

Recognizing Human Pose and Actions for Interactive Robots

Odest Chadwicke Jenkins¹, Germán González Serrano²
and Matthew M. Loper¹

¹*Brown University*, ²*Ecole Polytechnique Fédérale de Lausanne*
¹*USA*, ²*Switzerland*

1. Introduction

Perceiving human motion and non-verbal cues is an important aspect of human-robot interaction (Fong et al., 2002). For robots to become functional collaborators in society, they must be able to make decisions based on their perception of human state. Additionally, knowledge about human state is crucial for robots to learn control policies from direct observation of humans. Human state, however, encompasses a large and diverse set of variables, including kinematic, affective, and goal-oriented information, which has proved difficult to model and infer. Part of this problem is that the relationship between such decision-related variables and a robot's sensor readings is difficult to infer directly.

Our greater view is that socially interactive robots will need to maintain estimates, or beliefs as probabilistic distributions, about all of the components in a human's state in order to make effective decisions during interaction. Humans make decisions to drive their muscles and affect their environment. A robot can only sense limited information about this control process. This information is often partial observations about the human's kinematics and appearance over time, such as images from a robot's camera. To estimate a human's decision making policy, a robot must attempt to invert this partial information back through its model of the human control loop, maintaining beliefs about kinematic movement, actions performed, decision policy, and intentionality.

As a step in this direction, we present a method for inferring a human's kinematic and action state from monocular vision. Our method works in a bottom-up fashion by using a vocabulary of predictive dynamical primitives, learned from previous work (Jenkins & Mataric, 2004a) as "action filters" working in parallel. Motion tracking is performed by matching predicted and observed human movement, using particle filtering (Isard & Blake 1998, Thrun et al., 2005) to maintain nonparametric probabilistic beliefs. For quickly performed motion without temporal coherence, we propose a "bending cone" distribution for extended prediction of human pose over larger intervals of time. State estimates from the action filters are then used to infer the linear coefficients for combining behaviours. Inspired by neuroscience, these composition coefficients are related to the human's cognitively planned motion, or "virtual trajectory", providing a compact action space for linking decision making with observed motion.

Source: Human-Robot Interaction, Book edited by Nilanjan Sarkar,
ISBN 978-3-902613-13-4, pp.522, September 2007, Itech Education and Publishing, Vienna, Austria

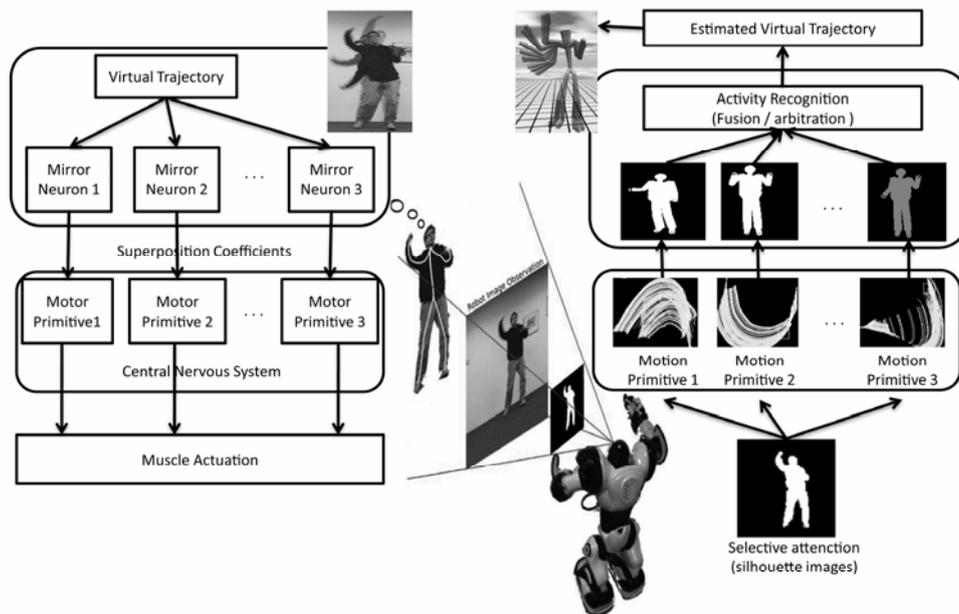


Figure 1. A “toy” example of our approach to human state estimation and movement imitation. The movement of a human demonstrator assumed to be generated by virtual trajectory executed as a weighted superposition of motor primitives, predictive low-dimensional dynamical systems. For movement imitation, a particle filter for each primitive performs kinematic state (or pose) estimation. Pose estimates across the vocabulary are fused at each timestep and concatenated over time to yield an estimate of the virtual trajectory for the robot to execute

We present results from evaluating our motion and action tracking system to human motion observed from a single robot camera. Presented results demonstrate our methods ability to track human motion and action, robust to performer speed and camera viewpoint with recovery from ambiguous situations, such as occlusion. A primary contribution of our work is interactive-time inference of human pose using sparse numbers of particles with dynamical predictions. We evaluate our prediction mechanism with respect to action classification over various numbers of particles and other prediction related variables. Our work has a broad scope of applications, ranging from robot learning to healthcare. Due to the predictive nature of the motion primitives, the methodology presented can be used to modify control policies according to predicted human actions. The combination of motion primitives and human tracking can also be used in healthcare, analyzing the performance of a given activity by a patient, seeing how much it deviates from a “natural” or “standard” performance due to an illness. An example would be gait-analysis, or rehabilitation

programmes. The analysis of human performance could be extended to sports training, analysing how much a sportsman deviates from the canonical performance described by the motion primitive and how much does that affect his performance.

We highlight the application of our tracking results to humanoid imitation. These results allow us to drive a virtual character, which could be used in videogames or computer animation. The player would be tracked, and his kinematics would be adapted to the closest known pose in the motion primitives. This way, we could correct for imperfect player's performance.

2. Background

2.1 Motor Primitives and Imitation Learning

This work is inspired by the hypotheses from neuroscience pertaining to models of motor control and sensory-motor integration. We ground basic concepts for imitation learning, as described in (Mataric, 2002), in specific computational mechanisms for humanoids. Mataric's model of imitation consists of: 1) a selective attention mechanism for extraction of observable features from a sensory stream, 2) mirror neurons that map sensory observations into a motor repertoire, 3) a repertoire of motor primitives as a basis for expressing a broad span of movement, and 4) a classification-based learning system that constructs new motor skills.

Illustrated in Figure 1, the core of this imitation model is the existence and development of computational mechanisms for mirror neurons and motor primitives. As proposed by (Mussa-Ivaldi & Bizzi, 2000), motor primitives are used by the central nervous system to solve the inverse dynamics problem in biological motor control. This theory is based on an equilibrium point hypothesis. The dynamics of the plant $D(x, \dot{x}, \ddot{x})$ is a linear combination of forces from a set of primitives, as configuration-dependent force fields (or attractors) $\phi_i(x, \dot{x}, \ddot{x})$:

$$D(x, \dot{x}, \ddot{x}) = \sum_{i=1}^K c_i \phi_i(x, \dot{x}, \ddot{x}) \quad (1)$$

where x is the kinematic configuration of the plant, c is a vector of scalar superposition coefficients, and K is the number of primitives. A specific set of values for c produces stable movement to a particular equilibrium configuration. A sequence of equilibrium points specifies a virtual trajectory (Hogan, 1985) that can be used for control, as desired motion for internal motor actuation, or perception, to understand the observed movement of an external performer.

Mataric's imitation model assumes the firing of mirror neurons specifies the coefficients for formation of virtual trajectories. Mirror neurons in primates (Rizzolatti et al., 1996) have been demonstrated to fire when a particular activity is executed, observed, or imagined. Assuming 1-1 correspondence between primitives and mirror neurons, the scalar firing rate of a given mirror neuron is the superposition coefficient for its associated primitive during equilibrium point control.

2.2 Motion Modeling

While Matarić's model has desirable properties, there remain several challenges in its computational realization for autonomous robots that we attempt to address. Namely, what are the set of primitives and how are they parameterized? How do mirror neurons recognize motion indicative of a particular primitive? What computational operators should be used to compose primitives to express a broader span of motion?

Our previous work (Jenkins & Matarić 2004a) address these computational issues through the unsupervised learning of motion vocabularies, which we now utilize within probabilistic inference. Our approach is close in spirit to work by (Kojo et al., 2006), who define a "proto-symbol" space describing the space of possible motion. Monocular human tracking is then cast as localizing the appropriate action in the proto-symbol space describing the observed motion using divergence metrics. (Ijspeert et al., 2001) encode each primitive to describe the nonlinear dynamics of a specific trajectory with a discrete or rhythmic pattern generator. New trajectories are formed by learning superposition coefficients through reinforcement learning. While this approach to primitive-based control may be more biologically faithful, our method provides greater motion variability within each primitive and facilitates partially observed movement perception (such as monocular tracking) as well as control applications. Work proposed by (Bentivegna & Atkeson, 2001) and (Grupen et al., 1995; Platt et al., 2004) approach robot control through sequencing and/or superposition of manually crafted behaviors.

Recent efforts by (Knoop et al., 2006) perform monocular kinematic tracking using iterative closest point and the latest Swissranger depth sensing devices, capable of precise depth measurements. We have chosen instead to use the more ubiquitous passive camera devices and also avoid modeling detailed human geometry.

Many other approaches to data-driven motion modeling have been proposed in computer vision, animation, and robotics. The reader is referred to other papers (Jenkins & Matarić, 2004a; Urtasun et al., 2005; Kovar & Gleicher, 2004; Elgammal A. M. and Lee Ch. S. 2004) for broader coverage of these methods.

2.3 Monocular Tracking

We pay particular attention to methods using motion models for kinematic tracking and action recognition in interactive-time. Particle filtering (Isard & Blake, 1998; Thrun et al., 2005) is a well established means for inferring kinematic pose from image observations. Yet, particle filtering often requires additional (often overly expensive) procedures, such as annealing (Deutscher et al., 2000), nonparametric belief propagation (Sigal et al., 2004; Sudderth et al., 2003), Gaussian process latent variable models (Urtasun et al., 2005), POMDP learning (Darrell & Pentland, 1996) or dynamic programming (Ramanan & Forsyth, 2003), to account for the high dimensionality and local extrema of kinematic joint angle space. These methods tradeoff real-time performance for greater inference accuracy. This speed-accuracy contrast is most notably seen in how we use our learned motion primitives (Jenkins & Matarić, 2004a) as compared to Gaussian process methods (Urtasun et al., 2005; Wang et al., 2005). Both approaches use motion capture as probabilistic priors on pose and dynamics. However, our method emphasizes temporally extended prediction to use fewer particles and enable fast inference, whereas Gaussian process models aim for accuracy through optimization. Further, unlike the single-action motion-sparse experiments with Gaussian process models, our work is capable of

inference of multiple actions, where each action has dense collections of motion. Such actions are generalized versions of the original training motion and allow us to track new variations of similar actions.

Similar to (Huber & Kortenkamp, 1998), our method aims for interactive-time inference on actions to enable incorporation into a robot control loop. Unlike (Huber and Kortenkamp, 1998), however, we focus on recognizing active motions, rather than static poses, robust to occlusion by developing fast action prediction procedures that enable online probabilistic inference. We also strive for robustness to motion speed by enabling extended look-ahead motion predictions using a “bending cone” distribution for dynamics. (Yang et al., 2007) define a discrete version with similar dynamics using Hidden Markov Model and vector quantization observations. However, such HMM-based transition dynamics are instantaneous with a limited prediction horizon, whereas our bending cone allows for further look-ahead with a soft probability distribution.

3. Dynamical Kinematic and Action Tracking

Kinematic tracking from silhouettes is performed via the steps in Figure 2, those are: 1) global localization of the human in the image, 2) primitive-based kinematic pose estimation and 3) action recognition. The human localization is kept as a unimodal distribution and estimated using the joint angle configuration derived in the previous time step.

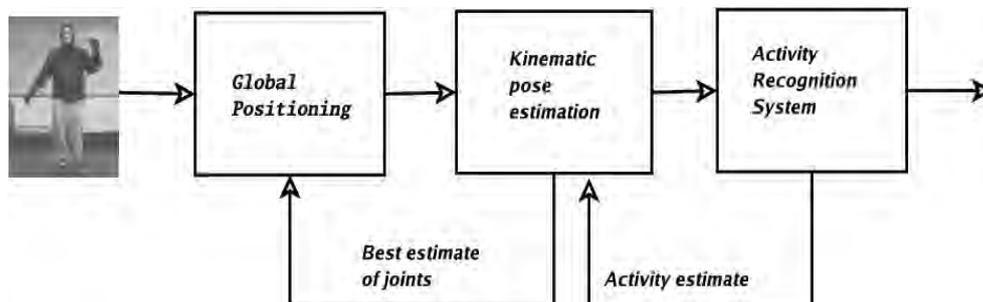


Figure 2. Illustration of the three stages in our approach to tracking: image observations are used to localize the person in 3D, then infer kinematic pose, and finally estimate of activity/action. Estimates at each stage are used to form priors for the previous stage at the next timestep

3.1 Dynamical Motion Vocabularies

The methodology of (Jenkins & Matarić, 2004a) is followed for learning dynamical vocabularies from human motion. We cover relevant details from this work and refer the reader to the citation for details. Motion capture data representative of natural human performance is used as input for the system. The data is partitioned into an ordered set of non-overlapping segments representative of “atomic” movements. Spatio-temporal Isomap (Jenkins & Matarić, 2004b) embed these motion trajectories into a lower dimensional space, establishing a separable clustering of movements into activities. Similar to (Rose et al., 1998),

each cluster is a group of motion examples that can be interpolated to produce new motion representative of the underlying action. Each cluster is speculatively evaluated to produce a dense collection of examples for each uncovered action. A primitive B_i is the manifold formed by the dense collections of poses X_i (and associated gradients) in joint angle space resulting from this interpolation.

We define each primitive B_i as a gradient (potential) field expressing the expected kinematic behaviour over time of the i^{th} action. In the context of dynamical systems, this gradient field $B_i(x)$ defines the predicted direction of displacement for a location in joint angle space $\hat{x}[t]$ at time t^l :

$$\hat{x}_i[t+1] = f_i(x[t], u[t]) = u[t]B_i(x) = u[t] \frac{\sum_{x \in \text{nbhd}(x[t])} w_x \Delta_x}{\left\| \sum_{x \in \text{nbhd}(x[t])} w_x \Delta_x \right\|} \quad (2)$$

where $u[t]$ is a fixed displacement magnitude, Δ_x is the gradient of pose x ², a motion example of primitive i , and w_x the weight³ of x with respect to $x[t]$. Figure 3 shows examples of learned predictive primitives.

Given results in motion latent space dimensionality (Urtasun et al., 2005; Jenkins & Mataric, 2004b), we construct a low dimensional latent space to provide parsimonious observables y_i of the joint angle space for primitive i . This latent space is constructed by applying Principal Components Analysis (PCA) to all of the poses X_i comprising primitive i and form the output equation of the dynamical system, such as in (Howe et al., 2000):

$$y_i[t] = g_i(x[t]) = A_i x[t] \quad (3)$$

Given the preservation of variance in A_i , it is assumed that latent space dynamics, governed by \bar{f}_i , can be computed in the same manner as f_i in joint angle space:

$$\frac{g_i^{-1}(\bar{f}_i(g_i(x[t]), u[t])) - x[t]}{\left\| g_i^{-1}(\bar{f}_i(g_i(x[t]), u[t])) - x[t] \right\|} \approx \frac{f_i(x[t], u[t]) - x[t]}{\left\| f_i(x[t], u[t]) - x[t] \right\|} \quad (4)$$

¹ nbhd() is used to identify the k-nearest neighbours in an arbitrary coordinate space, which we use both in joint angle space and the space of motion segments.

² The gradient is computed as the direction between y and its subsequent pose along its motion example.

³ Typically reciprocated Euclidean distance

3.2 Kinematic Pose Estimation

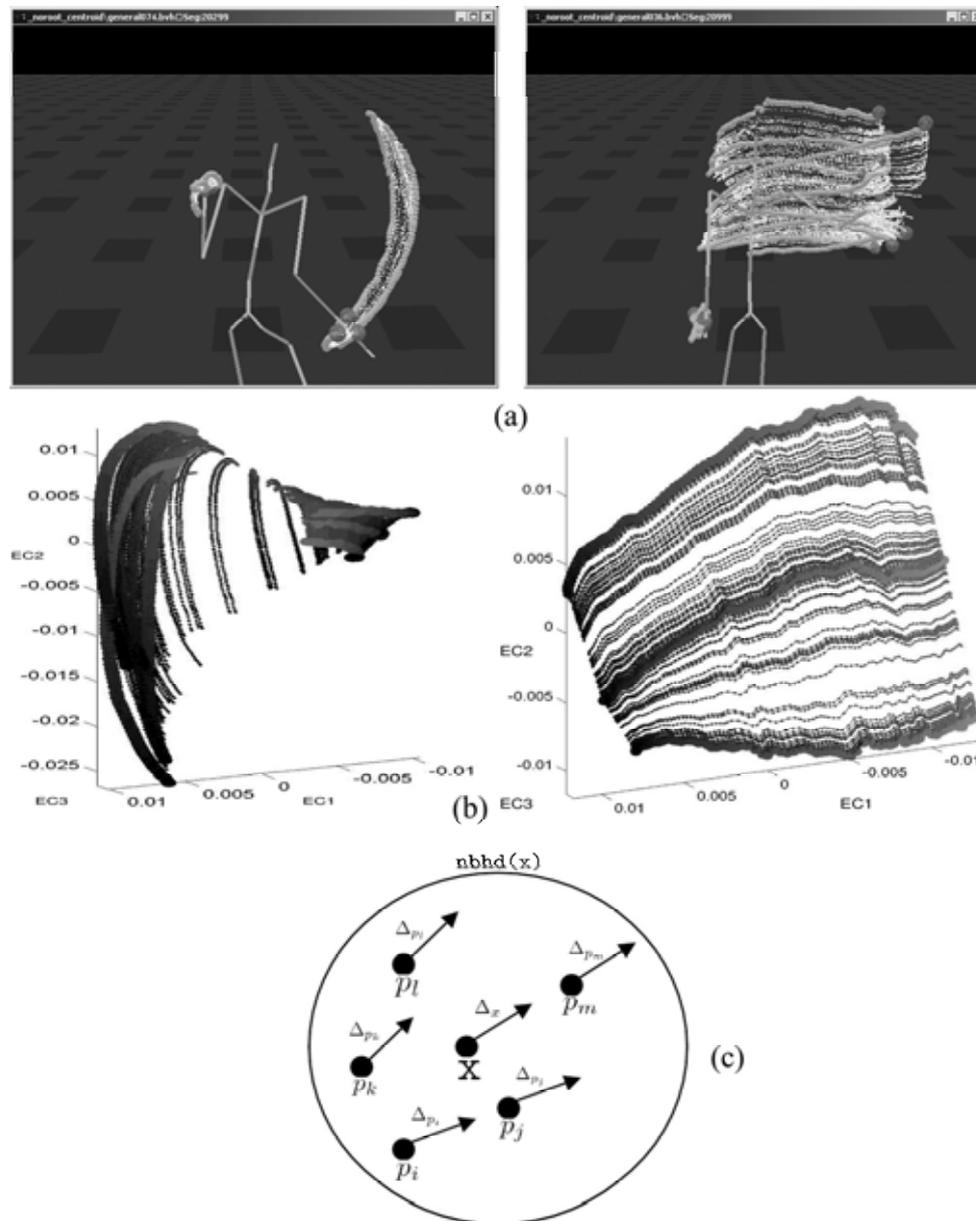


Figure 3. (a) Kinematic endpoint trajectories for learned primitive manifolds, (b) corresponding joint angle space primitive manifolds (view from first three principal components), and (c) an instantaneous prediction example (illustrated as a zoomed-in view on a primitive manifold)

Kinematic tracking is performed by particle filtering (Isard & Blake, 1998; Thrun et al., 2005) in the individual latent spaces created for each primitive in a motion vocabulary. We infer with each primitive individually and in parallel to avoid high-dimensional state spaces, encountered in (Deutscher et al., 2000). A particle filter of the following form is instantiated in the latent space of each primitive

$$p(y_i[1:t] | z_i[1:t]) \propto p(z[t] | g_i^{-1}(y_i[t])) \sum_{y_i} p(y_i[t] | y_i[t-1]) p(y_i[1:t-1] | z[1:t-1]) \quad (5)$$

where $z_i[t]$ are the observed sensory features at time t and g_i^{-1} is the transformation into joint angle space from the latent space of primitive i .

The likelihood function $p(z[t] | g_i^{-1}(y_i[t]))$ can be any reasonable choice for comparing the hypothesized observations from a latent space particle and the sensor observations. Ideally, this function will be monotonic with discrepancy in the joint angle space.

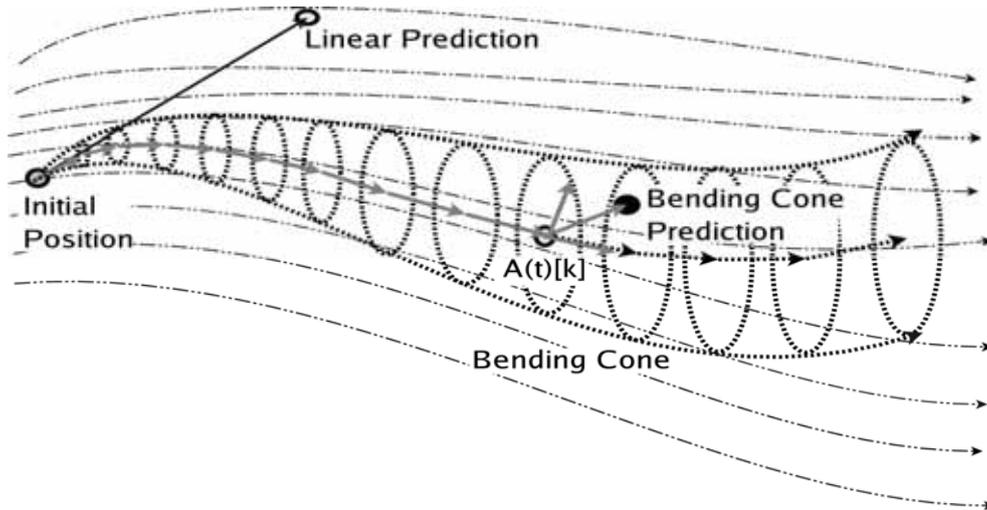


Figure 4. Illustration of the predictive bending cone distribution. The thin dashed black lines indicate the flow of a primitive's gradient field. Linear prediction from the current pose $y_i(t)$ will lead to divergence from the gradient field as the prediction magnitude increases. Instead, we use a bending cone (in bold) to provide an extended prediction horizon along the gradient field. Sampling a pose prediction $y_i(t+1)$ occurs by selecting a cross-section $A(t)[k]$ and adding cylindrical noise

At first glance, the motion distribution $p(z[t] | g_i^{-1}(y_i[t]))$ could be given by the instantaneous "flow", as proposed by (Ong et al., 2006), where a locally linear displacement with some noise is expected. However, such an assumption would require temporal coherence between the training set and the performance of the actor. Observations without temporal coherence cannot simply be accounted for by extending the magnitude of the displacement

vector because the expected motion will likely vary in a nonlinear fashion over time. To address this issue, a “bending cone” distribution is used (Figure 4) over the motion model. This distribution is formed with the structure of a generalized cylinder with a curved axis along the motion manifold and a variance cross-section that expands over time. The axis is derived from K successive predictions $\bar{y}_i[t]$ of the primitive from a current hypothesis $\mathbf{y}[t]$ as a piecewise linear curve. The cross-section is modelled as cylindrical noise $C(\mathbf{a}, \mathbf{b}, \boldsymbol{\sigma})$ with local axis $\mathbf{a}-\mathbf{b}$ and normally distributed variance $\boldsymbol{\sigma}$ orthogonal to the axis. The resulting parametric distribution, equation 6, is sampled by randomly selecting a step-ahead k and generating a random sample within its cylinder cross-section. Note that $f(k)$ is some monotonically increasing function of the distance from the cone origin; we used a linear function.

$$p(y_i[t] | y_i[t-1]) = \sum_{\bar{y}_i[t]}^k C(y_i[k+1], y_i[k], f(k)) \quad (6)$$

3.3 Action Recognition

For action recognition, a probability distribution across primitives of the vocabulary is created⁴. The likelihood of the pose estimate from each primitive is normalized into a probability distribution:

$$p(B_i[t] | z[t]) = \frac{p(z[t] | \bar{x}_i[t])}{\sum_B p(z[t] | \bar{x}_i[t])} \quad (7)$$

where $\bar{x}_i[t]$ is the pose estimate for primitive i . The primitive with the maximum probability is estimated as the action currently being performed. Temporal information can be used to improve this recognition mechanism by fully leveraging the latent space dynamics over time.

The manifold in latent space is essentially an attractor along a family of trajectories towards an equilibrium region. We consider *attractor progress* as a value that increases as kinematic state progresses towards a primitive's equilibrium. For an action being performed, we expect its attractor progress will monotonically increase as the action is executed. The attractor progress can be used as a feedback signal into the particle filters estimating pose for a primitive i in a form such as:

$$p(B_i[t] | z[t]) = \frac{p(z[t] | \bar{x}_i[t], w_i[1:t-1])}{\sum_B p(z[t] | \bar{x}_i[t], w_i[1:t-1])} \quad (8)$$

where $w_i[1:t-1]$ is the probability that primitive B_i has been performed over time.

⁴ We assume each primitive describes an action of interest.



Figure 5. Robot platform and camera used in our experiments

4. Results

For our experiments, we developed an interactive-time software system in C++ that tracks human motion and action from monocular silhouettes using a vocabulary of learned motion primitives. Shown in Figure 5, our system takes video input from a Fire-i webcam (15 frames per second, at a resolution of 120x160) mounted on an iRobot Roomba Discovery. Image silhouettes were computed with standard background modelling techniques for pixel statistics on colour images. Median and morphological filtering were used to remove noisy silhouette pixels. An implementation of spatio-temporal Isomap (Jenkins & Matarić, 2004b) was used to learn motion primitives for performing punching, hand circles, vertical hand waving, and horizontal hand waving.

We utilize a basic likelihood function, $p(z[t] | g^{-1}(y_i[t]))$, that returns the similarity $R(A,B)$ of a particle's hypothesized silhouette with the observed silhouette image. Silhouette hypotheses were rendered from a cylindrical 3D body model to an binary image buffer using OpenGL. A similarity metric, $R(A,B)$ for two silhouettes A and B , closely related to the inverse of the Generalized Hausdorff distance was used:

$$R(A,B) = \frac{1}{r(A,B) + r(B,A) + \epsilon} \quad (9)$$

$$r(A,B) = \sum_{a \in A} \left(\min_{b \in B} \|a - b\| \right)^2 \quad (10)$$

This measure is an intermediate between undirected and generalized Hausdorff distance. ϵ is used only to avoid divide-by-zero errors. An example Hausdorff map for a human

silhouette is shown in Figure 6. Due to this silhouetting procedure, the robot must be stationary (i.e., driven to a specific location) during the construction of the background model and tracking process. As we are exploring in future work, this limitation could be relaxed through the use of other sensor modalities, such as stereo vision or time-of-flight ranging cameras.

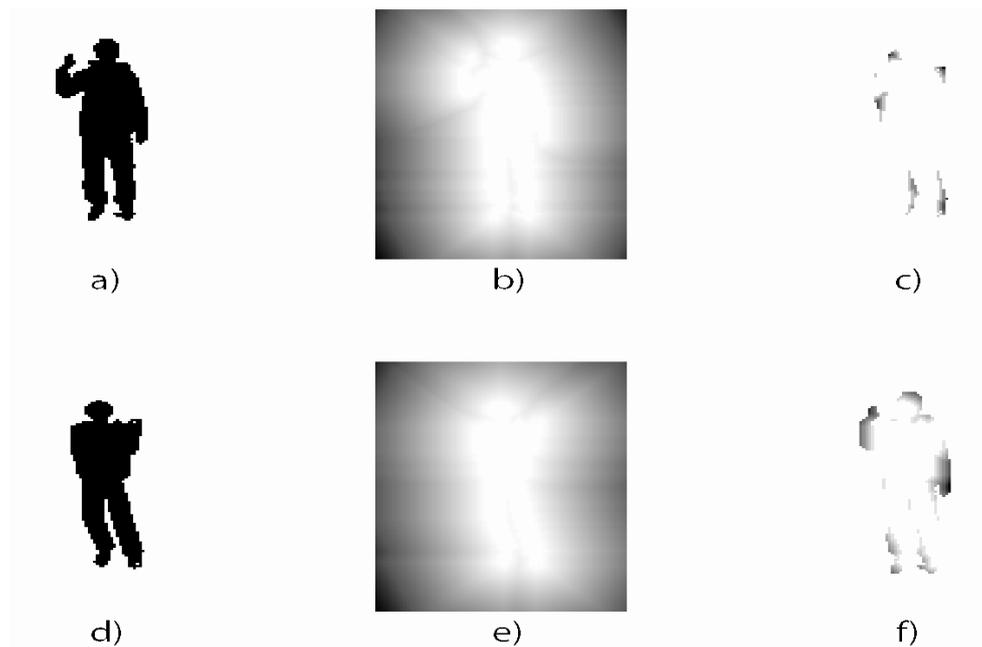


Figure 6. Likelihood function used in the system. (a) is the silhouette A extracted from the camera and (d) is the synthetic silhouette B generated from a pose hypothesis. (b) and (e) are the respective Hausdorff distance transforms, showing pixels with larger distances from the silhouette as dark. (c) and (f) illustrate the sums $r(A,B)$, how silhouette A relates B , and $r(B,A)$, silhouette B relates to A . These sums are added and reciprocated to assess the similarity of A and B

To enable fast monocular tracking, we applied our system with sparse distributions (6 particles per primitive) to three trial silhouette sequences. Each trial is designed to provide insight into different aspects of the performance of our tracking system.

In the first trial (termed multi-action), the actor performs multiple repetitions of three actions (hand circles, vertical hand waving, and horizontal hand waving) in sequence. As shown in Figures 7, reasonable tracking estimates can be generated from as few as six particles. As expected, we observed that the Euclidean distance between our estimates and the ground truth decreases with the number of particles used in the simulation, highlighting the tradeoffs between the number of particles and accuracy of the estimation.

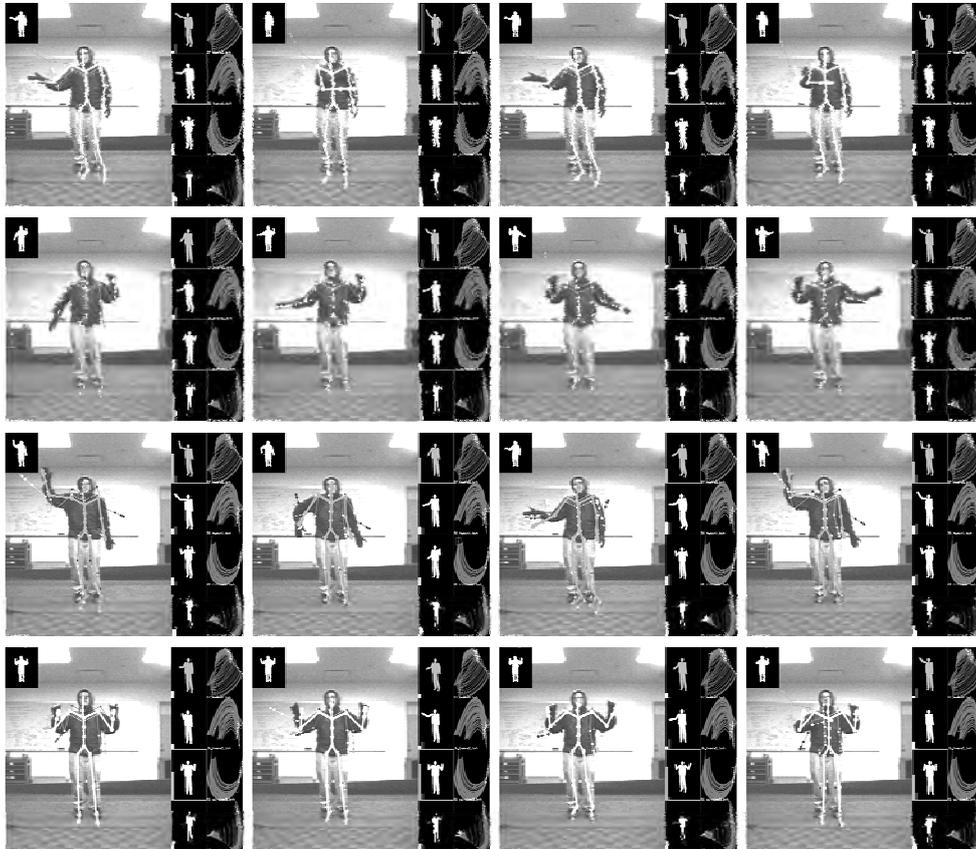


Figure 7. Tracking of a motion sequence containing three distinct actions performed in sequence without stopping. Each row shows the recognition of individual actions for waving a hand across the body (top row), bottom-to-top in a circular fashion (second and fourth row) and top-to-bottom (third row). The kinematic estimates are shown with a thick-lined stick figure; the color of the stick figures represents the action recognized. Each image contains a visualization of the dynamical systems and pose estimates for each action

To explore the effects of the number of particles, we ran our tracking system on the multi-action trial using powers-of-two number particles between 1 and 1024 for each action. The bending cone for these trials are generated using 20 predictions into the future and the noise aperture is $\pi/6$, which increases in steps of $20\pi/6$ per prediction. Shown in Figure 8, the action classification results from the system for each trial were plotted in the ROC plane for each action. The ROC plane plots each trial (shown as a labelled dot) in 2D coordinates where the horizontal axis is the “false positive rate”, percentage of frames incorrectly labelled as a given action, and the vertical axis is the “true positive rate”, percentage of correctly labelled frames. Visually, plots in the upper left-hand corner are indicative of good performance, points in the lower right-hand corner indicate bad performance, and points along the diagonal indicate random performance.

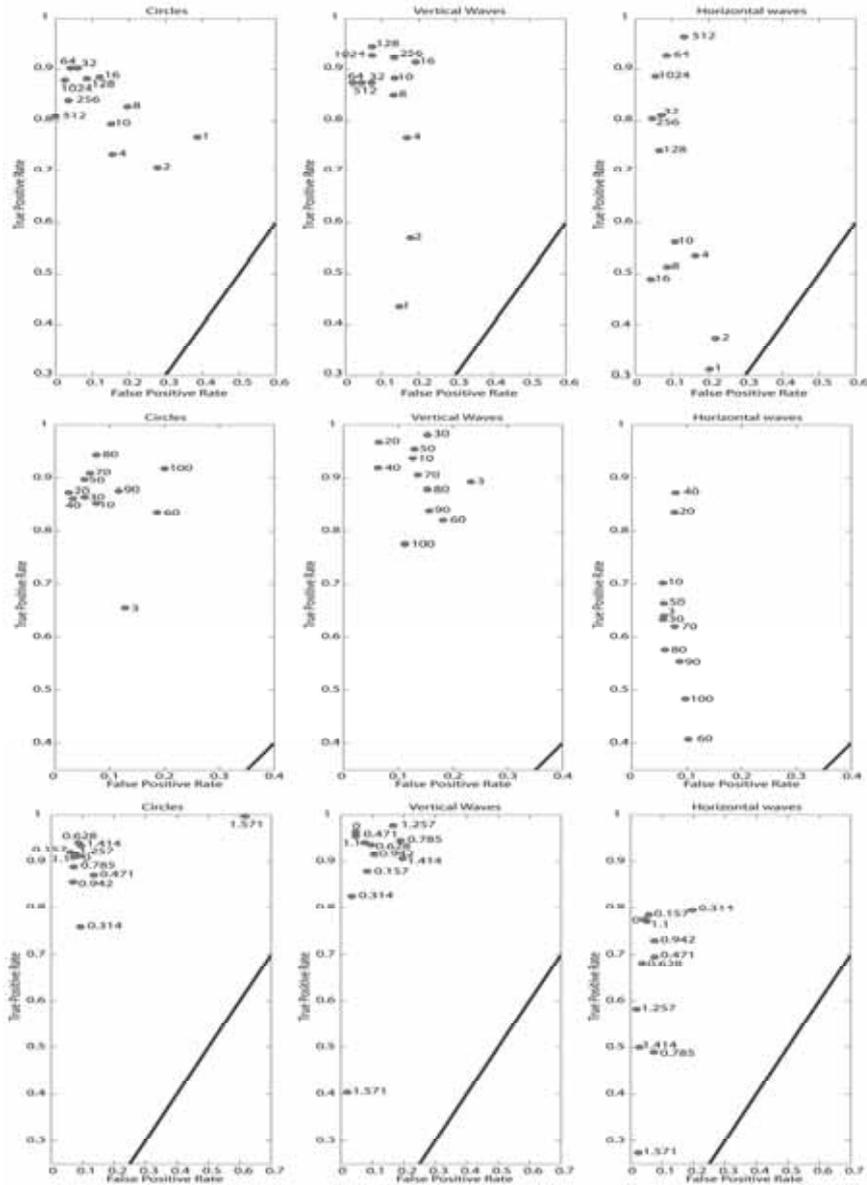


Figure 8. ROC plots of action classification on the multi-action sequence. Columns breakdown by the action performed: Circle action (left), Vertical Waving action (center), and Horizontal Waving action (right). Columns show the effect of varied numbers of particles (top, varied between 1 and 1024), bending cone prediction length (middle, varied between 3 and 100), and bending cone noise aperture (bottom, varied between 0 and $\Pi/2$). Each plot shows for each trial (plotted as a labeled point) the false positive rate (horizontal axis) and true positive rate (vertical axis). Note the difference in scales among rows

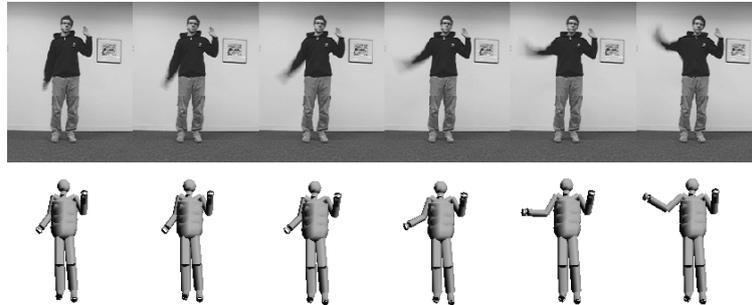


Figure 8. Tracking of a fast waving motion. Observed images (top) and pose estimates from the camera view (bottom)

The ROC plots for the numbers of particles indicate the classifier works better with more particles, but not always. Trials using 1024 particles per action are never the closest point to the upper-left corner. The most noticeable benefit to having more particles is when occlusion was present. In particular, the Circle action does not introduce occlusion when performed in profile, whereas the Horizontal Waving motion does require occlusion in its performance. Consequently, the Circle trials all plot near the ROC upper-left corner and the Horizontal Waving ROC trials fork bi-modally in classification performance.

ROC plots were also generated for varying bending cone prediction length and noise aperture, also shown in Figure 8. Plotted in the middle row, the variations in bending cone length were varied between 10 and 100 predictions in increments of 10, with an additional trial using 3 predictions. For these trials, the number of particles was fixed to 64 and the noise aperture to $\Pi/6$. Plotted in the bottom row, the variation in the bending cone noise aperture were varied between 0 and $\Pi/2$ in increments of $0.1 \cdot \Pi/2$. The number of particles and bending cone length were fixed at 64 and 20, respectively. These plots indicate variations similar to those in the numbers of particles plots. Specifically, good performance results regardless of the bending cone parameters when the action has little or no occlusion ambiguity, but performance drops off when such ambiguity is present. However, in these trials, increased prediction length and noise aperture does not necessarily improve performance. It is surprising that with 0 aperture (that is, staying fixed in the manifold), the classifier does not perform that bad. Instead, there are sweet spots between 20-40 predictions and under 0.15Π noise aperture for the bending cone parameters. Although including more particles is always increase accuracy, we have not yet explored how the sweet spots in the bending cone parameters could change as the numbers of particles vary.

In trial two (fast-wave motion), we analyzed the temporal robustness of the tracking system. The same action is performed at different speeds, ranging from slow (hand moving at ~ 3 cm/s) to fast motion (hand moving at ~ 6 m/s). The fast motion is accurately predicted as seen in Figure 9. Additionally, we were able to track a fast moving punching motion (Figure 10) and successfully execute the motion with our physics-based humanoid simulation. Our simulation system is described in (Wrotek et al., 2006).

In trial three (overhead-view), viewpoint invariance was tested with video from a trial with an overhead camera, shown in Figure 11. Even given limited cues from the silhouette, we are able to infer the horizontal waving of an arm. Notice that the arm estimates are consistent throughout the sequence.

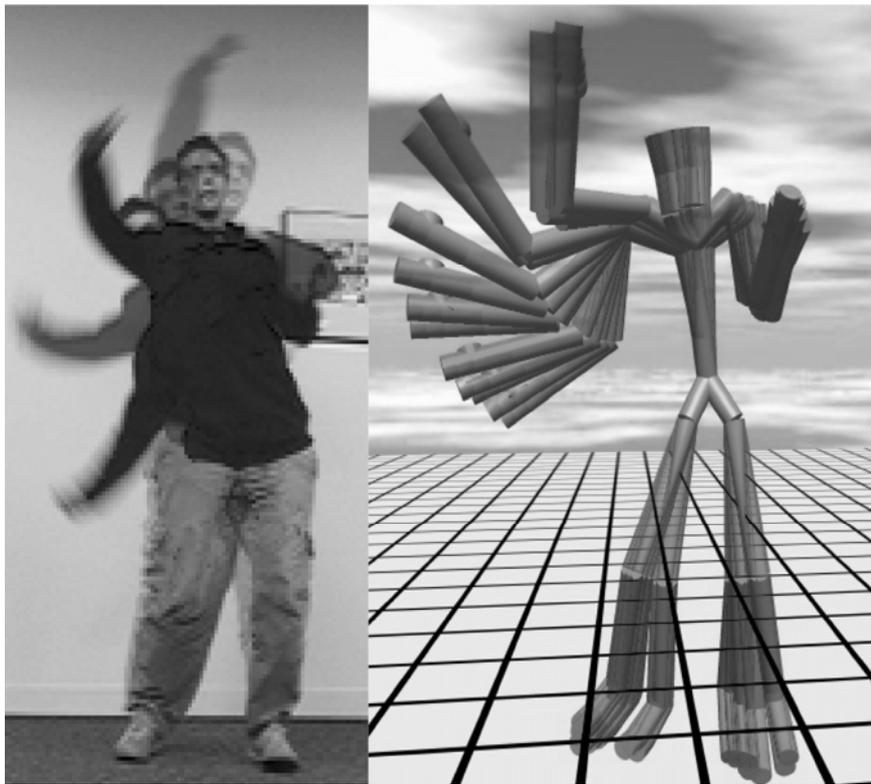


Figure 9. Illustrations of a demonstrated fast moving "punch" movement (left) and the estimated virtual trajectory (right) as traversed by our physically simulated humanoid simulation

Using the above test trials, we measured the ability of our system to recognize performed actions to provide responses similar to mirror neurons. In our current system, an action is recognized as the pose estimate likelihoods normalized over all of the primitives into a probability distribution, as shown in Figure 12.

Temporal information can be used to improve this recognition mechanism by fully leveraging the latent space dynamics over time. The manifold in latent space is essentially an attractor along a family of trajectories. A better estimator of action would consider *attractor progress*, monotonic progress towards the equilibrium region of an action's gradient

field. We have analyzed preliminary results from observing attractor progress in our trials, as shown in Figure 12. For an action being performed, its attractor progress is monotonically increasing. If the action is performed repeatedly, we can see a periodic signal emerge, as opposed to the noisier signals of the action not being performed. These results indicate that we can use attractor progress as a feedback signal to further improve an individual primitive's tracking performance

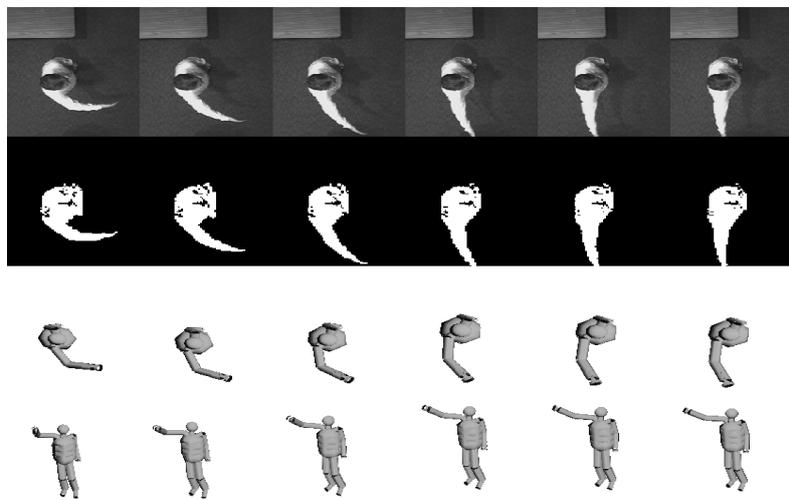


Figure 10. Illustrations of a demonstrated fast moving "punch" movement (left) and the estimated virtual trajectory (right) as traversed by our physically simulated humanoid simulation

Because of their attractor progress properties, we believe that we can analogize these action patterns into the firing of idealized mirror neurons. The firings of our artificial mirror neurons provide superposition coefficients, as in (Nicolescu et al., 2006). Given real-time pose estimation, online movement imitation could be performed by directly executing the robot's motor primitives weighted by these coefficients. Additionally, these superposition coefficients could serve as input into additional inference systems to estimate the human's emotional state for providing an affective robot response.

In our current system, we use the action firing to arbitrate between pose estimates for forming a virtual trajectory. While this is a simplification of the overall goal, our positive results for trajectory estimation demonstrate our approach is viable and has promise for achieving our greater objectives. As future work, we will extend the motion dynamics of the vocabulary into basis behaviours using our complementary work in learning behaviour fusion (Nicolescu et al., 2006).

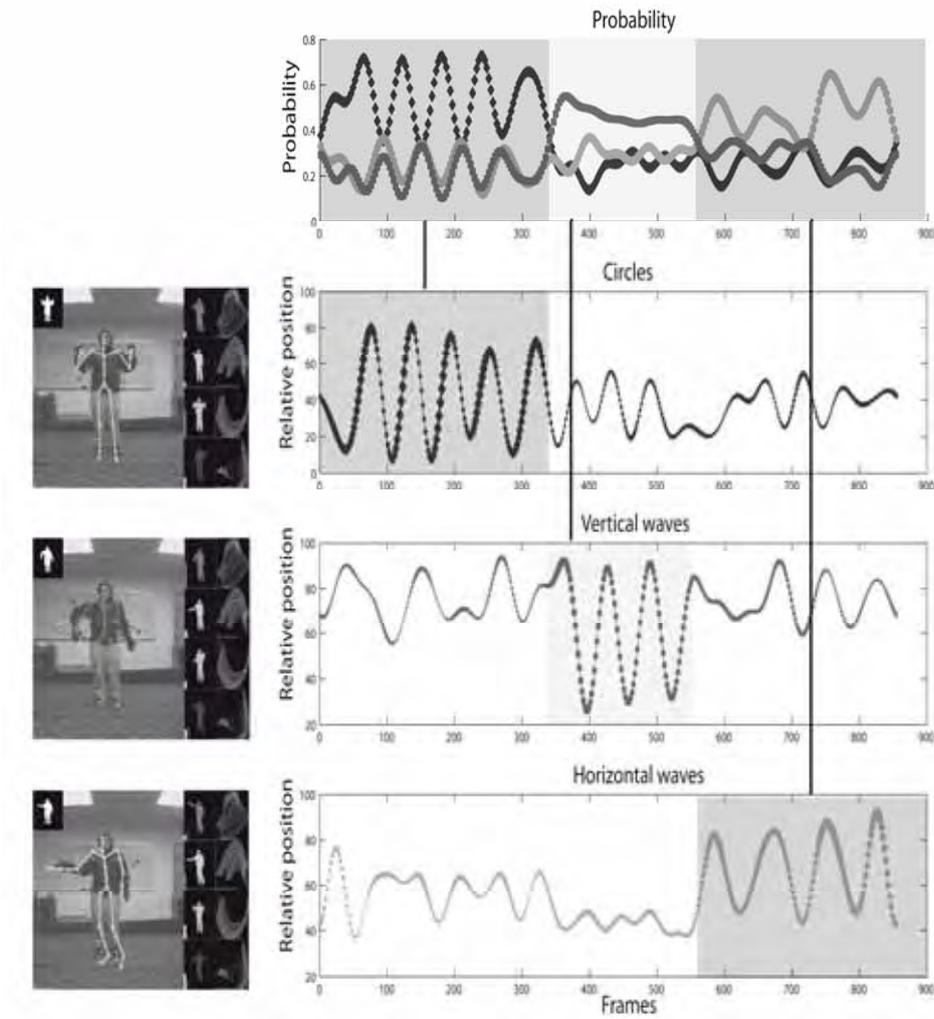


Figure 11. An evaluation of our action recognition system over time with a 3-action motion performing "hand circles", "horizontal waving", and "vertical waving" in sequence. The first row reflects the relative likelihood (idealized as mirror neuron firing) for each primitive with background sections indicating the boundary of each action. Each of the subsequent rows shows time on the x-axis, attractor progress on the y-axis, and the width of the plot marker indicates the likelihood of the pose estimate

5. Conclusion

We have presented a neuro-inspired method for monocular tracking and action recognition for movement imitation. Our approach combines vocabularies of kinematic motion learned offline with online estimation of a demonstrator's underlying virtual trajectory. A modular approach to pose estimation is taken for computational tractability and emulation of structures hypothesized in neuroscience. Our current results suggest our method can perform tracking and recognition from partial observations at interactive rates. Our current system demonstrates robustness with respect to the viewpoint of the camera, the speed of performance of the action, and recovery from ambiguous situations.

6. References

- Bentivegna, D. C. and Atkeson, C. G. (2001). Learning from observation using primitives. *In IEEE International Conference on Robotics and Automation*, pp 1988–1993, Seoul, Korea, May 2001, IEEE.
- Darrell, T. and Pentland, A. (1996). Active gesture recognition using learned visual attention. *Advances in Neural Information Processing Systems*, 8, pp 858–864, 0-262-20107-0, Denver, CO, USA, November 1995, The MIT Press.
- Deutscher J.; Blake, A. and Reid, I. (2000). Articulated body motion capture by annealed particle filtering. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2, pp 126–133, 0-7695-0662-3, Hilton Head, SC, USA, June 2000, IEEE Computer Society.
- Elgammal A. M. and Lee Ch. S. (2004). Inferring 3D body pose from silhouettes using activity manifold learning. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2, pp 681–688, 0-7695-0662-3, Washington DC, USA, July 2004, IEEE Computer Society.
- Fong, T.; Nourbaksh, I. and Daoutenhahn, K. (2002). A survey of socially interactive robots: Concepts, design and applications. Carnegie Mellon University Robotics Institute, Pittsburgh, PA Tech. Rep CMU-RI-TR02-29, November 2002.
- Gruppen, R. A.; Huber, M.; Coehlo Jr. J. A.; and Souccar, K. (1995). A basis for distributed control of manipulation tasks. *IEEE Expert*, 10, 2, (April 1995), pp. 9–14.
- Hogan, N. (1985) The mechanics of multi-joint posture and movement control. *Biological Cybernetics*, 52, (September 1985), pp. 315–331, 0340-1200.
- Howe, N. R.; Leventon, M. E. and Freeman, W. T. (2000). Bayesian reconstruction of 3D human motion from single-camera video. *Advances In Neural Information Processing Systems*, 12, Denver, CO, USA, 2000, The MIT Press.
- Huber, E. and Kortenkamp, D. (1998). A behavior-based approach to active stereo vision for mobile robots. *Engineering Applications of Artificial Intelligence*, 11, (December 1998), pp. 229–243, 0952-1976.
- Ijspeert, A. J.; Nakanishi, J. and Schaal, S. (2001). Trajectory formation for imitation with non-linear dynamical systems. *In IEEE Intelligent Robots and Systems*, pp 752–757, Maui, Hawaii, USA, October 2001, IEEE.
- Isard, M. and Blake, A. (1998). Condensation - conditional density propagation for visual tracking, *International Journal of Computer Vision*, 29, 1, (August 1998) , pp. 5-28, 0920-5691.

- Jenkins, O. C. and Matarić, M. J. (2004a). Performance-derived behavior vocabularies: Data-driven acquisition of skills from motion. *International Journal of Humanoid Robotics*, 1, 2, (June 2004) 237-288, 0219-8436.
- Jenkins, O. C. and Matarić, M. J. (2004b). A spatio-temporal extension to isomap non-linear dimension reduction, *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 441-448, Banff, Alberta, Canada, July 2004, Omnipress, Madison, WI, USA.
- Knoop, S.; Vacek, S. and Dillmann, R. (2006). Sensor fusion for 3D human body tracking with an articulated 3D body model. In *IEEE International Conference on Robotics and Automation*, pp 1686-1691, Orlando, FL, USA, May 2006, IEEE.
- Kojo, N.; Inamura, T.; Okada, K. and Inaba, M.(2006). Gesture recognition for humanoids using proto-symbol space. *Proceedings of the IEEE International Conference on Humanoid Robotics*. pp 76-81, Genova, Italy, December 2006, IEEE,
- Kovar, L. and Gleicher, M.(2004). Automated extraction and parameterization of motions in large data sets. *International Conference on Computer Graphics and Interactive Techniques, ACM Siggraph 2004*, pp 559-568, 0730-0301, Los Angeles, California, USA, 2004.
- Matarić, M. J. (2002). Sensory-motor primitives as a basis for imitation: Linking perception to action and biology to robotics. *Imitation in Animals and Artifacts*. MIT Press, 0-262-04203-7, Cambridge, Massachusetts, USA.
- Mussa-Ivaldi, F. and Bizzi, E.(2000). Motor learning through the combination of primitives. *Philosophical Transactions of the Royal Society: B: Biological Sciences*. 355, pp. 1755-1769 London, UK.
- Nicolescu, M.; Jenkins, O. C., and Olenderski A.(2006). Learning behavior fusion estimation from demonstration. In *IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN 2006)*, pp. 340-345, Hatfield, United Kingdom, September 2006, IEEE Computer Society.
- Ong, E.; Hilton, A. and Micilotta, A. (2006). Viewpoint invariant exemplar-based 3D human tracking. *Computer Vision and Image Understanding*, 104, 2, (November 2006), pp 178-189, ISSN:1077-3142.
- Platt R.; Fagg, A. H. and Grupen, R. R. (2004) Manipulation gaits: Sequences of grasp control tasks. In *IEEE Conference on Robotics and Automation*, pp 801-806, New Orleans, LA, USA, April 2004, IEEE.
- Sigal, L.; Bhatia, S.; Roth, S.; Black, M. J. and Isard, M. (2004). Tracking loose-limbed people. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp 421-428, 0-7695-2158-4, Washington, USA, July 2004., IEEE Computer Society .
- Sudderth, E. B.; Ihler, A. T.; Freeman, W. T. and Willsky, A. S. (2003). Nonparametric belief propagation. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp 605-612, 0-7695-1900-8, Madison, WI, USA, June 2003, IEEE Computer Society.
- Ramanan, D. and Forsyth, D. A. (2003). Automatic annotation of everyday movements. *Advances in Neural Information Processing Systems*, 16, 0-262-20152-6 , Vancouver, Canada, 2003, The MIT Press.
- Rizzolatti, G.; Fadiga, L.; Gallese, V. and Fogassi, L. (1996). Premotor cortex and the recognition of motor actions. *Cognitive Brain Research*, 3, 2, (March 1996), pp 131-141, 0006-8993.

- Rose, C.; Cohen, M. F. and Bodenheimer, B. (1998). Verbs and adverbs: Multidimensional motion interpolation. *IEEE Computer Graphics & Applications*, 18, 5, (September-October 1998), pp. 32-40, 0272-1716.
- Thrun, S.; Burgard, W. and Fox, D. (2005). *Probabilistic Robotics*. MIT Press, 0-262-20162-3, Cambridge, Massachusetts, USA.
- Urtasun, R.; Fleet, D. J.; Hertzmann, A. and Fua, P. (2005). Priors for people tracking from small training sets. In *International Conference in Computer Vision*, pp 403-410, 0-7695-2334-X, Beijing, China, October 2005, IEEE Computer Society.
- Wang, J.; Fleet, D. J. and Hertzmann, A. (2005). Gaussian Process Dynamical Models. *Advances in Neural Information Processing Systems*, 18, Vancouver, Canada, December 2005, The MIT Press.
- Wrotek, P.; Jenkins, O. C. and McGuire, M. (2006). Dynamo: Dynamic data-driven character control with adjustable balance. *Proceedings of the Sandbox Symposium on Video Games*, Boston, MA, USA, July 2006, ACM.
- Yang, H. D.; Park, A. Y. and Lee, S. W. (2007). Gesture Spotting and Recognition for Human-Robot Interaction. *IEEE transactions on Robotics*, 23, (April 2007), pp 256-270, 1552-3098.

7. Acknowledgments

This research was supported by grants from the Office of Naval Research (Award N000140710141) and National Science Foundation (Award IIS-0534858). The authors are grateful to Chris Jones and iRobot Corporation for support with the Roomba platform, Brian Gerkey the Player/Stage Project, RoboDynamics Corporation for Roomba interfaces, and Prof. Rod Beresford.