# Sim-to-Real in Reinforcement Learning Across Multiple Domains

TIM KNUDSEN

FH Wedel University of Applied Sciences
tim04knudsen@gmail.com

May 19, 2025

**Abstract**

*Reinforcement Learning is a modern and powerful approach to machine learning. It has achived remarkable success in simulated environments. But no matter how powerful your algorithm is, it is only as good as the deployment in the real world. Due to the complexity of the real world, the Sim-2-Real transfer is a major challenge in the field of Reinforcement Learning. In this paper, we will discuss the challenges of the Sim-2-Real transfer and how to overcome them. We will discuss different approaches to the Sim-2-Real transfer and how they can be used to improve the performance of Reinforcement Learning.*

## I. INTRODUCTION

Reinforcement Learning (RL) is a machine learning paradigm in which an agent learns to make decisions by interacting with an environment and receiving feedback in the form of rewards. Through trial-and-error, the agent improves its policy to maximize cumulative rewards. In recent years, deep RL has solved complex tasks in sumulation - from playing video games to controlling simulated robots - showcasing the protential of automomous learning. However, applying these learned policies on pysical systems presents a significant challenge known as the sim-to-real gap. Simulated training is appealing because collecting real-world data can be costly, slow, or unsafe, whereas simulations provide virtually unlimited, fast data in a safe setting. Yet differences between simulated and real physics, sensors, and visuals often cause a policy that works well in simulation to fail on the real hardware. Bridging this sim-to-real gap has become a crucial research focus in fields like robotics and atonomous vehicles.To address the reality gap, researchers have developed a spectrum of transfer techniques.

This paper revises the core sections of sim-to-real transfer in RL by organizing the discussion around the fundamental components of an RL problem, following the Observation–Action–Transition–Reward (OATR) framework proposed by recent work. We first introduce RL fundamentals in terms of these four components, highlighting how each contributes to the sim-to-real discrepancy. We then provide a structured review of sim-to-real transfer techniques, explicitly mapping each method to the OATR element it targets. We preserve key references and seminal examples from prior research – including CAD2RL for drone navigation, OpenAI's Dactyl and ADR (Automatic Domain Randomization) for robotic manipulation, and Tan et al.'s quadruped locomotion – to illustrate each category of method.

## II. BACKGROUND

### i. Definition

An RL task is typically formalized as a Markov Decision Process (MDP) defined by four core elements: Observation, Action, Transition, and

1

Reward. We briefly define each component and discuss how they factor into the sim-to-real context:

1. Observation(State): A feature representation $s$ that tries to describe the environment at a given time $t$. It can be discrete (like a chessboard) or continuous (like a robot's position in space). In the simulation environment, the state is often simplified and idealized, while in the real world, it can be noisy. Mismatches in the state space $S$ between the simulator and the real world can lead to a significant sim-to-real gap.

2. Action: The action space $A$ defines the set of actions the agent can take. Simulators often use discrete actions (like "left" or "right"), while real-world actions can be continuous and flexible, influenced by factors like physical constraints, noise, and delays.

3. Transition: The state transition function $T(s_t, a \rightarrow s_{t+1})$ defines how the world evolves from a given state with a action. In the simulation it's often deterministic, while in the real-world, it can be stochastic or different distributed. This could lead to a different real-world state $s_{t+1}^{real}$ than the $s_{t+1}^{sim}$.

4. Reward: The reward function is used to evaluate an selected action $a_t$ on the current state $s_t$. It's used to train the action-selection from a state (named policy $\pi$). In sim, engineers often create own functions for faster learning and using information that is easy to get in simulation. In reality, rewards may be delayed or doesn't fully capture the real task goals or constraints.

These four elements together determine the behaviour of an RL agent. The sim-to-real challenge can be viewed as discrepancy in one or more of these elements between the simulator MDP and the real-world MDP. Each technique can be used to mitigate one or more type of discrepancy.

## III.   Observation-Level Transfer

Observation discrepancies are the most visually apparent sim-to-real differences. These include divergences in visual appearance (textures, lighting, backgrounds), sensor noise patterns, and partial observability. There are two major families of methods address the observation (state) gap: domain randomization and domain adaptation, but also other approaches like sensor calibration or fusion can be used to align simulated observations with reality.

### i.   Domain Randomization

This technique randomizes the rendering of the simulator's observations so that the RL agent is trained on a wide variety of appearances, forcing it to learn features invariant to those visual details. Pioneering work by Tobin et al. (2017) demonstrated that randomizing textures, colors, lighting, and camera angles in simulation enabled a vision-based object detector to generalize to real images, even though the policy never saw a single real photograph during training. Sadeghi and Levine's CAD2RL experiment (2017) applied pure visual randomization to train a drone to fly in simulation using only RGB images; remarkably, the learned policy could control a real quadrotor through indoor corridors based solely on its onboard camera, achieving real-world collision avoidance without any real-image training data. In robotic manipulation, OpenAI's Dactyl project used extensive visual domain randomization to train a robot hand to reorient objects. They randomized lighting, textures, and backgrounds in the virtual scene such that the hand's camera observed many variations. This produced a policy that was resilient to real-world visual differences, allowing the real robot to manipulate a block and even solve a Rubik's Cube despite substantial changes in lighting or object appearance. The general principle is that by exposing the agent to diverse simulated observations, it learns a representation that is robust to the exact appearance of the state, thus coping with real-world vi-

suals it was never explicitly trained on. One caution is that unbounded randomization can make training unstable or overly difficult if the agent sees too much random variation too soon. To address this, researchers introduced curriculum strategies. For example, Automatic Domain Randomization (ADR) (OpenAI et al., 2019) begins with small random variations and progressively increases the randomness as the policy improves. ADR was used to great effect in training Dactyl's vision and control policies – starting with near-realistic visuals and gradually adding more extreme randomizations (e.g., wildly colored or textured objects) once the agent mastered easier variants. This curriculum ensures the agent is not overwhelmed initially, and it focuses training on closing the sim-to-real gap incrementally. Recent studies like Yuan et al. (2024) have further formalized curriculum randomization to stabilize RL training and achieve state-of-the-art vision-based policy transfer. In summary, visual domain randomization treats every rendered frame as a new "domain," so the agent essentially learns to ignore visual specifics and rely on task-relevant features, thereby mitigating observation mismatches between simulation and reality.

## ii. Domain Adaptation

While domain randomization modifies the simulator, domain adaptation typically uses data-driven post-processing or representation learning to bridge the gap from simulation outputs to real inputs (or vice versa). The idea is to align the feature distributions of simulated and real observations, so that a policy trained on simulation features will behave correctly on real features. One common approach is to use adversarial training: for instance, Bousmalis et al. (2018) employed a generative adversarial network to transform simulated camera images into a more realistic style before feeding them to a robotic grasping policy. By training a convolutional network to fool a discriminator into mistaking simulated observations for real, the feature representations become invariant to the origin (sim vs real). This enables the

policy to perform consistently across domains, as shown in experiments where a grasping policy learned in sim (augmented by an adaptation network) could directly grasp real objects with much higher success rates than an unadapted policy. Another example is Pixel-Level Domain Adaptation by Ganin and colleagues, and the SimGAN framework (2017) which refined simulator-rendered images to look more photorealistic using unpaired image translation. Domain adaptation isn't limited to vision: it can apply to any observation modality. For instance, aligning LiDAR point cloud distributions or calibrating simulator sensor noise models also fall under this category. Carlson et al. (2019) demonstrated sensor calibration for simulations of radar and cameras so that the statistical properties (e.g., noise covariance, artifact patterns) match those observed from real sensors. By minimizing the discrepancy in observations, the agent is less likely to be thrown off by unfamiliar inputs when deployed. Recent research has also explored self-supervised representation learning to find a shared latent state space for sim and real observations. Jeong et al. (2020) introduced a sequence-based autoencoder that learns latent states capturing the underlying environment dynamics, yielding improved sim-to-real performance in a robot manipulation task. The key takeaway is that domain adaptation techniques insert an adaptation layer or transform in the pipeline so that the observation fed into the policy looks as if it came from the same distribution as training data, even if it actually comes from a different domain. This can be achieved through adversarial loss, contrastive learning, or mapping functions between domains.

## iii. Other Observation-Level Strategies

In some cases, rather than rely on a single sensor modality, combining multiple observations can mitigate sim-to-real issues. For example, if a simulator provides perfect state information (like object positions) which are not available directly from sensors in reality, one could train

a policy that uses both vision and this privileged state in sim, then at deployment use real vision plus perhaps another sensor that helps approximate the missing state. This idea of sensor fusion and complementary observations is highlighted by Mahajan et al. (2024), who show that integrating data from multiple sensors (e.g., camera + depth sensor) yields more robust policies. The intuition is that even if one sensor's data differ between sim and real, the additional modalities provide redundancy. Another trend is leveraging large pretrained models – for instance, using a pretrained foundation model or pre-trained visual encoder that has seen many real images, to encode simulated observations into a representation that is already aligned with real-world visual features. By plugging such an encoder into the RL pipeline, the hope is that the agent's perception is grounded in real-world priors, reducing the gap in how it interprets simulated vs real observations. While these approaches are emerging (e.g., using CLIP or other vision transformers to encode state), they point toward the same goal: ensuring the agent "sees" the simulator and the real world in a consistent way.

In summary, observation-level transfer techniques focus on the state/observation space $S$ of the MDP. Methods like visual domain randomization widen the training distribution to encompass reality, whereas domain adaptation methods explicitly learn mappings or invariant features to align sim with real observations. Both approaches have enabled zero- or few-shot transfers in vision-based tasks that were once thought to require real data. Notably, the success of CAD2RL and Dactyl underscore that bridging the perception gap is often the first and pivotal step in sim-to-real reinforcement learning.

## IV. Action-Level Transfer

## V. Cross-Domain Case Studies

To ground the above concepts, we now highlight several case studies where sim-to-real techniques have enabled successful transfer in different application domains. These examples span robotic manipulation, drone flight, and legged locomotion – demonstrating both the commonalities and unique challenges of each domain. Each case study will mention the task, the approach used to bridge the gap, and key results achieved on real systems.

### i. Robotic Object Manipulation

One of the early triumphs of sim-to-real in RL was in robotic object manipulation – specifically, having robot arms or hands handle objects using policies learned in simulation. A representative example is object grasping and localization. In the 2017 study by Tobin et al. (discussed earlier in Domain Randomization), the task was for a robotic arm to pick up objects (toys of various shapes) on a table using input from an RGB camera. The challenge was primarily vision-based: the robot needed to perceive the objects' positions and plan grasps. The researchers trained a convolutional neural network to detect object locations by training entirely on synthetic images of a cluttered tabletop with random textures, lighting, and distractor objects (domain randomization). Even though the simulated images looked nothing like the real world individually, the neural network learned to focus on geometric cues and was robust to color/texture. When this network was used on the real robot's camera input, it localized objects with 1.5 cm accuracy and enabled reliable grasping in the real world – all without any real-image training. This was a landmark result showing that simulation-only training could produce a useful real-world perception system for manipulation.

Building on that, OpenAI's Dactyl (2018) tackled an even harder manipulation problem: dexterous in-hand manipulation. In that setup, a five-fingered robot hand had to reorient a block (and later a Rubik's Cube) in its fingers. They trained the policy in a MuJoCo simulator with extensive randomizations of both visuals and physics (frictions, object mass, joint gains, etc.). One especially important aspect was us-

ing an LSTM (recurrent network) in the policy to allow it to adapt online to the physics – essentially the policy would "feel" how heavy or slippery the object was and adjust its strategy. This can be seen as a form of implicit system identification happening inside the policy's hidden state (a clever trick). After training on on the order of $10^8$ timesteps (with massive parallelization), the policy was deployed to the real Shadow Hand. The result: the hand could successfully spin and reposition the block in precise ways, achieving the goal orientations consistently. Later, they even demonstrated the hand solving a Rubik's Cube – an extremely complex manipulation – using the same policy architecture (with some additional state estimation help). The success of Dactyl was a strong validation of the domain randomization approach for high-dimensional control. Figure 3 shows the real robotic hand from the Dactyl system manipulating a block, an outcome of this sim-to-real training.

Another manipulation example is assembly tasks, such as inserting pegs or screwing in bolts. These tasks require fine precision and often have contact-rich dynamics (which simulators struggle with). In 2019, researchers at NVIDIA (Mahler et al., 2019) showed that by using a high-fidelity simulator (Isaac Gym) and randomizing friction and part tolerances, a policy could learn to insert a peg into a hole in sim and then perform it on a real robot arm with high success. They also leveraged domain adaptation by using depth camera inputs (which are easier to match between sim and real than RGB images). Depth images from simulation are quite comparable to real depth sensor data, so the gap was smaller; in some sense, choosing the right sensing modality (depth over RGB) is also a strategy to ease sim-to-real.

To summarize the manipulation domain: Visual randomization has been critical for perception-driven tasks, and dynamics randomization plus recurrent policies or system ID has been important for contact-rich tasks. These strategies enabled policies trained in simulation to achieve real-world feats like multi-

object grasping and dexterous hand manipulation, which are significant milestones in robotics. Even so, challenges remain, such as more complex object interactions, generalizing to novel objects not seen in simulation, and long-horizon manipulation sequences. Current research is exploring combining imitation learning from a few real demos with simulation training to handle those cases (so that the policy gets a hint of real-world behavior to anchor it).

## ii. Drone Navigation (UAV Flight)

Unmanned aerial vehicles (UAVs or drones) present a compelling application for sim-to-real RL because crashing a real drone during learning is very costly, yet drones operate in highly dynamic environments where a robust policy is needed. The case study of CAD2RL by Sadeghi and Levine (2016) is instructive. They wanted a drone to perform collision avoidance in indoor hallways using only a single camera image as input (no fancy laser sensors). They trained a deep RL policy in a 3D simulator that generated random hallway textures and layouts – essentially the drone learned to react to visual input to avoid walls. Crucially, the network was never shown a real image during training, only the randomized simulated images. When they flew the policy on a real drone (a small quadrotor with an onboard camera), it was able to traverse real hallways and avoid obstacles for hundreds of meters of flight, all without collisions. This zero-shot transfer was achieved primarily by domain randomization of the visual domain (wall colors, pictures on the wall, floor patterns, lighting, etc.), so the drone didn't overfit to any particular appearance. The dynamics in that case were less of an issue because the policy was high-level (it output velocity commands); a low-level stabilizing controller was running on the drone to handle basic flight, which is common in such setups.

Another example in drone navigation is in outdoor flight or in simulated city environments. Researchers have used simulators like

Microsoft AirSim or Gazebo to train drone policies for tasks like following roads or avoiding moving obstacles. Domain adaptation has been applied here by taking simulated images (which might, for instance, not perfectly capture the complexity of natural outdoor scenes) and using neural networks to enhance realism. For instance, one could train a network to translate simulated aerial images to look like real aerial images using satellite photo datasets. By doing so, an RL agent trained on the translated images would be more familiar with real visuals when deployed.

One of the challenges unique to drones is the dynamics: unlike many ground robots, drones have very fast dynamics and are quite sensitive to aerodynamic effects (propeller wash, wind) which are hard to simulate exactly. Some recent works (2019–2021) have started to incorporate dynamics randomization for drones – e.g., randomizing the mass of the drone, motor thrust curves, latency in control signals, and wind disturbances during training. Combined with PID controllers or adaptive control on the drone, this has improved success in transferring agile flight maneuvers from simulation to real quadrotors. For example, if a simulated drone can learn to do a flip or a fast turn in the presence of random gusts of wind and motor power variations, the resulting policy is more likely to work on a physical drone where those factors occur naturally.

In summary, for drone navigation, vision-based tasks have benefited hugely from visual domain randomization (as in CAD2RL), while highly dynamic tasks require careful system identification (tuning simulators for drone dynamics) and possibly domain randomization of dynamics. Safety is a big concern – unlike a robot arm that might just miss a grasp, a bad drone policy can crash the vehicle. So often a conservative approach is used: the RL policy might be given to the real drone only when a safety controller can take over if it diverges, or training might happen progressively (low-speed flight policies then increasing speed). This relates to future directions in safe RL, but it's worth noting as a practical aspect of sim-to-

real for UAVs.

## iii. Legged Robot Locomotion

Legged robots, such as bipeds or quadrupeds, have high degrees of freedom and complex contact dynamics with the ground. They are an excellent testbed for sim-to-real because falls or instability in the real world are potentially damaging, yet simulation of legged locomotion can be made quite realistic with modern physics engines. One prominent case study is deep RL for quadruped locomotion by Tan et al. (2018). In that work, the task was to make a quadruped robot (with four legs, similar to MIT's Mini Cheetah or Unitree's Laikago) learn to trot and gallop – gaits that require dynamic balance. They used a combination of system identification and domain randomization: first, they improved their simulator's fidelity by identifying the robot's motor model and sensor latency (the time delay between sending a command and the motor responding, and between an event and the sensor reading). This alone significantly narrowed the gap because a lot of control policies are sensitive to latency. Next, they applied dynamics randomization during training: they varied parameters like friction coefficients of the feet and ground, small variations in motor torque output, and they even added random pushes to the robot in simulation to teach it to recover balance. They also constrained the observation space to be something the real robot could reliably sense (e.g., using on-board inertial measurements rather than assuming access to global orientation). With these measures, they trained a policy (using PPO) in simulation that could produce fast trotting and galloping gaits.

The result was that when this policy was deployed on the real quadruped robot, it could execute the learned gaits and remain stable, without any additional training on the robot. The robot achieved a trot at around 0.5 m/s and a galloping gait faster than that, which was quite an achievement at the time for learned controllers. The sim-to-real transfer here crucially depended on both identification (to get

the simulator close, e.g. matching the real motor dynamics) and randomization (to account for remaining discrepancies, like a slightly different friction on the lab floor versus the sim). If they had trained only on a nominal simulator, small unmodeled effects could have led to the robot tripping or oscillating. Because the policy had experience with those perturbations (random pushes, etc.), it was robust.

Following this, many other works came out applying similar ideas to different robots. For instance, Hwangbo et al. (2019) trained a bipedal robot to jump and run using an accurate simulator and randomizing ground compliance (how springy the ground is), which let the robot handle both hard flooring and softer mats in reality. The theme is consistent: identify what matters, randomize the rest. Another interesting case was by Muratore et al. (2021) on a ball-in-cup task with a robot arm (a task where the robot must swing a cup to catch a ball on a string). They used dynamics randomization with adaptive feedback, meaning the policy had parameters that could adjust based on observed behavior differences, effectively learning a form of online system ID (this blurs into meta-learning territory). The policy learned to compensate for slightly different string lengths and ball weights between sim and real, successfully performing the trick on real hardware.

For legged locomotion, another challenge is often the computational cost: simulation of many-legged robots with contact can be slow. This is where tools like Isaac Gym's massive parallelism have helped. Researchers have shown that you can train a policy for a simple biped to walk in just a few minutes using thousands of parallel simulations – a case of brute force enabled by GPUs. Once trained, transfer still needs the above techniques, but the training turnaround is faster.

To sum up the locomotion case: It highlights the importance of physics-side alignment. Visual realism is usually not an issue (these policies often don't rely on external cameras, but on internal sensors like joint encoders, IMUs, etc.). Instead, it's all about dynamics: masses,

friction, delays, and unexpected forces. By addressing those through system identification and domain randomization, researchers have achieved robust locomotion controllers via deep RL that operate on real legged robots – something that even a few years ago was very difficult. Still, open problems include handling environment changes (e.g., different terrains like slippery ice or tall grass – one can randomize terrain properties to an extent, but extreme changes may need on-the-fly adaptation) and multi-contact scenarios (like a humanoid climbing or using

## iv. (continued) Legged Locomotion (cont.)

...multi-contact scenarios. For example, getting a bipedal humanoid to not only walk but also climb or use its arms for support involves complex interactions that are hard to simulate perfectly. Current sim-to-real techniques have yet to fully conquer these extreme cases, and they remain active research problems.

## VI. Challenges and Future Directions

Despite significant progress, sim-to-real transfer in reinforcement learning still faces numerous challenges. We highlight some key limitations and promising directions for future research:

1. Simulator Fidelity vs. Generalization: Achieving a highly accurate simulator via system identification can be time-consuming and still incomplete – some aspects of reality (like subtle wear or complex contacts) resist exact modeling. On the other hand, relying on broad randomization without understanding can make training inefficient. A theoretical understanding of why domain randomization works and how to choose optimal randomization distributions is still lacking. Bridging this gap requires new insights. Future

work may involve differentiable simulators or learned simulators that can be directly tuned with gradient-based methods, providing a more automated way to improve fidelity.

2. Sample Efficiency and Real Data Usage: Many sim-to-real successes (e.g. OpenAI's results) required an enormous amount of simulated experience, which is computationally expensive. While simulation is cheaper than real life, training might still take days or weeks on clusters. An open challenge is improving sample efficiency so that policies can learn robustly with less data. Techniques like model-based RL, which learn a model of the environment and plan with it, or offline RL using logged data, could reduce the need for brute-force simulation. Additionally, using small amounts of real data in the training loop can dramatically help but needs careful integration. Future systems might do real-to-sim-to-real loops where occasional real trials are used to correct the simulator (via system ID) or fine-tune the policy, thus continually reducing the reality gap.

3. Domain Adaptation Limits: Current domain adaptation methods mostly assume a shared feature space or visual similarity between sim and real. In cases where the simulation and reality have fundamentally different observables (for instance, if a simulation provides full state info that a real robot can't sense directly), adaptation is non-trivial. Also, training GANs or adaptation models can be unstable. There is room for more robust adaptation techniques, possibly using architectures that incorporate physical constraints or semantic information (e.g., ensuring that "objects" in a scene remain consistent through the translation). Future research could explore one-shot adaptation, where the policy adapts on the fly from the first few real observations – a bridge towards meta-learning.

4. Safety and Exploration in the Real World: One reason we use simulation is to avoid unsafe explorations on real robots. However, when deploying the learned policy, there is still a risk if the policy encounters states it wasn't prepared for. Ensuring safety during deployment (and during any real-world fine-tuning) is paramount, especially in domains like autonomous driving or healthcare. This calls for safe reinforcement learning approaches to be combined with sim-to-real. For instance, one might enforce that the policy respects certain safety constraints by design (through reward shaping or a filter) so that even if the sim was imperfect, the real robot won't take catastrophic actions. Another idea is to have an intermediary training phase in a controlled real setting (like a lab with safety harnesses for a robot) to verify and adjust the policy before full deployment.

5. Combining Methods and Multimodal Transfer: As noted, the most successful strategies often mix techniques – e.g., calibrating some parameters (system ID), randomizing others, and possibly adapting visual inputs. Designing unified frameworks that can do all of the above automatically is an ongoing direction. Recent work on adaptive domain randomization tries to learn the randomization range itself (using algorithms that adjust the randomization to where the policy is weak, focusing training on those). Moreover, integrating sim-to-real with other forms of transfer learning like imitation learning or human demonstrations is promising. A human demonstration in the real world could be used to guide the simulation randomization (to ensure the simulator produces similar scenarios) or to initialize the policy via imitation, after which RL in sim fine-tunes it. This can dramatically cut down training time and yield more human-compatible behaviors.

6. Meta-Learning and Online Adaptation: One exciting avenue is meta-learning,

where the policy is trained to adapt quickly to new environments. In a sim-to-real context, one could meta-train a policy across many simulated variations so that it can infer the current environment's quirks (friction, weight, etc.) from a small amount of experience and adjust its behavior. This is related to what the recurrent policies (like the LSTM in Dactyl) are doing – implicitly identifying the environment. Future work could make this more explicit and efficient, perhaps by having a small network (an adaptation module) that updates some policy parameters using a few real-world trials. Such approaches blur the line between training and deployment, as the robot continues learning (or at least calibrating) once it's in the real world. This can make the transfer more robust to long-term changes (like battery aging or component wear).

7. Wider Domains and Generalization: So far, sim-to-real RL has mostly been demonstrated on specific tasks. A long-term goal is generalist agents that can transfer knowledge from simulation to reality across a variety of tasks without retraining from scratch for each one. This might involve training large RL policies or models (perhaps inspired by large-scale unsupervised learning) in simulation on many tasks and then adapting them to real tasks with minimal real data. In addition, applying sim-to-real in domains beyond robotics – for instance, in economics or internet simulations to real markets, or in training assistants in virtual environments before deploying to real human interaction – are frontiers that present their own challenges of fidelity and human unpredictability.

In summary, while sim-to-real in reinforcement learning has matured a lot, especially in robotics, there remains a rich landscape of research ahead. The combination of improving simulation fidelity, smarter training techniques, and clever adaptation algorithms will continue to push the boundary. The interplay of these methods – and ensuring theoretical understanding catches up with empirical success – will be key to making sim-to-real transfer more reliable and routine.

## VII. Conclusion

Simulation-to-real transfer enables reinforcement learning agents to leverage the best of both worlds: the abundance and safety of simulated training, and the authenticity of real-world deployment. In this paper, we reviewed how core RL principles and various transfer techniques come together to bridge the notorious reality gap. We discussed domain randomization, which broadens training conditions until reality is just another variation, domain adaptation, which tweaks observations or learned representations to eliminate visual or sensory discrepancies, and system identification, which tunes simulators to mirror reality as closely as possible. Through case studies in robotic manipulation, drone navigation, and legged locomotion, we saw these ideas in action – turning what were once simulation-bound successes into real-world achievements.

Sim-to-real in reinforcement learning has evolved from a niche endeavor to a flourishing area of research. Robots can now grasp, fly, and walk using brains baptized in virtual worlds. Yet, as we highlighted, challenges like achieving higher fidelity, ensuring safety, and reducing data requirements remain before sim-to-real becomes a plug-and-play tool. Encouragingly, many future directions are being actively explored: from algorithms that learn how to randomize or adapt (meta-learning, Bayesian randomization) to more encompassing frameworks that combine multiple techniques for maximum robustness. The continued development of better simulators (e.g., photo-realistic renderers, differentiable physics) and more powerful RL algorithms will further close the gap.

Ultimately, success in sim-to-real not only advances robotics and AI by saving time and resources, but it also opens the door to deploying learning agents in the physical world with

confidence. As our simulators become more like sandboxes for training "digital athletes," we can tackle increasingly complex real-world tasks with agents that have practiced extensively in simulation. The progress so far – from video-game trained policies now controlling real robots, to drones and quadrupeds mastering skills via virtual trial and error – offers a compelling preview of what's to come. By continuing to refine these methods, we move closer to a future where learning in simulation and performing in reality is a seamless pipeline, empowering RL to drive innovations in robotics, autonomous vehicles, and beyond.

## References