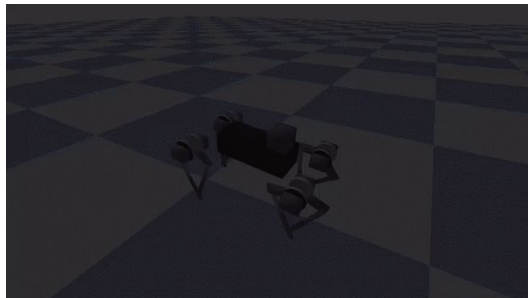# Sim-to-Real
## An Incomplete Overview

Jie Tan
CS 8803 Deep Reinforcement Learning for Intelligent Control
04/04/2022

# What's the sim-to-real gap?

Dynamics:



Perception:

# Goals

- Understand the causes of sim-to-real
- Review of the state-of-the-art methods
  - System Identification
  - Domain Randomization
  - Domain Adaptation
  - Meta Learning

# Brainstorming: Why Sim-to-Real?

### Real world

- Slow
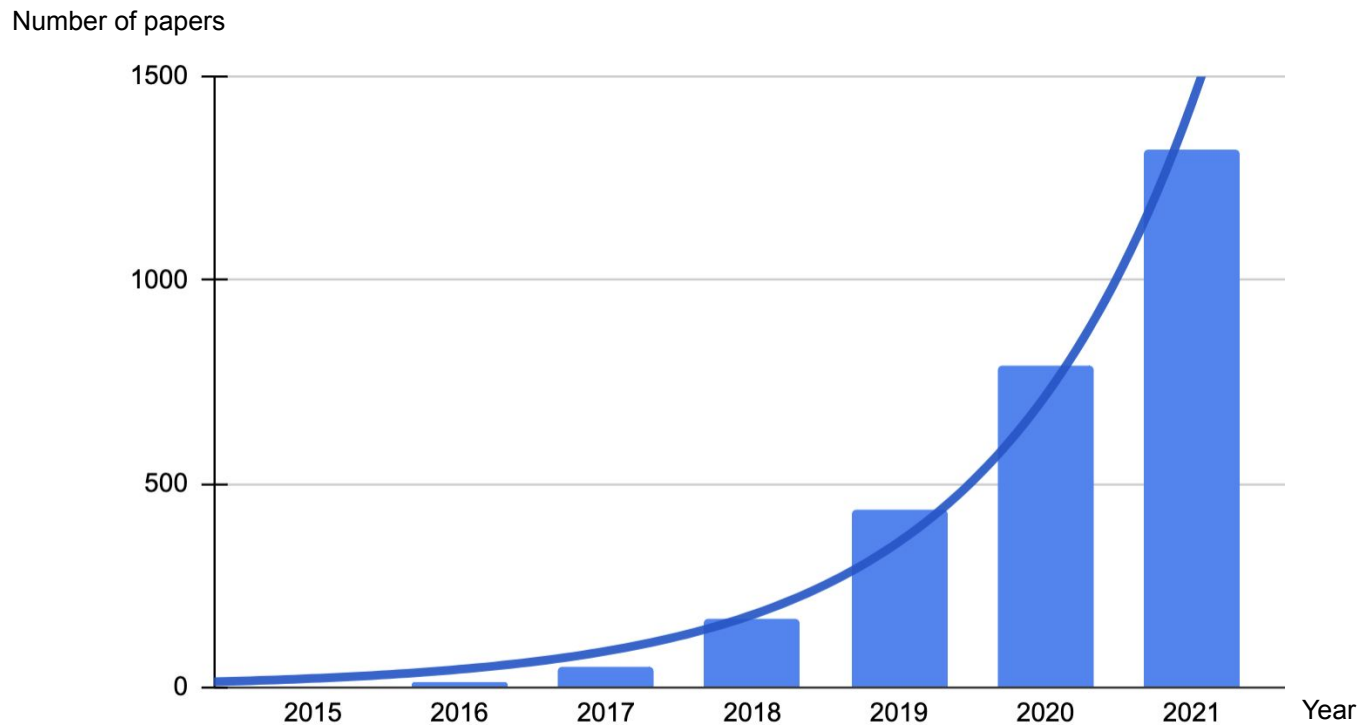- Unsafe
- Expensive
- Human supervision

### Simulation

- Fast
- Safe
- Cheap
- Scalable

# Why Sim-to-Real?

Transform hard robotic problems into large-scale computation problems

# Trend on Sim-to-Real

Number of papers

# How to overcome sim-to-real gap?

- Improve simulation
  - System identification
    - [Sim-to-Real: Learning Agile Locomotion For Quadruped Robots](#)
    - [Simulation-Based Design of Dynamic Controllers for Humanoid Balancing](#)
    - [Preparing for the Unknown: Learning a Universal Policy with Online System Identification](#)
- Improve policy
  - Domain randomization
    - [Sim-to-Real Transfer of Robotic Control with Dynamics Randomization](#)
    - [Closing the Sim-to-Real Loop: Adapting Simulation Randomization with Real World Experience](#)
  - Domain adaptation
    - [Learning Agile Robotic Locomotion Skills by Imitating Animals](#)
    - [Sim-to-Real Transfer for Biped Locomotion](#)
  - Meta learning
    - [Rapidly Adaptable Legged Robots via Evolutionary Meta-Learning](#)
    - [NoRML: No Reward Meta-Learning](#)
    - [Policy Transfer with Strategy Optimization](#)
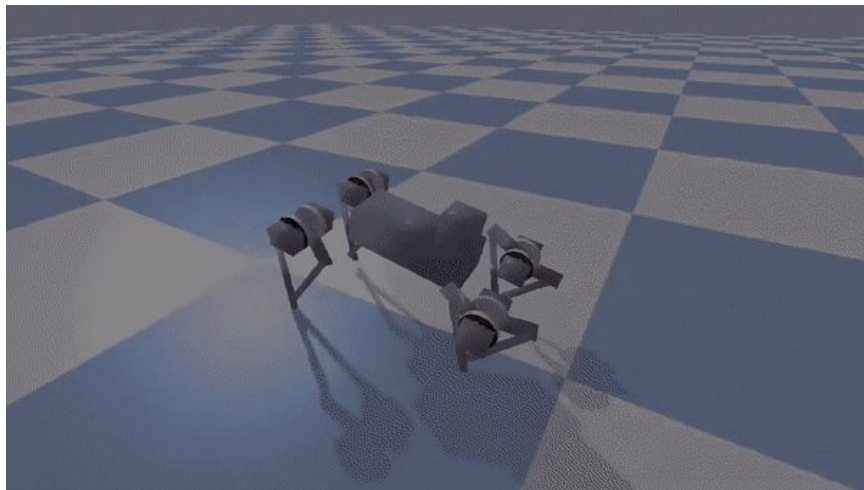
# System Identification

# What are the causes of sim-to-real gap?

- Unmodeled dynamics
- Wrong simulation parameters
- Inaccurate contact models
- Latency
- Actuator dynamics
- Noise
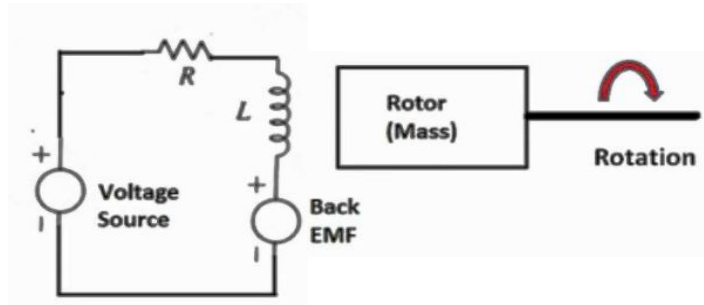- Stochastic real environment
- Numerical accuracy
- ...

**Actuator dynamics** and **latency** are two important causes of reality gap.



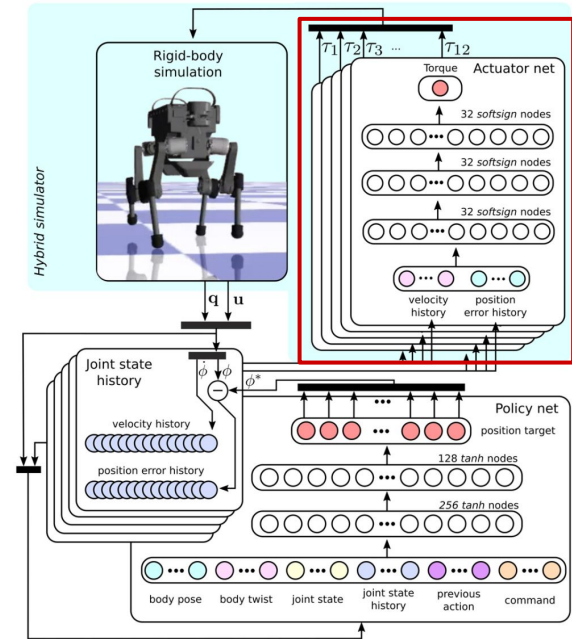[Sim-to-Real: Learning Agile Locomotion For Quadruped Robots, RSS 2018]
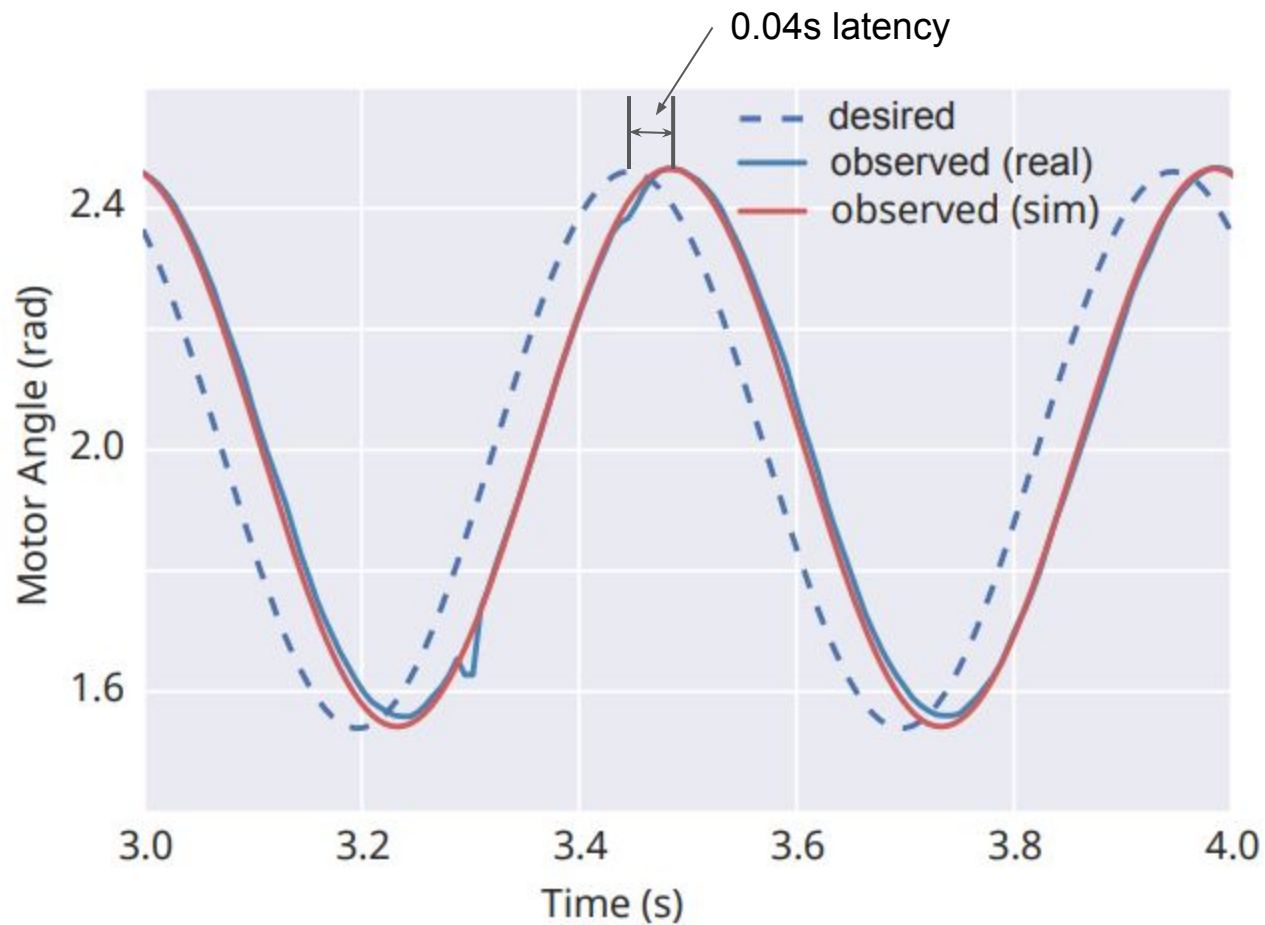
# Actuator Model



## Analytical models

$$\tau = f(I)$$
$$I = \frac{V * \mathrm{PWM} - V_{\mathrm{emf}}}{R}$$
$$V_{\mathrm{emf}} = K_t \dot{q}$$

[Sim-to-Real: Learning Agile Locomotion For Quadruped Robots, RSS 2018]

## Neural network models

[Learning agile and dynamic motor skills for legged robots, Science Robotics 2019]

0.04s latency

desired
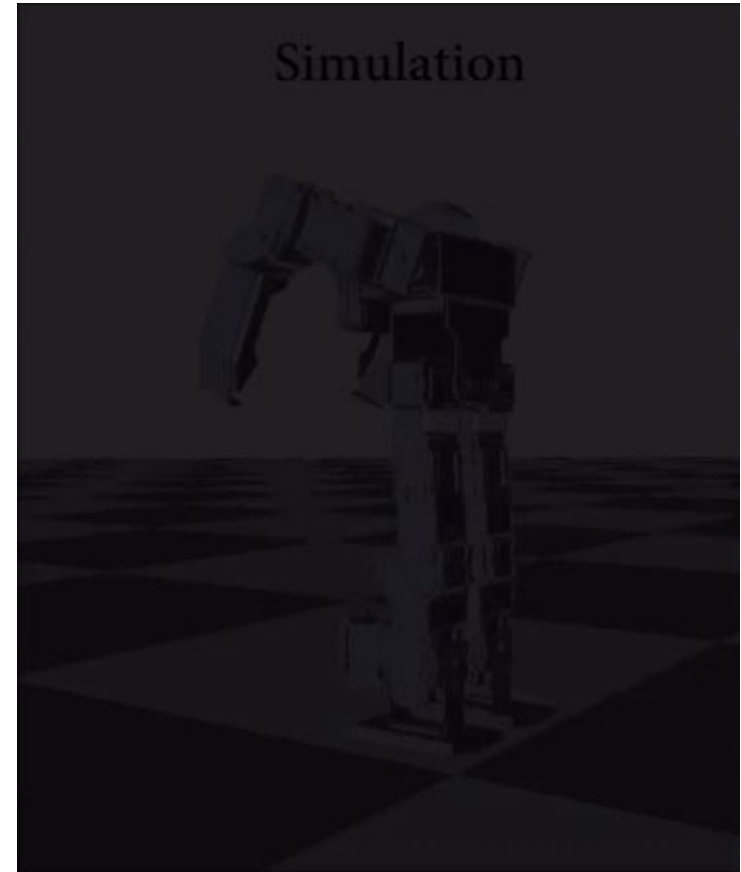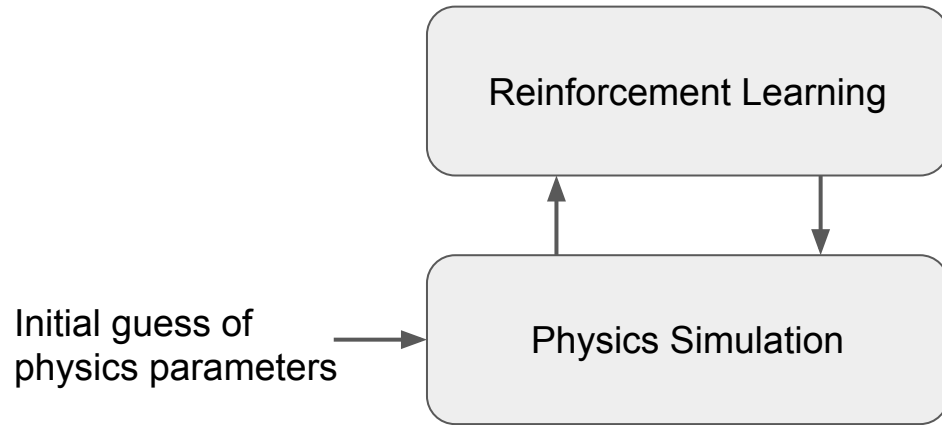observed (real)
observed (sim)

Motor Angle (rad)

Time (s)

[Sim-to-Real: Learning Agile Locomotion For Quadruped Robots, RSS 2018]

- Limitations
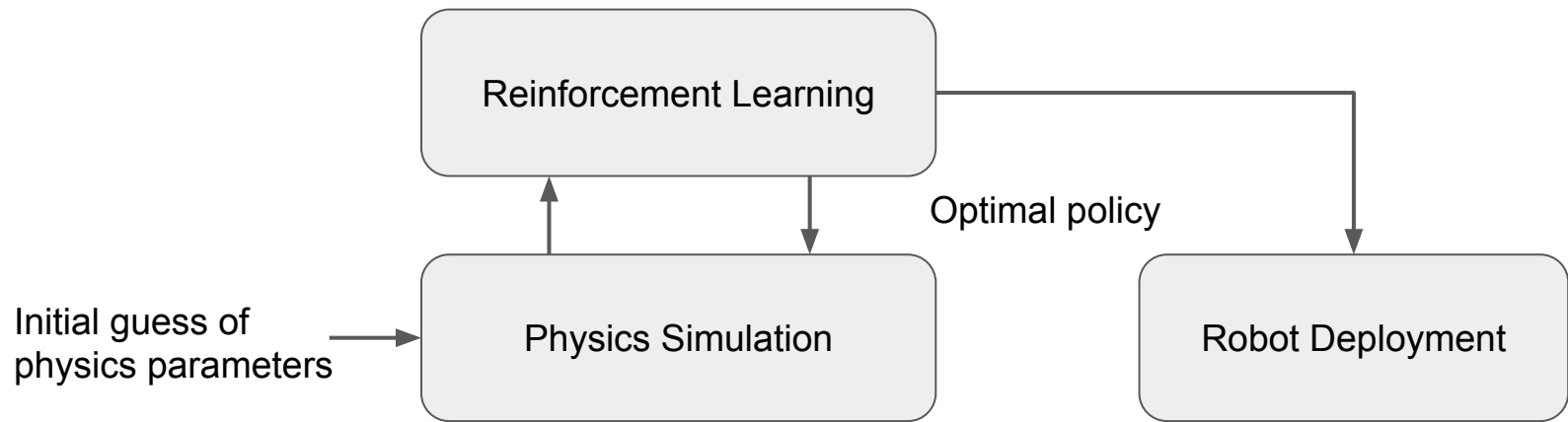  - Disassemble the robot
  - Decide what parameters to identify
  - Design experiments for individual parameters
  - Lots of manual work

Reinforcement Learning

Physics Simulation

Initial guess of physics parameters

Simulation

[Simulation-based design of dynamic controllers for humanoid balancing, IROS 2016]

[Simulation-based design of dynamic controllers for humanoid balancing, IROS 2016]

Simulation     Robot
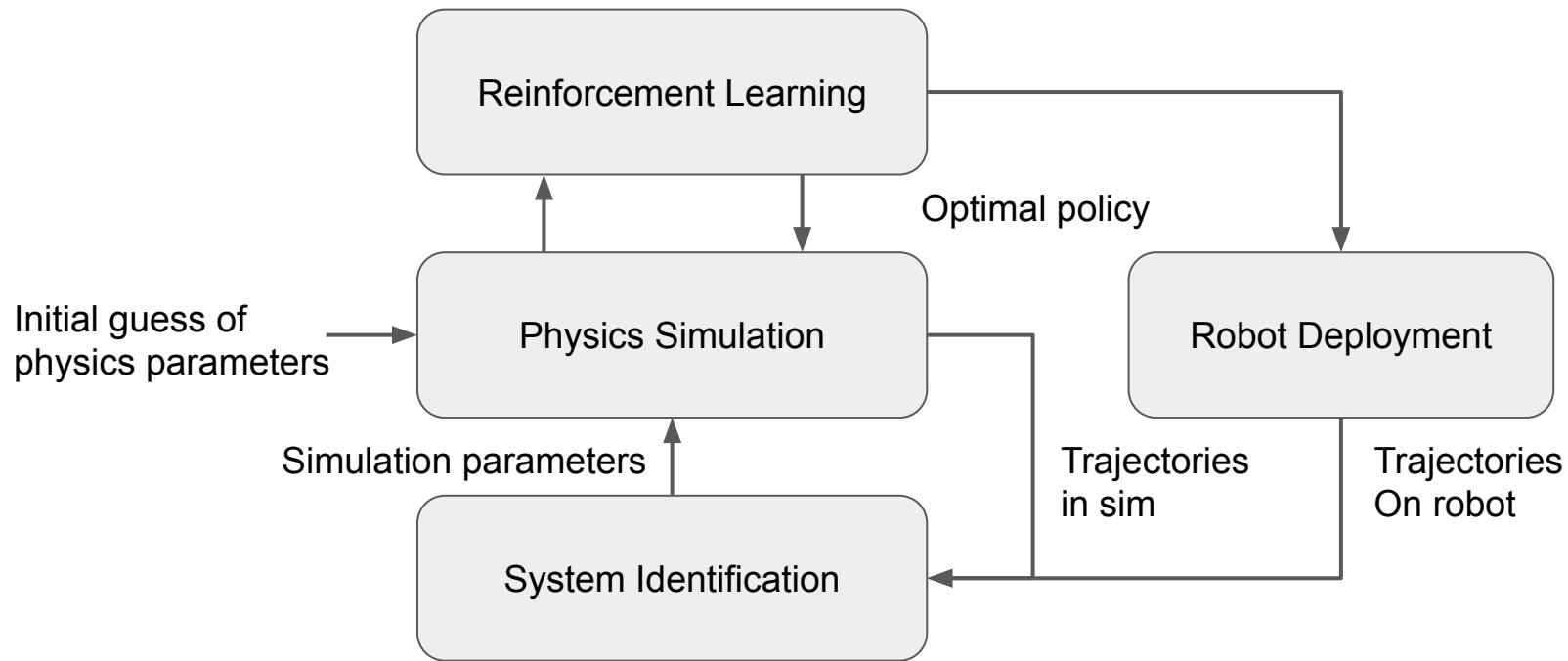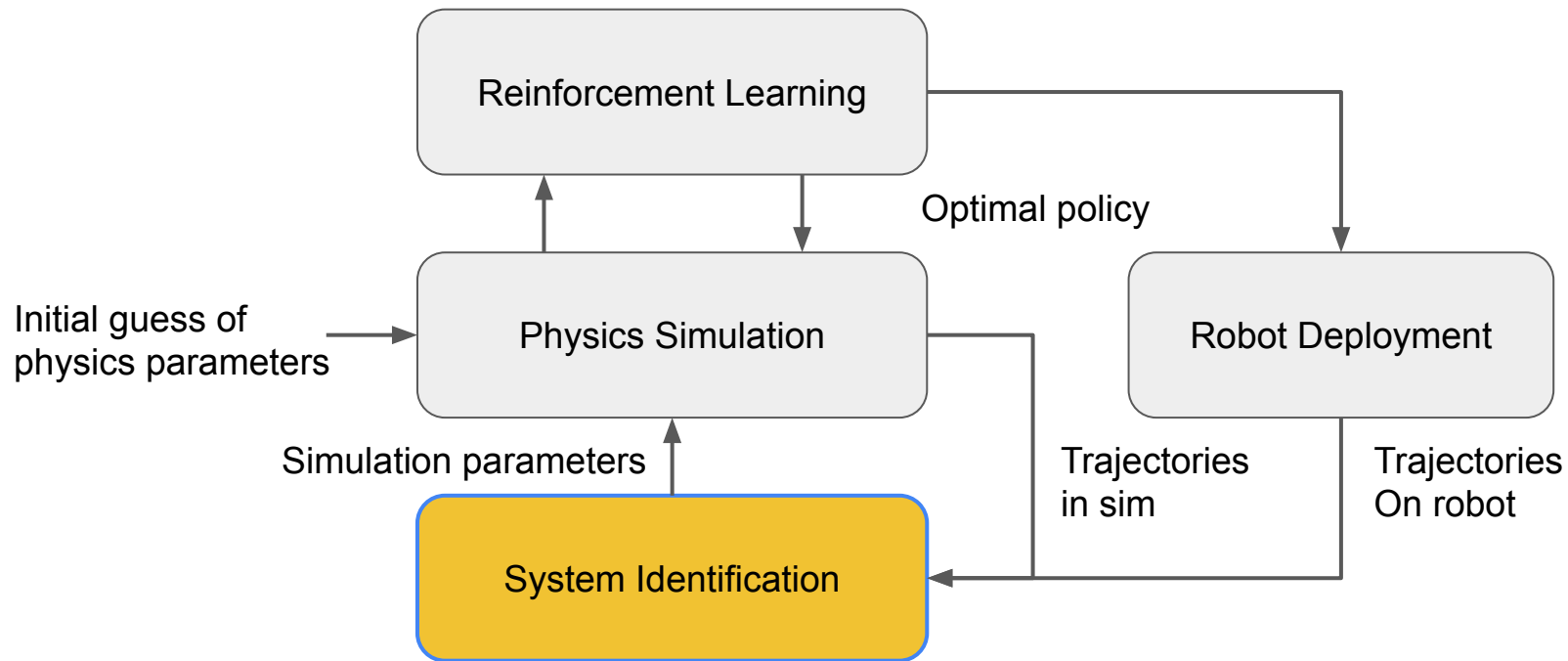
[Simulation-based design of dynamic controllers for humanoid balancing, IROS 2016]

[Simulation-based design of dynamic controllers for humanoid balancing, IROS 2016]

[Simulation-based design of dynamic controllers for humanoid balancing, IROS 2016]

# Automatic System Identification

- Measure sim-to-real discrepancy

$$\boldsymbol{\theta} = \arg\min \frac{1}{n} \sum_{i=1}^{n} \int_{0}^{T+1} ||\tilde{\mathbf{q}}_i(t) - \mathbf{q}_i(t; \boldsymbol{\theta})||_{\mathbf{W}}^2 \, dt$$

- Optimize the physics parameters
  - Covariance Matrix Adaptation-Evolution Strategy

Latency

**Ground truth physical
parameter**:
Latency = 5ms
Actuator strength = 10nm

Actuator
strength

**Latency**

**Randomly sampled physical parameter**:
Latency = 1ms
Actuator strength = 2Nm

**Sim trajectory**:

state

time

**Real trajectory**:

state

time

**Loss**: $\frac{1}{n}\sum_{i=1}^{n}\int_{0}^{T+1}||\tilde{\mathbf{q}}_i(t) - \mathbf{q}_i(t;\boldsymbol{\theta})||_{\mathbf{W}}^2 \mathrm{d}t$

Actuator strength

Robot

[Simulation-based design of dynamic controllers for humanoid balancing, IROS 2016]

# Automatic System Identification



- Limitations
  - Manual selection of physical parameters needed
  - Do not work if sim and real trajectory diverge too quickly
  - Not account for unmodeled dynamics
  - Physical parameters overfit

# Domain Randomization

# Domain Randomization

- Original objective: reward maximization

$$\mathbb{E}_{\tau \sim p(\tau|\pi)} \left[ \sum_{t=0}^{T-1} r(s_t, a_t) \right]$$



[Sim-to-Real Transfer of Robotic Control with Dynamics Randomization, ICRA 2018]

# Domain Randomization

- Original objective: reward maximization

$$\mathbb{E}_{\tau \sim p(\tau | \pi)} \left[ \sum_{t=0}^{T-1} r(s_t, a_t) \right]$$



- New objective with domain randomization

$$\mathbb{E}_{\mu \sim \rho_\mu} \left[ \mathbb{E}_{\tau \sim p(\tau | \pi, \mu)} \left[ \sum_{t=0}^{T-1} r(s_t, a_t) \right] \right]$$

Physical parameters

| Parameter | Range |
|---|---|
| Link Mass | $[0.25, 4] \times$ default mass of each link |
| Joint Damping | $[0.2, 20] \times$ default damping of each joint |
| Puck Mass | $[0.1, 0.4] kg$ |
| Puck Friction | $[0.1, 5]$ |
| Puck Damping | $[0.01, 0.2] Ns/m$ |
| Table Height | $[0.73, 0.77] m$ |
| Controller Gains | $[0.5, 2] \times$ default gains |
| Action Timestep $\lambda$ | $[125, 1000] s^{-1}$ |

[Sim-to-Real Transfer of Robotic Control with Dynamics Randomization, ICRA 2018]

our method

no randomization during training

[Sim-to-Real Transfer of Robotic Control with Dynamics Randomization, ICRA 2018]

# Memory (LSTM) in sim-to-real



[Sim-to-Real Transfer of Robotic Control with Dynamics Randomization, ICRA 2018]

# Memory (LSTM) in sim-to-real



[Sim-to-Real Transfer of Robotic Control with Dynamics Randomization, ICRA 2018]

our method

feedforward policy (no LSTM)

- Limitations
  - Trade optimality for robustness
  - Careful tuning needed for the range of randomization

[Sim-to-Real Transfer of Robotic Control with Dynamics Randomization, ICRA 2018]

Can we use a small amount of real data to mitigate these limitations?

- Limitations
    - Trade optimality for robustness
    - Careful tuning needed for the range of randomization

[Closing the Sim-to-Real Loop: Adapting Simulation Randomization with Real World Experience, ICRA 2019]

Swing-peg-in-hole
Simulated environment

# Domain Adaptation

# Domain Adaptation



[Sim-to-Real Transfer for Biped Locomotion, IROS 2019]

[Learning Agile Robotic Locomotion Skills by Imitating Animals, RSS 2020]

# Domain Adaptation

$$\mu = \begin{bmatrix} \text{Mass} \\ \text{Inertia} \\ \text{Motor Strength} \\ \text{Motor Friction} \\ \text{Latency} \\ \text{Lateral Friction} \\ \text{Etc...} \end{bmatrix}$$

[Sim-to-Real Transfer for Biped Locomotion, IROS 2019]

[Learning Agile Robotic Locomotion Skills by Imitating Animals, RSS 2020]

# Domain Adaptation



$\mu$ ⟹ Encoder ⟹ $\mathbf{Z}$

[Sim-to-Real Transfer for Biped Locomotion, IROS 2019]

[Learning Agile Robotic Locomotion Skills by Imitating Animals, RSS 2020]

# Domain Adaptation



[Sim-to-Real Transfer for Biped Locomotion, IROS 2019]

[Learning Agile Robotic Locomotion Skills by Imitating Animals, RSS 2020]

# Domain Adaptation



[Sim-to-Real Transfer for Biped Locomotion, IROS 2019]

[Learning Agile Robotic Locomotion Skills by Imitating Animals, RSS 2020]

# Domain Adaptation



$$z^* = \arg \max_{z} \; \mathbb{E}_{\tau \sim p^*(\tau | \pi, z)} \left[ \sum_{t=0}^{T-1} \gamma^t r_t \right]$$

[Sim-to-Real Transfer for Biped Locomotion, IROS 2019]

[Learning Agile Robotic Locomotion Skills by Imitating Animals, RSS 2020]

# Domain Adaptation vs. Domain Randomization

Dog Pace



No Randomization      Randomization      Domain Adaptation (Ours)

# Domain Adaptation vs. Domain Randomization

Dog Spin



| No Randomization | Randomization | Domain Adaptation (Ours) |

- Limitations
  - Policy is not updated
  - The latent space may not contain the optimal vector for the real world
  - Performance not necessarily improve with more real data

# Meta Learning

# MAML for Sim-to-Real Gap

Learn mostly in simulation; quickly adapt to the real world with few real rollouts.

# MAML for Sim-to-Real Gap

Learn mostly in simulation; quickly adapt to the real world with few real rollouts.

- Real robot: mass = unknown

$\boldsymbol{\theta}_{\text{real}}$

# MAML for Sim-to-Real Gap

Learn mostly in simulation; quickly adapt to the real world with few real rollouts.

- Real robot: mass = unknown
- Sim task 1: mass = 7.0kg
- Sim task 2: mass = 9.2kg
- Sim task 3: mass = 9.8kg

$\boldsymbol{\theta}_1$

$\boldsymbol{\theta}_2$

$\boldsymbol{\theta}_3$

$\boldsymbol{\theta}_{\text{real}}$

# MAML for Sim-to-Real Gap

Learn mostly in simulation; quickly adapt to the real world with few real rollouts.

- Real robot: mass = unknown
- Sim task 1: mass = 7.0kg
- Sim task 2: mass = 9.2kg
- Sim task 3: mass = 9.8kg

$\boldsymbol{\theta}$ ●

$\boldsymbol{\theta}_1$ ●

$\boldsymbol{\theta}_2$ ●

$\boldsymbol{\theta}_3$ ●

$\boldsymbol{\theta}_{real}$ ●

# MAML for Sim-to-Real Gap

Learn mostly in simulation; quickly adapt to the real world with few real rollouts.

- Real robot: mass = unknown
- Sim task 1: mass = 7.0kg
- Sim task 2: mass = 9.2kg
- Sim task 3: mass = 9.8kg

$\boldsymbol{\theta}$ •  ----  meta-training

$\boldsymbol{\theta}_{meta}$ •

$\boldsymbol{\theta}_1$ •

$\boldsymbol{\theta}_2$ •

$\boldsymbol{\theta}_3$ •

$\boldsymbol{\theta}_{real}$ •

# MAML for Sim-to-Real Gap

Learn mostly in simulation; quickly adapt to the real world with few real rollouts.

- Real robot: mass = unknown
- Sim task 1: mass = 7.0kg
- Sim task 2: mass = 9.2kg
- Sim task 3: mass = 9.8kg

$\theta$

$\theta_{meta}$

$-\nabla \mathcal{L}_{\mathcal{T}_1}$

$-\nabla \mathcal{L}_{\mathcal{T}_2}$

fine-tuning

$-\nabla \mathcal{L}_{\mathcal{T}_3}$

$\theta_1$

$\theta_2$

$\theta_3$

$\theta_{rea}$

l

# MAML for Sim-to-Real Gap

Learn mostly in simulation; quickly adapt to the real world with few real rollouts.

- Real robot: mass = unknown
- Sim task 1: mass = 7.0kg
- Sim task 2: mass = 9.2kg
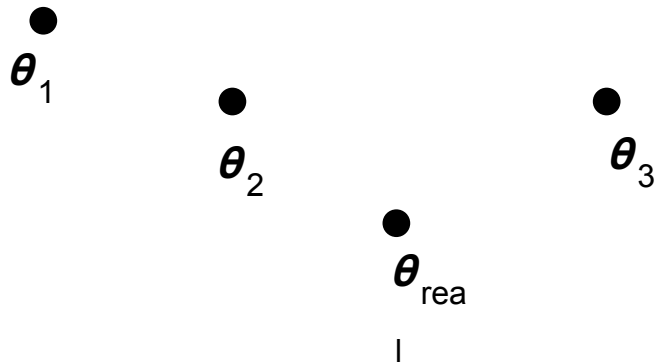- Sim task 3: mass = 9.8kg

$\boldsymbol{\theta}_{meta}$

fine-tuning

$\boldsymbol{\theta}_1$

$\boldsymbol{\theta}_2$

$\boldsymbol{\theta}_3$

$\boldsymbol{\theta}_{rea}$

l

# MAML for Sim-to-Real Gap

- Fine-tuning

$$\theta_i' = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$$

- Meta-training

$$\min_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i'})$$

---

**Algorithm 1** Model-Agnostic Meta-Learning

---

**Require:** $p(\mathcal{T})$: distribution over tasks
**Require:** $\alpha, \beta$: step size hyperparameters
1: randomly initialize $\theta$
2: **while** not done **do**
3:     Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
4:     **for all** $\mathcal{T}_i$ **do**
5:         Evaluate $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$ with respect to $K$ examples
6:         Compute adapted parameters with gradient descent: $\theta_i' = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$
7:     **end for**
8:     Update $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i'})$
9: **end while**

# MAML for Sim-to-Real Gap

- Fine-tuning

$$\theta_i' = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$$

- Meta-training

$$\min_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i'})$$

**Algorithm 1** Model-Agnostic Meta-Learning

**Require:** $p(\mathcal{T})$: distribution over tasks
**Require:** $\alpha, \beta$: step size hyperparameters
1: randomly initialize $\theta$
2: **while** not done **do**
3:      Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
4:      **for all** $\mathcal{T}_i$ **do**
5:          Evaluate $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$ with respect to $K$ examples
6:          Compute adapted parameters with gradient descent: $\theta_i' = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$
7:      **end for**
8:      Update $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i'})$
9: **end while**

# MAML for Sim-to-Real Gap

- Fine-tuning

$$\theta_i' = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$$

- Meta-training

$$\min_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i'})$$

**Algorithm 1** Model-Agnostic Meta-Learning

**Require:** $p(\mathcal{T})$: distribution over tasks
**Require:** $\alpha$, $\beta$: step size hyperparameters
1: randomly initialize $\theta$
2: **while** not done **do**
3:     Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
4:     **for all** $\mathcal{T}_i$ **do**
5:         Evaluate $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$ with respect to $K$ examples
6:         Compute adapted parameters with gradient descent: $\theta_i' = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$
7:     **end for**
8:     Update $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i'})$
9: **end while**

# Gradient Estimation using Evolution Strategy (ES)

$$\nabla_\theta \tilde{f}_\sigma(\theta) = \frac{1}{\sigma} \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d)} [f(\theta + \sigma \mathbf{g}) \mathbf{g}]$$

Perturbation
direction

Return of
perturbed policy

# Gradient Estimation using Evolution Strategy (ES)

$$\nabla_\theta \tilde{f}_\sigma(\theta) = \frac{1}{\sigma} \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d)} [f(\theta + \sigma \mathbf{g}) \mathbf{g}]$$

Estimated policy gradient

# ES-MAML for Reality Gap

**Algorithm 1** Model-Agnostic Meta-Learning

**Require:** $p(\mathcal{T})$: distribution over tasks
**Require:** $\alpha, \beta$: step size hyperparameters
1: randomly initialize $\theta$
2: **while** not done **do**
3:  Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
4:  **for all** $\mathcal{T}_i$ **do**
5:   Evaluate $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$ with respect to $K$ examples
6:   Compute adapted parameters with gradient descent: $\theta_i' = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$
7:  **end for**
8:  Update $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i'})$
9: **end while**

**Data:** initial policy $\theta_0$, adaptation step size $\alpha$, meta step size $\beta$, number of queries $K$
1 **for** $t = 0, 1, \dots$ **do**
2   Sample $n$ tasks $T_1, \dots, T_n$ and iid vectors $\mathbf{g}_1, \dots, \mathbf{g}_n \sim \mathcal{N}(0, \mathbf{I})$;
3   **foreach** $(T_i, \mathbf{g}_i)$ **do**
4     $\mathbf{d}^{(i)} \leftarrow \text{ESGRAD}(f^{T_i}, \theta_t + \sigma \mathbf{g}_i, K, \sigma)$;
5     $\theta_t^{(i)} \leftarrow \theta_t + \sigma \mathbf{g}_i + \alpha \mathbf{d}^{(i)}$;
6     $v_i \leftarrow f^{T_i}(\theta_t^{(i)})$;
7   **end**
8   $\theta_{t+1} \leftarrow \theta_t + \frac{\beta}{\sigma n} \sum_{i=1}^n v_i \mathbf{g}_i$;
9 **end**

# ES-MAML for Reality Gap

**Algorithm 1** Model-Agnostic Meta-Learning

**Require:** $p(\mathcal{T})$: distribution over tasks
**Require:** $\alpha, \beta$: step size hyperparameters
1: randomly initialize $\theta$
2: **while** not done **do**
3:     Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
4:     **for all** $\mathcal{T}_i$ **do**
5:         Evaluate $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$ with respect to $K$ examples
6:         Compute adapted parameters with gradient descent: $\theta_i' = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$
7:     **end for**
8:     Update $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i'})$
9: **end while**

**Data:** initial policy $\theta_0$, adaptation step size $\alpha$, meta step size $\beta$, number of queries $K$

1 **for** $t = 0, 1, \ldots$ **do**
2     Sample $n$ tasks $T_1, \ldots, T_n$ and iid vectors $\mathbf{g}_1, \ldots, \mathbf{g}_n \sim \mathcal{N}(0, \mathbf{I})$;
3     **foreach** $(T_i, \mathbf{g}_i)$ **do**
4         $\mathbf{d}^{(i)} \leftarrow \text{ESGRAD}(f^{T_i}, \theta_t + \sigma \mathbf{g}_i, K, \sigma)$;
5         $\theta_t^{(i)} \leftarrow \theta_t + \sigma \mathbf{g}_i + \alpha \mathbf{d}^{(i)}$;
6         $v_i \leftarrow f^{T_i}(\theta_t^{(i)})$;
7     **end**
8     $\theta_{t+1} \leftarrow \theta_t + \frac{\beta}{\sigma n} \sum_{i=1}^n v_i \mathbf{g}_i$;
9 **end**

Initial Policy | After 30 Episodes | After 50 Episodes

The initial policy shifts to the right.

Domain Randomization

PG-MAML

Our Method

# Summary

- Improve simulation
  - System identification
    - Always should be the first step for sim-to-real
    - Identified parameters can be reused
    - Require diverse data and careful experiment design
    - Can be tedious
- Improve policy
  - Domain randomization
    - Simple
    - Zero-shot transfer is often possible
    - Trade-off between robustness and optimality
  - Domain adaptation
    - Small amount of real-world data needed
    - Good experience so far in multiple locomotion projects
  - Meta learning
    - Novel and popular research direction

# Pipeline for Sim-to-Real

```
┌─────────────┐     ┌─────────────┐     ┌─────────────┐
│ Sanity check│ ──▶ │  Measure    │ ──▶ │   Model     │
│  on URDF    │     │  latency    │     │  actuators  │
└─────────────┘     └─────────────┘     └─────────────┘
```

# Pipeline for Sim-to-Real

```
┌──────────────┐     ┌──────────────┐     ┌──────────────┐
│ Sanity check │ ──→ │   Measure    │ ──→ │    Model     │
│   on URDF    │     │   latency    │     │  actuators   │
└──────────────┘     └──────────────┘     └──────────────┘
                                                 │
                                                 ↓
                                          ┌──────────────┐
                                          │    Domain    │
                                          │randomization │
                                          └──────────────┘
```

# Pipeline for Sim-to-Real

```
┌──────────────┐     ┌────────────┐     ┌────────────┐
│ Sanity check │ ──► │  Measure   │ ──► │   Model    │
│   on URDF    │     │  latency   │     │  actuators │
└──────────────┘     └────────────┘     └────────────┘
                                               │
                                               ▼
                     ┌────────────┐          ╱Work?╲        Yes
                     │   Domain   │ ──►    ◄  Work?  ►  ──────────►  **Done!**
                     │randomization│        ╲     ╱
                     └────────────┘          ╲  ╱
```
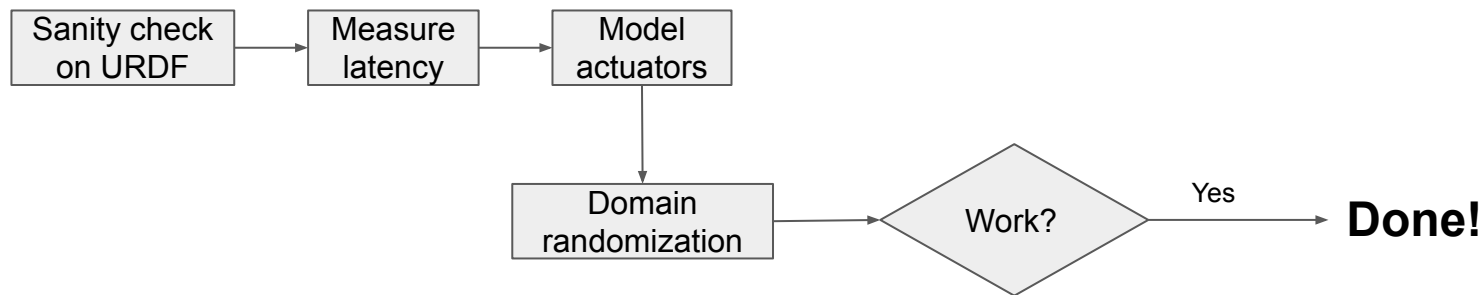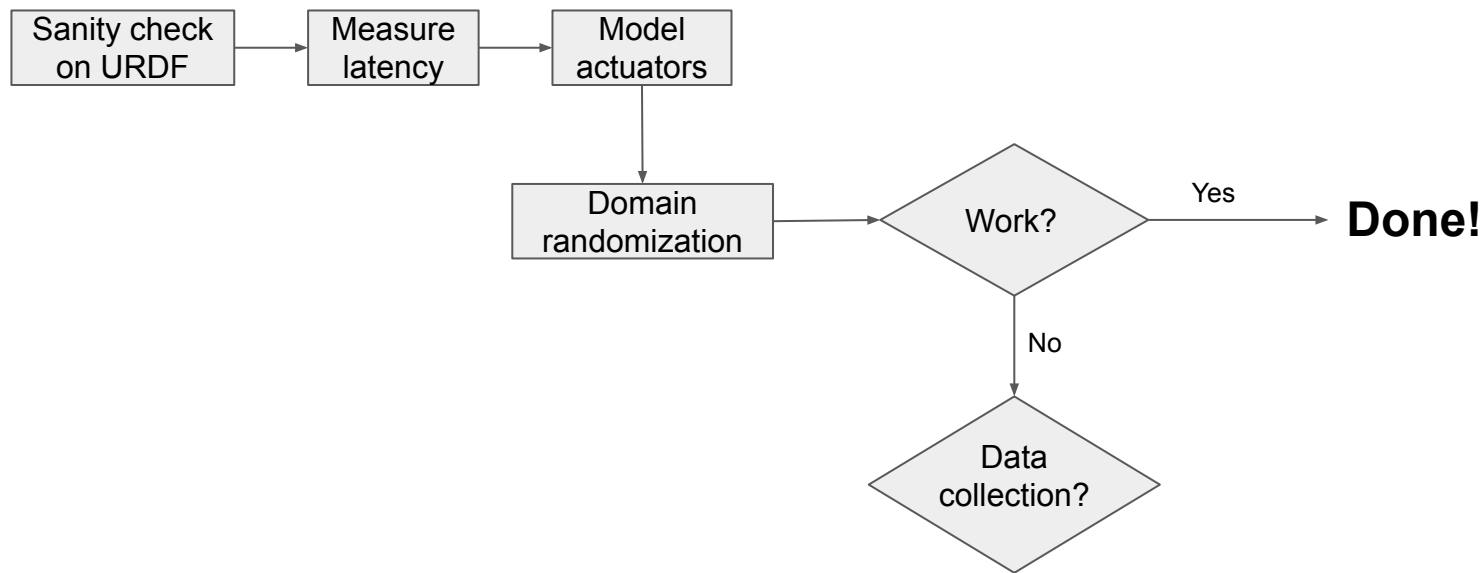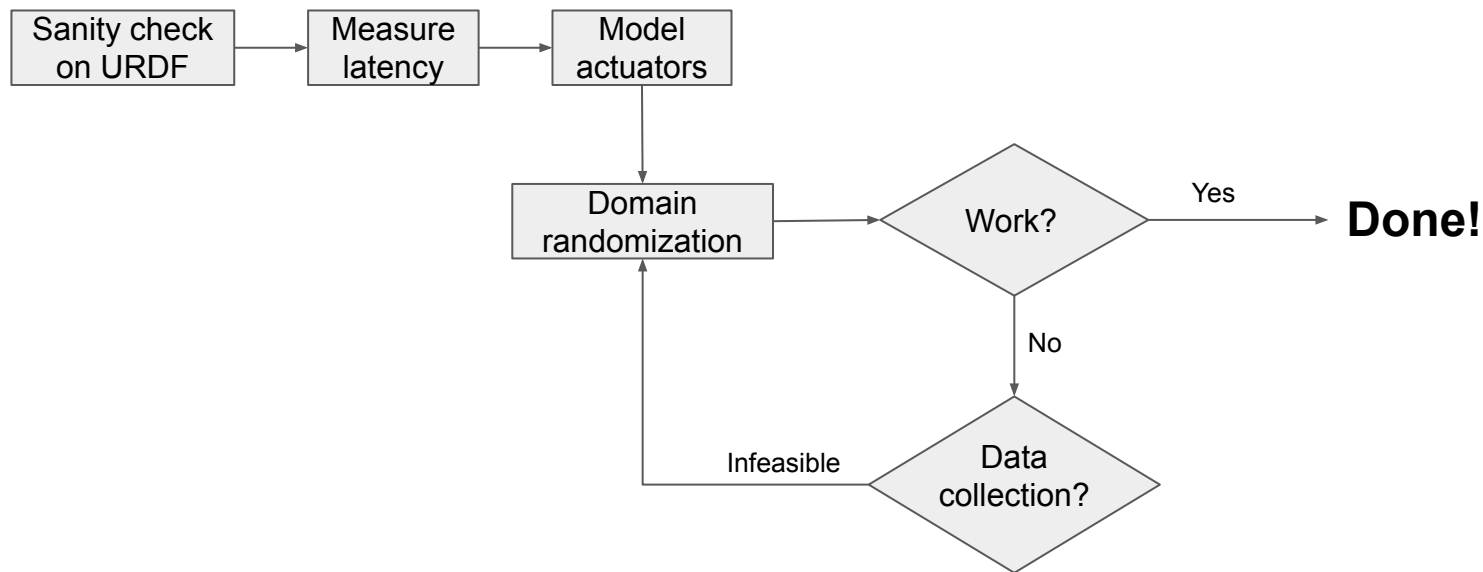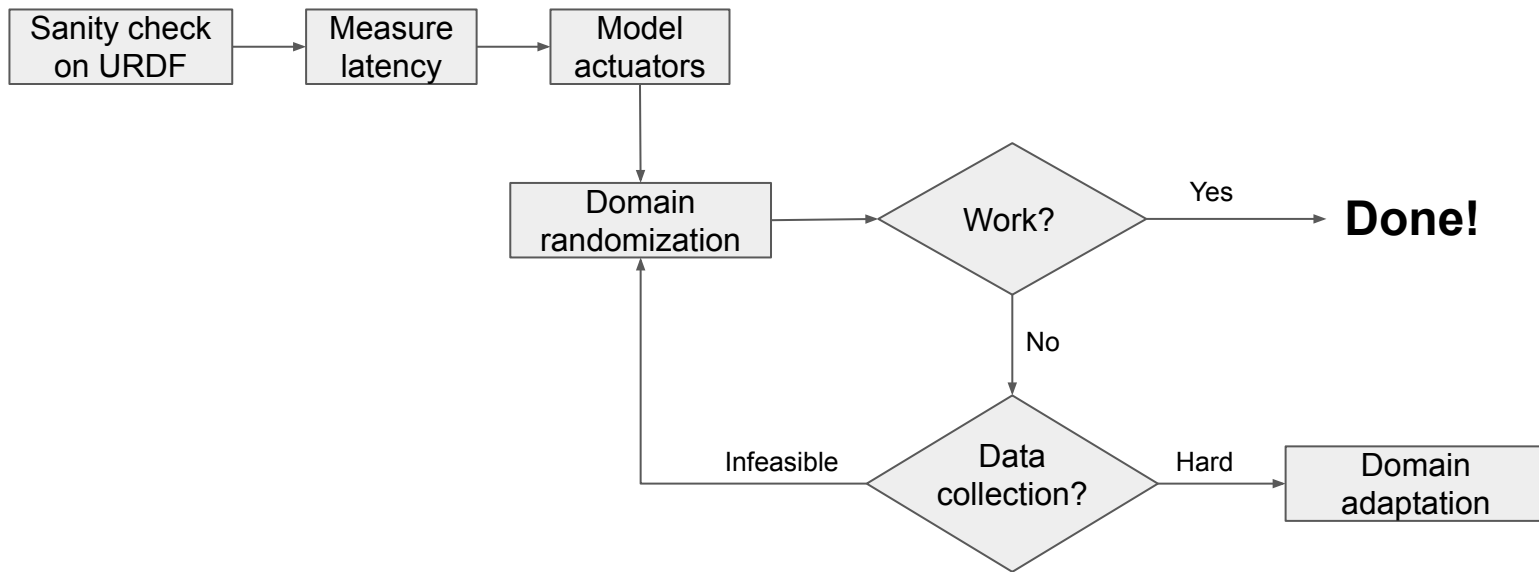
# Pipeline for Sim-to-Real

# Pipeline for Sim-to-Real

# Pipeline for Sim-to-Real

# Pipeline for Sim-to-Real