

FORGE: Force-Guided Exploration for Robust Contact-Rich Manipulation under Uncertainty

Michael Noseworthy¹, Bingjie Tang², Bowen Wen³, Ankur Handa³, Nicholas Roy¹
Dieter Fox³, Fabio Ramos³, Yashraj Narang³, Ireteyo Akinola³

Abstract—We present FORGE, a method that enables sim-to-real transfer of contact-rich manipulation policies in the presence of significant pose uncertainty. FORGE combines a *force threshold* mechanism with a *dynamics randomization* scheme during policy learning in simulation, to enable the robust transfer of the learned policies to the real robot. At deployment, FORGE policies, conditioned on a maximum allowable force, adaptively perform contact-rich tasks while respecting the specified force threshold, regardless of the controller gains. Additionally, FORGE autonomously predicts a termination action once the task has succeeded. We demonstrate that FORGE can be used to learn a variety of robust contact-rich policies, enabling multi-stage assembly of a planetary gear system, which requires success across three assembly tasks: nut-threading, insertion, and gear meshing. Project website: <https://noseworm.github.io/forge/>

I. INTRODUCTION

We are interested in developing *sim-to-real* techniques for learning assembly primitives (e.g., low-clearance insertion or nut-threading). Over the past decade, work in simulation and sim-to-real techniques has led to advances in challenging areas such as dexterous manipulation and legged locomotion [1], [2], [3], [4]. However, similar results have only recently been achieved for robotic assembly, which requires efficient and accurate simulation of both the robot and the detailed, low-clearance parts [5], [6], [7], [8], [9], [10].

Even with these advances, successfully deploying sim-to-real policies for assembly tasks remains challenging. Previous approaches typically only consider small amounts of perceptual noise relative to part size. This assumption aligns with industrial-style robot workcells where uncertainty is typically engineered away. Strategies to deal with uncertainty include careful mechanical design of fixtures and adapters, extensive calibration processes, and the use of high-precision sensing. As a result, the time and money required to set up a workcell with different parts, poses, and tasks can be prohibitive for small-scale enterprises, or in less-structured environments. We aim to develop control methods that are robust to higher levels of pose estimation error, which is unavoidable in less structured environments.

When there is pose uncertainty, behaviours that search for and rely on contact can be used to ensure success [11], [12]. However, the required contact between the parts can lead to undesirable outcomes if the force is too high. Parts can slip or become damaged, making the task difficult or impossible to complete. Previous heuristic approaches, such as spiral

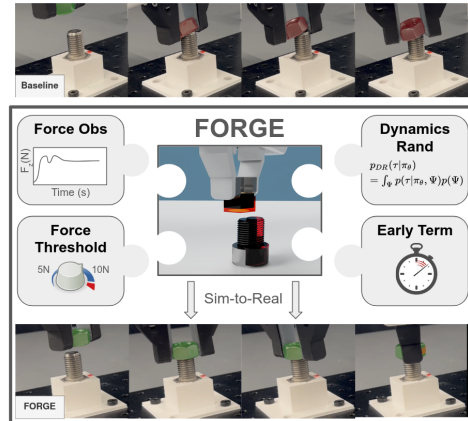


Fig. 1. FORGE uses force feedback to learn search behaviours for contact-rich tasks with pose estimation uncertainty. It combines *dynamics randomization*, a *force threshold*, and *early termination* for robust sim-to-real transfer. The resulting policies are *safe* and *efficient* (bottom) compared to aggressive baseline policies that cause parts to slip (top).

search [11], [13], can limit the applied force but are task-specific and can be inefficient. RL offers a general paradigm for developing more flexible search behaviours. However, the sim-to-real gap makes it difficult to transfer policies learned in simulation to the real world. In particular, it is challenging to simulate a similar contact distribution to what the real robot would experience. Even if the simulator has an accurate robot model (itself a time-consuming calibration procedure), it is difficult to know *a priori* the material and inertial properties of the parts the robot will interact with.

In this work, we propose FORGE: a framework for developing sim-to-real policies that safely and efficiently perform assembly tasks in the presence of significant pose uncertainty. FORGE trains policies in simulation that are *robust* to a wide range of contact interactions. Additionally, policies are trained without precise knowledge of part poses, leading to emergent search behaviours.

FORGE has two complementary components to ensure policies are robust to contact. First, we propose to condition policies on a *force threshold* that should not be exceeded during task execution. Second, policies are trained to maintain this threshold under a wide range of dynamics randomizations (we randomize *robot*, *controller*, and *part* properties). The dynamics randomization ensures that the search behaviour is robust to simulator mis-specification and has the added benefit that the resulting policies can be deployed without extensive controller tuning.

As assembly policies become more performant, we em-

Correspondence: mnosew@mit.edu
¹MIT ²USC ³NVIDIA.

phasize the importance of reporting metrics beyond *success rate*, such as *mean force* or *time to success*. Standard practice in *sim-to-real* assembly is to execute policies for a fixed duration [7]. However, due to task variation, this will often lead to premature termination or leave a robot “waiting” after the task is finished. Instead, the policy can determine when to terminate. This is itself a difficult task that can benefit from contact (e.g., a successfully inserted peg cannot move laterally). FORGE proposes a method for early termination that expands the action space so that the policy learns to predict task success. We show that early termination, also trained in simulation, robustly transfers to the real world.

In summary, our contributions are:

- 1) **A method to specify maximum allowable contact-force** during policy execution. This results in policies that exhibit safe search behaviour even with significant levels of pose estimation error (up to $5mm$).
- 2) **A dynamics randomization scheme** that enables robust sim-to-real transfer, and minimizes the need to tune controller gains.
- 3) **A method for early termination prediction** that allows efficient policy execution.
- 4) **A demonstration of multi-part assembly** of a planetary gearbox requiring a diverse set of skills, including the challenging task of fastening nuts and bolts.

Results are shown over > 500 real-world trials and multiple tasks. We plan to release the code with the paper.

II. RL FOR CONTACT-RICH ASSEMBLY

We want to learn policies for tasks with tight tolerances and detailed geometry. We first describe the problem formulation before introducing FORGE in the next section.

A. Assembly Tasks

Each task involves mating two parts: one held in the gripper and another fixed to the workspace. We consider all three tasks from *Factory* [5] and demonstrate the first sim-to-real transfer for threading a small M16 nut (see Fig. 2).

Peg Insertion: A small round peg with $8mm$ diameter needs to be inserted into a corresponding socket with $0.5mm$ diametrical clearance. There is position uncertainty such that a successful search behaviour requires lateral exploration.

Gear Meshing: Gears need to be inserted onto pegs with $0.5mm$ clearance. Other gears are present and the teeth of adjacent gears must be aligned for successful meshing. In addition to lateral exploration, rotational exploration may be necessary to mesh the teeth.

Nut Threading: Instead of fully lowering a nut onto a bolt as in *Factory*, we define the *nut threading* task as successfully threading the nut such that it cannot be lifted by a vertical motion (we find lowering by a quarter-thread is sufficient). Because our robot has joint limits, and to prevent the need to regrasp, we assume the nut and bolt are initially oriented¹ such that success can be achieved with a

¹We leave the more challenging, yet realistic, scenario involving completely unobserved thread orientation to future work.

single revolution of the wrist joint. We consider nuts with a relatively small size (M16) compared to previous sim-to-real work (M48) [14]. A successful search behaviour will resolve lateral uncertainty and place the nut on the bolt before rotating the wrist (otherwise the threads may not mesh).

B. POMDP Formulation

We formulate our problem as a *Partially Observable Markov Decision Process (POMDP)* [15], [16] to reflect the partial observability of most contact-rich manipulation setups. A POMDP is defined by the tuple: $(\mathcal{S}, \mathcal{A}, \Omega, T, O, R, \gamma)$. The goal is to learn a parameterized policy, $\pi_\theta(a_t|o_1, \dots, o_t)$, that maximizes the expected return:

$$J(\pi_\theta) = \mathbb{E}_{\tau \sim p(\tau, |\pi_\theta, \Psi)}[\sum_{t=0}^{\infty} \gamma^t r_t] \quad (1)$$

where $\tau = (s_0, a_0, o_0, s_1, a_1, o_1, \dots)$ is the trajectory of states, actions, and observations resulting from the robot following policy π_θ . Below, we further specify the components of the POMDP for contact-rich tasks.

States (\mathcal{S}): A state, $s_t \in \mathcal{S}$ consists of the pose and velocities of the end-effector (EE), fixed part, and held part: $p^{ee}, p^{fixed}, p^{held} \in SE(3)$ and $v^{ee}, v^{held} \in \mathbb{R}^6$.² We also include the contact force experienced by the end-effector, $F^{ee} \in \mathbb{R}^3$, and time-invariant information about the dynamics properties of the robot, controller, and parts (e.g., mass or joint-friction): $\Psi = (\psi_{robot}, \psi_{control}, \psi_{parts})$.

Observations (Ω): It is difficult to accurately estimate the full state for contact-rich manipulation. Instead, all our policies use the following observations:

- Noisy EE pose and velocity: $\hat{p}^{ee} \in SE(3)$, $\hat{v}^{ee} \in \mathbb{R}^6$
- Estimated contact force: $\hat{F}^{ee} \in \mathbb{R}^3$
- Noisy estimate of the fixed part’s pose: $\hat{p}^{fixed} \in SE(3)$

We do not include pose or velocity of the held part because it can move in the gripper and be difficult to track without tactile sensing. Likewise, we do not observe Ψ , but include the previous action, a_{t-1} , to help infer unknown dynamics.

Actions (\mathcal{A}): Control targets for a task-space impedance controller [17], [7]. As in previous work [5], [7], we assume all parts are in an upright orientation. Thus it is sufficient for the policy to only have control authority over the (x, y, z, yaw) -dimensions: $a_t \in \mathcal{A} = \mathbb{R}^4$.

Transition Function (T): T determines the next state and is parameterized by the dynamics parameters, Ψ : $T_\Psi : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$. T_Ψ is often specified using a simulator (in our case *IsaacGym* [18]) with the corresponding set of simulation parameters, Ψ^{sim} . The sim-to-real gap comes from the mismatch between Ψ^{sim} and Ψ^{real} .

Observation Function (O): The observation function generates noisy observations from state: $O : \mathcal{S} \times \mathcal{A} \rightarrow \Omega$. The position of the fixed part is assumed to have up to $5mm$ error. Independent Gaussian noise is added to each of the other observations at every timestep (except velocity, where positional noise is propagated through finite differencing). Noise values can be found on the website.

² v^{held} may be different from v^{ee} if the part slips in the gripper.

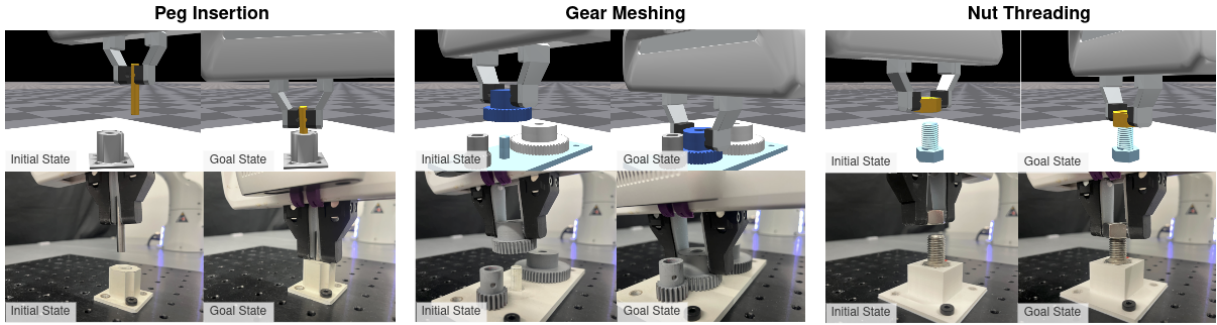


Fig. 2. FORGE is evaluated on three tasks proposed in the *Factory* work [5]: Peg Insertion, Gear Meshing, and Nut Threading. Each task is trained solely in simulation (top) and transferred directly to the real robot (bottom).

Reward (R): The reward function, $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, uses a keypoint formulation as its main component: $R_{kp}(p^{fixed}, p_t^{held})$. The target keypoints, k_t^{targ} , represent the desired position of the held part, while k_t^{held} represent its current position. We modify the keypoint reward from previous work [5], [19] to account for small, threaded geometries (see the website for more details). We also add two discrete *bonus* rewards that are given when important phases of the tasks are reached: once the held part is centered on top of the fixed part and once the task is successful:

$$R_{bonus}(p^{fixed}, p_t^{held}) = \mathbb{I}_{place} + \mathbb{I}_{success}. \quad (2)$$

The relative z -position of bottom of the held part to the top of the fixed part is used to check each condition (see website). We found the bonuses led to more robust learning when there is significant pose uncertainty.

III. FORGE: ROBUST SEARCH UNDER UNCERTAINTY

FORGE uses on-policy RL to learn exploratory behaviours in simulation. A *force threshold* (Section III-A) and *dynamics randomization* (Section III-B) are introduced to achieve robust search behaviours. FORGE also introduces an *early termination* procedure (Section III-C) for efficient execution.

A. Force Threshold

During policy execution, excessive force can cause parts to slip or become damaged (e.g., electronic components with fragile pins). Although it may be possible to recover from small amounts of slip with the right sensors (e.g., wrist camera or tactile), we prefer to avoid these scenarios.

To develop *safe* policies, we propose to condition the policy on a *force threshold*, $F_{th}: \pi(a|o, F_{th})$. During training, the policy is penalized if the contact force, F_t^{ee} , experienced by the arm exceeds the threshold. Concretely, we add an additional term to the reward function:

$$R_{contact_pen}(F_t^{ee}) = -\beta * \max(0, \|F_t^{ee}\| - F_{th}). \quad (3)$$

In simulation, the true contact force can be measured. Note that this penalty can be used during training whether the policy has access to the force observation or not.

B. Dynamics Randomization

To successfully deploy policies trained in simulation, it is important that the trajectory distribution experienced during training is similar to what it would be when deployed:

$p(\tau^{real}|\pi_\theta, \Psi^{real}) \approx p(\tau^{sim}|\pi_\theta, \Psi^{sim})$. The difference between these distributions is usually referred to as the *sim-to-real gap*. To gain insight into why minimizing the gap is important, specifically for contact-rich tasks, we can look more deeply into how trajectories are sampled:

$$\tau \sim p(\tau|\pi, \Psi) = p(s_0) \prod_{t=1}^T \left[\pi(a_t|o_{1:t}) p(o_t|s_t, a_{t-1}, \Psi) p(s_t|s_{t-1}, a_{t-1}, \Psi) \right]. \quad (4)$$

From this equation, we see the dynamics parameters can impact both the next-state and observation distributions. For the same action, different dynamics parameters can lead to parts being in different locations. Further, similar actions may lead to different observed contact forces.

The *sim-to-real gap* is usually handled by (1) system identification (Sys-ID) [20] or (2) dynamics randomization (DR) [21], [10]. The goal of Sys-ID is to tune Ψ^{sim} to be close to Ψ^{real} . This itself is a complicated tuning procedure that may need to be redone for every new set of parts. Instead, we follow the DR approach which learns policies that are *robust* to a wide range of dynamics parameters. Concretely, we optimize a version of Eq. 1 where:

$$\tau \sim p_{DR}(\tau|\pi_\theta) = \int p(\tau|\pi_\theta, \Psi) p(\Psi) d\Psi. \quad (5)$$

The integral is approximated with Monte Carlo samples from a randomization distribution (see website for values). The next subsections describe the variables that are randomized (at the beginning of each episode).

Controller Randomization: The controller has a large impact on what force will be experienced. This work uses impedance-control where applied forces are computed as:

$$p_t^{targ} = clip(combine(a_t, p^{fixed}), \lambda), \quad (6)$$

$$F^{targ} = k_p(p_t^{targ} - p_t^{ee}) - k_d v_t^{ee}. \quad (7)$$

First, the policy outputs a relative-pose, a_t , which is applied to the fixed part's pose to get an absolute target pose, p_t^{targ} . This pose is clipped by an action scale, λ , to ensure that the target is not too far from the EE's current pose. As in previous work, we use critically damped gains to ensure stable controllers: $k_d = 2\sqrt{k_p}$ [10], [22], [23]. The controller thus depends on two parameters which govern how much force can be commanded: $\lambda \times k_p$. We randomize both quantities so that the range of maximum commandable forces

is in $[6.4, 20.0]N$. Note that the control parameters are not included in the observations, so the policy must adjust its behavior based on force measurements. This reduces the policy’s dependence on a particular controller implementation.

Controller tuning [7] or optimization [23] is a costly and often complex procedure. Randomization has the additional benefit that the policy is robust to a range of control parameters, greatly simplifying deployment.

Part Friction Randomization: As parts slide against each other, the material friction determines how much lateral force the sensor will experience. To ensure policies can work across a range of materials, we randomize part friction.

Robot Dynamics Randomization: Due to phenomena such as joint friction, the applied force may be smaller than the commanded force. We implement a simple way to account for this: inducing a randomized *dead-zone* in simulation. Each episode, a dead-zone is selected for each dimension, F_i^{DZ} , where commanded forces below this value are clamped to zero: $|F_i^{applied}| = \max(0, |F_i^{target}| - F_i^{DZ})$. This enables the policy to increase its target which can help apply more force when needed or reduce steady-state error.

These randomizations lead to a policy that is robust to a wide range of dynamics parameters. Combined with the force threshold, the policy can modulate its actions to achieve safe interaction. For example, with higher gains, the policy will output smaller actions to limit the contact force.

C. Early Termination

Ideally, we want the policy to terminate as soon as the task has succeeded and no sooner. Although success is clearly defined in simulation where we have access to the positions of each part, it is much more difficult to reliably predict success in the real world [24]. Consider the nut-threading task, where the distance between a successfully threaded nut and a loose nut is a fraction of a millimeter.

We propose to train a success predictor which can robustly transfer from sim-to-real and be used to make early termination decisions. Concretely, we share the weights of the policy network with the success predictor by expanding the action space of the policy to include an early termination action: $a_t^{ET} \in [0, 1]$. To train the policy to output the correct action, we include an early termination reward, R_t^{ET} , which penalizes incorrect success predictions:

$$R_t^{ET}(a_t, y_t) = -|a_t^{ET} - y_t|, \quad (8)$$

where y_t is the true success label at time t . During training, episodes are always executed for the maximum length.

At deployment time, a confidence threshold, p_{term} , can be used to terminate the episode as soon as the policy believes it has succeeded: $a_t^{ET} > p_{term}$. This allows us to behave efficiently, a desirable property for industrial applications where it is important to reduce cycle times.

IV. EXPERIMENT SETUP

A. Robot System

For all experiments, we use a *Franka Panda* robot and the *FrankaPy* [25] library for the impedance controller. All

policies send targets to the controller at a rate of $15Hz$ while the controller operates at $1000Hz$. The Panda has joint-torque sensing, which is projected to EE-frame force values when needed by the policy: $\hat{F}^{ee} = J^\dagger(q)\tau^{ext}$, where J^\dagger is the Jacobian pseudo-inverse and τ^{ext} are the estimated external joint torques [26]. Alternatively, a force-torque sensor can be used to estimate the contact force.

For the majority of our experiments, we calibrate the poses of each fixed object and artificially add noise. This allows us to analyze performance under known levels of pose estimation error. The calibration is done by guiding the arm to a successful pose for the respective task from which a nominal initial pose can be backed out. Unless otherwise reported, our real experiments use the same initial state randomization as in simulation (see website).

For our last experiment, we assemble a planetary gear box (Section V-E) using the perception system from *IndustReal* [7] (retrained using data we collected). This model assumes the z -position of parts are known and uses a Mask-RCNN model [27] to estimate bounding boxes from which planar locations can be backed out. The perception errors in this system are largely caused by extrinsics calibration errors and minor bounding box prediction errors.

B. Policy Training

Simulator: All policies are trained using the *Factory* simulation methods within IsaacGym [5]. In simulation, we estimate the external contact force experienced by the end-effector (akin to attaching a force-torque sensor on the robot, or projection from joint sensing of external torques as done on Franka robots). Noisy sensor values are used as policy input, whereas ground-truth sensor values are used to compute the excessive-force penalty. For all RL, we use recurrent PPO [28] with asymmetric actor-critic [29] to handle partial observability. Details on initial state and observation randomization can be found on the website.

Checkpoint Selection: For all tasks and models, we train three policies with separate random seeds. For the peg-insertion and gear-meshing tasks, all policies are deployed on the real-robot and reported results are averaged across the three policies. For the nut-threading task, we found that not all policies transferred reliably to the real world, even when high success rates are achieved in simulation. As such, for this task, we report results for the *best* of the three checkpoints (determined using 18 test runs each).

Observation and Action Frames: To allow efficient generalization across the workspace, we assume actions and observations are defined relative to the tip of the fixed part. Specifically, the policy outputs a $4D$ relative transform from the tip of the fixed part (we assume upright parts). The control target is computed from the fixed part’s pose estimate and the relative pose from the policy. The policy output is bounded, which limits the operational volume of the end-effector (we allow targets to be up to $5cm$ away in all directions). Similar to the action space, all position observations are relative to the tip of the fixed part.

V. RESULTS AND DISCUSSION

A. Baseline Comparisons

We first compare FORGE to a baseline method that does not include any FORGE components. Specifically, the *Baseline* method was not trained using force observations, an excessive-force penalty, or dynamics randomization; however, it was trained with the early termination procedure so meaningful episode durations could be reported. We also considered a version of FORGE that does not use force observations. The questions we seek to answer are: **(Q1) Does FORGE lead to more robust sim-to-real transfer?** **(Q2) Does FORGE lead to policies with more desirable behavioural properties?**

We considered one variation of each task from Section II-A: the *8mm Peg*, the *Medium Gear* (with two abutting gears), and the *M16 Nut*. Along with policy success rate, used to measure robustness for Q1, the following metrics are reported to answer Q2:

- Duration (s): The average trial length when using an early termination threshold of $p_{term} = 0.9$.
- $F_{mean}, F_{max}(N)$: Forces experienced by the robot.
- Early Term. Precision: The fraction of early-terminated trials that were actually successful.
- Early Term. Recall: The fraction of successful trials which were terminated correctly with $a^{ET} > p_{term}$.

Each reported metric represents 45 trials spread across 5 workspace locations for the fixed part, and 3 pose-estimation error levels ranging from 0 – 5mm (see Fig. 3). Similar randomization ranges were used as in simulation except for the in-hand part randomization where the part was placed centrally in the gripper. Results are reported in Table I.

One conclusion for Q1 is that FORGE outperformed the *Baseline* method for all tasks whether force is included or excluded from the observation space. This suggests that the primary benefits of FORGE come from the dynamics randomization and excessive-force penalty. For FORGE, comparison to the policy without force observation shows that although using force sensing was useful for the easier insertion and gear meshing tasks, it harmed performance for the nut-threading task. We hypothesize that this is because nut-threading policies rely more heavily on force observations in simulation; these policies would therefore be more sensitive to any sim-to-real gap for this observation modality.

Examining the behavioural metrics for Q2, we notice that FORGE used less force than the baseline and had minor improvements in trial duration. During experiments, we observed FORGE led to gentler interactions between the parts (see accompanying video). The reduced force produced by this policy was especially helpful for the M16 Nut which was more susceptible to slipping than the peg or gear.

B. Noise Analysis

We next aim to answer **(Q3): How is policy performance, in terms of success rate, affected by pose-estimation error?** We use the same trials from the previous section, but show a breakdown of the results across different error

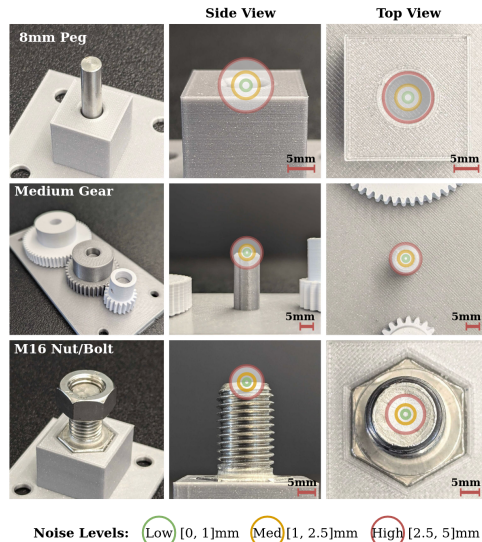


Fig. 3. **Perception Error** For each task, we visualize what the different pose estimation errors look like overlaid on the fixed part.

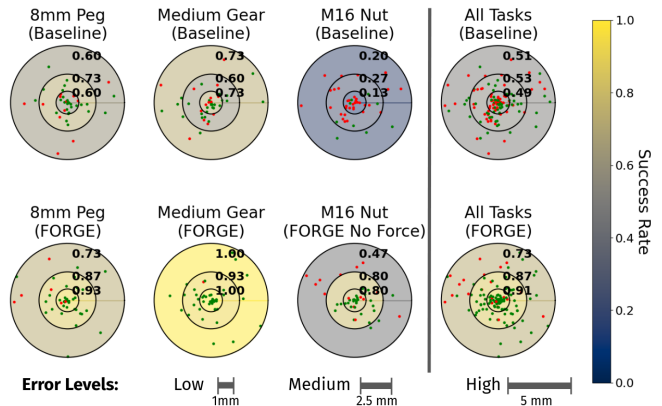


Fig. 4. **Noise Analysis** Performance broken down by level of pose error. Each subplot is a planar representation of the error levels where each ring corresponds to low (0-1mm), medium (1-2.5mm), and high (2.5-5mm) error. Success rate, stated in black text, is also represented by the shade of the corresponding ring. Dots represent x-y noise samples for successful (green) and failed (red) trials. FORGE results in good performance across tasks even with high error levels.

levels. During each trial, artificial perception error was added to the fixed part’s pose (calibrated as described in Section IV-A³). A third of the trials fell in each of the three considered error levels (see Fig. 3): Low (0-1mm), Medium (1-2.5mm), and High (2.5-5mm). We considered 3D position error by sampling a perturbation vector with a radius uniformly sampled in the desired error range and a direction uniformly sampled from the unit-sphere.

Figure 4 visualizes the performance of the baseline policy vs. FORGE at different noise levels. Each subplot is a 2D representation of how much x-y error there was for each trial (z-dimension error not visualized). Each point corresponds to either a successful (green) or unsuccessful (red) trial. The color of the ring represents the success rate at the corresponding error levels (increasing outwards). For the

³Adding artificial noise allows us to better characterize performance across error level compared to a perception system whose bias and variance can be difficult to estimate and control.

	Episode		Force		Early Termination	
	Success Rate \uparrow	Duration (s) \downarrow	F_{mean} (N) \downarrow	F_{max} (N) \downarrow	Precision \uparrow	Recall \uparrow
8mm Peg						
FORGE	0.84 (0.05)	5.01 (0.17)	5.51 (0.24)	12.84 (0.37)	1.00 (0.0)	1.00 (0.0)
FORGE (No Force)	0.82 (0.06)	7.30 (0.42)	7.09 (0.35)	14.16 (0.39)	0.59 (0.08)	0.81 (0.07)
No FP (400kp)	0.64 (0.07)	6.06 (0.40)	6.94 (0.13)	11.94 (0.24)	0.83 (0.07)	0.92 (0.05)
No FP (600kp)	0.71 (0.07)	5.28 (0.35)	10.66 (0.15)	16.58 (0.32)	0.91 (0.05)	0.97 (0.03)
Baseline	0.64 (0.07)	5.09 (0.30)	11.81 (0.21)	17.93 (0.41)	0.97 (0.03)	1.00 (0.0)
Medium Gear						
FORGE	0.98 (0.02)	6.34 (0.42)	7.95 (0.11)	15.10 (0.45)	0.95 (0.03)	1.00 (0.0)
FORGE (No Force)	0.93 (0.04)	9.02 (0.79)	8.49 (0.23)	14.68 (0.39)	0.60 (0.08)	1.00 (0.0)
No FP (400kp)	0.82 (0.06)	6.44 (0.23)	6.52 (0.14)	10.97 (0.24)	1.00 (0.0)	1.00 (0.0)
No FP (600kp)	0.73 (0.07)	6.99 (0.46)	9.48 (0.23)	15.73 (0.30)	0.94 (0.04)	0.97 (0.03)
Baseline	0.69 (0.07)	7.57 (0.50)	11.67 (0.45)	18.29 (0.40)	0.90 (0.05)	0.97 (0.03)
M16 Nut						
FORGE	0.44 (0.07)	24.50 (1.35)	6.88 (0.14)	13.34 (0.17)	0.50 (0.11)	1.00 (0.00)
FORGE (No Force)	0.69 (0.07)	13.16 (0.66)	7.78 (0.17)	14.41 (0.28)	1.00 (0.0)	1.00 (0.0)
Baseline	0.20 (0.06)	27.89 (2.22)	7.32 (0.23)	16.73 (0.29)	0.11 (0.10)	1.00 (0.0)

TABLE I

BASILINE COMPARISON FORGE (WITH AND WITHOUT FORCE OBSERVATIONS) IS COMPARED TO BASELINES THAT DO NOT INCLUDE ROBUST SIM-TO-REAL COMPONENTS. IT IS ADDITIONALLY COMPARED TO ABLATIONS THAT DO NOT USE AN EXCESSIVE-FORCE PENALTY. EVALUATIONS ARE PERFORMED OVER A TOTAL OF 585 TRIALS ON THE REAL ROBOT (45 PER ROW). STANDARD ERRORS ARE INCLUDED IN PARENTHESES.

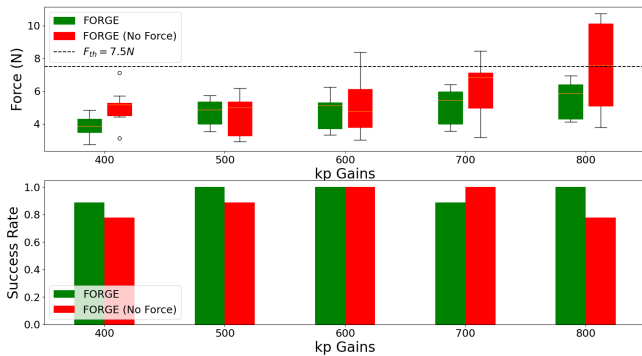


Fig. 5. **Gains Analysis** (90 trials, 8mm Peg) With force sensing (+Force), FORGE can achieve robust success rates (bottom) across varying controller gains at deployment time. Even with different gains, force sensing allows the policy to modulate its actions to achieve low contact forces (top).

M16 Nut task, we include results for the *No Force* ablation as this was the most robust policy that used dynamics randomization and a force threshold.

FORGE achieved high success rates (> 0.8) for all tasks at low and medium error levels, even for the M16 nut. Although performance degraded with error $> 2.5mm$, FORGE still significantly outperformed the baseline. With high error, the effects of contact are more pronounced because the robot may need to search longer before the task is complete.

C. Force Analysis

Next, we investigate how FORGE limits forceful interactions. **(Q4) How important is the excessive-force penalty for safe interactions? (Q5) Can FORGE limit the applied force without extensive controller tuning?**

Excessive-Force Penalty (Q4): In Table I, we compare to an ablation, *No FP*, that was trained without the excessive-force penalty of FORGE (but still used force observations and dynamics randomization). We used the same evaluation procedure as for FORGE but deployed with two different controller gains (we chose values at the lower and middle of the gain randomization range). Note that ablation results are not reported for nut-threading as we found that the nut always

slipped out of the gripper. We found that policies deployed with the lower gains achieved similar average forces to FORGE while those deployed with higher gains naturally experienced more force. Both policies had lower success rates than FORGE which was deployed with controller gains at the middle of the randomization range.

Gains Robustness (Q5): To measure how robust FORGE is to controller gains, we performed an additional experiment where we varied the gains at deployment time and measured success rate. We compared FORGE and the *No Force* ablation to gain insight into how important force sensing is to limiting applied forces. The experiment was carried out for the 8mm peg task at a single workspace location, with medium pose estimation error and limited initial-state randomization. We considered 5 proportional gain levels across the randomization range (corresponding to an 8N range in the maximum force the controller could apply) and each condition was evaluated 9 times (3 runs per checkpoint).

In Fig. 5 (bottom), we see that FORGE achieves high success rates across a wide range of controller gains. However, performance is less consistent without force observations. In Fig. 5 (top), we use a box plot to show the spread of F_{mean} across the 9 trials of each condition. The dotted line shows the deployment force-threshold: $F_{th} = 7.5N$. We see that when the force observation was included, contact force was consistently low across gains. However, without force observations, the spread of forces across episodes was high, often exceeding the threshold at higher gains. This highlights the importance of force sensing to enable the policy to effectively modulate the contact force.

D. Success Prediction Analysis

To evaluate the early termination procedure, we ask: **(Q6) How much efficiency is gained when the policy determines when to terminate? (Q7) Does success prediction, trained in simulation, transfer to the real world?**

To answer Q6, we use *Delay Time* (s) to capture efficiency (lower values are better). Delay time measures the time

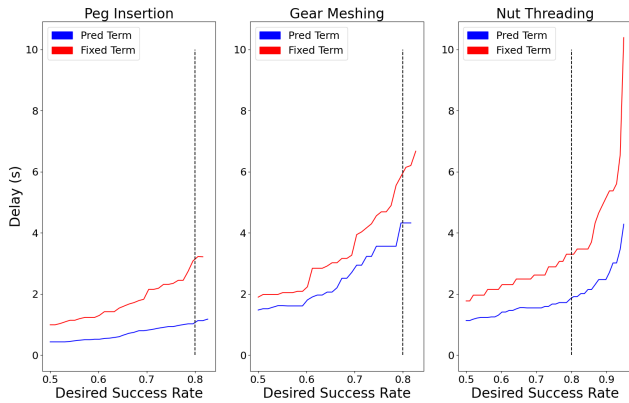


Fig. 6. **Success Prediction Analysis** Relationship between Delay Time and Success Rate for two early termination methods (generated by varying each method’s respective parameter: T or p_{term}). The *Pred Term* method leads to lower delays than the *Fixed Term* method, especially at higher success rates. The vertical line shows a 0.8 success rate.

between when success occurred and when the episode was terminated. We compare the proposed method (*Pred Term*) to a standard termination method that stops the policy after a fixed duration, T (*Fixed Term*). Each method has a parameter that can be tuned to produce a different success rate (fraction of episodes that are successful when terminated). However, this will introduce a trade-off with delay time:

- *Fixed Term* (T): Waiting too long is inefficient while terminating too early will harm success rates.
- *Pred Term* (p_{term}): A high threshold can cause extra delay while a low threshold can affect success rate.

Fig. 6 is a simulated analysis that shows the relationship between *Delay Time* and *Success Rate* for each method. Each line was generated by measuring the success rate and corresponding delay time across a fine discretization of each method’s termination parameter. These were then sorted by success rate and plotted.⁴ As a practitioner, one could choose a desired success rate and find the resulting delay.

Across all tasks, we see that the *Fixed Term* method leads to longer delays, especially at higher success rates (we plot a vertical line to show the 0.8 success rate). The early termination action, a^{ET} , allows for dynamic episode lengths, leading to high success rates with smaller delay times.

Finally to answer Q7, results in Table I show that early termination prediction transferred well to the real world. While the termination method tended to correctly identify successes for all models (high and often perfect recall), we see that precision was best when using force observations for the gear and peg tasks. This shows the benefit of force for sensing task completion: when the bottom of the socket has been reached, or the gear has been fully meshed.

E. Multi-Stage Assembly

To culminate this work, we show that FORGE enables the multi-stage assembly of a planetary gearbox using a simple perception system (see Fig. 7 for the initial and final states). We assume the assembly sequence is known *a priori* and

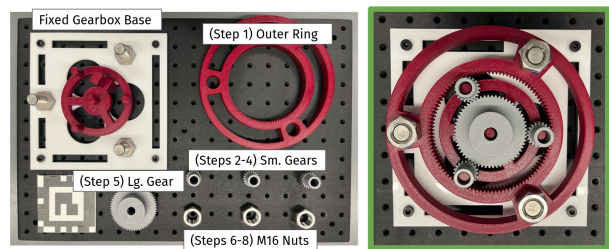


Fig. 7. FORGE policies enable a robot to complete long-horizon tasks such as assembling a planetary gearbox (from initial state [left] to goal state [right, enlarged]).

train FORGE policies for *Small Gear*, *Large Gear*, and *M16 Nut* tasks. We additionally introduce a new *Ring Insertion* task, which must also be robust to orientation estimation noise such that the three bolts align with the holes in the outer ring. Successfully assembling the planetary gearbox requires executing 8 contact-rich primitives.

We ran 5 trials resulting in the following success rates: Ring Insertion (5/5), Small Gear (15/15), Large Gear (3/5), M16 Nut (15/15). Early terminations saved on average 65s in a single trial compared to executing policies for a fixed duration. Overall, the complete assembly succeeded in 3/5 trials where the failures correspond to the large gear insertion (which has to align the teeth of three already inserted small gears). Please see the accompanying video for a demonstration of the multi-stage assembly and the website for more experimental details.

VI. RELATED WORK

Assembly tasks typically involve mating parts with tight clearances and detailed geometries [30], [31]. Various approaches have been proposed to handle pose uncertainty in such tasks. Mechanically, *remote centers of compliance* [32] or chamfers can mitigate small misalignments. Compliant control [33] and strategies such as spiral search [11], [34] have also been used for insertion. These strategies typically consider low noise levels and are task-specific.

Sim-to-Real Transfer: System identification [20] is often time-consuming and difficult to apply to contact-rich tasks. Instead, *dynamics randomization* randomizes parameters such as part friction/stiffness [10], [21], [22], [35], [36], controller gains [12], [21], or F/T observation scale [12], [37]. Even with randomization, excessive forces can occur when deployed. An expert can tune the controller gains at deployment or choose an action-space that is safe by design [7], [25], [38]. Gains can also be adapted online via optimization [23] or an explicit *gain-tuning* model [37].

Similar to FORGE, other works have proposed to use a force-threshold [10], [17], [36]. These works have a fixed threshold during training which is often very large to primarily prevent damage (e.g., 40N). However, especially with small parts, slip can occur with much lower contact forces. Most similar to FORGE, [10] introduces a method to specify the *desired* interaction force at deployment time.

Most prior work focus on insertion-style tasks. We show how the combined application of a force-threshold and dynamics randomization can lead to robust sim-to-real transfer

⁴Similar to an ROC plot, but higher areas above the curve are better.

for a range of tasks, including the complicated nut-threading task. Prior work on sim-to-real for nut-threading [14] focused on large parts (*M48* nuts) that were fixed to the gripper. In addition, we show these techniques are applicable for sim-to-real transfer of early termination procedures. Additional discussion of related works can be found on the website.

Early-Termination: Previous *sim-to-real* approaches execute policies for a fixed duration [7]. Instead, we would like to terminate once success is achieved. For some tasks, success can be manually specified from sensor data [39], [40]. For others, a classifier can be learned from visual data [41], [42]. However, for contact-rich tasks, visual and proprioceptive data alone may be insufficient to determine success [43]. In such cases, the robot can execute actions to verify success [44]. Previous work learns a separate policy to check success *after* task execution [24]. Instead, we jointly trained a policy to predict success *during* task execution.

VII. CONCLUSION

In conclusion, we present FORGE, a method to train sim-to-real policies for robust execution with pose estimation uncertainty. FORGE uses a force threshold and dynamics randomization to learn *safe* exploration behaviours, enabling successful policy execution with up to 5mm of position estimation error. In addition, FORGE can predict early termination, allowing efficient policy execution. In future work, we plan to investigate torque sensing to help develop more efficient search strategies. We also believe research in *real-to-sim* will help automatically tune simulation models for robust transfer in complicated tasks such as nut-threading.

ACKNOWLEDGMENT

The authors thank Chad Kessens, the Seattle Robotics Lab, and the Robust Robotics Group for their valuable feedback.

REFERENCES

- [1] I. Akkaya *et al.*, “Solving rubik’s cube with a robot hand,” *arXiv:1910.07113*, 2019.
- [2] A. Handa *et al.*, “DeXtreme: Transfer of Agile In-hand Manipulation from Simulation to Reality,” in *ICRA*. IEEE, 2023.
- [3] J. Tan *et al.*, “Sim-to-Real: Learning Agile Locomotion For Quadruped Robots,” in *RSS*, 2018.
- [4] J. Hwangbo *et al.*, “Learning agile and dynamic motor skills for legged robots,” *Science Robotics*, 2019.
- [5] Y. Narang *et al.*, “Factory: Fast Contact for Robotic Assembly,” in *RSS*, 2022.
- [6] J. Yoon, M. Lee, D. Son, and D. Lee, “Fast and Accurate Data-Driven Simulation Framework for Contact-Intensive Tight-Tolerance Robotic Assembly Tasks,” *arXiv:2202.13098*, 2022.
- [7] B. Tang *et al.*, “IndustReal: Transferring Contact-Rich Assembly Tasks from Simulation to Reality,” in *RSS*, 2023.
- [8] G. Schoettler and *et al.*, “Meta-reinforcement learning for robotic industrial insertion tasks,” in *IROS*. IEEE, 2020.
- [9] S. Kozlovsky, E. Newman, and M. Zacksenhouse, “Reinforcement Learning of Impedance Policies for Peg-in-Hole Tasks: Role of Asymmetric Matrices,” *IEEE RA-L*, 2022.
- [10] C. Beltran-Hernandez, D. Petit, I. Ramirez-Alpizar, and K. Harada, “Variable compliance control for robotic peg-in-hole assembly: A deep-reinforcement-learning approach,” *Applied Sciences*, 2020.
- [11] S. Chhatpar and M. Branicky, “Search strategies for peg-in-hole assemblies with position uncertainty,” in *IROS*. IEEE, 2001.
- [12] S. Jin, X. Zhu, C. Wang, and M. Tomizuka, “Contact Pose Identification for Peg-in-Hole Assembly under Uncertainties,” in *ACC*. IEEE, 2021.

- [13] K. Van Wyk, M. Culleton, J. Falco, and K. Kelly, “Comparative peg-in-hole testing of a force-based manipulation controlled robotic hand,” *IEEE T-RO*, 2018.
- [14] D. Son, H. Yang, and D. Lee, “Sim-to-Real Transfer of Bolting Tasks with Tight Tolerance,” in *IROS*. IEEE, 2020.
- [15] L. Kaelbling, M. Littman, and A. Cassandra, “Planning and acting in partially observable stochastic domains,” *Artificial intelligence*, 1998.
- [16] Y. Jiang, C. Wang, R. Zhang, J. Wu, and L. Fei-Fei, “TRANSIC: Sim-to-Real Policy Transfer by Learning from Online Correction,” *arXiv:2405.10315*, 2024.
- [17] R. Martín-Martín, M. Lee, R. Gardner, S. Savarese, J. Bohg, and A. Garg, “Variable impedance control in end-effector space: An action space for reinforcement learning,” in *IROS*. IEEE, 2019.
- [18] V. Makoviychuk *et al.*, “Isaac gym: High performance gpu-based physics simulation for robot learning,” *arXiv:2108.10470*, 2021.
- [19] A. Allshire *et al.*, “Transferring dexterous manipulation from gpu simulation to a remote real-world trifinger,” in *IROS*. IEEE, 2022.
- [20] L. Ljung, “System identification,” in *Signal analysis and prediction*. Springer, 1998, pp. 163–173.
- [21] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, “Sim-to-Real Transfer of Robotic Control with Dynamics Randomization,” in *ICRA*. IEEE, 2018.
- [22] O. Spector and M. Zacksenhouse, “Learning Contact-Rich Assembly Skills Using Residual Admittance Policy,” in *IROS*. IEEE, 2021.
- [23] X. Zhang, C. Wang, L. Sun, Z. Wu, X. Zhu, and M. Tomizuka, “Efficient Sim-to-real Transfer of Contact-Rich Manipulation Skills with Online Admittance Residual Learning,” in *CORL*, 2023.
- [24] K. Huang, E. Hu, and D. Jayaraman, “Training Robots to Evaluate Robots: Example-Based Interactive Reward Functions for Policy Learning,” in *CORL*, 2022.
- [25] K. Zhang, M. Sharma, J. Liang, and O. Kroemer, “A modular robotic arm control stack for research,” *arXiv:2011.02398*, 2020.
- [26] R. Petrea, M. Bertoni, and R. Oboe, “On the Interaction Force Sensing Accuracy Of Franka Emika Panda Robot,” in *IECON*. IEEE, 2021.
- [27] K. He, G. Gkioxari, P. Dollár, and R. Girshick, “Mask R-CNN,” in *ICCV*. IEEE, 2017.
- [28] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv:1707.06347*, 2017.
- [29] L. Pinto, M. Andrychowicz, P. Welinder, W. Zaremba, and P. Abbeel, “Asymmetric actor critic for image-based robot learning,” in *RSS*, 2018.
- [30] J. Xu, Z. Hou, Z. Liu, and H. Qiao, “Compare contact model-based control and contact model-free learning,” *arXiv:1904.05240*, 2019.
- [31] Z. Jia, A. Bhatia, R. Aronson, D. Bourne, and M. Mason, “A survey of automated threaded fastening,” *IEEE T-ASE*, 2018.
- [32] S. H. Drake, “Using compliance in lieu of sensory feedback for automatic assembly,” Ph.D. dissertation, MIT, 1978.
- [33] T. Lozano-Perez, M. Mason, and R. Taylor, “Automatic synthesis of fine-motion strategies for robots,” *IJRR*, 1984.
- [34] W. Newman, Y. Zhao, and Y. Pao, “Interpretation of force and moment signals for compliant peg-in-hole assembly,” in *ICRA*. IEEE, 2001.
- [35] A. Apolinarska *et al.*, “Robotic assembly of timber joints using reinforcement learning,” *Automation in Construction*, 2021.
- [36] M. Hebecker, J. Lambrecht, and M. Schmitz, “Towards Real-World Force-Sensitive Robotic Assembly through Deep Reinforcement Learning in Simulations,” in *AIM*. IEEE, 2021.
- [37] X. Zhang, M. Tomizuka, and H. Li, “Bridging the Sim-to-Real Gap with Dynamic Compliance Tuning for Industrial Insertion,” in *ICRA*. IEEE, 2024.
- [38] N. Vuong, H. Pham, and Q. Pham, “Learning Sequences of Manipulation Primitives for Robotic Assembly,” in *ICRA*. IEEE, 2021.
- [39] L. Pinto and A. Gupta, “Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours,” in *ICRA*. IEEE, 2016.
- [40] B. Wen and *et al.*, “You only demonstrate once: Category-level manipulation from single visual demonstration,” *RSS*, 2022.
- [41] Z. Su, O. Kroemer, G. Loeb, G. Sukhatme, and S. Schaal, “Learning manipulation graphs from demonstrations using multimodal sensory signals,” in *ICRA*. IEEE, 2018.
- [42] J. Fu, A. Singh, D. Ghosh, L. Yang, and S. Levine, “Variational inverse control with events: A general framework for data-driven reward definition,” *NeurIPS*, 2018.
- [43] A. Rodriguez and *et al.*, “Failure detection in assembly: Force signature analysis,” in *IEEE CASE*, 2010.
- [44] O. Kroemer, S. Niekum, and G. Konidaris, “A review of robot learning for manipulation,” *JMLR*, 2021.