

***Stefan Schaal***

*Computer Science & Neuroscience  
University of Southern California, Los Angeles  
&  
ATR Computational Neuroscience Laboratory  
Kyoto, Japan*

sschaal@usc.edu  
<http://www-clmc.usc.edu>

## Scaling Reinforcement Learning to Complex Motor Systems

A Lecture Series at the Okinawa Computational Neuroscience Course, June 2005

1



## Joint Work With:

- *Auke Ijspeert*
- *Aaron D'Souza*
- *Jun Nakanishi*
- *Jan Peters*
- *Sethu Vijayakumar*
- *Dimitris Pongas*



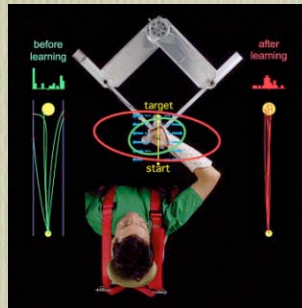
2





# How are Motor Skills Generated?

A Question Shared by Biological and Robotics Research

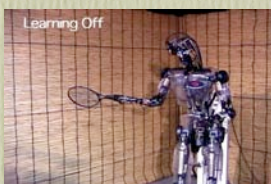


3



# How are Motor Skills Generated?

A Question Shared by Biological and Robotics Research



Movies from collaborations with C. Akteson, S. Kotosaka, S. Vijayakumar

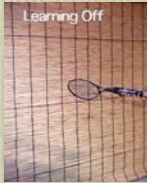
4





# How are Motor Skills Generated?

A Question Shared by Biological and Robotics Research



Unfortunately, each of these skills required manual generation of representations control policies, and learning mechanisms



Movies from collaborations with C. Akteson, S. Kotosaka, S. Vijayakumar

5

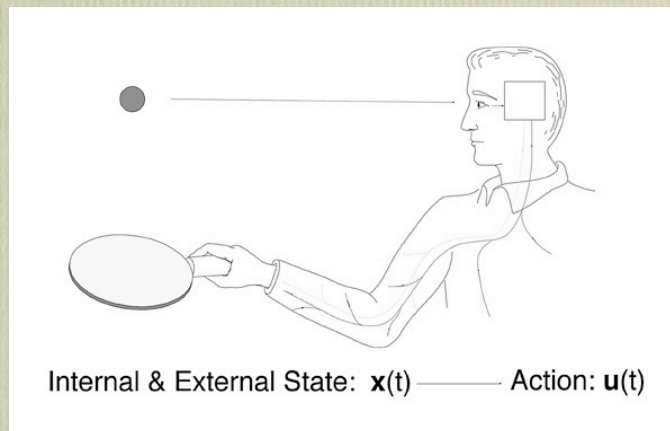


# A Computational Approach to Motor Control: Control Policies

- The General Goal of Motor Learning:

Control Policies

$$\mathbf{u}(t) = \mathbf{p}(\mathbf{x}(t), t, \alpha)$$

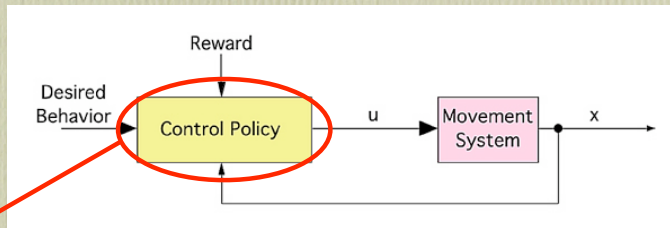


6



# Where Do Control Policies Come From?

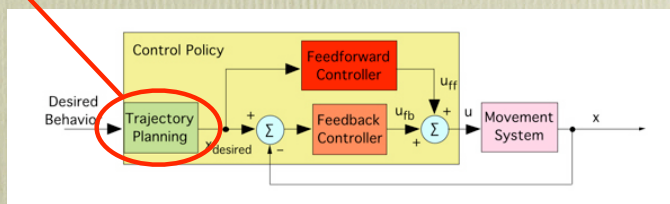
## Direct Control (Model Free)



Possible learning methods include:

- Reinforcement Learning
- Dynamic Programming
- Adaptive Control
- Supervised Learning
- Evolutionary Methods

## Indirect Control (Model-Based)



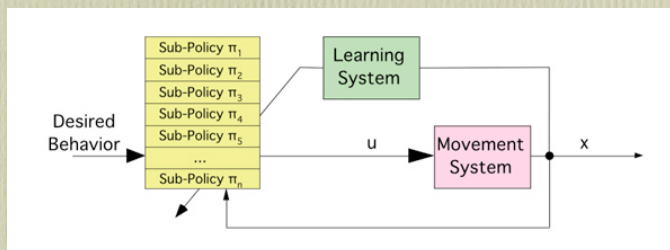
Unstructured Learning of Policies is infeasible

7



# Making Motor Learning Feasible

## • Modular Motor Control



## • Parameterization

$$\mathbf{u}(t) = p(\mathbf{x}(t), t, \alpha)$$

What is needed is research on  
**Motor Primitives**

8





# Motor Primitives

The Search for a Common Building Block in Motor Control

## • *Previous Suggestions Included:*

- Organizational Principles
  - *2/3 Power Law*
  - *Piecewise Planarity*
  - *Speed-Accuracy Tradeoff*
  - *Optimization of Energy, Jerk, Torque Change, Motor Command Change, Task Variance, Stochastic Feedback Control, Effort, etc.*
- Equilibrium Point/Trajectory Hypotheses
- Force Fields
- Pattern Generators and Dynamics System Theory
  - *Focusing mostly on coupling phenomena (e.g. inter-limb, perception-action, intra-limb) and the necessary interaction of control and musculoskeletal dynamics*
- Contraction Theory
  - *A version of control theory for modular control*
- ... and many more

9



# (At Least) Two Different Approaches to Motor Primitives

## *The Self-organizing View*

- *Regularities of motor coordination are the “emergent” results of (well-tuned) self-organizing dynamic systems*

### *Examples:*

- *Interlimb coordination*
- *Locomotion*
- *Perception-Action coupling*
- *etc.*

## *The Optimizing View*

- *Regularities of motor coordination are the results of explicit or implicit learning/optimization processes based on inherent organizational criteria*

### *Examples:*

- *Visually guided reaching*
- *Eye movement*
- *Motor learning*
- *etc.*

10





# Optimization vs. Self-organization

## *The Self-organizing View*

### *Pros:*

- *Independent of initial conditions (generalization)*
- *Inherent stability due to attractor dynamics*
- *Coupling with external signals is relatively straightforward*
- *etc.*

### *Cons:*

- *Hard to analyze*
- *Hard to design in general*
- *Hard to apply learning*
- *etc.*

## *The Optimizing View*

### *Pros:*

- *Learning is relatively easy*
- *Potential existence of general optimization criteria*
- *Established numerical tools to perform optimization*

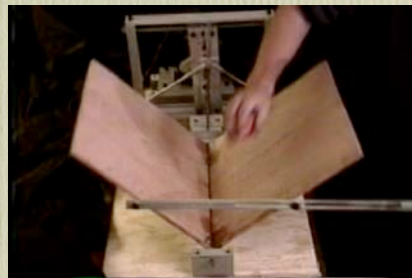
### *Cons:*

- *Optimization is time consuming and often complex*
- *Dependence of initial conditions*
- *Often explicit time dependence*
- *Problems with generalization*
- *How to express complex tasks*
- *etc.*

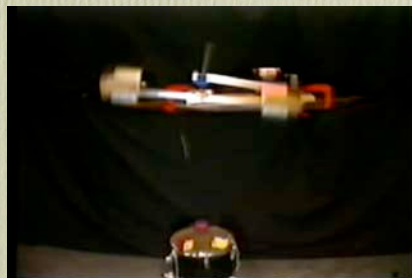
11



# Investigating Principles of Self-organization



Start with  
simple  
systems ...



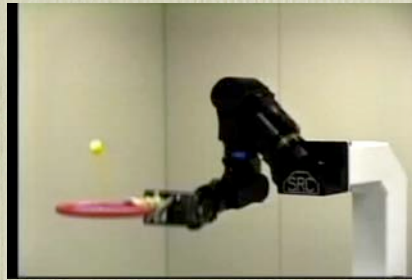
12





## Investigating Principles of Self-organization

...move on  
to more  
complex  
systems ...

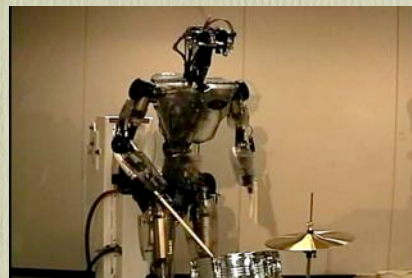
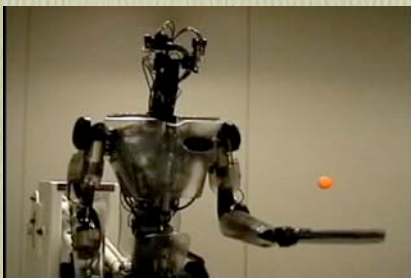


13



## Investigating Principles of Self-organization

...even  
more  
complex  
systems ...



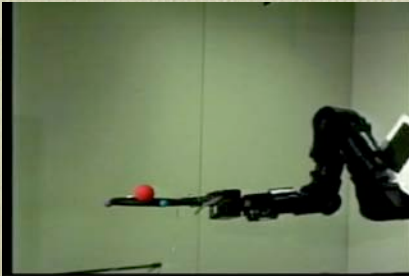
14





## Investigating Principles of Self-organization

...and add additional feedback control.



15



## Investigating Principles of Self-organization

- However, the amount of human insight and manual tuning remained very significant in all these examples.
- General principles of designing and adjusting dynamic systems were missing.

16





## Goals of this Lecture Series

- *Introduce Dynamic Motor Primitives (DMPs)*
- *Discuss Some Evidence of DMPs in Behavioral Science and Neuroscience*
- *Introduce the Formal Framework of DMPs*
- *Imitation Learning with DMPs*
- *Movement Recognition with DMPs*
- *Reinforcement Learning with DMPs*

17



## Outline

- *Part I: Dynamic Movement Primitives as a Computational Model for Human Movement?*
  - Some behavioral and fMRI data ...
- *Part II: The Formal Framework of Dynamic Motor Primitives*
  - Algorithms, imitation learning, and movement recognition
- *Part III: Reinforcement Learning with DMPs*
  - Optimization, skill learning, and other applications

18



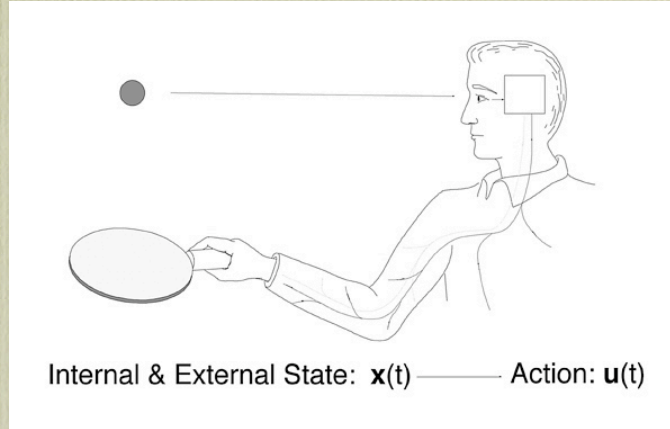


## Reminder: Control Policies

- *The General Goal of Motor Learning:*

**Control Policies**

$$\mathbf{u}(t) = p(\mathbf{x}(t), t, \alpha)$$



19



## Parameterized Dynamic Movement Primitives

- *Note the similarity between a generic control policy*

$$\mathbf{u}(t) = p(\mathbf{x}(t), t, \alpha)$$

*and nonlinear differential equations*

$$\mathbf{u}(t) = \dot{\mathbf{x}}_{desired}(t) = p(\mathbf{x}_{desired}(t), goal, \alpha)$$

*From the “Self-Organizing View”, this creates a natural distinction between two major movement classes:*

- **Rhythmic Movement**
- **Discrete Movement**

20





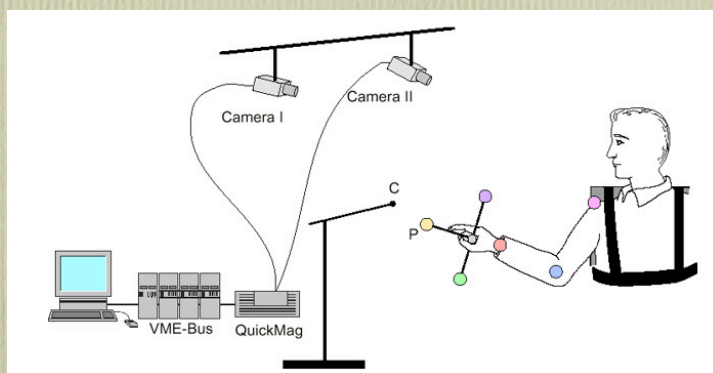
## Are Rhythmic and Discrete Movement Different Primitives?

- *Two major hypotheses in the behavioral and neurophysiological literature:*
  - All Movements are Discrete (i.e., stroke-based)
    - *Morasso, Lacquaniti, Terzuolo, Soechting, Viviani, Equilibrium-Point supporters, Optimization-supporters (via-points), and many others (1981-today)*
    - *Derived from monkey and human experiments*
  - All Movement are Oscillatory Dynamical Systems (i.e., complete and incomplete limit cycles)
    - *Dynamic System's Approach to Motor Control (Turvey, Kelso, Kugler, Schöner, et.c, 1978-today)*
    - *CPG Research (Grillner, Selverston, Cohen, etc., 1970-today)*
    - *Derived from invertebrate, lower vertebrate, and human (infant) experiments*

21



## Behavioral Experiments with Humans and Robots



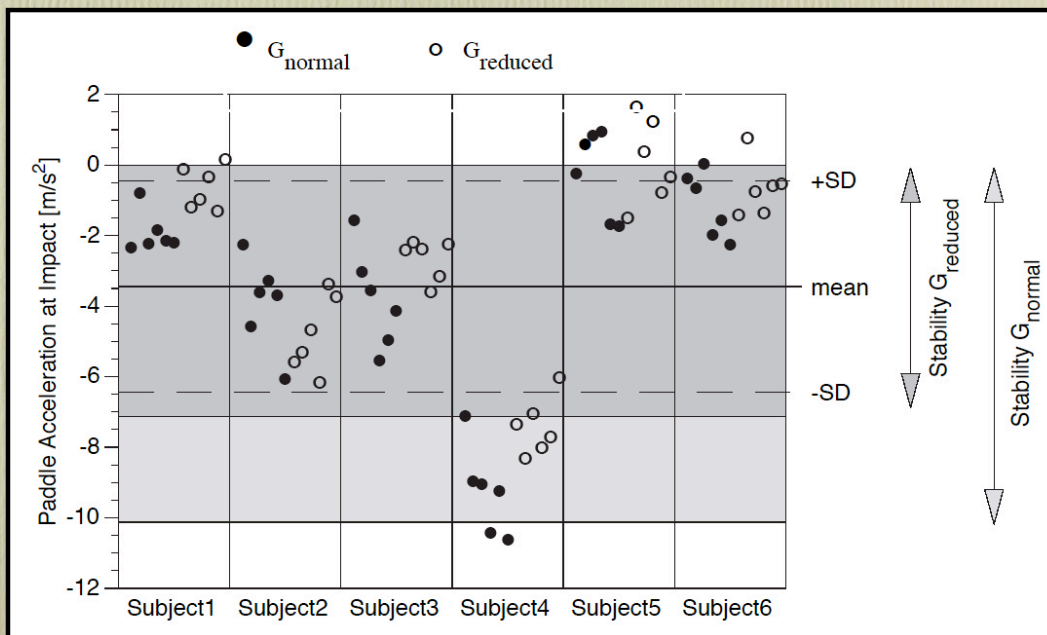
Are all arm movements stroke-based, i.e., discrete?

22





## Ball-Bouncing Seems to Be a Coupled Oscillator System

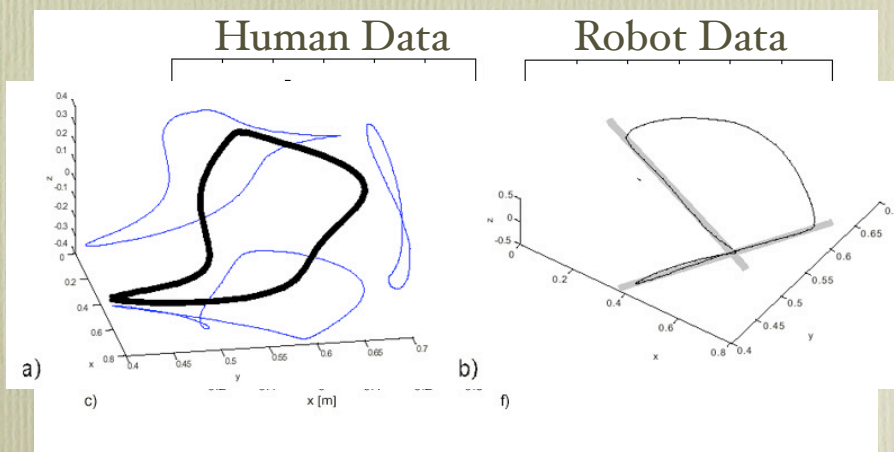


23



## The Phenomenon of Piecewise Planarity is an Artifact

- *Human hand movements in rhythmic movement exhibit piecewise planarity*



**Piecewise planarity can be explained by joint space rhythmic pattern generators**

24

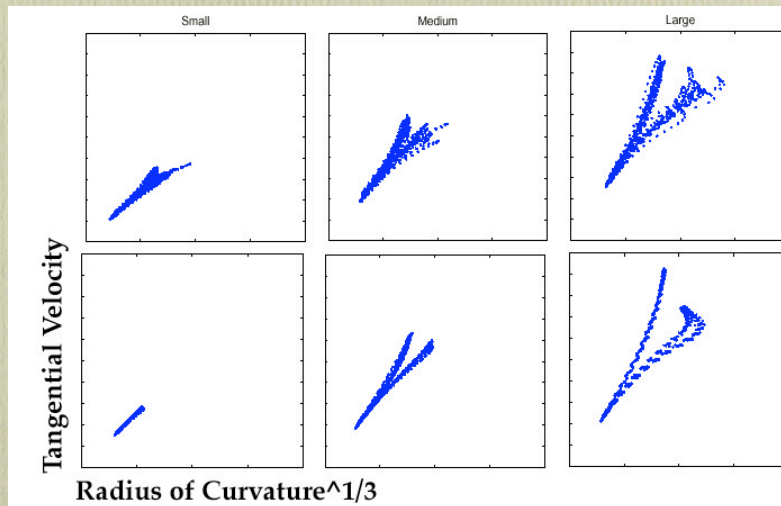




## Segmentation Based on the $2/3$ Power Law is an Artifact

Human Data

Robot Data



**Movement segmentation based on the  $2/3$  Power Law can be explain by joint space rhythmic pattern generators**

25

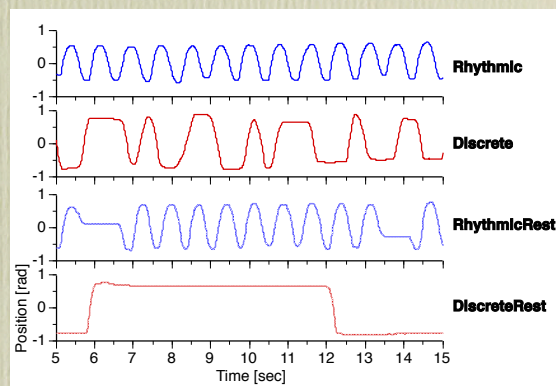


## Rhythmic and Discrete Movement in fMRI



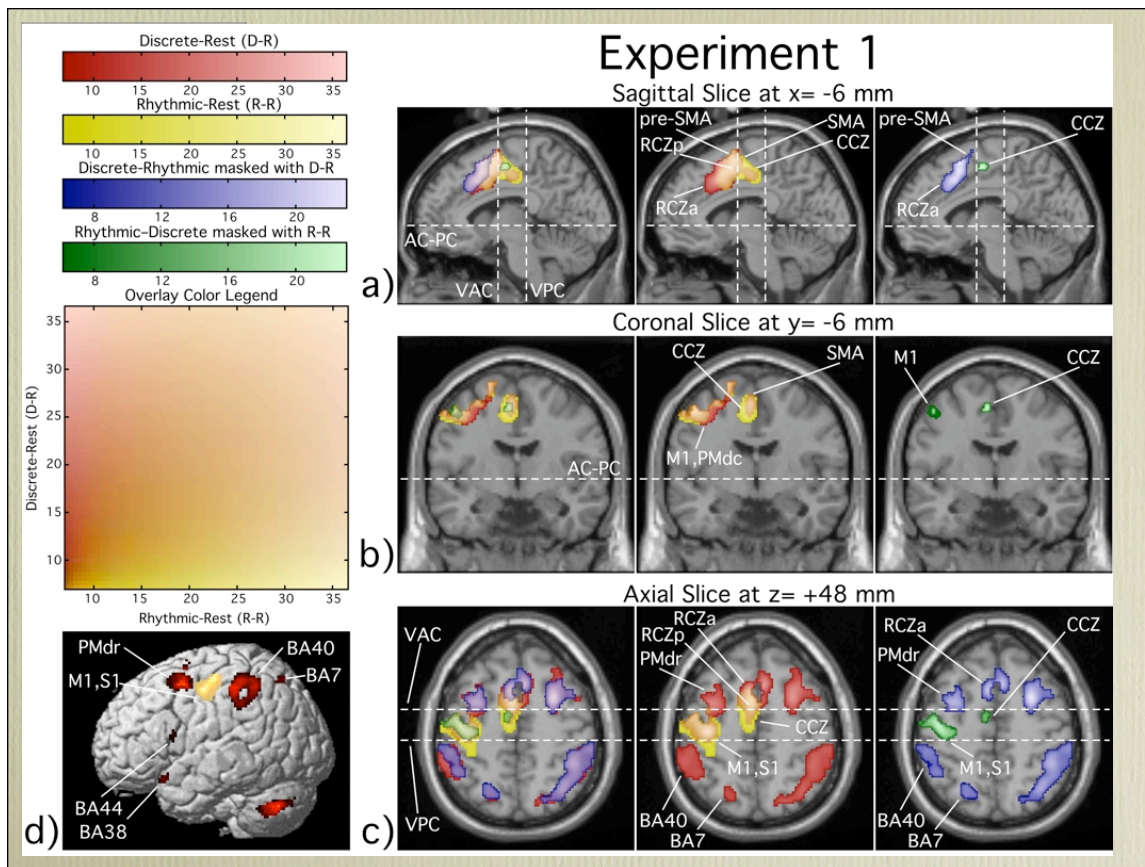
Subject in fMRI Scanner

### Typical Wrist Trajectories

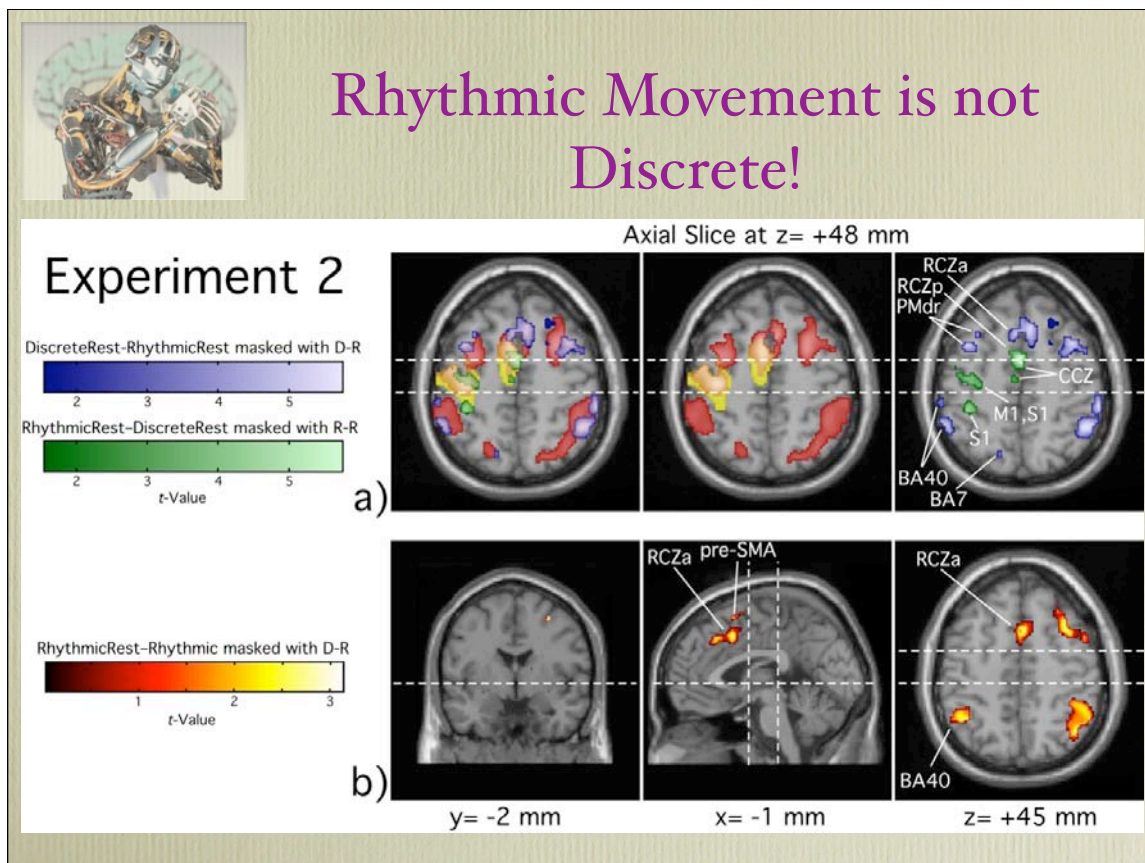


26





27



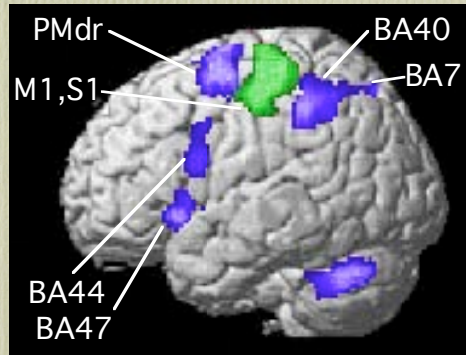
28





## Rhythmic Movement is not Discrete!

*Summary  
Results*



DISCRETE-RHYTHMIC  
RHYTHMIC-DISCRETE

29



## Summary of Behavioral Results

- *Rhythmic Movement is not generated from discrete strokes*
- *Rhythmic and Discrete Movements seem to be two different motor primitives in human motor control*
- *The hypothesis of Dynamic Movement Primitives may be viable*

*...but how would one model motor control with dynamic systems?*

30





## Outline

- *Part I: Dynamic Movement Primitives as a Computational Model for Human Movement?*
  - Some behavioral and fMRI data ...
- *Part II: The Formal Framework of Dynamic Motor Primitives*
  - Algorithms, imitation learning, and movement recognition
- *Part III: Reinforcement Learning with DMPs*
  - Optimization, skill learning, and other applications

31



## How to Model Motor Primitives? Some Inspiration from Hodgkin-Huxley Models

Total current flow :

$$I_m = C_m \dot{V} + I_{Na} + I_K + I_{leak}$$

where

$$I_{Na} = (V - V_{Na}) \bar{g}_{Na} m^3 h$$

$$I_K = (V - V_k) \bar{g}_K n^4$$

$$I_{leak} = (V - V_{leak}) g_{leak}$$

Channel Dynamics :

$$\dot{m} = \alpha_m (1 - m) - \beta_m m$$

$$\dot{h} = \alpha_h (1 - h) - \beta_h h$$

$$\dot{n} = \alpha_n (1 - n) - \beta_n n$$

Can one build (or learn) motor primitives from similar equations?

32





## Computational Goals

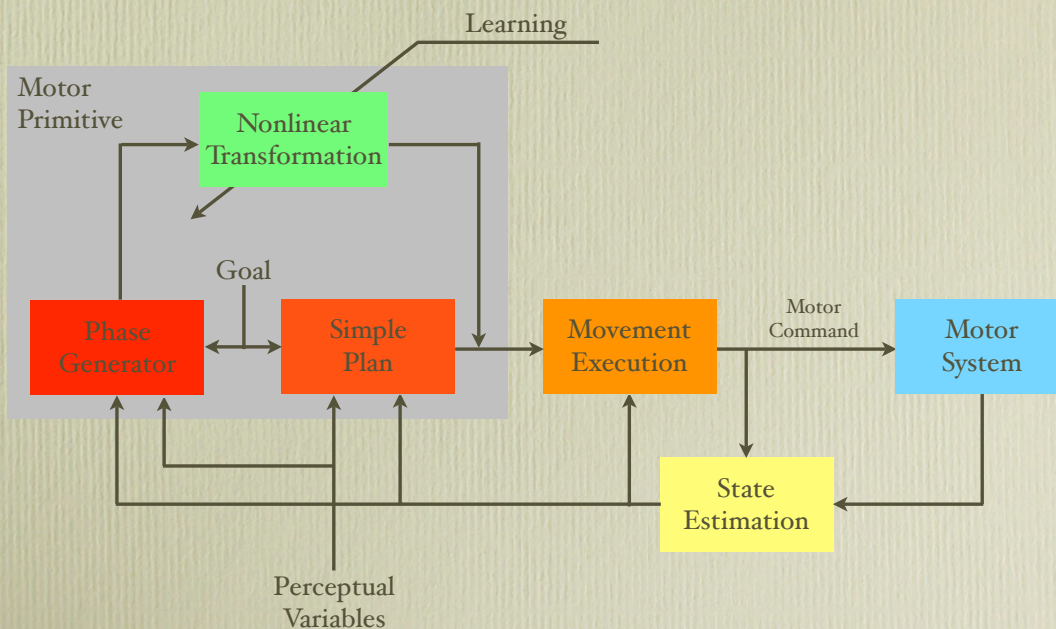
$$\dot{\mathbf{x}} = f(\mathbf{x}, \text{goal})$$

- *A Class of Dynamic Systems that Can Code:*
  - Point-to-point and periodic behavior as their attractor
  - Multi-dimensional systems that required phase locking
  - Attractors that have rather complex shape (e.g., complex phase relationships, movement reversals)
  - Learning and optimization
  - Coupling phenomena
  - Timing (without requiring explicit time)
  - Generalization (structural equivalence for parameter changes)
  - Robustness to disturbances and interactions with the environment
  - Stability guarantees

33



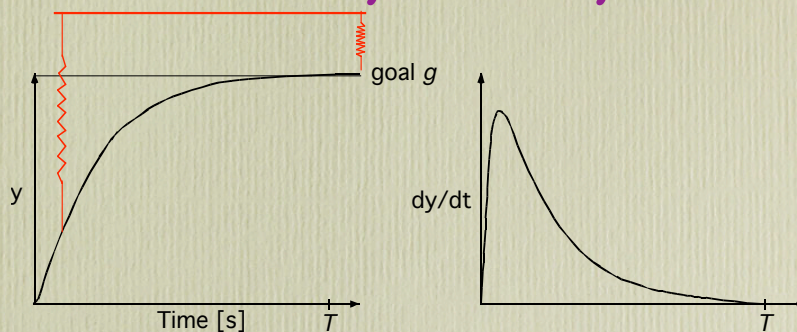
## A Sketch of a Control Diagram Using Motor Primitives



34



## Point-to-Point Movements as Dynamic Systems



E.g., for a one degree-of-freedom movement, start  
with a simple damped spring model

$$\text{---} \left[ \begin{array}{c} \text{---} \text{---} \text{---} \\ \text{---} \end{array} \right] \text{---} \equiv \ddot{y} + b\dot{y} + k(g - y) = 0 \equiv \begin{array}{l} \dot{z} = \alpha_z(\beta_z(g - y) - z) \\ \dot{y} = z \end{array}$$

35

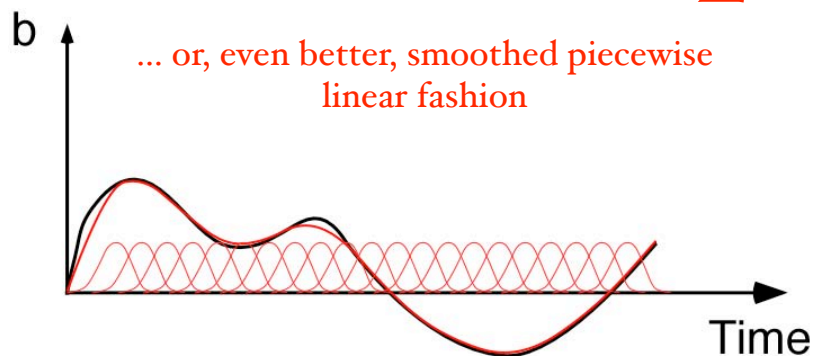


## Point-to-Point Movements as Dynamic Systems

### How to Create a Very Nonlinear Spring?

For instance, create a time varying damping term:

$$f(?) = \sum b_i \varphi(?)$$

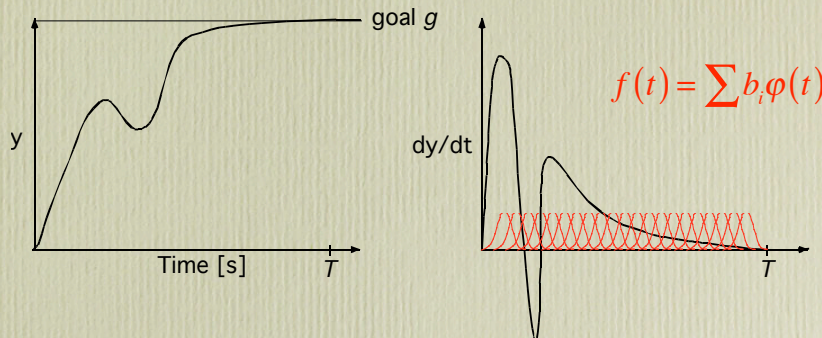


36

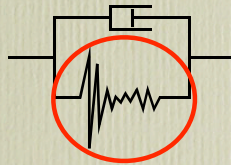




## Point-to-Point Movements as Dynamic Systems



Can one create more complex dynamics with  
a very nonlinear spring?



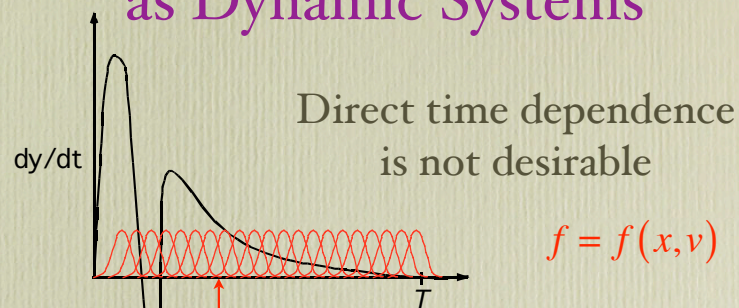
$$\dot{z} = \alpha_z (\beta_z (g - y) - z)$$

$$\dot{y} = f(?) - z$$

37

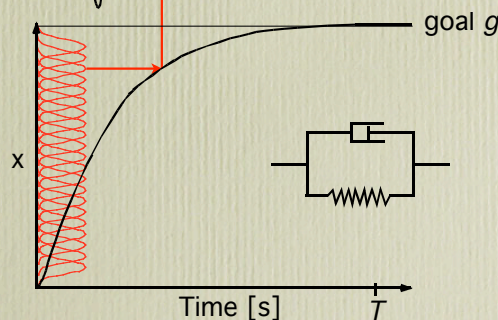


## Point-to-Point Movements as Dynamic Systems



Direct time dependence  
is not desirable

Introduce a  
behavioral  
phase dynamics



$$\dot{v} = \alpha_v (\beta_v (g - x) - v)$$

$$\dot{x} = \alpha_x v$$

38





## A Dynamic Systems Model for Discrete Movement

- *A learnable nonlinear point attractor with guaranteed stability properties*

Behavioral Phase

$$\dot{v} = \alpha_v (\beta_v (g - x) - v)$$

$$\dot{x} = \alpha_x v$$

Nonlinear Function

$$f(x, v)$$

$$\dot{z} = \alpha_z (\beta_z (g - y) - z)$$

$$\dot{y} = \alpha_y (f(x, v) + z)$$

Nonlinear Spring-Damper



39



## A Dynamic Systems Model for Discrete Movement

- *A learnable nonlinear point attractor with guaranteed stability properties*

Trajectory Plan  
Dynamics

$$\begin{cases} \dot{z} = \alpha_z (\beta_z (g - y) - z) \\ \dot{y} = \alpha_y (f(x, v) + z) \end{cases}$$

where

Canonical  
Dynamics

$$\begin{cases} \dot{v} = \alpha_v (\beta_v (g - x) - v) \\ \dot{x} = \alpha_x v \end{cases}$$

Local Linear  
Model Approx.

$$\begin{cases} f(x, v) = \frac{\sum_{i=1}^k w_i b_i v}{\sum_{i=1}^k w_i} \\ w_i = \exp\left(-\frac{1}{2} d_i (\bar{x} - c_i)^2\right) \text{ and } \bar{x} = \frac{x - x_0}{g - x_0} \end{cases}$$

“Phase Velocity”  
for Nonlinear  
Amplification

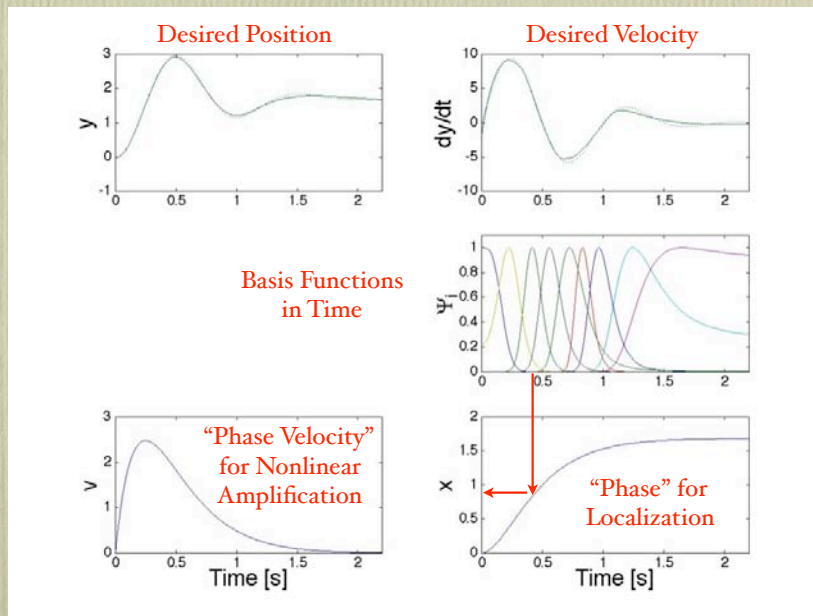
Learnable  
Weights

40





## An Example



41



## Extension to Periodic Systems

- *A learnable nonlinear limit cycle attractor with guaranteed stability properties*

Behavioral Phase

Nonlinear Function

$$\dot{r} = \alpha_r (A - r)$$

$$f(r, \varphi)$$

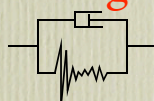
$$\dot{\varphi} = \omega$$

Phase Oscillator  
with amplitude A

$$\dot{z} = \alpha_z (\beta_z (g - y) - z)$$

$$\dot{y} = \alpha_y (f(r, \varphi) + z)$$

Nonlinear Oscillating Spring-Damper



42



## Extension to Periodic Systems

- *Use van Mises Basis Function for Local Linear Models*

Trajectory Plan Dynamics 
$$\begin{cases} \dot{z} = \alpha_z (\beta_z (g - y_m) - z) \\ \dot{y} = \alpha_y (f(r, \varphi) + z) \end{cases}$$

where

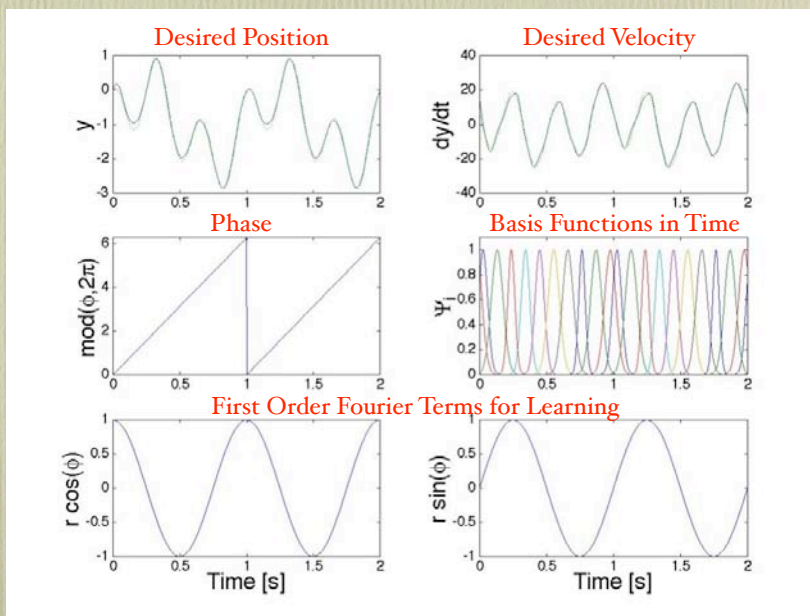
Canonical System 
$$\begin{cases} \dot{r} = \alpha_r (A - r) \\ \dot{\varphi} = \omega \end{cases}$$

Local Linear Models using van Mises bases 
$$\begin{cases} f(x, v) = \frac{\sum_{i=1}^k w_i \mathbf{b}_i^T \mathbf{x}}{\sum_{i=1}^k w_i} \text{ where } \mathbf{x} = \begin{bmatrix} r \cos \varphi \\ r \sin \varphi \end{bmatrix} \\ w_i = \exp(d_i (\cos(\varphi - c_i) - 1)) \end{cases}$$

43



## An Example



44





## Real Demos

Pointer Attractor  
Demo

Limit Cycle  
Demo

45



## How Did These Demos Work?

- *Given*

- A goal
- A desired trajectory

- *Algorithm*

- Extract features from the desired trajectory
- Adjust time constants of canonical dynamics to movement duration
- Using Locally Weighted Learning to solve nonlinear function approximation problem

$$\dot{y}_{\text{target}} = \frac{\dot{y}_{\text{demo}}}{\alpha_y} - z = f(x, v)$$

- where  $z$  can be calculated from desired trajectory

Trajectory Plan  
Dynamics

$$\begin{cases} \dot{z} = \alpha_z (\beta_z (g - y) - z) \\ \dot{y} = \alpha_y (f(x, v) + z) \end{cases}$$

where

Canonical  
Dynamics

$$\begin{cases} \dot{v} = \alpha_v (\beta_v (g - x) - v) \\ \dot{x} = \alpha_x v \end{cases}$$

Local Linear  
Model Approx.

$$\begin{cases} f(x, v) = \frac{\sum_{i=1}^k w_i b_i v}{\sum_{i=1}^k w_i} \\ w_i = \exp\left(-\frac{1}{2} d_i (\bar{x} - c_i)^2\right) \text{ and } \bar{x} = \frac{x - x_0}{g - x_0} \end{cases}$$

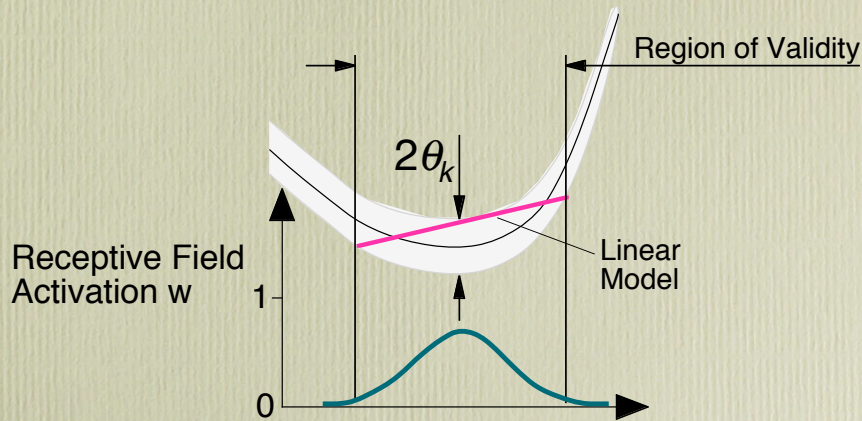
Note: This is a one-shot  
learning problem, i.e.,  
no iterations!

46





## Locally Weighted Learning with locally linear models

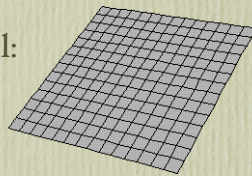


47



## Locally Weighted Learning with locally linear models

Linear Model:



learned with

Recursive weighted least squares:

$$\beta_k^{n+1} = \beta_k^n + w \mathbf{P}_k^{n+1} \mathbf{x} (\mathbf{y} - \tilde{\mathbf{x}}^T \beta_k^n)^T$$

$$\mathbf{P}_k^{n+1} = \frac{1}{\lambda} \left( \mathbf{P}_k^n - \frac{\mathbf{P}_k^n \tilde{\mathbf{x}} \tilde{\mathbf{x}}^T \mathbf{P}_k^n}{\frac{\lambda}{w} + \tilde{\mathbf{x}}^T \mathbf{P}_k^n \tilde{\mathbf{x}}} \right)$$

$$\mathbf{y} = \beta_x^T \mathbf{x} + \beta_0 = \beta^T \tilde{\mathbf{x}} \quad \text{where} \quad \tilde{\mathbf{x}} = [\mathbf{x}^T \ 1]^T$$

Weighting Kernel:



learned with

Gradient descent in penalized leave-one-out local cross-validation (PRESS) cost function:

$$\mathbf{M}_k^{n+1} = \mathbf{M}_k^n - \alpha \frac{\partial J}{\partial \mathbf{M}}$$

$$J = \frac{1}{\sum_{i=1}^N w_{k,i}} \sum_{i=1}^N w_{k,i} \|\mathbf{y}_i - \hat{\mathbf{y}}_{k,i,-i}\|^2 + \gamma \sum_{i=1, j=1}^n D_{k,ij}^2$$

$$w = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{c})^T \mathbf{D}(\mathbf{x} - \mathbf{c})\right) \quad \text{where} \quad \mathbf{D} = \mathbf{M}^T \mathbf{M}$$

Combined Prediction:

$$\mathbf{y} = \frac{\sum_{i=1}^K w_i \mathbf{y}_i}{\sum_{i=1}^K w_i}$$

add model when

$$\text{if } \min_k(w_k) < w_{gen}$$

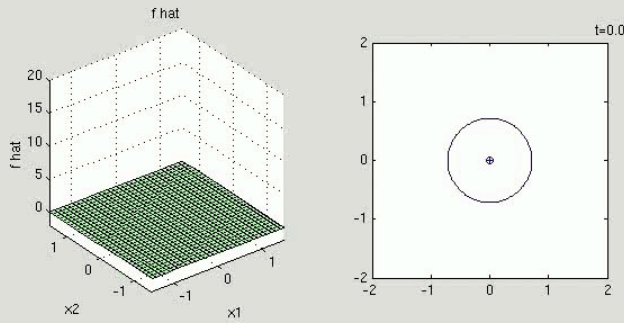
create new RF at  $\mathbf{c}_{K+1} = \mathbf{x}$

48





## Visualization of On-line Function Approximation



49



## Additional Issues to Be Addressed

- *Coordinating multiple Degrees-of-Freedom*
- *Movement Recognition*
- *Dealing with Perturbations*

50





## Coordinating Multiple Degrees-of-Freedom

### Behavioral Phase

$$\dot{r} = \alpha_r (A - r)$$

$$\dot{\phi} = \omega$$

### Nonlinear Function

$$f(r, \phi)$$

$$\dot{z} = \alpha_z (\beta_z (g - y) - z)$$

$$\dot{y} = \alpha_y (f(r, \phi) + z)$$

51



## Coordinating Multiple Degrees-of-Freedom

### Behavioral Phase

$$\dot{r} = \alpha_r (A - r)$$

$$\dot{\phi} = \omega$$

### Nonlinear Function

$$\dot{z}_1 = \alpha_z (\beta_z (g_1 - y_1) - z_1)$$

$$\dot{y}_1 = \alpha_y (f_1(r, \phi) + z_1)$$

$$f_1(r, \phi) \text{ DOF-1}$$

$$\dot{z}_2 = \alpha_z (\beta_z (g_2 - y_2) - z_2)$$

$$\dot{y}_2 = \alpha_y (f_2(r, \phi) + z_2)$$

$$f_2(r, \phi) \text{ DOF-2}$$

$$\dot{z}_3 = \alpha_z (\beta_z (g_3 - y_3) - z_3)$$

$$\dot{y}_3 = \alpha_y (f_3(r, \phi) + z_3)$$

$$f_3(r, \phi) \text{ DOF-3}$$

$\vdots$

$$\dot{z}_N = \alpha_z (\beta_z (g_N - y_N) - z_N)$$

$$\dot{y}_N = \alpha_y (f_N(r, \phi) + z_N)$$

$$f_N(r, \phi) \text{ DOF-N}$$

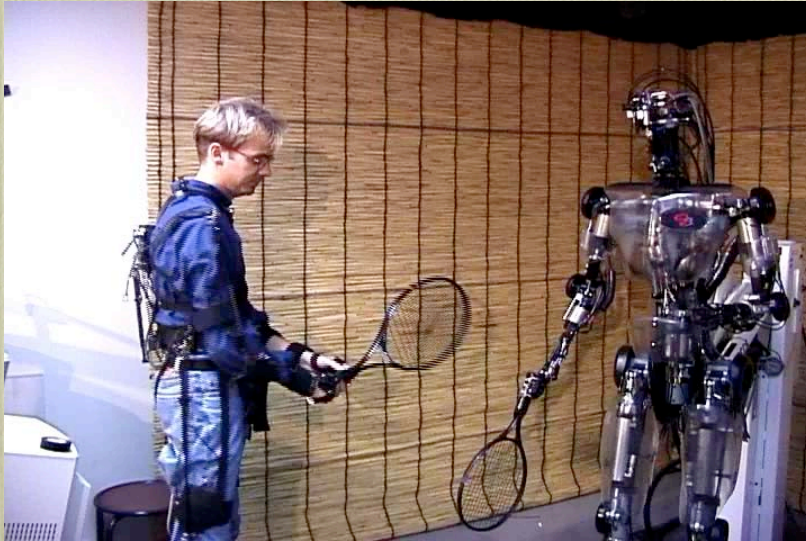
A Central  
Pattern  
Generator?

52





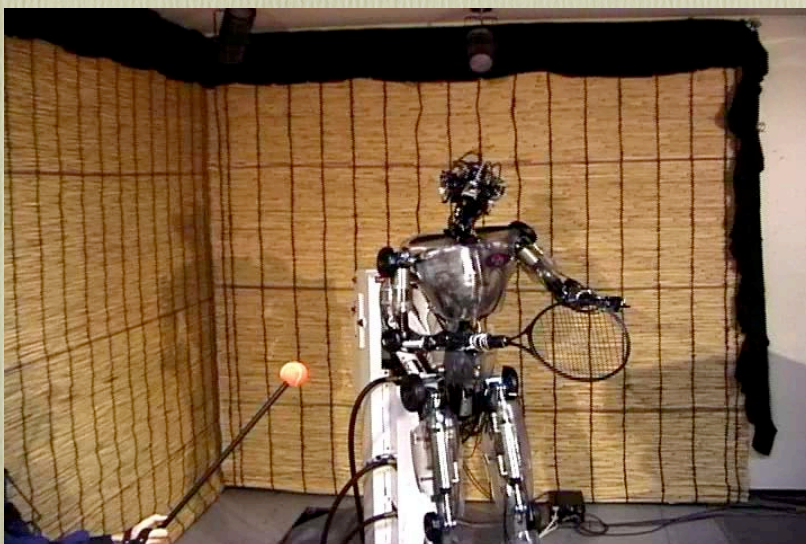
## Example: A Tennis Backhand as a Dynamic System



53



## Example: A Tennis Backhand as a Dynamic System

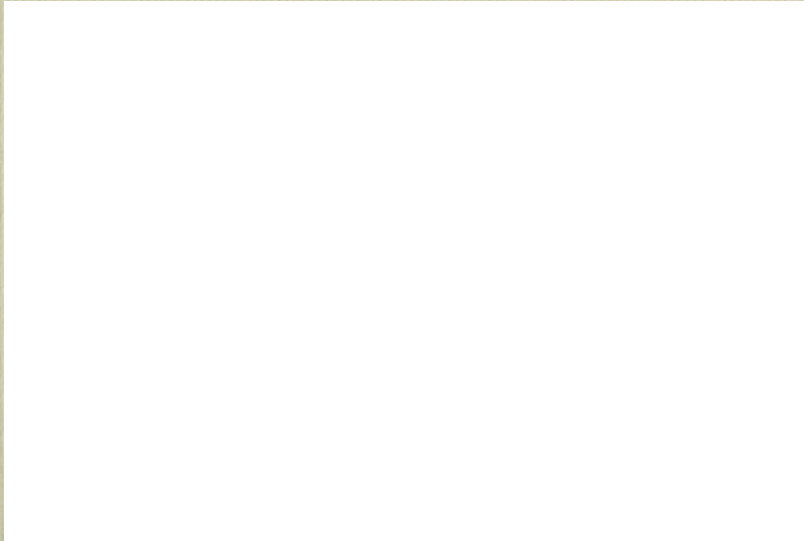


54





## Example: Various Periodic Movement Patterns



55



## Dealing with Perturbations: Adding On-line Modification

- *Among the most interesting properties of the DMP approach is the on-line modification with coupling terms*

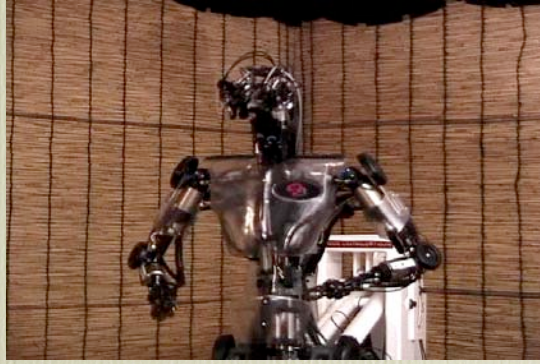
Trajectory Plan Dynamics	$\begin{cases} \dot{z} = \alpha_z (\beta_z (g - y) - z) \\ \dot{y} = \alpha_y (f(x, v) + z) + \alpha_p (y_{real} - y) \end{cases}$ <p>where</p>
Canonical Dynamics	$\begin{cases} \dot{v} = \alpha_v (\beta_v (g - x) - v) \\ \dot{x} = \alpha_x v \left( 1 + \alpha_p (y_{real} - y)^2 \right)^{-1} \end{cases}$
Local Linear Model Approx.	$\begin{cases} f(x, v) = \frac{\sum_{i=1}^k w_i b_i v}{\sum_{i=1}^k w_i} \\ w_i = \exp\left(-\frac{1}{2} d_i (\bar{x} - c_i)^2\right) \text{ and } \bar{x} = \frac{x - x_0}{g - x_0} \end{cases}$

56





## Example: A Periodic Movement with Perturbation



57



## Movement Recognition

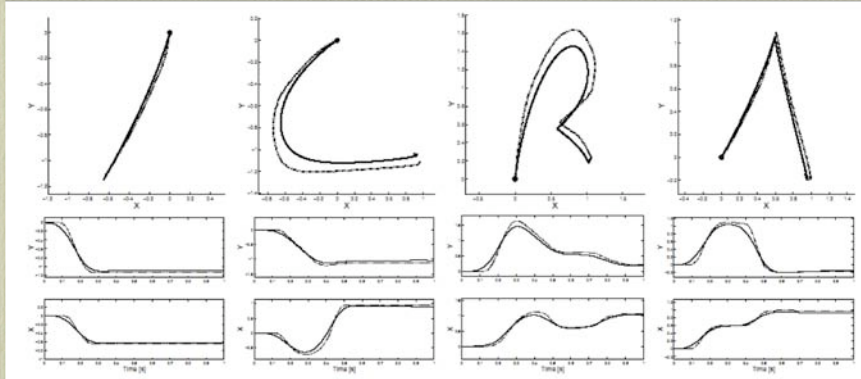
- *An Important Property:*
  - By design, the dynamic systems are structurally equivalent under scaling the distance to the goal for point attractor systems, and the amplitude for limit cycles
  - Structural equivalence also holds for a uniform scaling of the time constants
  - Thus, the parameters of the nonlinear function are invariant under spatial and temporal scaling of a movement and can be used to classify a movement pattern

58





## Example: Recognizing Graffiti Characters



Characters are fit with 2DOF  
dynamic point attractors

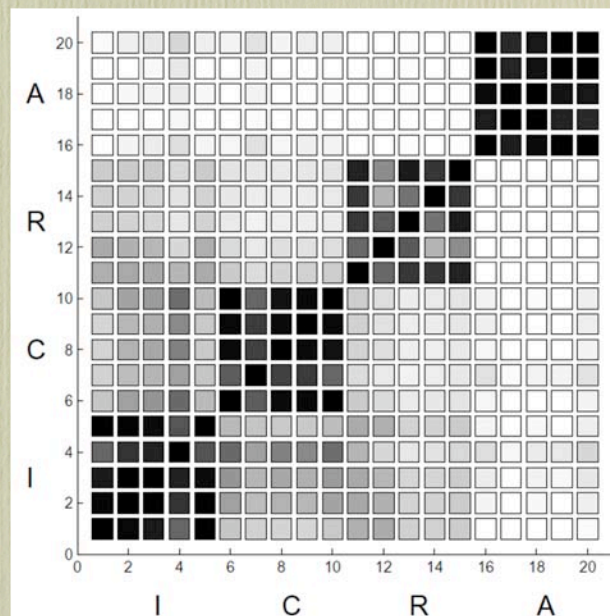
59



## Example: Recognizing Graffiti Characters

About 87%  
correct  
classification  
using  
correlation  
coefficient

$$Score = \frac{\mathbf{b}_{Template}^T \mathbf{b}_{Observed}}{\|\mathbf{b}_{Template}\| \|\mathbf{b}_{Observed}\|}$$

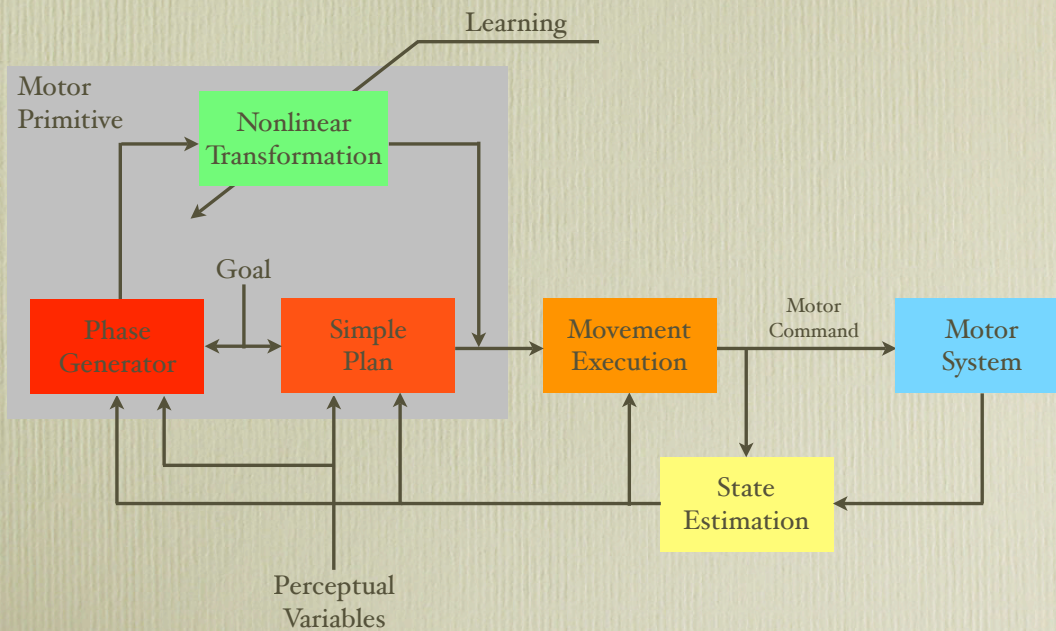


60





## A Sketch of a Control Diagram Using Motor Primitives



61



## Computational Goals: What Did We Accomplish?

$$\dot{\mathbf{x}} = f(\mathbf{x}, \text{goal})$$

- *A Class of Dynamic Systems that Can Code:*
  - ✓ • Point-to-point and periodic behavior as their attractor
  - ✓ • Multi-dimensional systems that required phase locking
  - ✓ • Attractors that have rather complex shape (e.g., complex phase relationships, movement reversals)
  - ✓ • Learning and optimization
  - ✓ • Coupling phenomena
  - ✓ • Timing (without requiring explicit time)
  - ✓ • Generalization (structural equivalence for parameter changes)
  - ✓ • Robustness to disturbances and interactions with the environment
  - ✓ • Stability guarantees

62





## Outline

- *Part I: Dynamic Movement Primitives as a Computational Model for Human Movement?*
  - Some behavioral and fMRI data ...
- *Part II: The Formal Framework of Dynamic Motor Primitives*
  - Algorithms, imitation learning, and movement recognition
- *Part III: Reinforcement Learning with DMPs*
  - Optimization, skill learning, and other applications

63



## Reinforcement Learning with Movement Primitives

- *Given:*
  - An acceleration-based motor primitives with parameters  $\theta$

$$\frac{1}{\tau} \dot{z} = \alpha_z (\beta_z (s - y) - z) + f \quad \frac{1}{\tau} \dot{v} = \alpha_v (\beta_v (g - x) - v)$$

$$\frac{1}{\tau} \dot{y} = z \quad \frac{1}{\tau} \dot{x} = v$$

$$\frac{1}{\tau} \dot{s} = \alpha_g (g - s)$$

such that

$y, \dot{y}, \ddot{y} = \dot{z}$  is the desired trajectory

$$f(x, v, g) = \frac{\sum_{i=1}^N \psi_i b_i v}{\sum_{i=1}^N \psi_i}, \text{ where } \psi_i = \exp \left( -h_i \left( \frac{x - x_0}{g - x_0} - c_i \right)^2 \right)$$

64





## Reinforcement Learning with Movement Primitives

- *Given (cont'd):*

- A stochastic realization of the motor primitive

$$u = \ddot{y} + \varepsilon$$

$$\pi(u | \mathbf{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(u - \ddot{y}(\mathbf{x}))^2\right)$$

$$\mathbf{x} = [x, v, z, y, s]^T$$

- Several n-step roll-outs (a.k.a. trajectory, sample path, etc.)

$$\xi = \{\mathbf{x}_1, u_1, \mathbf{x}_2, u_2, \dots, \mathbf{x}_n, u_n, \mathbf{x}_{n+1}\}$$

- An immediate reward

$$r_i = r(\mathbf{x}_i, u_i)$$

- A discounted long term reward of a roll-out

$$R(\xi) = \sum_{i=1}^n \gamma^{i-1} r_i$$

65



## Reinforcement Learning with Movement Primitives

- *Given (cont'd):*

- A cost criterion to be optimized

$$J(\theta) = E\{R(\xi)\}_{\xi}$$

- *Approach*

- Gradient descent in the primitive parameter using stochastic policy gradient methods

$$\theta = [\mathbf{b}^T \quad \sigma^2]^T$$

$$\theta^{n+1} = \theta^n + \alpha \frac{\partial J}{\partial \theta}$$

66





## Computing the Policy Gradient

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \\
 &= \nabla_{\theta} E\{R(\xi)\} \\
 &= \nabla_{\theta} \int_{\Xi} p_{\theta}(\xi) R(\xi) d\xi \\
 &= \int_{\Xi} \nabla_{\theta} p_{\theta}(\xi) R(\xi) d\xi \\
 &= \int_{\Xi} p_{\theta}(\xi) \nabla_{\theta} \log p_{\theta}(\xi) R(\xi) d\xi \\
 &= \langle \nabla_{\theta} \log p_{\theta}(\xi) R(\xi) \rangle_{\xi}
 \end{aligned}$$

$$p_{\theta}(\xi) = p(\mathbf{x}_1) \prod_{i=1}^n \pi_{\theta}(u_i | \mathbf{x}_i) p(\mathbf{x}_{i+1} | \mathbf{x}_i, u_i)$$

$$\nabla_{\theta} \log p_{\theta}(\xi) = \sum_{i=1}^n \nabla_{\theta} \log \pi_{\theta}(u_i | \mathbf{x}_i)$$

67



## Computing the Policy Gradient

- *A useful observation*

$$\int_{\Xi} p_{\theta}(\xi) d\xi = 1$$

Thus

$$\nabla_{\theta} \int_{\Xi} p_{\theta}(\xi) d\xi = \int_{\Xi} p_{\theta}(\xi) \nabla_{\theta} \log p_{\theta}(\xi) d\xi = 0$$

and

$$\beta \int_{\Xi} p_{\theta}(\xi) \nabla_{\theta} \log p_{\theta}(\xi) d\xi = 0$$

for any parameter  $\beta$

68





# Policy Gradient I

## Classic Policy Gradient

- Thus, the policy gradient with baseline becomes

$$\nabla_{\theta} J(\theta) = \left\langle \left( \sum_{i=1}^n \nabla_{\theta} \log \pi_{\theta}(u_i | \mathbf{x}_i) \right) (R(\xi) - \beta) \right\rangle_{\xi}$$

Episodic REINFORCE

- The baseline that minimized the variance of the policy gradient can be shown to be:

$$\beta_i = \frac{\left\langle (\nabla_{\theta_i} \log p_{\theta}(\xi))^2 R(\xi) \right\rangle_{\xi}}{\left\langle (\nabla_{\theta_i} \log p_{\theta}(\xi))^2 \right\rangle_{\xi}}.$$

William, 1992; Lawrence et al. 2004

69



# Improving the Gradient

- Inserting the definition of the roll-out reward

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \left\langle \left( \sum_{i=1}^n \nabla_{\theta} \log \pi_{\theta}(u_i | \mathbf{x}_i) \right) (R(\xi) - \beta) \right\rangle_{\xi} \\ &= \left\langle \left( \sum_{i=1}^n \nabla_{\theta} \log \pi_{\theta}(u_i | \mathbf{x}_i) \right) \left( \sum_{i=1}^n \gamma^{i-1} r_i - \beta \right) \right\rangle_{\xi} \end{aligned}$$

- Realizing that rewards at time  $k$  cannot be affected by actions at time  $i > k$ , the gradient can be rewritten

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \left\langle \sum_{k=1}^n \left( \sum_{j=1}^k \nabla_{\theta} \log \pi_{\theta}(u_j | x_j, u_j) \right) (\gamma^{k-1} r(x_k, u_k) - \beta(k)) \right\rangle_{\xi} \\ \beta_i(k) &= \frac{\left\langle (\nabla_{\theta_i} \log p_{\theta}(\xi_{1:k}))^2 r(x_k, u_k) \right\rangle_{\xi_{1:k}}}{\left\langle (\nabla_{\theta_i} \log p_{\theta}(\xi_{1:k}))^2 \right\rangle_{\xi_{1:k}}} \end{aligned}$$

GPOMDP

Baxter et al. 1999

70





## Adding Function Approximation

- *Replace the roll-out reward with a function approximator*

$$\nabla_{\theta} J(\theta) = \left\langle \left( \sum_{i=1}^n \nabla_{\theta} \log \pi_{\theta}(u_i | \mathbf{x}_i) \right) (R(\xi) - \beta) \right\rangle_{\xi}$$

$$R(\xi) - \beta = f(\xi)$$

- *It can be shown (Sutton et al, Konda et al, 2000) that in order to avoid biasing the gradient, the function approximator needs to be of the form*

$$f(\xi) = [\nabla_{\theta} \log p_{\theta}(\xi)^T \quad 1] \mathbf{w}$$

71



## The Natural Gradient

- *Inserting the function approximator into the gradient, and canceling all irrelevant terms results in*

$$\nabla_{\theta} J(\theta) = \left\langle \left( \sum_{i=1}^n \nabla_{\theta} \log \pi_{\theta}(u_i | \mathbf{x}_i) \right) \left( \sum_{i=1}^n \nabla_{\theta} \log \pi_{\theta}(u_i | \mathbf{x}_i) \right)^T \right\rangle_{\xi} \tilde{\mathbf{w}}$$

$$= \mathbf{F} \tilde{\mathbf{w}}$$

where  $\tilde{\mathbf{w}}$  is the  $\mathbf{w}$  vector without the constant coefficient ,  
and  $\mathbf{F}$  is the Fisher Information Matrix

- *Amari (1999) demonstrated that a more efficient gradient in stochastic optimization is*

$$\nabla_{\theta} J(\theta)_{nat} = \mathbf{F}^{-1} \nabla_{\theta} J(\theta)$$

Thus:

$$\nabla_{\theta} J(\theta)_{nat} = \mathbf{F}^{-1} \mathbf{F} \tilde{\mathbf{w}} = \tilde{\mathbf{w}}$$

72





# Policy Gradient II

## Natural Policy Gradient

$$\nabla_{\theta} J(\theta)_{nat} = \tilde{\mathbf{w}}$$

where

$$\mathbf{w} = [\tilde{\mathbf{w}}^T \ w_0]^T$$

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X}) \mathbf{X}^T \mathbf{Y}$$

$$\mathbf{X} = \begin{bmatrix} \nabla_{\theta} \log p_{\theta}(\xi_1)^T & 1 \\ \nabla_{\theta} \log p_{\theta}(\xi_2)^T & 1 \\ \vdots & \vdots \\ \nabla_{\theta} \log p_{\theta}(\xi_K)^T & 1 \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} R(\xi_1) \\ R(\xi_2) \\ \vdots \\ R(\xi_K) \end{bmatrix}$$

Note: There is also a “natural” gradient based on GPOMDP

73



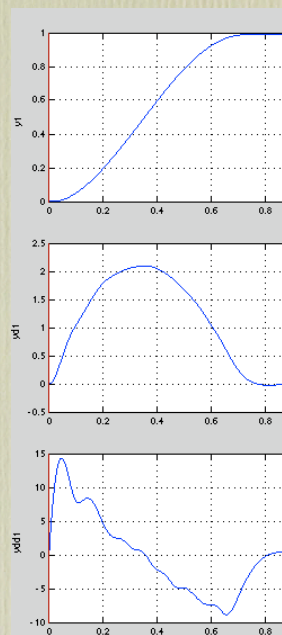
## Example: Learning a Minimum Motor Command Trajectory

- Assume that goal and duration of a movement are given, learn the trajectory that minimizes

$$r_i = cu_i^2$$

accumulated over the entire trajectory

- Compare different stochastic policy algorithms

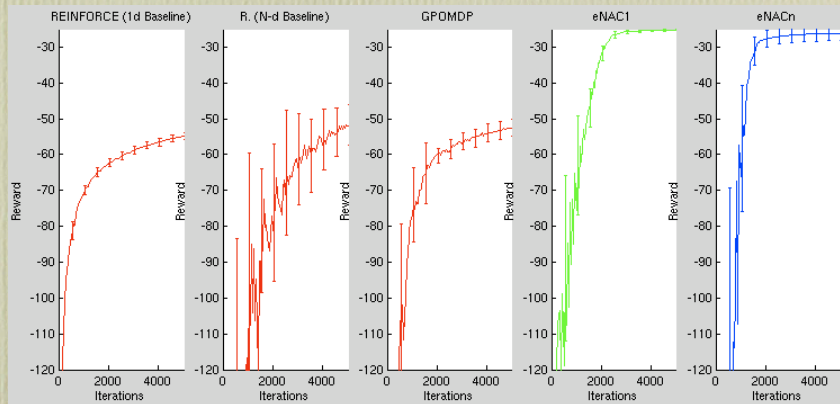


74





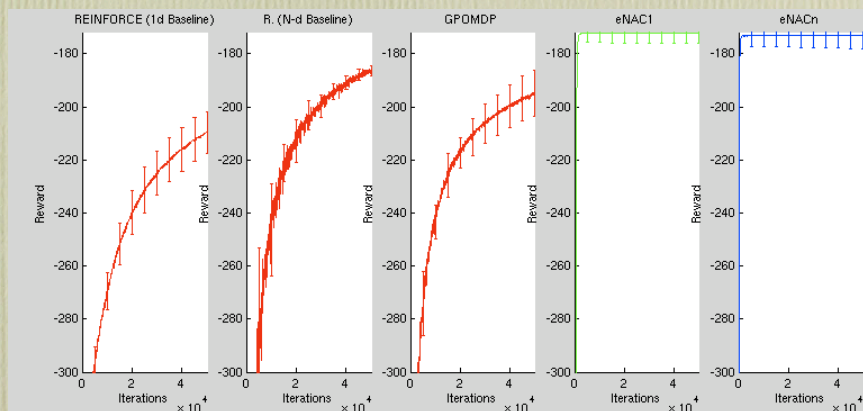
## Example: Learning a Minimum Motor Command Trajectory



75



## Example: Learning a Minimum Motor Command Trajectory



These results used a cubic spline representation.

76





## Example: Imitation Learning with Self-Improvement



Goal: Hit ball precisely      Note: about 150 trials are needed.

77



## Coupling of External Beat and Control



78





## Discussion

- *Evidence from behavioral and fMRI data supports the idea of motor primitives, in particular in a dynamic systems framework*
- *Formulating motor primitives as kinematic dynamic systems for movement planning offers a model of movement generation, which can address many issues, including:*
  - Optimization
  - Reinforcement learning, supervised learning, imitation learning
  - Perception-Action Coupling
  - Motor Primitives
  - Generalization
- *The suggested approach is more of a design principle rather than fixed formalism*
- *The necessary computations of this approach remain (so far) manageable and potentially biologically plausible, and may thus serve as a tool to model primate motor control phenomena*