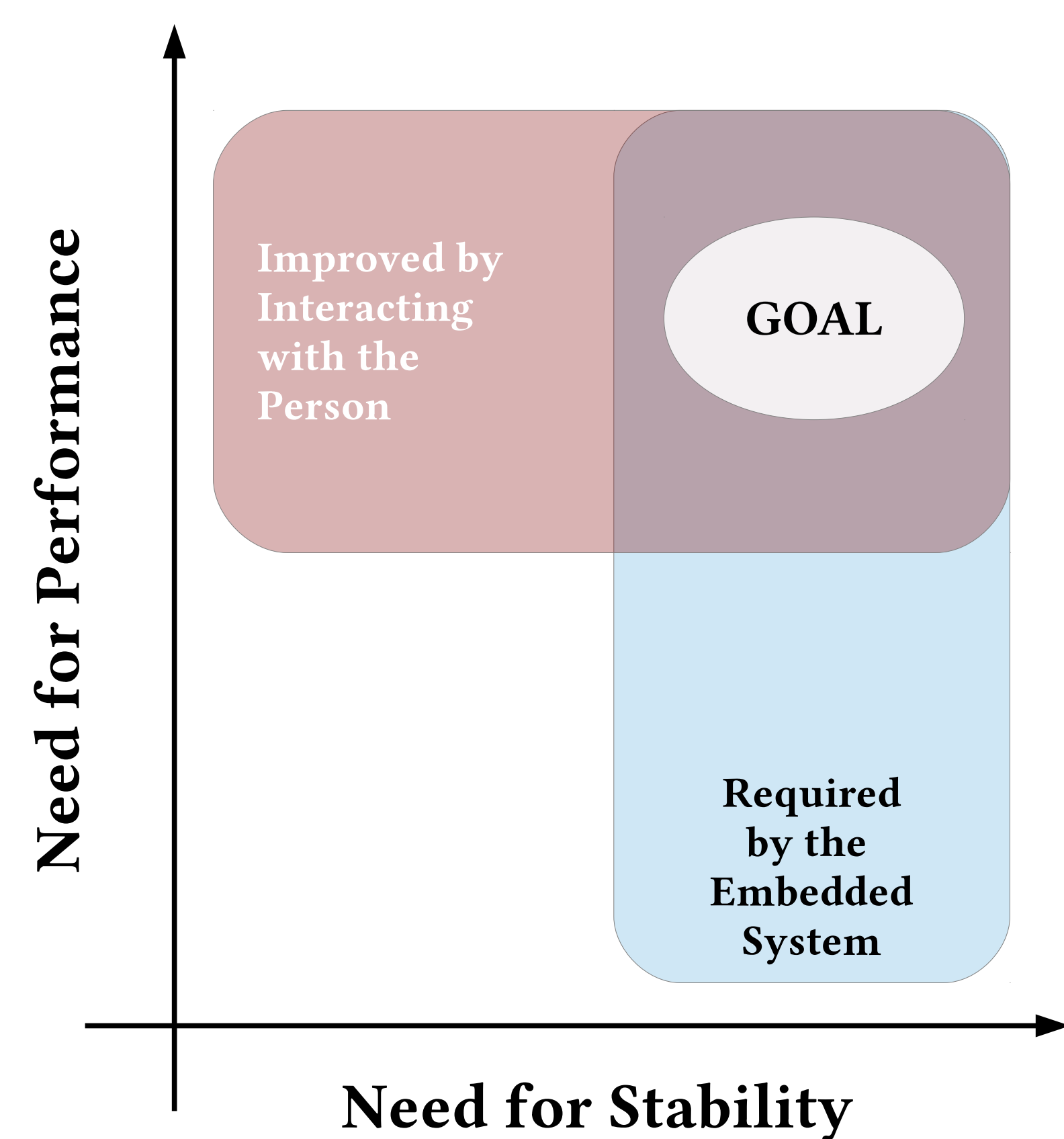# Mutually Stabilized Correction in Physical Demonstration

**PI (lead): Todd Murphey,** *Northwestern University*
**Co-PI: Brenna Argall,** *Northwestern University*
**PI: Magnus Egerstedt,** *Georgia Institute of Technology*

## How much should a person be allowed to interact with a controlled machine?



**Figure 1:** Our goal is to algorithmically resolve the tension between the need for stability and the need for performance.

**Aim:** Balancing the ability of a person to direct a cyber-physical system, against the system's representation of its own capabilities and limitations.

**Why?** Physical interactions with cyber-physical systems need to be understood in terms of **shared autonomy**, where the embedded software and the human together have to interface directly with the system dynamics.

**How?** Develop a **science of trust**, that bridges human operator capabilities and physical system safety.

For cyber-physical systems, an understanding of *each by the other* is of crucial importance. The *human operator* needs to understand the automated system order to provide good shaping guidance and sound control input. The *automated system* needs to understand the quality limitations of the guidance and controls provided by the operator.

## Mutually Controlled Motion

**Idea:** Derive control behaviors via **optimal control**, while....

Engaging the human operator for **corrective** demonstrations via **physical guidance**.

**Challenge:** The operator may destabilize the system.

This risk changes from operator to operator.



**Figure 2:** Trust (T) is based on both the Cyber and Physical components of the system, and informs how both the computer (C) and the operator (O) control the machine (M).

## A Science of Trust

**Solution:** **Verify** derived controllers for **stability** and **robustness**.

**Compute** a formal **measure of trust** in the operator.

This trust measure decides **how much control to cede** to the operator during physical correction.

**Result:** A computable notion of trust → The system assesses the safety of the instruction.
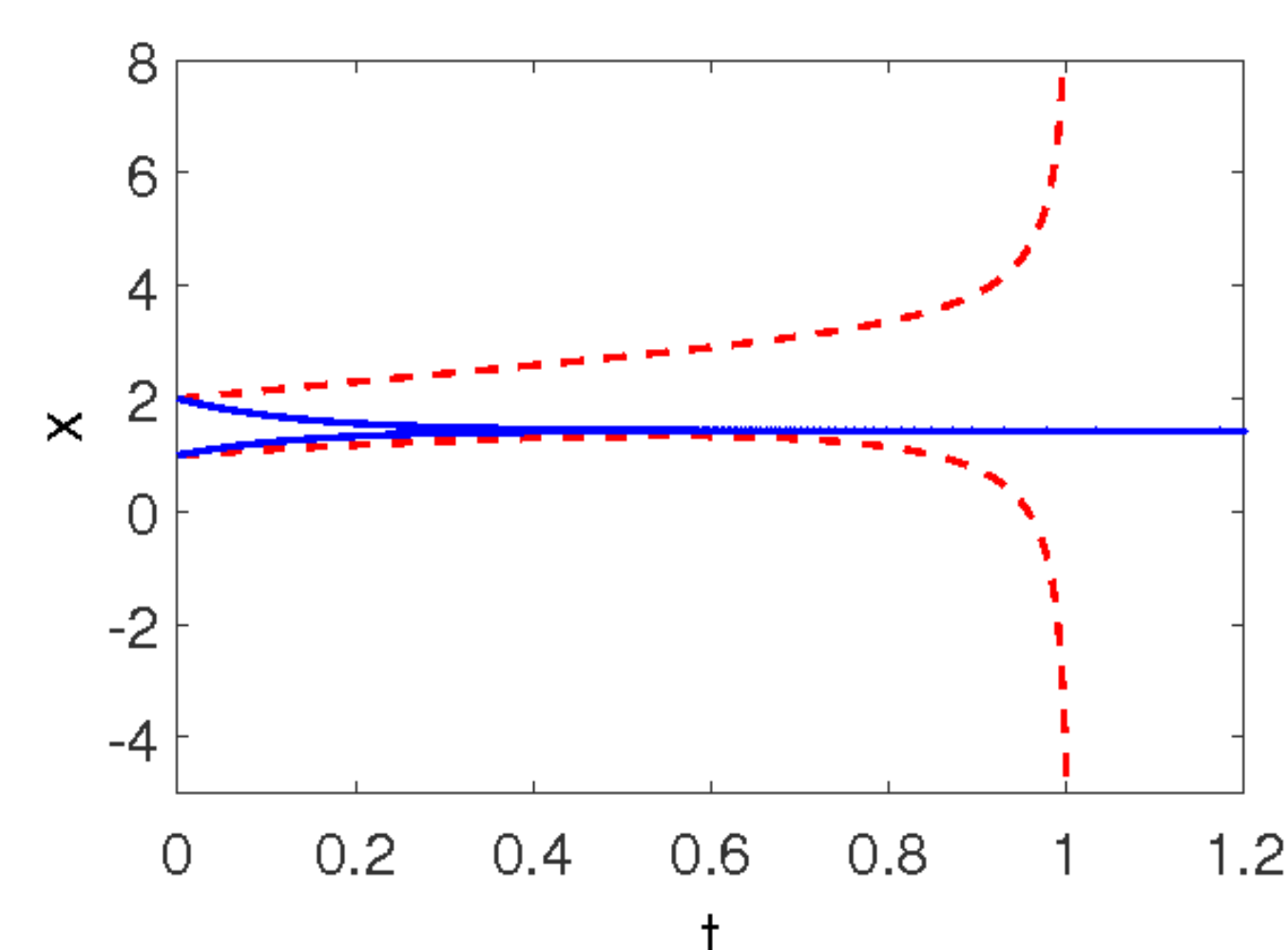


**Figure 3:** Trust-based interactions in multi-agent teams. Trust evolves according to how much an agent's neighbor is assisting in achieving their inter-agent goal. The state dynamics and trust dynamics are coupled, such that the trust values weight the standard gradient-descent based state update laws. Plots depict a two-agent rendezvous task. The solid trajectories show the agents' states when the sum of the initial trust values is positive, causing the two agents to reach an agreement asymptotically. The dashed trajectories correspond to a negative initial trust sum, resulting in diverging states in finite time.
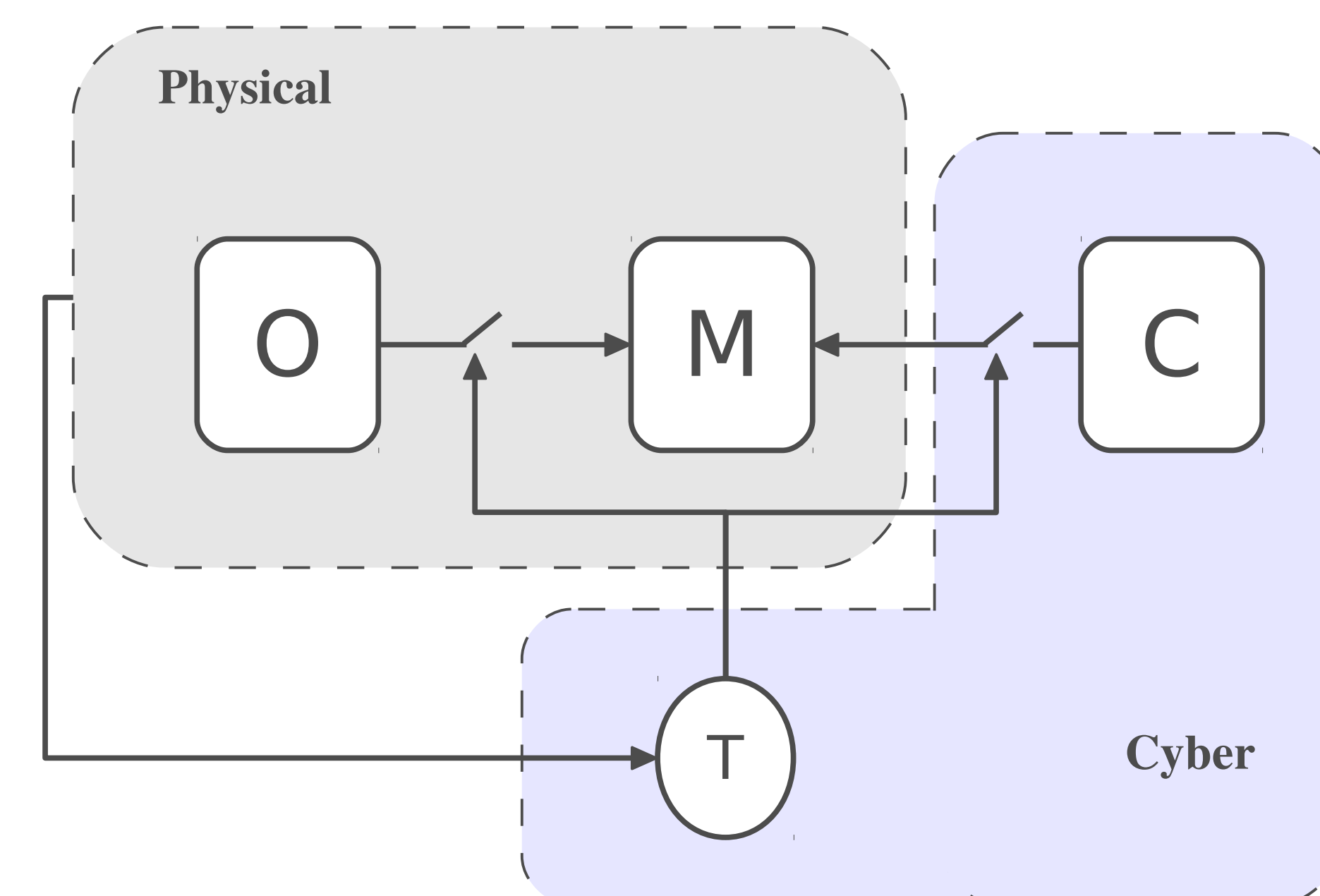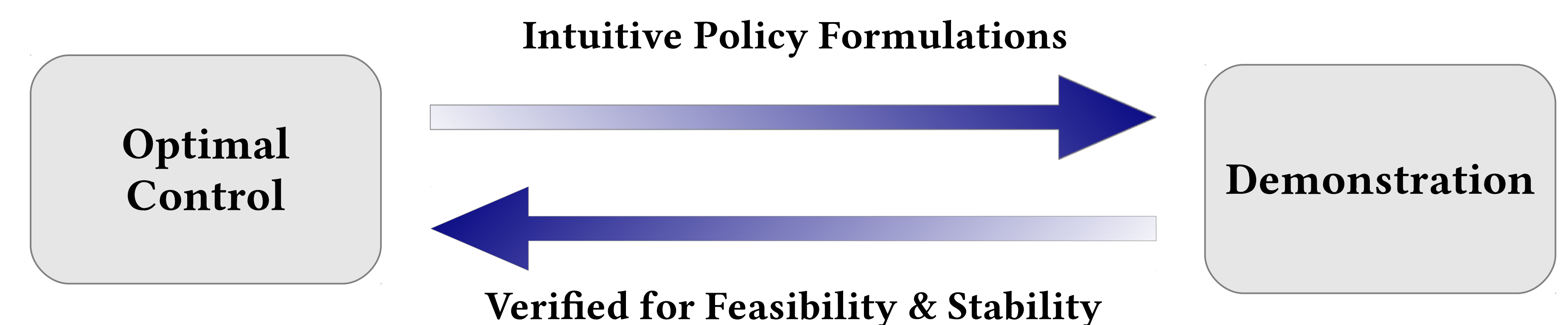


**Figure 4:** Respective benefits of two procedures for deriving controllers: Optimal Control and Learning from Demonstration.

### Assessments

**Platforms:** Underactuated articulated rigid bodies → easily destabilized.

**Goal:** By the end of this project, an operator will be able to physically manipulate any of the rigid body configurations of the experimental testbed to produce a desired motion, while getting feedback from the automatic control system about the stability of that motion.
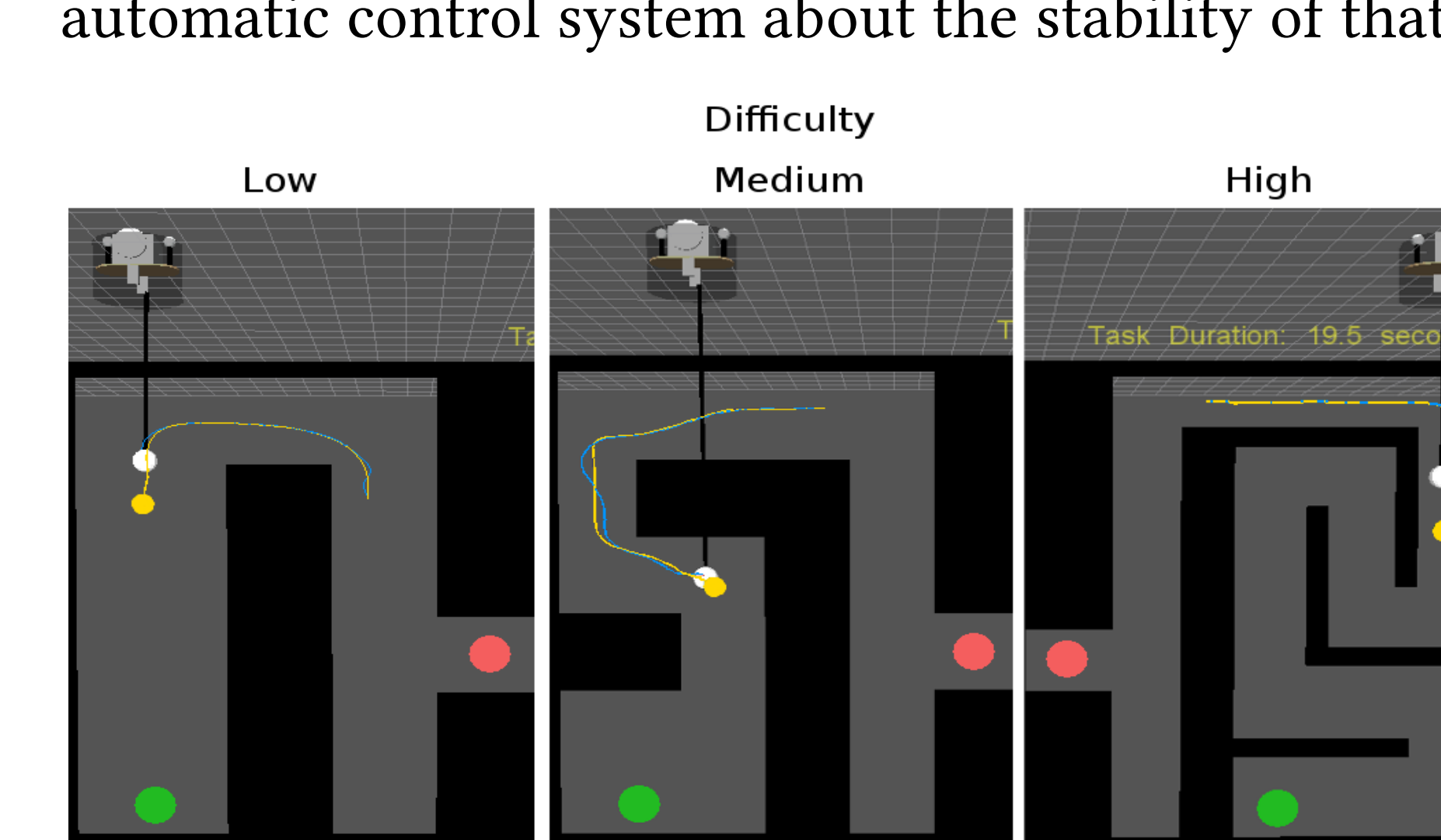


**Figure 5:** The simulated crane system (above) can be controlled in real-time, with an operator in the loop, and the control system in the Robot Operating System (ROS) automatically updates a trust measure and modulates the amount of control authority provided to the operator. This system was used in a subject study with 22 participants to evaluate adaptive versus static methods for computing the trust metric (results in Figure 6).
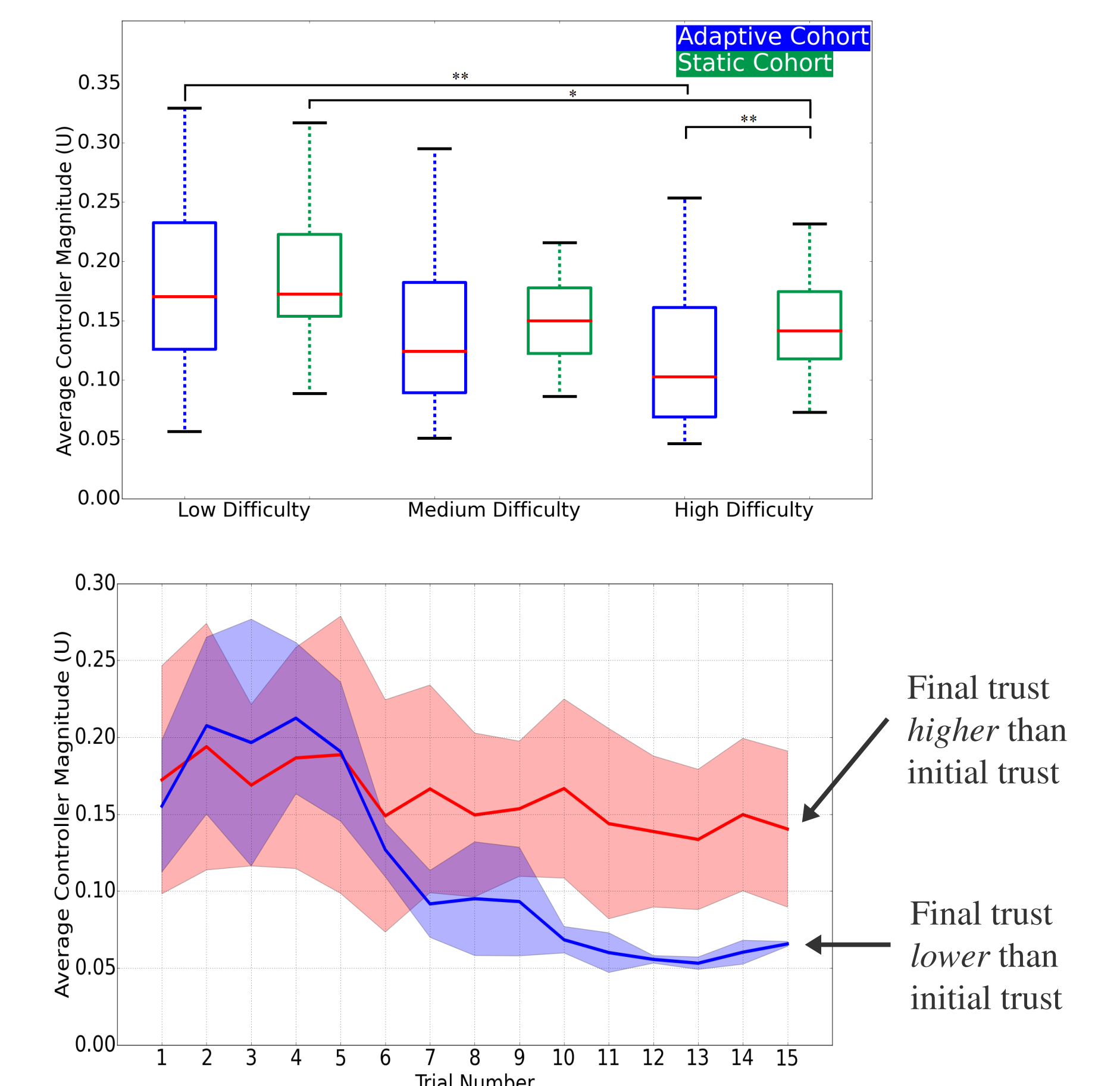


**Figure 6:** *Top:* Average controller magnitude for the static (green) and adaptive (blue) trust cohorts. We see that the adaptive cohort requires a significantly ($p < 0.01$) diminished average controller magnitude than the static cohort in the final maze configuration. Key : * $p < 0.05$ and ** $p < 0.01$. *Bottom:* Evolution of the average controller magnitude per trial. We see a significant decrease in the required average controller magnitude both in users whose final trust value was lower ($p < 0.01$) and higher ($p < 0.05$) than the initial estimate, demonstrating that the results hold regardless of whether the initial control authority allocation is an over- or under-estimate of the user's expertise

### Relevance to Cyber-Physical Systems

**Impact:** Cyber-physical systems for which (i) control authority is shared between the human and machine, (ii) the machine automation is adaptable by and able to receive instruction from a human who is not an automation expert, (iii) there are physical, possibly destabilizing, interactions between the human and machine.

**Domains:** Immediately impacted: Rehabilitation, assistive devices, and human augmentation. Near-term impact: Manufacturing, which will soon involve skilled workers working side-by-side with robots and teaching robots tasks. More broadly: Non-mechanical but highly interconnected systems, such as air traffic control and power grid management. Such systems are often too complex to understand completely, yet the operator still must provide instruction that is feasible for the system to reliably execute.

**Relationship to CPS Needs:** *Interaction and potential interference among CPS and humans*, by explicitly reasoning about when to cede control authority to a human operator, and when to request instruction for stability assistance. *Cross-disciplinary collaborative research*, by building a synergy between the areas of data-driven machine learning and formal control theory. *Jointly modeling the interaction of both cyber and physical components*, by taking steps to quantify the level of understanding needed by the human to provide effective corrections, and by explicitly computing the system's understanding of the consequences of physical or interaction during instruction. *Incorporating CPS science into education*, by incorporating CPS-centric coverage in the Control of Mobile Robotics MOOC taught by co-PI Egerstedt.