

ICRT: In-Context Imitation Learning via Next-Token Prediction

Max (Letian) Fu^{1*}, Huang Huang^{1*}, Gaurav Datta^{1*}, Lawrence Yunliang Chen¹
Will Panitch¹, Fangchen Liu¹, Hui Li², Ken Goldberg¹

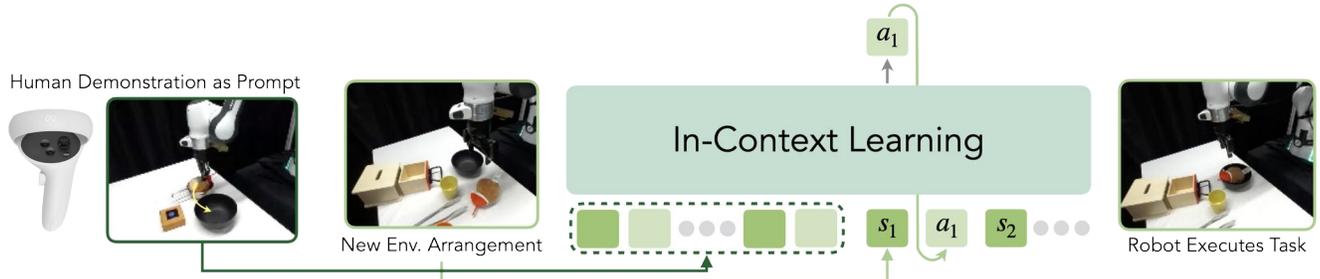


Fig. 1: In-Context Robot Transformer (ICRT) is a novel robot foundation model designed for in-context imitation learning without fine-tuning. In ICRT, the prompt includes 2-3 image-action sequences of teleoperated demonstration trajectories with the object of interest. Given a previously unseen task demonstration, ICRT performs next-token prediction to execute the task with real-time continuous control.

Abstract—In-context imitation learning is the capability to perform novel tasks when prompted with task demonstration examples. In-Context Robot Transformer (ICRT) is a causal transformer that performs autoregressive prediction on sensorimotor trajectories, which include images, proprioceptive states, and actions. This approach supports flexible and training-free execution of new tasks at test time. Experiments with a Franka Emika robot demonstrate that ICRT can adapt to new environment configurations that differ from both the prompt and the training data. In a multi-task environment setup, ICRT significantly outperforms current state-of-the-art robot foundation models on generalization to unseen tasks. Code, data, and appendix are available on <https://icrt.dev>.

I. INTRODUCTION

Learning-based single and multi-task robot policies have become increasingly capable [1–10]. This improvement in robot capabilities can largely be attributed to progress in related fields, particularly in vision and language modeling. Inspired by the recent development of large language models (LLMs) and large vision models (LVMs) [11–13], which formulate natural language processing and vision problems as next-token-prediction, recent works also have formulated robot learning as next-token-prediction problems and achieved state-of-the-art performance [7, 8, 14, 15]. Concurrently, there has been a surge in collecting large-scale generalist robot datasets [16–23] and pre-training models on them [15, 24–27].

Despite being pre-trained on large datasets and showing some generalization ability, it is still challenging to teach these models to perform unseen tasks in different environments without additional training. New human demonstrations via teleoperation or new data collected from hand-crafted motion primitives, as well as subsequent fine-tuning of pretrained models, are often needed to complete the new tasks. This process adds time and complexity to the workflow, making it challenging to apply these methods in real-world

environments. Ideally, given some demonstrations, the robot should be able to perform a new task *immediately*. In their respective domains, LLMs and LVMs [11–13] have exhibited a similar ability, using *in-context learning*, where the model to adapt to the task corresponding to the prompt provided at inference time without additional fine-tuning.

Is the in-context learning capability of next-token prediction models limited to vision and language domains? In this paper, we introduce the In-Context Robot Transformer (ICRT) to explore how next-token prediction models can be extended to perform real-robot in-context learning. For ICRT, the context is provided as robot trajectories corresponding to a new task. The model infers from this context to perform the task in a different environment configuration without requiring additional training. Each robot trajectory is a sequence of image observations, robot proprioceptive states, and actions. This trajectory implicitly encodes task primitives and the objects the robot needs to interact with. The model extracts this information from the prompt and then executes actions following a similar pattern in the new environment.

Compared to existing few-shot imitation learning approaches, ICRT offers a framework that avoids complicated loss functions, prior knowledge, and the need to identify key points or key frames, and operates directly on raw robot trajectories for continuous control. Additionally, unlike existing next-token prediction models for robot learning, ICRT features a long context window, allowing it to train on multiple sensorimotor trajectories from the same task and use one or more sensorimotor trajectories as prompts during inference.

Importantly, we find that certain properties of the pre-training dataset are crucial for enabling in-context learning on real robots. Specifically, datasets that allow multiple tasks to be performed from the same initial observation are particularly beneficial. Unlike existing single-task datasets or many multi-task datasets where each environment uses one unique object for robot interaction, complex environments require the model to rely on the prompt to correctly identify the task and

*Equal contribution, ¹University of California, Berkeley, ²Autodesk Research

determine the appropriate object for interaction.

We make the following contributions:

- 1) We introduce ICRT, a novel robot foundation model that performs in-context learning on a real robot, which can effectively learn from prompt trajectories and perform unseen tasks without fine-tuning.
- 2) We provide a new multi-task robot dataset and a training paradigm for fostering multi-task and in-context capability at inference time.
- 3) Physical experiments on a Franka Emika robot suggest that ICRT can learn from the provided context and perform unseen tasks at various generalization levels.

II. RELATED WORKS

A. Multi-Task Imitation Learning for Robotics

Imitation learning is an effective paradigm for equipping robots with various skills. The simplest algorithm in this domain, behavior cloning, has been successful across a wide range of tasks [28–30]. In recent years, alternative architectures such as energy-based models [31] and diffusion models [1] have also been proposed. Typically, these approaches require training a *separate* model for each task, although multi-task policies can be distilled from these task-specific models after training [32].

Recent advancements have shown that using transformers for next-token prediction in sequence modeling has been particularly effective in both language and vision domains, especially for *multi-task learning* [12, 33, 34]. In pursuit of developing generalist agents and multi-task robot policies, robot action planning is framed as a next-token prediction task using transformer-based architectures trained on large multi-task robot datasets [5, 7, 8, 14, 15, 35–40]. Octo [15] and OpenVLA [14] represent the state-of-the-art among multi-task robotic policies. Octo [15] conditions on both goal images and language instructions, utilizing a transformer architecture with a diffusion head that fuses these inputs with current image observations to predict robot actions. OpenVLA [14] conditions solely on language instructions, fine-tuning a pre-trained vision-language model to predict robot actions based on visual observations and language inputs.

B. In-Context Robot Learning

Despite training on large datasets, multi-task policies often struggle when faced with new objects, tasks, or environments, frequently requiring fine-tuning. Meta-learning has been shown to increase fine-tuning efficiency for generalization to new tasks [41–43], which has led to progress in few-shot imitation learning. To simplify the application of learned policies in the real world, recent approaches focus on methods that avoid fine-tuning model parameters for task generalization. Instead, these methods inform the model by providing demonstrations of tasks in the prompt [44, 45]. Brown *et al.* [33] refers to this as “in-context learning”, distinguishing it from approaches that rely on subsequent stages of parameter fine-tuning.

Prior in-context robot learning methods often employ contrastive learning to train context encoders, which identify

the most similar training tasks to the test task in the latent space [37, 46]. However, how to effectively integrate these methods within the next-token-prediction framework remains unclear. Valassakis *et al.* [47] achieved one-shot in-context learning by training a visual servoing network to align the robot’s end-effector with the object’s relative pose during the demonstration, but this approach requires an additional object segmentation model. Di Palo *et al.* [44] introduced Keypoint Action Tokens, demonstrating in-context imitation learning using a large language model by representing demonstration trajectories as 3D coordinates with few-shot prompting. Unlike these approaches, ICRT operates without additional perception modules, processing raw image observations directly. Additionally, Vid2Robot [45] developed an encoder-decoder transformer that uses a demonstration video of a human and the current robot state as the prompt to generate robot actions. However, this method requires many auxiliary losses while ICRT uses a simple next-token prediction loss. More recently, methods like MOKA [48] have explored prompting closed-source VLMs with in-context examples to generate waypoints for manipulation tasks. However, relying on off-the-shelf VLMs can introduce significant latency in control, potentially slowing task execution and increasing overall completion time.

In this paper, we focus on enhancing next-token-prediction models to perform real-world in-context imitation learning with robots. ICRT bypasses the need for additional context encoders by directly using robot sensorimotor trajectories from new tasks as prompts for the transformer-based model. ICRT is closely related to the seminal works, One-Shot Imitation Learning [49] and Prompting Decision Transformer [50]. [49] predicts the next action by applying cross-attention between a demonstration sequence on a new task and the current environment’s state, while [50] employs a short trajectory prompt to encode task-specific information for guiding policy generation in offline reinforcement learning, using full state information and known reward functions. While both of these methods show their effectiveness in simulation, it is hard to have full-state information and known reward functions for all real-world robot manipulation tasks. To address these challenges, ICRT does not model rewards, utilizes a significantly longer context window, and demonstrates in-context learning capabilities in physical experiments using image observations.

III. PROBLEM STATEMENT

We consider in-context imitation learning in a real-robot manipulation setting. The first step is to train a model with in-context learning capabilities using a multi-task robotic dataset. At test time, the model can handle unseen tasks in novel environment configurations by using a few new human-teleoperated robot demonstrations as prompts. Here, environment configuration refers to the objects present in the scene and their spatial arrangement. Notably, this process is achieved *without any additional fine-tuning* on the new demonstrations.

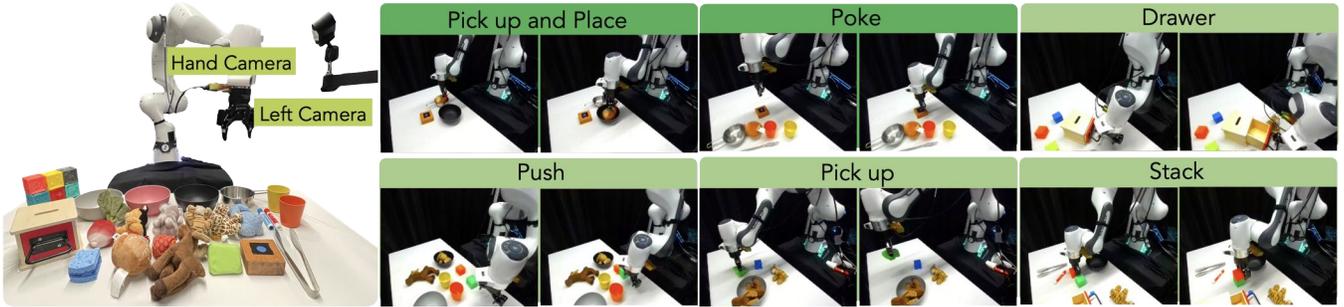


Fig. 2: Our physical setup with the Franka Emika robot, the wrist and side camera and the objects used in training and evaluation. We consider 6 primitives for training and choose “pick up and place” and “poke” as the primitives for evaluation (dark green).

We define *motion primitives* as distinct robot actions utilized to accomplish various *tasks*. Each task is characterized by 1) a specific motion primitive and 2) the set of objects the robot interacts with using that primitive. By altering the environment configuration at test time compared to the one in the prompt, we assess the model’s ability to select the appropriate motion primitive and identify the correct object for interaction. In this work, we consider new tasks to be tasks involving unseen objects but using motion primitives from the training data (for example, training on picking up a tiger toy and testing on picking up a cube).

We make the following assumptions for ICRT experiments:

- 1) The model is trained on a diverse multi-task demonstration dataset. Each trajectory contains RGB observations from a fixed camera and a wrist-mounted camera, proprioception and action.
- 2) Tested primitives are within distribution and the reachable workspace of the robot.

IV. APPROACH

In this section, we first introduce the data composition to facilitate in-context imitation learning. We then introduce the architecture and training objective for the transformer-based policy to effectively leverage the data.

A. Data Formulation

For model training, we consider a dataset \mathcal{D} of visuomotor trajectories \mathcal{T} . Each trajectory of length t is a sequence of camera images i_t , proprioceptive robot states s_t , and actions a_t : $\mathcal{T} = (i_1, s_1, a_1, \dots, i_t, s_t, a_t)$. We use the absolute end-effector pose as the robot’s proprioceptive state and the delta robot end-effector pose between time steps as the action, which consists of delta translation, delta rotation and the continuous gripper action. We assume a known grouping of the trajectories so that the dataset can be partitioned into disjoint sets of tasks $\mathcal{D} = \bigcup_{k=1}^K \mathcal{S}_k$, with $\mathcal{S}_k \cap \mathcal{S}_\ell = \emptyset$, $k \neq \ell$, where $\mathcal{S}_k = \{\mathcal{T}_{k_1}, \dots, \mathcal{T}_{k_n}\}$. In practice, this grouping can be retrieved from the semantic labels of the dataset. In this work, we utilize the existing large robotic dataset DROID [51] and a multi-task dataset manually collected in our robot setup, which we name ICRT-Multi-Task (ICRT-MT).

DROID [51] is a joint effort from different organizations and contains 76k real-world demonstrations. We randomly sample 10k demonstrations from DROID after filtering

out demonstrations shorter than 30 steps and longer than 450 steps. DROID dataset labels the task through human-specified language instructions, which may be different for the same task. We organized the DROID data by grouping demonstrations based on their language instructions CLIP text embedding cosine similarity. Specifically, we use a threshold of 0.9 for grouping demonstrations. To further facilitate in-context learning, we make sure that each task group contains at least 4 trajectories so that there are sufficient trajectories to serve as prompts for each other. This results in roughly 2k trajectories that we use for pre-training ICRT.

Many trajectories in the DROID dataset are collected in a single-task setup, where only one task can be completed in a scene (e.g. only one object is presented). In such a setup, the model can learn a shortcut solution to perform the task, by only focusing on the current observation but not the prompt. Therefore, multi-task data is crucial for the model to learn from the prompt. We manually collected a multi-task dataset ICRT-Multi-Task (ICRT-MT) using the DROID setup (Figure. 2). This dataset has 1098 trajectories in total, and contains 29 tasks with 6 primitives: picking, pick-and-place, stacking, pushing, poking, opening and closing drawers. Objects used in the data collection and examples of the primitives are shown in Figure. 2. In ICRT-MT, each environment is set so that there exist more than 2 possible tasks for the current observation so that the model has to distinguish and learn the motion from the prompt.

During the training, for each trajectory, we independently apply vision augmentation on the image observations by augmenting the brightness and contrast. We additionally apply random crops and scaling to the side camera observation. We also apply proprioception noise sampled from a normal Gaussian distribution $\mathcal{N}(0, 0.005)$. For each epoch, we randomly shuffle the order of trajectories from each task and concatenate them to form the training sequence. For each batch, we sample a subsequence of length $L = 512$ as the input to the model, where L is the sequence length defined as the number of observation, state, and action tuples. In practice, 512 steps usually contain up to 5 trajectories from the same task. We randomly select the first k trajectories and label them as the prompt within the sequence. At least one complete trajectory is included in the prompt. This data grouping aims to capture inter-trajectory patterns, encouraging the model to generate action conditioned on the prompt

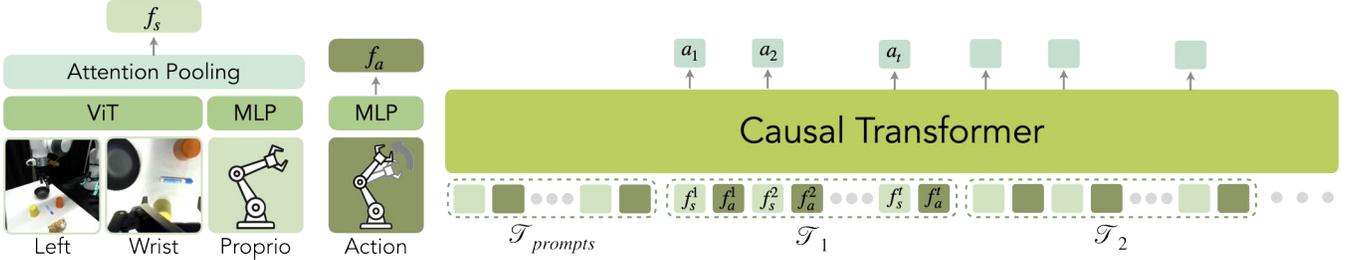


Fig. 3: Method Overview: (Left) We encode camera observations with a pre-trained vision transformer. Additionally, we encode proprioception with an MLP. We concatenate the visual latent and the proprioception’s latent and use attention pooling to extract a feature f_s as the current state representation. We encode the current action with an MLP to get f_a . (Right) We concatenate multiple trajectories of the same task and randomly sample the first k trajectories as the prompt. A causal transformer autoregressively predicts the next series of tokens. We decode the tokens that are at the position of the state features to generate the next $h = 16$ action via an MLP.

trajectories. This approach differs from traditional behavior cloning methods, which typically use short input sequences that focus on modeling intra-trajectory behaviors.

B. Model Architecture

We construct the ICRT model with three parts: a pre-trained vision encoder, a series of projectors for each input modality, and a causal transformer backbone (Figure 3).

Vision Encoder The model processes multi-view image observations through a pre-trained vision transformer. However, most visual pre-trained networks are trained on ImageNet or human videos [24, 27, 52, 53], which exhibit a significant domain gap when compared to typical images from robot datasets, where the images frequently include robots or grippers. To minimize the domain gap, we pre-train a vision transformer [54] (ViT-Base) on an equal mix of ImageNet [55] and Open X-Embodiment [40] data, using CrossMAE as an efficient pre-training method [56]. During the training of the ICRT model, we freeze the vision encoder for efficiency. The vision encoder outputs the entire feature map for each of the cameras and is then fed into the proprioception projector (Figure 3 left).

Modality-Specific Projectors To project image observations, the robot’s proprioceptive state, and actions into a shared latent space for sequence modeling, we design modality-specific projectors. At each timestep, the model takes as input a token representing either an observation or an action. To produce a single state token that captures fine-grained visual information and the proprioceptive state of the robot, we use attention pooling [57] between all visual tokens from a single camera’s observation and a proprioception embedding produced by a multi-layer perceptron (MLP). The resulting embeddings for each camera are concatenated to produce a single state token f_s^t of dimension equal to the transformer latent dimension. Similar to proprioception, the action is embedded with an MLP into an action token f_a^t . This process produces a sequence of state and action tokens that are passed into the transformer.

Transformer Model The encoded sequence of state and actions is passed into a Transformer model [58], following the design of Llama2 [12]. The transformer takes as input the sequence of state and action features $(f_s^1, f_a^1, \dots, f_s^t, f_a^t)$ that are produced by the modality-specific projectors. We

add MLP decoders to produce state and action outputs from the last layer of the transformer at the appropriate positions. We denote the transformer with the decoder heads as g_θ . Therefore, the desired outputs are the shifted sequence of proprioceptive states and actions $(a^1, s^2, a^2, \dots, a^t, s^{t+1})$. This naturally forms a next token prediction problem, as $g_\theta(f_s^1)$ predicts a^1 and $g_\theta(f_s^1, f_a^1, \dots, f_s^n)$ predicts a^{n+1} . In practice, we find it beneficial to predict the next h actions at each time step, and use temporal ensembling [2] to execute the final action.

Inspired by Octo [15] and vision transformers [54], we consider a randomly initialized Llama2 model of 12 layers with a latent dimension of 768, which we name *Llama2-Base*. In addition, multiple works have shown that multimodal inputs can be aligned to large-language models [8, 34, 59–61]. Multimodal language model, Palm-E [10] has shown success in enhancing generalization when being directly incorporated into robotic control [8]. Therefore, we also investigate the effectiveness of using a large-language model for in-context robot learning by initializing the transformer with a pre-trained Llama2-7B. Due to the large domain gap between natural language and robot trajectories, a frozen language model may not be sufficient. Therefore, similar to prior work in multimodal alignment, we fine-tune the language model with LoRA [62], with a rank of 32. Due to compute resource limitations, we are unable to fully fine-tune the model.

Loss Function To provide more supervision signals so that the model can better respond to the trajectory “prompt” we provide at test time, we reference works in training multi-turn conversation chatbots [34, 63], where they only compute loss on the response generated by the chatbot, instead of the prompt. Recall that in Section IV-A, we randomly sample the subsequence of the concatenated trajectories as the prompt trajectory. Analogously, we only compute action prediction with L1-loss for the actions after the prompt trajectories.

Inference The simplicity of the next-token prediction objective makes inferencing with ICRT straightforward at test time. As shown in Figure. 4, we provide one or more human-teleoperated demonstrations in the form of robot sensorimotor trajectories (formatted identically to the training data), along with the current image observations and the robot’s proprioceptive state as inputs. The model then predicts the next action, which is executed by the robot. After each

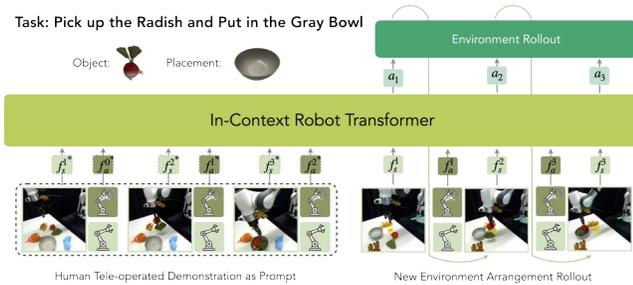


Fig. 4: Example inference pipeline of ICRT on the task of picking up the radish and putting in the gray bowl. A human teleoperated demonstration trajectory consisting of image observations, proprioception and actions are provided as the prompt. ICRT takes the prompt and the current observation in a different environment to accomplish the task.

action, the policy receives updated image observations and proprioceptive state, allowing it to iteratively predict and execute subsequent actions.

A key advantage of this framework is its use of the transformer’s sequential processing capability. Instead of reprocessing the entire sequence history for each model evaluation, as seen in previous works [7, 8, 14, 15], the model employs a key-value (KV) caching mechanism, as discussed in [12]. This mechanism stores previous outputs, allowing the model to compute only the outputs for the new token. This approach significantly reduces computational overhead, lowering the complexity from quadratic to linear relative to the sequence length.

C. Training and Inference Details

All models are trained on 4 NVIDIA A100 80GB GPUs. ICRT pre-training on DROID takes 56 minutes and fine-tuning on ICRT-MT takes 18 hours. ICRT-Llama7B takes roughly 28 hours to finetune. All tests are performed on a workstation with NVIDIA RTX 3090Ti and Intel i5-12400F with 64GB memory. By using KV cache, ICRT can inference at 39.6 Hz and ICRT-Llama2 can inference at 10.7 Hz, allowing them to perform real-time close-loop control.

V. EXPERIMENTS

In this section, we design an experimental setup to evaluate the in-context learning capabilities of the proposed models and compare them against several baselines. Instead of focusing on the difficulty of learning a specific task primitive, we design the experiments to assess the policy’s ability to accomplish novel tasks based on the provided prompt trajectories.

Experiment Design We consider two action primitives: a *pick-and-place* primitive and a *poking* primitive. For each action primitive, we design *six unseen tasks* (as defined in Section III), with three tasks evaluating *in-domain* object generalization (selected from *yellow cube*, *red cube*, *black cube*, *pink bowl*, and *blue bear* and three objects *unseen* during training (selected from *radish*, *blue sponge*, *grey dog*, and *black dog*).

Each task has five difficulty tiers. In the pick-and-place task, the model must identify the correct object to grasp and where to place it in a multi-object or multi-placement scenario. The tiers include: 1) no distractors, 2) one distractor

	Pick and Place	Poke	Average
Goal Condition	33.3 (± 6.5)	6.7 (± 4.6)	20.0 (± 4.3)
Octo [15]	5.0 (± 2.7)	13.3 (± 6.2)	9.2 (± 3.5)
OpenVLA [14]	11.7 (± 4.6)	3.3 (± 3.3)	7.5 (± 2.9)
ICRT	65.0 (± 7.3)	93.3 (± 4.6)	79.2 (± 4.6)

TABLE I: Success Rates. ICRT outperforms two state-of-the-art robot foundation models that are conditioned on goals or language in both pick-and-place and poking tasks. We evaluated each task primitive using six tasks not seen during training, conducting five trials per task for a total of 30 trials per primitive and 60 trials overall to calculate average performance. For each model, we report the mean success rate for each task, the overall success rate, and the corresponding standard error in parentheses.

object, 3) two distractors, 4) three distractors, and 5) one distractor placement position. In the poking task, the robot closes the gripper, pokes the object, lifts the end-effector, and opens the gripper, with tiers involving 0-4 distractor objects.

The pick-and-place task is scored with 0.5 for a correct pick and 1 for a successful placement. In the poking task, failure is marked if the wrong object is poked. The model has 25 seconds (375 steps) for retries. Each difficulty level is attempted once, and we report the average success rate per task, along with the average success rate and standard deviation across the six tasks for each action primitive.

Models The default **ICRT** is a randomly initialized Llama2-Base model pretrained on DROID and fully fine-tuned on ICRT-MT. We evaluate the impact of model initialization and training datasets by introducing the following three variants: 1) **ICRT-Llama2**, a pre-trained Llama2-7B language model fine-tuned on ICRT-MT with LoRA; 2) **ICRT (DROID)**, a randomly initialized Llama2-Base model trained only on the DROID dataset; and 3) **ICRT (MT)**, a randomly initialized Llama2-Base model trained only on the ICRT-MT dataset.

We consider 3 baseline models. We train a goal-conditioned policy, where the goal observations are always prepended to the sequence, and each sequence is from one trajectory. This resembles the normal goal-conditioned imitation learning setup. Additionally, we finetune Octo [15], the state-of-the-art goal-image and language conditioned policy, and OpenVLA [14], the state-of-the-art language conditioned multi-task imitation learning policy. Octo is fine-tuned using their official fine-tuning recipe. We incorporate action chunking into OpenVLA by asking it to predict the next 16 actions, which performs better than vanilla OpenVLA which predicts only the next step. Both of these methods are representative of robot policies that use next-token prediction objectives.

Prompt Generation For each task, we collect 3 demonstrations (with zero, one distractor object, a distractor placement for pick-and-place, or two distractor objects for poking) as the prompt in total before running the experiment. Please refer to the website for a visual example. During testing, a random demonstration is drawn as a prompt to assess the model’s ability to generalize to different prompts. It’s important to note that the environment setup during policy rollout *differs* from the prompts’ setup, ensuring that the evaluation measures the model’s understanding of task-relevant information from the prompt, rather than simply copying actions from it.

Results We present the results in Table I. For the pick-and-place primitive, we observe that the goal-conditioned

	Pick and Place	Poke	Average
ICRT-Llama2	43.3 (± 7.9)	73.3 (± 8.2)	58.3 (± 6.0)
ICRT (DROID only)	0.0 (± 0.0)	0.0 (± 0.0)	0.0 (± 0.0)
ICRT (MT only)	76.7 (± 7.1)	70.0 (± 8.5)	73.3 (± 5.5)
ICRT +Prompt Loss	21.7 (± 6.2)	23.3 (± 7.9)	22.5 (± 5.0)
ICRT	65.0 (± 7.3)	93.3 (± 4.6)	79.2 (± 4.6)

TABLE II: Ablation Study. We ablate three key design choices: fine-tuning ICRT using a language model, training solely on the relevant dataset or the fine-tuning dataset, and the impact of including prompt loss in the training process. We report the mean success rates and their corresponding standard errors for the pick-and-place and poke tasks, as well as the overall average performance.

policy generally succeeds in identifying the correct object to grasp when no distractor objects are present. However, its performance degrades significantly as the number of distractors increases. When the goal image only specifies the task but not the specific way to achieve it in the current environment, goal-conditioned policies often fail to execute the task effectively.

Octo struggles with determining which object to interact with and where it should be placed, highlighting the challenges posed by our experimental setup for multi-task policies. OpenVLA, while often moving towards the correct object, frequently fails in grasping the object or mistakenly performs the wrong task (e.g., grasping instead of poking, and vice versa). This indicates that OpenVLA may require more demonstrations (more than 50) per task to achieve better performance, and that relying solely on language conditioning may not be sufficient for generalization to new tasks.

The results suggest that ICRT outperforms the goal-conditioned policy in identifying the correct object to pick up and the appropriate placement location. The poking task presents a significant challenge for the goal-conditioned policies, as the goal position often closely resembles the start configuration. However, after conditioning on the prompt trajectory, ICRT is able to correctly identify the task as poking, and the results indicate that it consistently reaches the correct target object while ignoring distractors. Despite this, we do observe some failure modes with ICRT, such as missing the grasp of the target object, grasping the wrong object, or placing objects in incorrect locations. Specifically, when a distractor object shares the same color but has a different shape, the model struggles to accurately determine which object to grasp. This implies that additional fine-tuning of the vision encoder might be required to enhance model performance, a conclusion also reached by OpenVLA [14].

VI. ABLATIONS

In this section, we provide additional experiments presented Table II that ablate on a few core design choices. We provide additional ablation studies on the official website.

A. Model Initialization

We conducted ablation studies to examine the impact of using a pretrained Llama2 on language data and fine-tune it for robot sensorimotor sequence modeling. The results, presented in Table II, show that although ICRT-Llama2-7B achieves a lower training loss, its performance is worse compared to its smaller counterparts. This discrepancy may be attributed

to a lower inference frequency of ICRT-Llama2 (10.7 Hz vs 39.6 Hz). Future work can focus on optimizing the inference speed of ICRT-Llama2 to improve performance.

B. Training Dataset

We find that training on the DROID subset (see Section IV-A) is insufficient for completing any of the test tasks; the policy (ICRT (DROID)) shows no progress across all tasks. This suggests that although the DROID subset may offer greater visual diversity, the unique structure of ICRT-MT—where multiple tasks are performed from the same initial observation—is particularly beneficial in developing the in-context learning capabilities of a next-token prediction robot model.

ICRT (MT) shows similar performance to ICRT that is pre-trained on DROID, especially for the pick-up and place primitive, even surpassing ICRT on the *put radish in grey bowl* task. However, ICRT (MT) does not perform as well on the poking primitive. The results suggest that it may be beneficial to pre-train the autoregressive model on a large dataset, as a diverse dataset may help the transformer to perform better alignment between visual features and control.

C. No Prompt Loss

Following the design of many multi-turn conversation large language models or vision language model fine-tuning works [34, 63–65], we do not calculate the loss for the predicted action in the prompt trajectories but only do so on the predictions after the prompt trajectories. We mark the model that calculates loss on the prompt as **ICRT +Prompt Loss** and the default model as **ICRT**. The results are shown in Table II. We find that only predicting the trajectories after the designated prompt trajectories can significantly improve the model’s performance. We hypothesize that in the situation where there is a loss on the prompt trajectories, the model is forced to do unconditional generation based on current observations for those prompts. This may cause the model to stop paying attention to the prompt, especially when there are multiple possible tasks available.

VII. LIMITATIONS AND CONCLUSION

While results suggest that ICRT learns from the prompt trajectories and generalizes to unseen objects, tasks, and certain primitives that resemble the ones in training¹, it is still unclear how to generalize to completely unseen action primitives. Future works should investigate how scaling model capacity and scaling dataset can help with primitive-level generalization. In addition, ICRT assumes a fixed robot morphology with a fixed impedance controller. Future works can also investigate how to facilitate transfer between different robot morphologies by learning a unified policy on different robots. ICRT-Llama2 has a low inference frequency which may contribute to its low performance. We hope to speed up ICRT-Llama2 at inference time in the future.

¹Please refer to the appendix and videos on the website for more detail.

REFERENCES

- [1] C. Chi *et al.*, “Diffusion policy: Visuomotor policy learning via action diffusion,” *arXiv preprint arXiv:2303.04137*, 2023.
- [2] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn, “Learning fine-grained bimanual manipulation with low-cost hardware,” *arXiv preprint arXiv:2304.13705*, 2023.
- [3] C. Lynch *et al.*, “Interactive language: Talking to robots in real time,” *IEEE Robotics and Automation Letters*, 2023.
- [4] S. Reed *et al.*, “A generalist agent,” *arXiv:2205.06175*, 2022.
- [5] D. Shah *et al.*, “ViNT: A Foundation Model for Visual Navigation,” in *7th Annual Conference on Robot Learning (CoRL)*, 2023.
- [6] H. Bharadhwaj, J. Vakil, M. Sharma, A. Gupta, S. Tulsiani, and V. Kumar, “RoboAgent: Towards sample efficient robot manipulation with semantic augmentations and action chunking,” *arxiv*, 2023.
- [7] A. Brohan *et al.*, “Rt-1: Robotics transformer for real-world control at scale,” *arXiv:2212.06817*, 2022.
- [8] A. Brohan *et al.*, “RT-2: Vision-language-action models transfer web knowledge to robotic control,” *arXiv preprint arXiv:2307.15818*, 2023.
- [9] X. Chen *et al.*, *Pali-x: On scaling up a multilingual vision and language model*, 2023. arXiv: 2305.18565 [cs.CV].
- [10] D. Driess *et al.*, “Palm-e: An embodied multimodal language model,” *arXiv:2303.03378*, 2023.
- [11] J. Achiam *et al.*, “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [12] H. Touvron *et al.*, “Llama 2: Open foundation and fine-tuned chat models,” *arXiv preprint arXiv:2307.09288*, 2023.
- [13] Y. Bai *et al.*, “Sequential modeling enables scalable learning for large vision models,” *arXiv preprint arXiv:2312.00785*, 2023.
- [14] M. J. Kim *et al.*, *Openvla: An open-source vision-language-action model*, 2024. arXiv: 2406.09246 [cs.RO].
- [15] Octo Model Team *et al.*, “Octo: An open-source generalist robot policy,” in *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, 2024.
- [16] A. Depierre, E. Dellandrea, and L. Chen, “Jacquard: A large scale dataset for robotic grasp detection,” in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, IEEE, 2018, pp. 3511–3516.
- [17] D. Kalashnikov *et al.*, “Scalable deep reinforcement learning for vision-based robotic manipulation,” in *CoRL*, 2018.
- [18] S. Levine, P. Pastor, A. Krizhevsky, J. Ibarz, and D. Quillen, “Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection,” *IJRR*, 2018.
- [19] C. Eppner, A. Mousavian, and D. Fox, “ACRONYM: A large-scale grasp dataset based on simulation,” in *2021 IEEE Int. Conf. on Robotics and Automation, ICRA*, 2020.
- [20] N. M. M. Shafiqullah *et al.*, *On bringing robots home*, 2023. arXiv: 2311.16098 [cs.RO].
- [21] H.-S. Fang *et al.*, “RH20T: A robotic dataset for learning diverse skills in one-shot,” in *RSS 2023 Workshop on Learning for Task and Motion Planning*, 2023.
- [22] F. Ebert *et al.*, “Bridge data: Boosting generalization of robotic skills with cross-domain datasets,” *arXiv:2109.13396*, 2021.
- [23] H. Walke *et al.*, *Bridgedata v2: A dataset for robot learning at scale*, 2023. arXiv: 2308.12952 [cs.RO].
- [24] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta, “R3m: A universal visual representation for robot manipulation,” *arXiv:2203.12601*, 2022.
- [25] T. Xiao, I. Radosavovic, T. Darrell, and J. Malik, “Masked visual pre-training for motor control,” *arXiv:2203.06173*, 2022.
- [26] Y. J. Ma, S. Sodhani, D. Jayaraman, O. Bastani, V. Kumar, and A. Zhang, “Vip: Towards universal visual reward and representation via value-implicit pre-training,” *arXiv preprint arXiv:2210.00030*, 2022.
- [27] I. Radosavovic, T. Xiao, S. James, P. Abbeel, J. Malik, and T. Darrell, “Real-world robot learning with masked visual pre-training,” *arXiv:2210.03109*, 2022.
- [28] D. A. Pomerleau, “Alvin: An autonomous land vehicle in a neural network,” in *NeurIPS*, D. Touretzky, Ed., vol. 1, Morgan-Kaufmann, 1988.
- [29] B. D. Argall, S. Chernova, M. Veloso, and B. Browning, “A survey of robot learning from demonstration,” *Robotics and autonomous systems*, vol. 57, no. 5, pp. 469–483, 2009.
- [30] S. Levine, C. Finn, T. Darrell, and P. Abbeel, “End-to-end training of deep visuomotor policies,” *JMLR*, 2016.
- [31] P. R. Florence *et al.*, “Implicit behavioral cloning,” in *CoRL*, 2021.
- [32] H. Ha, P. Florence, and S. Song, “Scaling up and distilling down: Language-guided robot skill acquisition,” in *Conference on Robot Learning*, PMLR, 2023, pp. 3766–3777.
- [33] T. B. Brown *et al.*, “Language models are few-shot learners,” *NeurIPS*, 2020.
- [34] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” in *NeurIPS*, 2023.
- [35] I. Radosavovic, B. Shi, L. Fu, K. Goldberg, T. Darrell, and J. Malik, “Robot learning with sensorimotor pre-training,” *arXiv:2306.10007*, 2023.
- [36] I. Radosavovic *et al.*, *Humanoid locomotion as next token prediction*, 2024. arXiv: 2402.19469 [cs.RO].
- [37] E. Jang *et al.*, “Bc-z: Zero-shot task generalization with robotic imitation learning,” in *Conference on Robot Learning*, 2022.
- [38] Y. Jiang *et al.*, “VIMA: General robot manipulation with multimodal prompts,” *International Conference on Machine Learning (ICML)*, 2023.
- [39] S. Reed *et al.*, “A generalist agent,” *arXiv preprint arXiv:2205.06175*, 2022.
- [40] E. Collaboration *et al.*, *Open x-embodiment: Robotic learning datasets and rt-x models*, 2024. arXiv: 2310.08864 [cs.RO].
- [41] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” in *International conference on machine learning*, PMLR, 2017, pp. 1126–1135.
- [42] C. Finn, T. Yu, T. Zhang, P. Abbeel, and S. Levine, “One-shot visual imitation learning via meta-learning,” in *Conference on robot learning*, PMLR, 2017, pp. 357–368.
- [43] M. Xu, Y. Lu, Y. Shen, S. Zhang, D. Zhao, and C. Gan, “Hyperdecision transformer for efficient online policy adaptation,” *arXiv preprint arXiv:2304.08487*, 2023.
- [44] N. Di Palo and E. Johns, “Keypoint action tokens enable in-context imitation learning in robotics,” *arXiv preprint arXiv:2403.19578*, 2024.
- [45] V. Jain *et al.*, “Vid2robot: End-to-end video-conditioned policy learning with cross-attention transformers,” *arXiv preprint arXiv:2403.12943*, 2024.
- [46] Z. Mandi, F. Liu, K. Lee, and P. Abbeel, *Towards more generalizable one-shot visual imitation learning*, 2022. arXiv: 2110.13423 [cs.RO].
- [47] E. Valassakis, G. Papagiannis, N. Di Palo, and E. Johns, “Demonstrate once, imitate immediately (dome): Learning visual servoing for one-shot imitation learning,” in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, IEEE, 2022, pp. 8614–8621.
- [48] K. Fang, F. Liu, P. Abbeel, and S. Levine, “Moka: Open-world robotic manipulation through mark-based visual prompting,” *Robotics: Science and Systems (RSS)*, 2024.
- [49] Y. Duan *et al.*, “One-shot imitation learning,” *Advances in neural information processing systems*, vol. 30, 2017.
- [50] M. Xu *et al.*, “Prompting decision transformer for few-shot policy generalization,” in *international conference on machine learning*, PMLR, 2022, pp. 24 631–24 645.
- [51] A. Khazatsky *et al.*, *Droid: A large-scale in-the-wild robot manipulation dataset*, 2024. arXiv: 2403.12945 [cs.RO].
- [52] A. Majumdar *et al.*, “Where are we in the search for an artificial visual cortex for embodied intelligence?” *arXiv preprint arXiv:2303.18240*, 2023. arXiv: 2303.18240 [cs.CV].
- [53] S. Chen, R. Garcia, I. Laptev, and C. Schmid, “Sugar: Pre-training 3d visual representations for robotics,” *arXiv preprint arXiv:2404.01491*, 2024.
- [54] A. Dosovitskiy *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *ICLR*, 2020.
- [55] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *CVPR*, 2009.
- [56] L. Fu *et al.*, “Rethinking patch dependence for masked autoencoders,” *arXiv preprint arXiv:2401.14391*, 2024.
- [57] J. Lee, Y. Lee, J. Kim, A. Kosiorek, S. Choi, and Y. W. Teh, “Set transformer: A framework for attention-based permutation-invariant neural networks,” in *International conference on machine learning*, PMLR, 2019, pp. 3744–3753.
- [58] A. Vaswani *et al.*, “Attention is all you need,” in *NeurIPS*, 2017.
- [59] J. Han *et al.*, *Imagebind-llm: Multi-modality instruction tuning*, 2023. arXiv: 2309.03905 [cs.MM].

- [60] L. Fu *et al.*, “A touch, vision, and language dataset for multimodal alignment,” *arXiv preprint arXiv:2402.13232*, 2024.
- [61] S. Mirchandani *et al.*, *Large language models as general pattern machines*, 2023. arXiv: 2307.04721 [cs.AI].
- [62] E. J. Hu *et al.*, “Lora: Low-rank adaptation of large language models,” *arXiv preprint arXiv:2106.09685*, 2021.
- [63] W.-L. Chiang *et al.*, “Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality,” See <https://vicuna.lmsys.org> (accessed 14 April 2023), vol. 2, no. 3, p. 6, 2023.
- [64] H. Liu, C. Li, Y. Li, and Y. J. Lee, *Improved baselines with visual instruction tuning*, 2023.
- [65] W. Dai *et al.*, “Instructblip: Towards general-purpose vision-language models with instruction tuning,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.