# Learning Interpretable and Safe Control Policies: Interface Between Model-free Learning and Model-based Control
### (Code 7a85i)

**Organizers:**

- **Nathan Lawrence,** nplawrence@berkeley.edu
  University of California, Berkeley, USA
- **Ali Mesbah,** mesbah@berkeley.edu
  University of California, Berkeley, USA

**Abstract:**

Learning-based control promises to provide powerful algorithms for synthesizing control policies from data. Model-free reinforcement learning has proved to be a particularly versatile approach. At the same time, there is a critical need to explain the behavior of learning-based control strategies. One avenue is to incorporate well-understood structures into the learning pipeline, for example, through physics-informed learning, built-in safety guarantees of control policies, or control-oriented constraints, such as those based on Lyapunov arguments. This Open Invited Track serves as a forum to grapple with this tension between recent advances in learning-based control and the need for interpretable control policies. In this vein, interpretability can be critical for the success of learning-based control, as it encompasses safety issues like stability, robustness, and explainability through prior system and domain knowledge.

**Relevant IFAC Technical Committee:**

TC 1.2. Adaptive and Learning Systems
TC 2.4. Optimal Control
TC 3.2. Computational Intelligence in Control
TC 6.1. Chemical Process Control

**Description:**

The potential to learn control policies purely from data is an exciting prospect, poised to have numerous benefits in practical control settings such as improved performance and reduced operational costs. Meanwhile, real-world control systems operate on principles of safety and robustness, usually backed with guarantees based on models of the underlying dynamics.
Therefore, as interest in deploying learning-based control strategies grows, there is also increasing urgency to balance performance and data-driven machine learning methods with classical and modern control-theoretic structures.

This Open Invited Track is a call for works on learning-based control that focus on learning interpretable policies. A good reference point is the notion of a "model-based agent" [1]. Such an agent contains a physics-based model to guide its search for safe control actions, while at the same time being able to learn from its interactions with an uncertain environment. To give this notion of a model-based agent more texture, one may consider a policy based on model predictive control (MPC) in which a policy is based on optimizing a cost subject to a dynamics model and constraints. While such a structure contains rich theory on stability and robustness, a practical

problem is configuring these ingredients—cost, model, constraints—towards good performance. This opens the door to many machine learning approaches for tuning MPC policies. One fruitful direction is to incorporate reinforcement learning (RL) algorithms into the design of MPC, drawing from their shared connections to dynamic programming and value functions [2,3].

This Open Invited Track welcomes practical and theoretical works at the intersection of RL and model-based control. However, this is only one possible framework for learning interpretable control policies and we look to foster fruitful discussions on the broader theme of model-based agents. Possible topics for the session include, but are not limited to:
- Computational aspects of learning MPC policies in RL settings
- Theoretical work on the interface between RL and model-based policies
- Model-free or derivative-free optimization (such as Bayesian optimization) for controller tuning
- Safety aspects of online learning of interpretable policies
- Design of RL algorithms or policy architectures with stability, robustness, or other control-theoretic guarantees, such as those based on Lyapunov methods or IQCs
- Real-world case studies of model-based agents

**About the Proposers:**

Ali Mesbah is Associate Professor of Chemical and Biomolecular Engineering at the University of California at Berkeley. He is a Senior Member of the IEEE and AIChE. His research interests include learning-based analysis and control of uncertain systems, with applications to materials processing and manufacturing systems. He has organized over 20 invited sessions and pre-conference workshops at major IFAC and IEEE conferences.

Nathan Lawrence is a Postdoctoral Scholar at the University of California at Berkeley. His research interests include the design of learning-based algorithms and architectures based on reinforcement learning and model predictive control. He has experience organizing workshops at AdCONIP 2022 and delivering tutorial sessions at Upper Bound 2024 and ACC 2025 on reinforcement learning and control theory.

**References:**
1. Banker T, Lawrence NP, Mesbah A, "Local-Global Learning of Interpretable Control Policies: The Interface between MPC and Reinforcement Learning," American Control Conference, Denver, 2025.
2. Lawrence NP, Loewen PD, Forbes MG, Gopaluni RB, Mesbah A, "A view on learning robust goal-conditioned value functions: Interplay between RL and MPC," arXiv preprint arXiv:2502.06996, 2025.
3. Lawrence NP, Banker T, Mesbah A, "MPCritic: A plug-and-play MPC architecture for reinforcement learning," Conference on Decision and Control, Rio de Janeiro, 2025. (forthcoming)