and paired forward-inverse models [6]. Related work for learning vocabularies of motion, called *movemes*, is presented in [1]. The representation for structuring motion vocabularies presented in this paper is closely related to [12], an approach for manually creating motion vocabularies. [12] use a “Verbs and Adverbs” representation for representing a set of adverb-parameterized verbs, similar to our behavior primitives. Our approach differs in that we aim to learn many of the tasks an animator is asked to perform manually. We use non-linear DR to find the intrinsic structure of the motion data and derive primitives. In particular, we use and extend Isomap [18] to find spatio-temporal structure. Other methods for IDR that we explored with less success include Locally Linear Embedding [13], Kernel PCA [15], and non-linear blind source separation [5].

## 3 Intrinsic Dimensionality Reduction

As stated in [18], traditional DR techniques (e.g., PCA) can only find the actual structure of data on

or near a linear subspace of the input space. The underlying (or intrinsic) structure of the data can be assumed to be on a manifold, a non-linear lower-dimensional subspace embedded in the input space. Using this assumption, [18] define the Isomap algorithm to perform non-linear DR by:

1. building spatially local neighborhoods (using a radius threshold or nearest neighbors)
2. computing a matrix,  $D$ , of all-pairs shortest paths distances
3. performing PCA using  $D$

### 3.1 Accounting for Spatio-Temporal Dependencies

Previous methods for intrinsic DR assume that the provided data is independent identically distributed (i.e., spatial with no temporal ordering). Motion data, however, have temporal dependencies. Thus, the need for a spatio-temporal version of Isomap arises. The extension of Isomap must achieve a good balance between spatial and temporal information (Figure 3.1). Spatio-temporal Isomap is demonstrated for the simple case of the figure-eight in Figures 3.1 and 3.1.

We used two approaches to spatio-temporal intrinsic DR: spatio-temporal neighborhood weighting and phase space weighting. In spatio-temporal neighborhood weighting (Figure 4), weighted spatial neighborhoods are computed as in Isomap, except for the inclusion of temporally adjacent data. These spatial weights are then strengthened or weakened based on common temporal neighbors. Two data points,  $x$  and  $y$ , such that  $y \in nbhd(x)$  (neighborhood of  $x$ ), have common temporal neighbors if their temporally adjacent data points,  $x_1$  and  $y_1$ , share a spatial neighborhood,  $y_1 \in nbhd(x_1)$ . If common temporal neighbors occur, the spatial weighting between the data pair is strengthened by dividing by a constant  $c_t n$ ; else, the spatial weighting is weakened or unchanged. Adjacent temporal neighbors are then strengthened by a constant  $c_a t n$ . PCA is then applied to the weight matrix, as in Isomap. By varying  $c_a t n$ , the amount of flattening from the spatial to the temporal 1D manifold can be controlled. By varying  $c_t n$ , the amount of variation from the temporal manifold can be controlled. In the phase space view of spatio-temporal Isomap, this weighting scheme could be seen as just introducing the velocity of the data into the weight matrix.

In the phase space view of spatio-temporal Isomap, velocity information is explicitly included into the vector of each data point. Regular spatial Isomap is performed on the phase space data with a weighted

distance metric. The weighted distance between two points,  $x$  and  $y$ , is expressed as:

$$d^2 = z^T W z \quad (1)$$

where  $d$  is the distance between the data pair,  $W$  is the weighting matrix, and  $z$  is the difference vector between phase space vectors for  $x$  and  $y$ ,

$$z = \begin{bmatrix} x - y \\ \dot{x} - \dot{y} \end{bmatrix} \quad (2)$$

The weighting matrix can be split into four blocks:

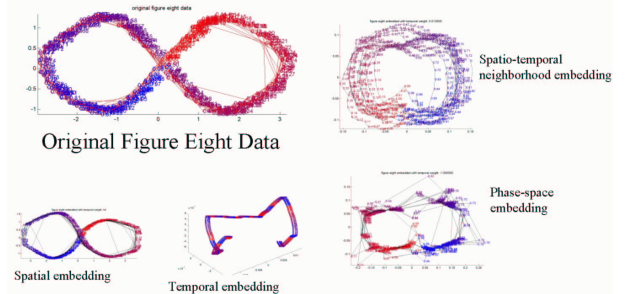
$$W = \begin{bmatrix} W_{11} & W_{21} \\ W_{12} & W_{22} \end{bmatrix} \quad (3)$$

such that  $W_{11}$  represents the weighting of spatial information and  $W_{22}$  weights velocity information, and the blocks  $W_{12}$  and  $W_{21}$  weight spatial and velocity information. Phase space weighting provides a continuous weighting function, in contrast to the common temporal neighbors step function.

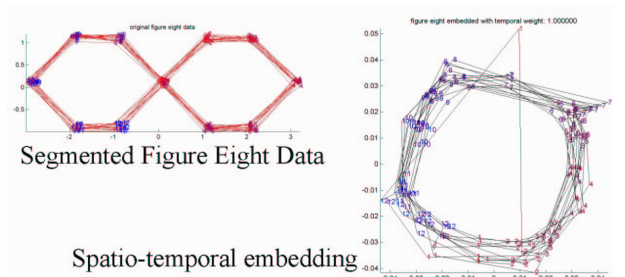
Both viewpoints have strengths and weaknesses. With spatio-temporal neighborhoods, explicit control is provided for how the underlying structure is represented on a continuum between completely spatial and 1D temporal manifolds. Also, the common temporal neighbors step function can produce cleaner, more aesthetic embeddings. However, spatio-temporal neighbors requires much more parameter tuning than the phase space weighting method. The latter is simpler to implement and easier to represent mathematically. Furthermore, a reasonable guess of the elements of the weight matrix  $W$  can provide useful embeddings faster than parameter tuning spatio-temporal neighborhoods.

#### 4 Deriving Action Units from Motion Capture

We applied our method for spatio-temporal intrinsic dimensionality reduction to the set of motion segments described in Section 6. For this work, degrees of freedom (DOFs) were separated into three groups representing the significant kinematic substructures (left arm, right arm, and torso). Spatio-temporal Isomap was performed separately for motion data of each substructure. We initially assume limbs are working independently, even though coordination of limbs is possible. The data for each kinematic substructure were reduced to 5 dimensions. Motion segments produced by the same underlying action will form proximal groups in the resulting embedding. By applying clustering, segments of motion are grouped into units representative of the same action, or action



**Figure 2:** Application of variations of Isomap to figure-eight data: (1) plot of points on a figure-eight with a small amount of noise; (2) embeddings resulting from original (spatial) Isomap and temporal Isomap; (3) embeddings resulting from spatio-temporal Isomap with spatio-temporal neighborhood and phase space weightings.

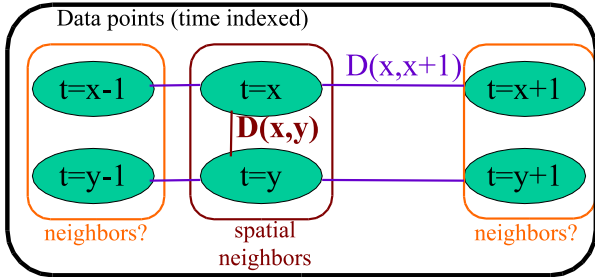


**Figure 3:** Application of spatio-temporal Isomap to pre-segmented figure-eight data: (1) pre-segmented figure-eight data; (2) embedding result, clusters remain most of their integrity.

units. Action units could be computed for each kinematic substructure separately, furthering the limb independence assumption. In this work, however, body-level action units are formed by concatenating the dimensions from each substructure before clustering. A bounding box clustering method was used based on the one-dimensional sweep-and-prune technique described in [2].

##### 4.1 Parameterized Action Primitives

The action units arising from the intrinsic embedding and clustering can be generalized into action primitives using motion interpolation. Each action unit is a group of motions that are instances of the same underlying action being performed. Interpolation between these instances allows for the generation of motions within the span of the instances. Assuming any motion in this span is an instance of the same action, we generalize an action unit into a



**Figure 4:** Illustration of weighting data pairs for spatio-temporal neighborhoods.  $t$  is the temporal index of a data point and  $D(a, b)$  is the weight between two data points  $a$  and  $b$ . If  $D(x, y) \gg D(x, x + 1)$ , the embedding will reflect spatial dependencies. If  $D(x, y) \ll D(x, x + 1)$ , the embedding will reflect the temporal 1D manifold. Balanced weights for  $D(x, y)$  and  $D(x, x + 1)$  accounts for both spatial and temporal dependencies.

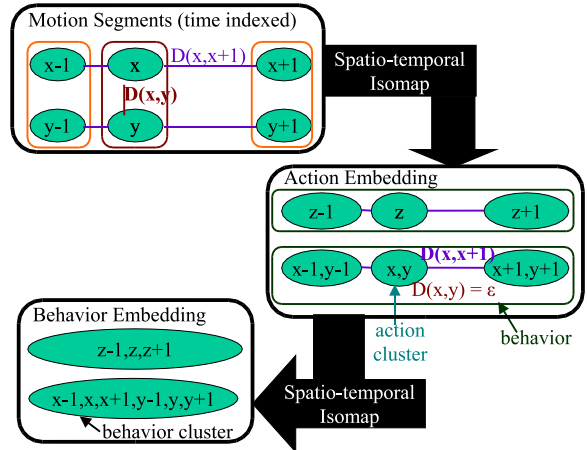
parameterized action primitive by interpolation. We have applied methods based on Shepard’s interpolation [16] and interpolation for sparse examples [12] [17]. The action primitives are considered atomic-level primitives because they represent movement at the level of atomic units of motion.

## 5 Deriving Behavior Units from Actions

By applying further iterations of spatio-temporal Isomap, behavior units representing extended-duration meta-level movement can be derived. Starting from the initial embedding for actions, another iteration of spatio-temporal Isomap is performed (Figure 3.1). As in the case for action units, we expect that segments generated from the same behavior will become proximal in the second embedding. More specifically, segments from actions that form an underlying behavior when sequenced will group together. Clustering is applied in the new embedding to find behavior units. Further embedding iterations can be used to find behaviors at higher levels.

### 5.1 Behavior Primitives

Unlike generalization of action units into primitives, we already have a lower-level primitive representation that needs to connect appropriately to represent the higher-level behavior. In order to generalize a behavior unit into a behavior primitive, the actions associated with each behavior must be determined, linked in the proper sequence, and matched for pre-conditions and post-conditions. Actions associated with behaviors and their sequencing can be found trivially from the action unit indexing into the seg-



**Figure 5:** Flowchart for embedding actions and behaviors. The first embedding groups motion segments from the same action, actions are performed in sequence for a behavior. Further embeddings group motions from the same behavior.

ments. However, linking adjacent actions is more difficult. The motion represented by an action primitive could potentially span a wide area in the space of joint angle trajectories. If one such action must be performed in sequence with another action, then there exist instances of these actions such that a large jump in joint angle space must be taken. To account for this problem, pre-conditions and post-conditions can be computed for each action primitive based on where instances of the action begin and end in joint angle space, respectively. Thus, a mapping between linked actions such that post-conditions for the current action can be used to find valid parameters for the next action in the sequence.

## 6 Experimental Validation

Our system for deriving primitives was validated experimentally through the application to collected motion data. Derived action primitives were used to synthesize desired joint-space trajectories for a humanoid robot. The collected motion and synthesized variations of those motions were actuated on our dynamical humanoid simulation, Adonis.

### 6.1 Motion Data

The motion data used in this work were captured by a Vicon optical motion capture system, which provides the trajectories of the markers in 3D. The marker data were converted into configuration information relative to a kinematic model (joint angles) through a post-process. Joint angle information was

then manually segmented into atomic units of motion based on user criteria. Irrelevant DOFs with respect to kinematic structure of the robot are omitted. The specific motion data used for this work consisted of 22,549 frames, split into 250 atomic segments and included performances of various types of arm waving, dancing, punching, semaphore-spelling, and circle-trajectory motions.

By applying our methodology, we automatically derived distinct behavior primitives for every type of performed motion, except for the semaphore motions. The semaphores, used to spell “animate”, expectedly did not form a coherent behavior because these motions contained a single execution without common temporal neighbors. As example of the results from deriving actions and behaviors, the punching performances from the input motion formed two behavior primitives for “left-arm punching” and “right-arm punching”. The left-arm punching behavior contained 4 component actions corresponding to “Assume punching posture” (hands in front of chest), “Reach back with left arm”, “Punch with left arm”, and “Return left arm to fighting posture”, as shown in Figure 7.

## 6.2 Robot Platform

The humanoid robot platform used to actuate motions produced by our action and behavior modules was Adonis [11], shown in figure 8 performing a slapping punch. Adonis is a rigid-body dynamic simulation of a human torso, with static legs. Adonis consists of eight rigid links connected in a tree structure by a combination of hinge, gimbal, and ball joints. The simulator has a total of 20 DOF and 7 DOF in each arm. Equations of motion are calculated using a commercial solver, SD/Fast [7]. The simulation acts under gravity and accepts other external forces from the environment. Collision detection implemented on Adonis was not necessary for the motion collected or synthesized. Any joint-space controller could be used to actuate our system’s output. We used PD servos for simplicity.

## 7 Continuing Work and Conclusion

We have presented a method for deriving action and behavior primitives from pre-segmented motion data. We extended the Isomap algorithm to find spatio-temporal structure in a set of data. By applying this extended algorithm to motion data, we demonstrated that atomic-level actions and meta-level behaviors can be derived.

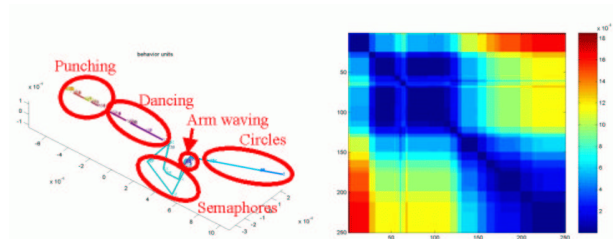
Given that we have derived a set of action and behavior primitives, algorithms for human activity classification, modular motor control, and imitation become implementable within our framework. For activity

classification, new observed movement can be classified into the actions by computing neighborhoods as in spatio-temporal Isomap and using those distances to find the closest action. Over the stream of actions, meta-level behaviors can be classified using pattern matching. For motor control, an action instantiated with a valid set of parameters can produce motion that serves as a set of desired for the PD controller. Behaviors can be instantiated by random walks through or constrained traversal of the actions in the behavior. Imitation can be performed by classifying new observed movement into the action and behavior set and creating motion desired from instantiating actions or behaviors.

As future work, we will be investigating using our primitive derivation methodology with automatic segmentation and without segmentation.

## Acknowledgments

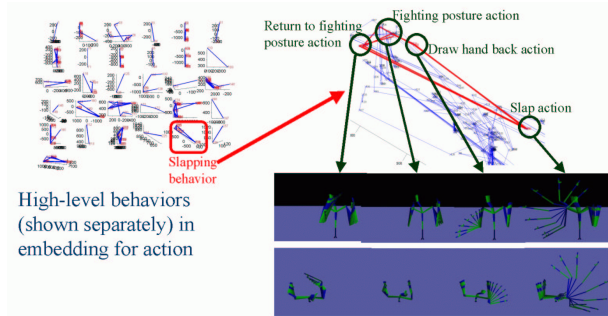
This work is supported in part by the DARPA MARS Program grant DABT63-99-1-0015 and ONR MURI grant N00014-01-1-0890. The authors would also like to thank Jessica Hodgins for providing motion capture data, Evan Drumwright for motion interpolation assistance, and Andrew Howard and Kazunori Okada for various insights.



**Figure 6:** (Left) The fifth embedding of the motion segment data. Behaviors have clustered into knots along a near 1D manifold describing their performance sequence. (Right) The data pair weight matrix during the fifth embedding. Behavior intervals become apparent by dark squares along the diagonal, with dark indicating a strong dependency. The large dark square in the upper left represent the performance of a variety of arm waving instances interrupted by the performance of another behavior.

## References

- [1] C. Bregler. Learning and recognizing human dynamics in video sequences. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 568–574, 1997.
- [2] Jonathan D. Cohen, Ming C. Lin, Dinesh Manocha, and Madhav K. Ponamgi. I-COLLIDE: An interac-



**Figure 7:** Illustration of action and behavior units for motion. (Above) a view of the action embedding, highlighting component actions units of the derived “slapping punch” behavior highlighted. The behavior consists of four sequential actions, as labeled in the image. Below is a sequence of visualization of motions synthesized (using interpolation) for each of the component action primitives. Front and top views are provided for each motion.



**Figure 8:** Snapshots of Adonis performing a slapping punch.

tive 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.

- [3] J. Demiris and G. Hayes. Active and passive routes to imitation. In *Proceedings of the AISB symposium on imitation in animals and artifacts.*, 1999.
- [4] A. Fod, M. Matarić, and O. Jenkins. Automated derivation of primitives for movement classification. *Autonomous Robots*, 12(1):39–54, January 2002.
- [5] S. Harmeling, A. Ziehe, M. Kawanabe, B. Blankertz, and K. Muller. Nonlinear blind source separation using kernel feature spaces. In *Third International Conference on Independent Component Analysis and Signal Separation*, 2001.
- [6] M. Haruno, D. M. Wolpert, and M. Kawato. Multiple forward-inverse models for human motor learning and control. In *Advances in Neural Information Processing Systems 11*. MIT Press, 1999.
- [7] M. Hollars, D. Rosenthal, and M. Sherman. Sd/fast user’s manual. Technical report, Symbolic Dynamics, Inc., 1991.
- [8] 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.
- [9] O. C. Jenkins, M. Matarić, and S. Weber. Primitive-based movement classification for humanoid imitation. In *Proc. of First IEEE-RAS International Conference on Humanoid Robots*, Cambridge, MA, USA, 2000.
- [10] Maja J Matarić. Sensory-motor primitives as a basis for imitation: Linking perception to action and biology to robotics. In C. Nehaniv & K. Dautenhahn, editor, *Imitation in Animals and Artifacts*. The MIT Press, 2001.
- [11] Maja J Matarić, Victor B. Zordan, and Z. Mason. Movement control methods for complex, dynamically simulated agents: Adonis dances the macarena. In *Autonomous Agents*, pages 317–324, 1998.
- [12] Charles Rose, Michael F. Cohen, and Bobby Bodenheimer. Verbs and adverbs: Multidimensional motion interpolation. *IEEE Computer Graphics & Applications*, 18(5):32–40, September - October 1998. ISSN 0272-1716.
- [13] S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000.
- [14] S. Schaal and D. Sternad. Programmable pattern generators. In *Third International Conference on Computational Intelligence in Neuroscience*, pages 48–51, 1998.
- [15] Bernhard Scholkopf, Alex J. Smola, and Klaus-Robert Muller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10(5):1299–1319, 1998.
- [16] D. Shepard. A two-dimensional interpolation function for irregularly-spaced data. In *Proceedings of the ACM national conference*, pages 517–524. ACM Press, 1968.
- [17] P. Sloan, C. F. Rose, and M. F. Cohen. Shape and animation by example. Technical Report MSR-TR-2000-79, Microsoft Research, 2000.
- [18] J. B. Tenenbaum, V. de Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.