

# Episode-Based Reinforcement Learning using Movement Primitives

appliedAI Seminar - Reinforcement Learning

---

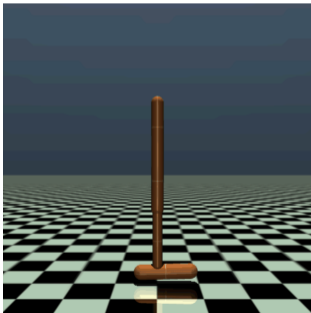
Maximilian Hüttenrauch

TransferLab@appliedAI

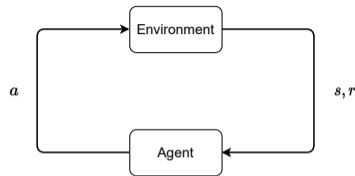
# Introduction

## Deep RL

- Widely used in sequential decision making problems
- Works well in broad set of tasks (e.g., forward motion)
- Step-based paradigm  $\pi(\mathbf{a} \mid \mathbf{s})$



Source: <https://gymnasium.farama.org/environments/mujoco/hopper/>



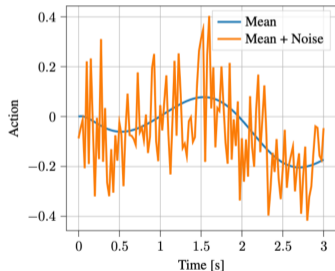
## MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$

- State space  $\mathcal{S}$
- Action space  $\mathcal{A}$
- Transition dynamics  
 $P = p(s' \mid s, a)$
- Reward function  $R = r(s, a)$
- Discount factor  $\gamma$

# Introduction

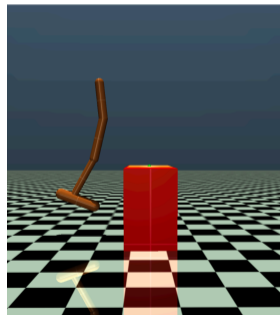
## Why episode-based RL?

- Action noise bad for robotics
- Learned policies are not energy efficient [Otto et al., 2022]



Source: [Li et al., 2024]

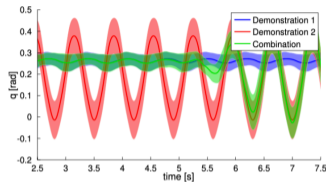
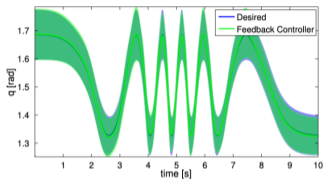
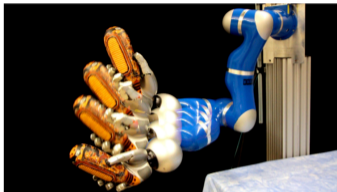
- Not all tasks can be modelled with dense/Markovian reward



Source: [Hüttenrauch and Neumann, 2024]

**Goal: Provide a different approach to learn episodic tasks**

# Movement Primitives



Source: [Paraschos et al., 2013]

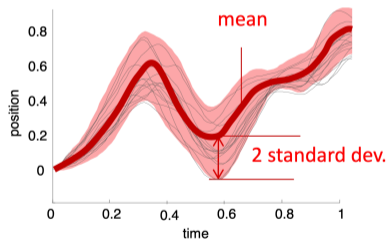
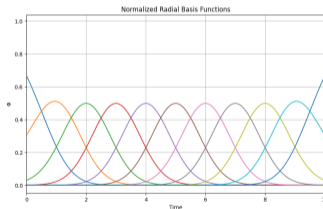
- Parameterized trajectory  $\boldsymbol{\tau} = \mathbf{y}_{1:T} = \begin{bmatrix} q_{1:T} \\ \dot{q}_{1:T} \end{bmatrix} = f(\mathbf{w})$
- Modular movement representation and generation
- Linear basis function model as function approximator
- Execution: Trajectory tracking controller

# Probabilistic Movement Primitives [Paraschos et al., 2013]

- Distribution over trajectories  $p(\boldsymbol{\tau} | \boldsymbol{w}) = \prod_t \mathcal{N}(\boldsymbol{y}_t | \boldsymbol{\Phi}_t \boldsymbol{w}, \boldsymbol{\Sigma}_y)$
- Trajectory defined directly in terms of position
- Start position cannot be enforced

ProMPs (mean trajectory)

$$\boldsymbol{y}_t = \begin{bmatrix} q_t \\ \dot{q}_t \end{bmatrix} = \boldsymbol{\Phi}_t \boldsymbol{w} + \boldsymbol{\epsilon}_y$$



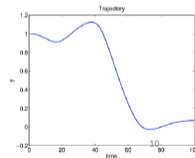
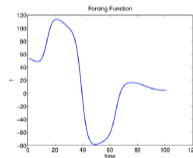
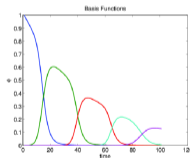
Source (right): [Neumann, 2018]

- Spring damper second order dynamical system resulting in smooth trajectory
- Trajectory defined in terms of acceleration profile
- Starting position built in, final position can be adapted

## DMPs

$$\tau \ddot{y} = \alpha_z (\beta_z (g - y) - \dot{y}) + f(t)$$

$$f(t) = \Phi_t w$$



Source: [Neumann, 2018]

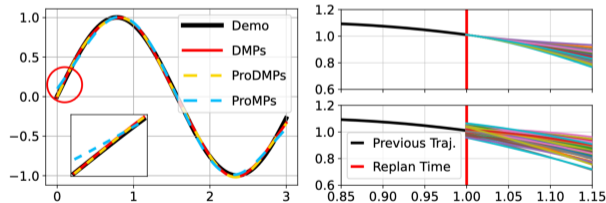
# Probabilistic Dynamic Movement Primitives [Li et al., 2023]

- Combine advantages of ProMPs and DMPs
- Offline numerical integration
- Start position can be enforced, enables smooth replanning

## ProDMPs

$$y(t) = c_1 y_1(t) + c_2 y_2(t) + \Phi(t)w$$

$$\dot{y}(t) = c_1 \dot{y}_1(t) + c_2 \dot{y}_2(t) + \dot{\Phi}(t)w$$



Source: [Li et al., 2023]

A ProMP can be learned using linear regression

- Record trajectories of successful examples (kinesthetic teaching, tele-op)
- For each trajectory  $\tau_i$ , obtain  $w_i$

$$w_i = (\Phi^T \Phi + \sigma I)^{-1} \Phi^T \tau_i$$

- Compute mean and variance

$$\mu_w = \frac{1}{N} \sum_i w_i$$

$$\Sigma_w = \frac{\sum_i (w_i - \mu_i)(w_i - \mu_i)^T}{N - 1}$$



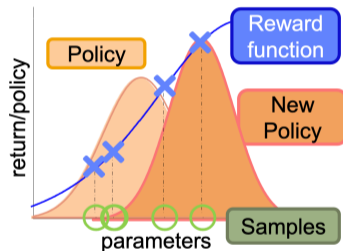
# Reinforcement Learning by Direct Policy Search

What if we don't have demonstrations?

- Find  $\pi(\mathbf{w})$  with maximum return

$$\pi^* = \arg \max_{\pi} \int \pi(\mathbf{w}) \mathcal{R}(\mathbf{w}) d\mathbf{w}$$
$$\mathcal{R}(\mathbf{w}) = \mathbb{E} \left[ \sum_{t=1}^T r_t \mid \mathbf{w} \right]$$

- High level policy  $\pi(\mathbf{w})$ 
  - Mean: Current best estimate of the maximum
  - Variance: Direction to explore
- Low level policy: Use open loop control



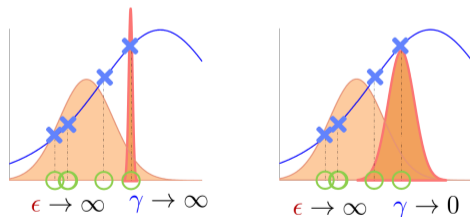
Source: [Neumann, 2017]

Towards an algorithm:

- Unconstrained optimization leads to pre-mature convergence
- Introduce constraints to tackle exploration-exploitation trade-off
  - Kullback Leibler Divergence:  
 $KL(p||q) = \int q(\mathbf{w}) \log \frac{q(\mathbf{w})}{p(\mathbf{w})} d\mathbf{w}$
  - Bound entropy loss:  
 $H(p) = - \int p(\mathbf{w}) \log p(\mathbf{w}) d\mathbf{w}$

## Information-Theoretic Policy Update

$$\arg \max_{\pi} \int \pi(\mathbf{w}) \mathcal{R}(\mathbf{w}) d\mathbf{w}$$
$$\text{s.t.} \quad KL(\pi || \pi_{old}) \leq \epsilon,$$
$$H(\pi_{old}) - H(\pi) \leq \gamma$$



Source: [Neumann, 2017]

Idea:

- Replace return function with learned model
- Compatible function approximation  
 $\mathcal{R}(\mathbf{w}) \approx \hat{\mathcal{R}}(\mathbf{w}) = -\frac{1}{2}\mathbf{w}^T\mathbf{A}\mathbf{w} + \mathbf{w}^T\mathbf{a} + a$

Iterative algorithm:

- Estimate  $\mathbf{A}$ ,  $\mathbf{a}$ ,  $a$  from sampled returns
- Solve using method of Lagrange multipliers
  - Minimize Lagrange dual (closed form)
  - Update rules in terms of optimal L. multiplier
  - Natural gradient direction

MORE [Abdolmaleki et al., 2015]

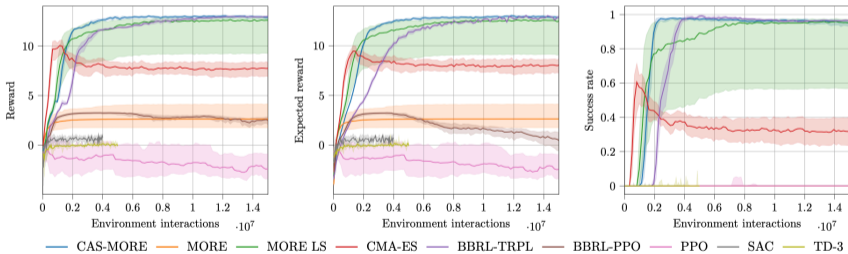
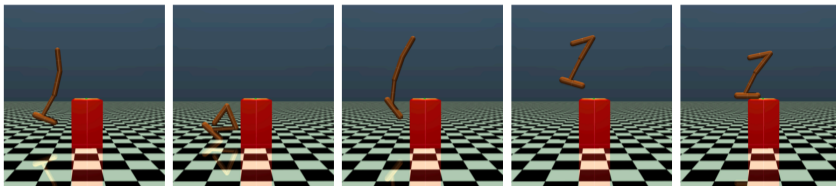
$$\begin{aligned} \arg \max_{\pi} \quad & \int \pi(\mathbf{w})\hat{\mathcal{R}}(\mathbf{w})d\mathbf{w} \\ \text{s.t.} \quad & KL(\pi\|\pi_{old}) \leq \epsilon, \\ & H(\pi) \geq \beta \end{aligned}$$

CAS-MORE

[Hüttenrauch and Neumann, 2024]

$$\begin{aligned} \arg \max_{\boldsymbol{\mu}} \quad & \int \pi(\mathbf{w})\hat{\mathcal{R}}(\mathbf{w})d\mathbf{w} \Big|_{\boldsymbol{\Sigma}=\boldsymbol{\Sigma}_{old}} \\ \text{s.t.} \quad & KL(\pi\|\pi_{old}) \Big|_{\boldsymbol{\Sigma}=\boldsymbol{\Sigma}_{old}} \leq \epsilon_{\boldsymbol{\mu}} \\ \arg \max_{\boldsymbol{\Sigma}} \quad & \int \pi(\mathbf{w})\hat{\mathcal{R}}(\mathbf{w})d\mathbf{w} \Big|_{\boldsymbol{\mu}=\boldsymbol{\mu}_{old}} \\ \text{s.t.} \quad & KL(\pi\|\pi_{old}) \Big|_{\boldsymbol{\mu}=\boldsymbol{\mu}_{old}} \leq \epsilon_{\boldsymbol{\Sigma}} \end{aligned}$$

# Model-Based Relative Entropy Stochastic Search



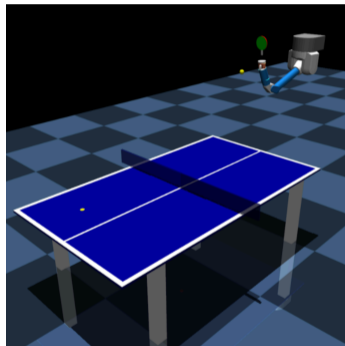
Source: [Hüttenrauch and Neumann, 2024]

# Contextual Episode-Based Policy Search

- Conditions described by context  $\mathbf{c}$
- Policy:  $\pi(\mathbf{w} \mid \mathbf{c}) = \mathcal{N}(\mathbf{w} \mid \boldsymbol{\mu}(\mathbf{c}), \boldsymbol{\Sigma}(\mathbf{c}))$
- Find  $\mathbf{w}$  that maximizes return given a context  $\mathbf{c}$

$$\pi^* = \arg \max_{\pi} \int p(\mathbf{c}) \int \pi(\mathbf{w} \mid \mathbf{c}) \mathcal{R}(\mathbf{w}, \mathbf{c}) d\mathbf{w} d\mathbf{c}$$

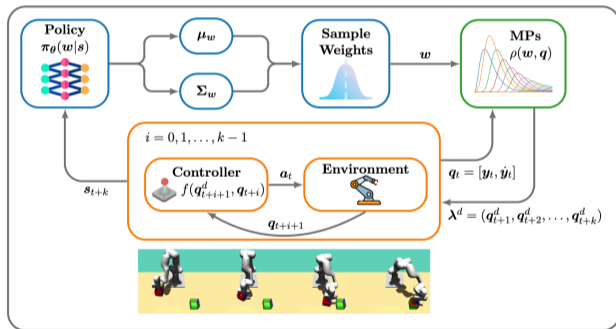
$$\mathcal{R}(\mathbf{c}, \mathbf{w}) = \mathbb{E} \left[ \sum_{t=1}^T r_t \mid \mathbf{c}, \mathbf{w} \right]$$



Source: [Otto et al., 2023]

# Deep Black-Box Reinforcement Learning [Otto et al., 2022]

- Neural network mapping from context to distribution parameters
- Differentiable trust region layers [Otto et al., 2020] for implicit constraint satisfaction

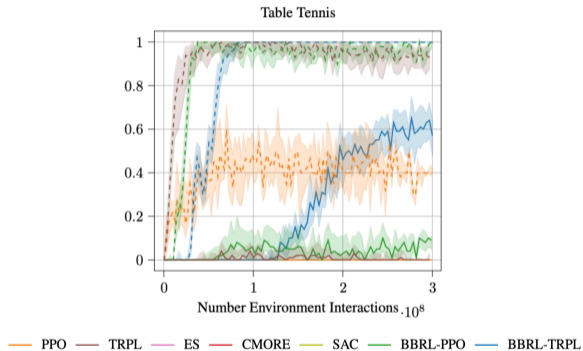


Source: [Otto et al., 2023]

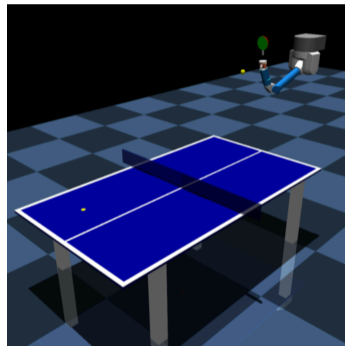
## BBRL TRPL

$$\arg \max_{\pi} \mathbb{E} \left[ \frac{\pi(w | c)}{\pi_{old}(w | c)} A^{\pi_{old}}(c, w) \right]$$
$$A^{\pi}(c, w) = \mathcal{R}(c, w) - V_{\phi}^{\pi}(c)$$

# Deep Black-Box Reinforcement Learning



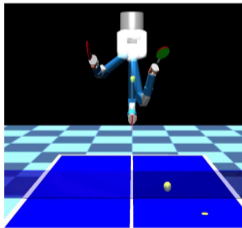
Source: [Otto et al., 2022]



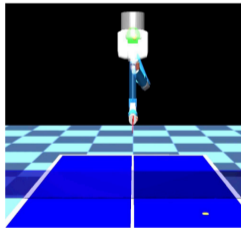
Source: [Otto et al., 2023]

# Diverse Skill Learning

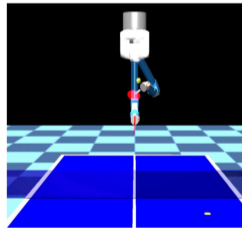
- Learn multiple solutions for a given  $\mathbf{c}$
- Mixture of Experts policy  $\pi(\mathbf{w} | \mathbf{c}) = \sum_o p(o | \mathbf{c})p(\mathbf{w} | \mathbf{c}, o)$ 
  - Gating policy  $p(o | \mathbf{c})$  assigns an expert  $o$  to a given context  $\mathbf{c}$
  - Expert  $p(\mathbf{w} | \mathbf{c}, o)$  movement primitive (Gaussian)



(a) Strike from left to right.



(b) Forehand counter strike.



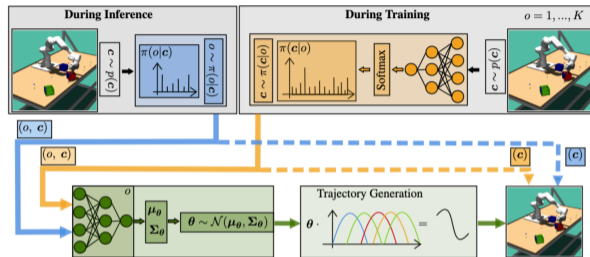
(c) Backhand backspin strike.

Source: [Celik et al., 2023]



## Curriculum learning

- Generate preferable contexts during learning
- Introduce learned context distribution  $\rho(\mathbf{c}) = \sum_o p(\mathbf{c} | o)p(o)$



Source: [Celik et al., 2023]

## Di-SkilL

$$\max_{\pi(\mathbf{w} | \mathbf{c}), \rho(\mathbf{c})} \mathbb{E}_{\rho(\mathbf{c})} [\mathbb{E}_{\pi(\mathbf{w}|\mathbf{c})} [\mathcal{R}(\mathbf{c}, \mathbf{w})] + \alpha H(\pi)] - \beta \text{KL}(\rho || p)$$

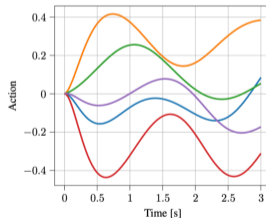
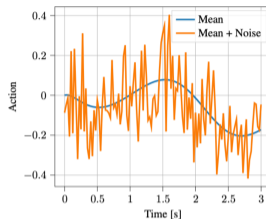
Depending on the setting different methods exist

- Non-contextual: Black-box methods (MORE [Abdolmaleki et al., 2015, Hüttenrauch and Neumann, 2024], CMA-ES [Hansen, 2023])
- Contextual: Linear (C-MORE [Tangkaratt et al., 2017]) and non-linear (BBRL TRPL [Otto et al., 2022]) mapping from context to MP
- Diverse Skill Learning: Linear MoE (LADIPS [End et al., 2017])
- Diverse Skill and Curriculum Learning: Linear MoE (SVSL [Celik et al., 2022]) and non-linear MoE (Di-Skill [Celik et al., 2023])

# Episode-Based vs Step-Based RL

- Step-based
  - Standard random exploration unsuitable for robotics
  - Each transition  $(s, a, r, s')$  is a data point
- Episode-based
  - Smooth exploration
  - Each episode  $(\tau, R)$  is a data point

Can we be more sample efficient than pure episode-based RL while keeping smooth exploration?



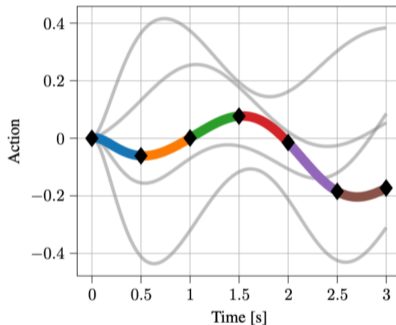
Source: [Li et al., 2024]

# Temporally Correlated Episodic RL [Li et al., 2024]

- Similar setup to previous BBRL approach
- Trajectory segments as data points
- Policy update based on segment likelihood

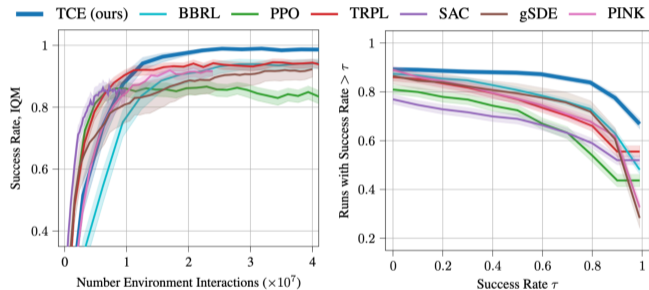
TCE

$$J(\theta) = \mathbb{E}_{\pi_{old}} \left[ \frac{1}{K} \sum_{k=1}^K \frac{\rho([\mathbf{y}_t]_{t=t_k:t'_k} | \mathbf{s})}{\rho_{old}([\mathbf{y}_t]_{t=t_k:t'_k} | \mathbf{s})} A^{\pi_{old}}(\mathbf{s}_{t_k}, [\mathbf{y}_t]_{t=t_k:t'_k}) \right]$$

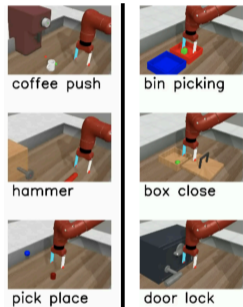


Source: [Li et al., 2024]

# Temporally Correlated Episodic RL



Source: [Li et al, 2024]



Source: [Yu et al, 2021]

- ERL better suited for robotics due to smooth exploration
- ERL good for sparse and non Markovian rewards
- Policy search objective robust to noise compared to classical BB optimizers
- Trade-off expressiveness vs complexity
- MP design may need domain knowledge and tuning

[Abdolmaleki et al., 2015] Abdolmaleki, A., Lioutikov, R., Peters, J. R., Lau, N., Pualo Reis, L., and Neumann, G. (2015).

**Model-based relative entropy stochastic search.**

In Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.

[Celik et al., 2023] Celik, O., Taranovic, A., and Neumann, G. (2023).

**Reinforcement Learning of Diverse Skills using Mixture of Deep Experts.**

In *Intrinsically-Motivated and Open-Ended Learning Workshop @NeurIPS2023*.

[Celik et al., 2022] Celik, O., Zhou, D., Li, G., Becker, P., and Neumann, G. (2022).

**Specializing Versatile Skill Libraries using Local Mixture of Experts.**

In *Proceedings of the 5th Conference on Robot Learning*, pages 1423–1433. PMLR.

[End et al., 2017] End, F., Akrou, R., Peters, J., and Neumann, G. (2017).

**Layered direct policy search for learning hierarchical skills.**

*In 2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6442–6448.

[Hansen, 2023] Hansen, N. (2023).

**The CMA Evolution Strategy: A Tutorial.**

[Hüttenrauch and Neumann, 2024] Hüttenrauch, M. and Neumann, G. (2024).

**Robust Black-Box Optimization for Stochastic Search and Episodic Reinforcement Learning.**

*JMLR (in submission).*



- [Ijspeert et al., 2013] Ijspeert, A. J., Nakanishi, J., Hoffmann, H., Pastor, P., and Schaal, S. (2013).  
**Dynamical Movement Primitives: Learning Attractor Models for Motor Behaviors.**  
*Neural Computation*, 25(2):328–373.
- [Li et al., 2023] Li, G., Jin, Z., Volpp, M., Otto, F., Lioutikov, R., and Neumann, G. (2023).  
**ProDMP: A Unified Perspective on Dynamic and Probabilistic Movement Primitives.**  
*IEEE Robotics and Automation Letters*, 8(4):2325–2332.
- [Li et al., 2024] Li, G., Zhou, H., Roth, D., Thilges, S., Otto, F., Lioutikov, R., and Neumann, G. (2024).  
**Open the Black Box: Step-based Policy Updates for Temporally-Correlated Episodic Reinforcement Learning.**  
*In The Twelfth International Conference on Learning Representations.*

[Neumann, 2017] Neumann, G. (2017).

**Seminar talk: Information-theoretic policy search methods for learning versatile, reusable skills.**

*[https://www.robots.ox.ac.uk/~seminars/seminars/Extra/2017\\_01\\_23\\_GerhardNeumann.pdf](https://www.robots.ox.ac.uk/~seminars/seminars/Extra/2017_01_23_GerhardNeumann.pdf)*

[Neumann, 2018] Neumann, G. (2018).

**Iros tutorial: Movement Primitives 2: Time-Dependent Primitives.**

*<https://www.idiap.ch/project/iros2018-tutorial/iros2018-tutorial-part3.pdf>*

[Otto et al., 2020] Otto, F., Becker, P., Ngo, V. A., Ziesche, H. C. M., and Neumann, G. (2020).

**Differentiable Trust Region Layers for Deep Reinforcement Learning.**

*In International Conference on Learning Representations.*

- [Otto et al., 2022] Otto, F., Celik, O., Zhou, H., Ziesche, H., Ngo, V. A., and Neumann, G. (2022).  
**Deep black-box reinforcement learning with movement primitives.**  
*In Conference on Robot Learning*, pages 1244–1265. PMLR.
- [Otto et al., 2023] Otto, F., Zhou, H., Celik, O., Li, G., Lioutikov, R., and Neumann, G. (2023).  
**MP3: Movement Primitive-Based (Re-)Planning Policy.**  
(arXiv:2306.12729).
- [Paraschos et al., 2013] Paraschos, A., Daniel, C., Peters, J. R., and Neumann, G. (2013).  
**Probabilistic movement primitives.**  
*Advances in neural information processing systems*, 26.

[Tangkaratt et al., 2017] Tangkaratt, V., Van Hoof, H., Parisi, S., Neumann, G., Peters, J., and Sugiyama, M. (2017).

**Policy Search with High-Dimensional Context Variables.**

*Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1).

[Yu et al., 2021] Yu, T., Quillen, D., He, Z., Julian, R., Narayan, A., Shively, H., Bellathur, A., Hausman, K., Finn, C., and Levine, S. (2021).

**Meta-World: A Benchmark and Evaluation for Multi-Task and Meta Reinforcement Learning.**

(arXiv:1910.10897).