

Autonomous Human-Robot Interaction via Operator Imitation

Sammy Christen¹, David Müller¹, Agon Serifi¹, Ruben Grandia¹,
Georg Wiedebach², Michael A. Hopkins², Espen Knoop¹, Moritz Bächer¹

Abstract—Teleoperated robotic characters can perform expressive interactions with humans, relying on the operators’ experience and social intuition. In this work, we propose to create autonomous interactive robots, by training a model to imitate operator data. Our model is trained on a dataset of human-robot interactions, where an expert operator is asked to vary the interactions and mood of the robot, while the operator commands as well as the pose of the human and robot are recorded. Our approach learns to predict continuous operator commands through a diffusion process and discrete commands through a classifier, all unified within a single transformer architecture. We evaluate the resulting model in simulation and with a user study on the real system. We show that our method enables simple autonomous human-robot interactions that are comparable to the expert-operator baseline, and that users can recognize the different robot moods as generated by our model. Finally, we demonstrate a zero-shot transfer of our model onto a different robotic platform with the same operator interface.

I. INTRODUCTION

Today, a wide variety of robotic systems are capable of expressing a rich set of behaviors, including quadrupeds [1], [2], humanoids [3], and non-anthropomorphic robots [4]. These systems are suitable for human-robot interaction (HRI), in particular when controlled by a skilled operator.

Operators can assess the environment, interpret the behavior of the human, and initiate appropriate interactions effectively. A key challenge is achieving autonomous HRI without an operator in the loop.

Full autonomy in HRI combines decision-making, motion control, and social interactions [5], [6]. Furthermore, it requires incorporating aspects of the theory of mind [7], such as understanding intent, beliefs, emotional states, and desires—both of the robot itself and the human. While such capabilities remain an open challenge, we observe that human operators already enable robots to convey expressive and lifelike behaviors through their experience and simple social cues such as human movements. For example, a shy robot might follow the human while keeping a safe distance and looking towards the ground, or a joyful robot might run closely behind the human and initiate dancing motions from time to time. Inspired by this, we aim to develop a framework that allows robots to autonomously engage in interactions—without contact—by approximating the human state through the robot-relative human pose, and express different moods at a level comparable to a human operator.

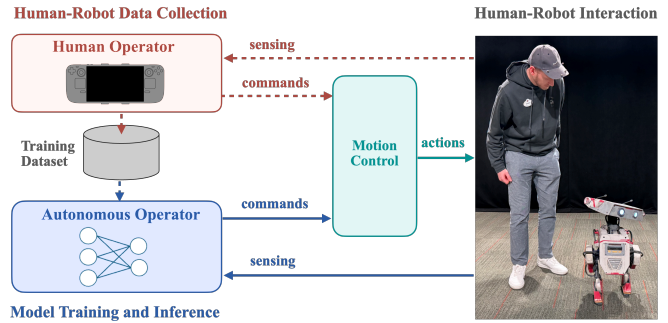


Fig. 1. **Overview of our approach.** We first collect a human-robot interaction dataset with a remote operator (red, top left). A motion control module takes operator commands and maps them to robot actions (cyan, middle). We then train a diffusion-based model that learns to imitate the operator (blue, bottom left). Our system learns to perform simple interactions with a human and express different moods. Dashed lines indicate components only required for data collection and training.

Developing an expressive and autonomous system for HRI requires methods that are flexible, expressive, and scalable. Heuristic-based approaches [8], [9] are inherently limited in scalability, require expert knowledge, and may need to be redefined for each robot. Similar challenges arise when defining rewards for a reinforcement learning policy [10], with the additional difficulty of realistically simulating human motions [11] or long training times in the real world [12]. An alternative is to collect real-world data and train supervised models, a popular approach in robotic manipulation [13], [14]. In the context of HRI, collecting such data can be time-consuming. Recent advances in user-controllable robots via gamepads offer a promising approach to simplifying data collection in such settings.

In this work, we introduce a novel framework for HRI that learns to imitate operators. As shown in Fig. 1, we first collect a human-robot interaction dataset, where an external operator controls the robot to interact with a human and express varying moods. The dataset includes the poses of the robot and the human, along with the operator’s teleoperation commands. Crucially, the problem statement focuses on imitating the operator rather than the robot’s actions, offering several advantages over relearning low-level control. Leveraging existing motion control eliminates the need to relearn the complex and data-intensive process of the robot’s dynamics. Furthermore, it ensures the robot maintains its motion style, robustness, and safety constraints.

In our method, we use diffusion models to imitate operators and predict diverse interactions. We introduce several simple yet effective techniques to adapt these models for real robotic platforms and human-robot interactions. As opposed

¹Disney Research, Switzerland. ²Disney Research, USA

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

to previous motion diffusion frameworks, which primarily focus on continuous predictions [15], [16], we propose a model that can predict both continuous commands and discrete events, matching the layout of gamepads. Our method consists of a unified transformer backbone that integrates a diffusion-based module for continuous signal prediction and a classifier for discrete event prediction. To enable our model to react dynamically to human movement, we condition on the time-varying robot-relative human pose. Additionally, we propose masking input signals after encoding rather than before [15], since zero input signals can be a valid input, such as indicating the proximity of the human to the robot.

Our experiments show that with less than one hour of data, our framework learns to perform autonomous interactions, exhibit multiple moods, and switch between different control modes, such as walking and standing. In a user study, we let 20 users interact with our system and show that they struggle to distinguish between autonomous and operated behavior for simple interactions. Furthermore, we show that users can successfully recognize the different moods of the robot. Lastly, we demonstrate the zero-shot transfer of our models trained with interaction data collected with a non-anthropomorphic robot to a humanoid robot.

In summary, our contributions are:

- a novel system enabling simple autonomous human-robot interactions and expressing different moods using a small dataset of human operator data for training,
- a model to imitate operator commands based on human pose, predicting both continuous and discrete signals,
- and empirical evidence through simulated and real-world experiments that our model generates realistic human-robot interactions and transfers across robotic platforms with the same operator interface.

II. RELATED WORK

We split related work into character control and learning-based HRI. Note that we consider our work as building up the machinery required for effective autonomous HRI, and there is extensive orthogonal research in social HRI that looks into the effects of robot personalities [17] and emotions [18]. This opens the door for interesting future studies. We refer to [6], [19], [20] for a detailed overview of social HRI.

A. Character Control

Learning-based controllers have shown remarkable success in controlling simulated and real-world characters or robots. Broadly, these controllers can be categorized into two types: *Imitation-based* controllers, which rely on dense kinematic reference motions as input [21], [22], [23], [24], [25], and *goal-driven* controllers [4], [15], [26], [27], which take high-level commands—such as joystick inputs—to control the character. These two paradigms are often closely intertwined. Many controllers leverage dense motion references during training, but exploit control through sparse user input at inference time [28].

With the advances of diffusion models [29] and their applications in motion generation [16], recent work investigated

their application in control settings. CAMDM [15] generates kinematic motion based on control signals in an autoregressive fashion. Closer to our work is the more direct intersection between text-conditioned diffusion models and physics-based controllers. RobotMDM [30] generates motions tailored for a pretrained controller, whereas CLoSD [31] integrates a controller into the generation process. Our model builds on such physics-based approaches, but instead of dense references, we learn to imitate user inputs provided during data collection for human-robot interaction, which reduces the data requirement.

Diffusion models have also been explored in the context of HRI. Yoneda *et al.* [32] propose a system for shared autonomy, where an autonomous agent assists a user in teleoperating a robot. Through the use of a diffusion model, the authors learn a model that trades off user autonomy for optimal behavior. Diffusion Co-Policy [33] is a method for a collaborative table moving task. Like them, we base our modeling on a diffusion model. However, in contrast to their approach, we support an action space with more than 2D continuous dimensions, supporting a combination of discrete *and* continuous actions. Moreover, our model controls the robot fully autonomously instead of controlling a subsystem that is attached to a robot commanded by a human.

B. Learning-Based Human-Robot Interactions

The application of deep learning techniques to autonomous HRI has recently gained popularity. Often, task-specific and isolated behaviors are investigated, such as human-robot handshakes [34], [35], human-robot handovers [36], [37], or table carrying [38]. For a broader overview that includes non-learning-based techniques, we refer to [5], [6].

A large body of work has focused on using learning-based methods for physical human-robot interaction. A common approach is to leverage human-human interaction demonstrations. Nikolaidis *et al.* [39] propose an unsupervised learning algorithm to cluster human behavior and learn robot policies that align with the user. In [40], the authors introduce a deep generative model representing the joint distribution of interactions in latent space via variational auto-encoders. Similarly, some works focus on learning joint latent representations between the human and robot from demonstrations, such as for human-robot hand interactions [41] or social motion forecasting [42]. In Co-GAIL [43], a human policy evolves along with a robotic policy via adversarial imitation learning. To adhere to the diversity of human behavior, a latent representation that represents the human’s intent is used, similar to [44] for multi-agent interactions. While there are similarities to these methods, our setting has different assumptions. In particular, we do not assume access to human-human demonstrations, as remotely operated robotic systems can be substantially different from the human anatomy, rendering a mapping between human demonstration and robot challenging. Furthermore, these methods typically focus on collaborative tasks with physical interactions, whereas we focus on non-physical interactions, such as following humans, expressing moods, and initiating pre-defined behaviors.

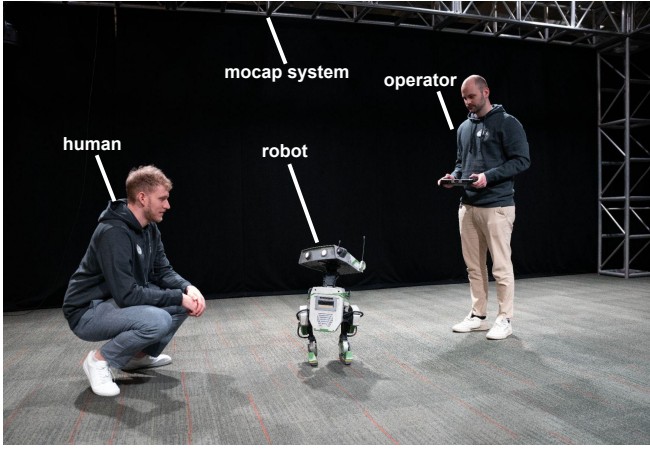


Fig. 2. **Our capture setup.** An operator controls the robot to interact with a human participant. The poses of the human and robot, as well as the operator commands, are recorded in a motion capture studio.

III. BACKGROUND DIFFUSION MODELS

Diffusion models [29] have become popular in various domains such as the generation of images [45], videos [46] or motion sequences [16]. They are a family of generative models particularly well suited for modeling complex data distributions. Diffusion models comprise a forward and a backward process. The forward process is a discretization of Langevin dynamics, where in each step, Gaussian noise ϵ_t is added to a data sample \mathbf{x}_0 , leading to progressively noisier versions of that sample \mathbf{x}_t :

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \epsilon_t, \epsilon_t \sim \mathcal{N}(0, \mathbf{I}). \quad (1)$$

Recursively applying this update yields

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}), \quad (2)$$

$$\bar{\alpha}_t = \prod_{i=1}^t (1 - \beta_i), \quad (3)$$

where β_i is a variance scheduler for the noise, with i or t indicating the diffusion step. In the reverse process, the data is gradually denoised into a clean sample, which is modeled by the probability distribution $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$. This denoising step is learned via a neural network parameterized by θ . We follow MDM [16] and directly predict the clean output $\hat{\mathbf{x}}_0$. Additionally, in conditional diffusion models, the probability distribution is extended by conditioning signals \mathbf{c} , yielding $p_\theta(\mathbf{x}_0 | \mathbf{x}_t, \mathbf{c})$. During inference, the model can generate diverse outputs from randomly sampled Gaussian noise and the conditioning signals.

IV. PROBLEM SETTING

Given a sequence of past human poses $\mathbf{p} \in \mathbb{R}^{M \times 7}$ and robot poses $\mathbf{r} \in \mathbb{R}^{M \times 7}$, which each consist of M frames of global 3D positions and 4D orientations (quaternions), our goal is to predict a future sequence of continuous operator commands $\mathbf{x}^{1:N} \in \mathbb{R}^{N \times j}$, where N indicates the prediction horizon and j is the number of commands per frame. In addition to the continuous command predictions, such as

joysticks on a gamepad, we are also interested in predicting discrete events, which are typically triggered through button presses on operator interfaces. We differentiate between discrete events that trigger a certain behavior \mathbf{d}_b , such as a pre-defined jumping or dancing motion (see Sec. VI-A), and the mode \mathbf{d}_m of the robot, e.g., whether it should be in standing or walking mode. Each prediction window has one prediction per discrete input, rather than a separate prediction per frame. Both the continuous and discrete signals are passed as conditioning to a motion control module that controls the robot (see Sec. VI-A). Please note that we use subscripts to indicate the diffusion step and superscripts to describe the sequential prediction.

V. METHOD

Our framework consists of two main components: human-robot interaction data collection and a model to train autonomous interactions.

A. Human-Robot Interaction Data Collection

Our capture setup is shown in Fig. 2 and consists of a robot, an operator controlling the robot, and a human interacting with it. We aim to capture the human poses \mathbf{p} , the robot poses \mathbf{r} , the continuous operator commands \mathbf{x} as well as the discrete commands \mathbf{d}_b and \mathbf{d}_m . We record operator commands directly on the input interface. To capture human and robot poses, we leverage the OptiTrack motion capture system [47], placing markers on the robot and a hat worn by the human.

To collect data, we ask an expert operator to perform interactions that are responsive to human pose, such as following the human and varying between expressing different moods through motions. For instance, in the shy mood, the robot tends to look at the ground while interacting with the human, whereas in the angry mood, the robot frequently shakes its head and refuses interactions when the human approaches. See Tab. I for an overview of the collected data and a more detailed description of the moods. In our data collection, the robot is operated by a single expert operator and two humans, one after another, interact with it.

B. Method Architecture

Our method architecture is illustrated in Fig. 3. Inspired by the CAMDM framework for character control in simulation [15], we propose a framework for real-world human-robot interaction. To address the challenges of human-robot interactions, we contribute several domain-specific adjustments and extensions to the architecture. Our model has a transformer backbone and is conditioned on multiple components; the history of human poses \mathbf{c}_p and operator commands \mathbf{c}_x , and the diffusion step t . Before passing the history of human poses to the model, we transform them into a robot-relative frame using the history of robot poses \mathbf{r} .

In summary, our model \mathcal{G} is defined as follows:

$$(\hat{\mathbf{x}}_0, \mathbf{d}_b, \mathbf{d}_m) = \mathcal{G}(\mathbf{c}_p, \mathbf{c}_x, t, \mathbf{x}_t, \mathbf{q}_b, \mathbf{q}_m). \quad (4)$$

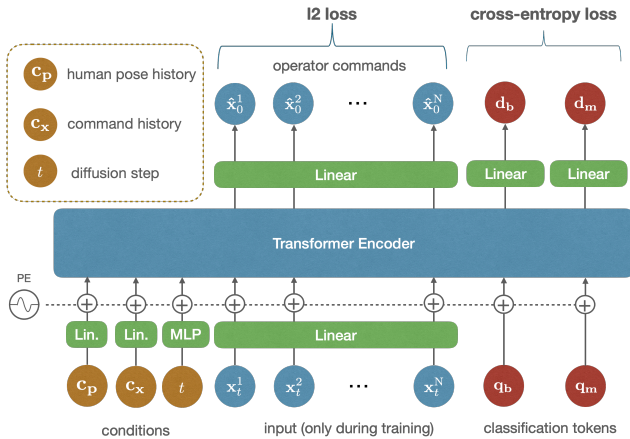


Fig. 3. **Overview of our method architecture.** We use conditions in the form of past human poses and commands (yellow). A diffusion model predicts operator commands to control the robot model (blue). The transformer also outputs discrete predictions for different behavior and the mode (red).

1) *Continuous Command Prediction via Diffusion:* To predict continuous operator commands, we use a diffusion model (see Sec. III). More specifically, during training, we alternate between noising clean data samples \mathbf{x}_0 to noisy signals \mathbf{x}_t according to Eq. 2 and denoising steps to predict the clean signals $\hat{\mathbf{x}}_0$.

We apply the separation of conditioning tokens and classifier-free guidance on past commands as proposed by [15]. In contrast to their model, we found that dropout in the positional encoding leads to noisy outputs, hence do not adopt it in our model. Furthermore, the conditioning signals in our case can be close to zero, for example, when the human is close to the robot. Therefore, we propose to apply masking to the conditioning *after* the encoding layers instead of to the raw conditions [15], because otherwise zero-masking could be mistaken by the model for actual zero conditions. To account for the diversity in human height, we augment the available human pose data by adding a random offset (± 0.3 m) in the negative gravity direction during training.

2) *Discrete Event Prediction:* It is non-trivial to mix continuous and discrete signals in a diffusion process. We propose to add the discrete event predictions as auxiliary tasks to our transformer model. In particular, we add classification query tokens, specifically \mathbf{q}_b for discrete behavior and \mathbf{q}_m for mode, to the transformer architecture [48]. The output of the transformer encoder has a head to predict the mode \mathbf{d}_m and pre-defined discrete behavior \mathbf{d}_b , respectively. The model learns to select from several classes of behaviors and a default class (i.e., no discrete event occurring). Note that these predictions use the same transformer encoder, but no diffusion is applied. To train the classification heads, we use weighted cross-entropy losses, due to the imbalance in the data between the rare discrete events and the default class.

TABLE I
MOTION TYPES, LENGTHS, AND DESCRIPTIONS

Category	Length (min)	Description
Default	8	Follows the human, retreats if the human walks toward it, looks at the human when in standing mode.
Angry	6	Ignores the human, walks away if approached, shakes its head, occasionally triggers angry animations.
Sad	8	Walks away from the human, shakes its head, looks mostly down at the ground.
Shy	7	Approaches the human carefully, occasionally stops, mostly looks at the ground, and tries to avoid eye contact.
Happy	8	Turns in circles, runs after the human, and often expresses positivity through animations such as dancing or jumping.

VI. EXPERIMENTS

In Sec. VI-A, we describe the robotic platform used in our experiments and provide implementation details about model parameters and runtime of our approach in Sec. VI-B. We conduct experiments in simulation in Sec. VI-C. In Sec. VI-D, we perform a user study in the real world and demonstrate zero-shot transfer to a different robot in Sec. VI-E. See our video¹ for more qualitative results.

A. Robotic Platform

We build on the platform developed by Grandia *et al.* [4], which proposes a new bipedal character design and control approach for entertainment applications. A reinforcement learning based control architecture is used to robustly imitate artistic motions conditioned on continuous and discrete command signals. During runtime, these command signals are generated by an animation engine that fuses user-inputs with predefined animation content. An intuitive operator interface, implemented on the hand-held operator controller, enables expressive real-time show performances with the robot.

The robot, 0.66 m tall, with a total mass of 15.4 kg, has 5 degrees of freedom (DoF) per leg and a 4 DoF neck. Moreover, the robot is equipped with a set of show functions: a pair of actuated antennas and illuminated eyes, and a headlamp. Additionally, the robot has a stereo pair of loudspeakers in both the body and the head. These components provide additional means to express moods. Their behavior is synchronized with the motion of the robot through animation signals from the animation engine and state feedback. We refer the reader to [4] for more details.

B. Implementation Details

Our transformer encoder has a latent dimension of 128, a feedforward size of 256, and two attention heads across 2-layers. We train for 5000 epochs with a batch size of 128 and a learning rate of $1e-4$. We select a past window of 15 frames and predict 25 future frames (M and N in Sec. IV). The number of diffusion steps is set to 8. We predict 10 continuous signals, up to 8 discrete behaviors (depending on

¹<https://youtu.be/4U4etupwzhQ?si=NNHM7NqjAKWoo32B>

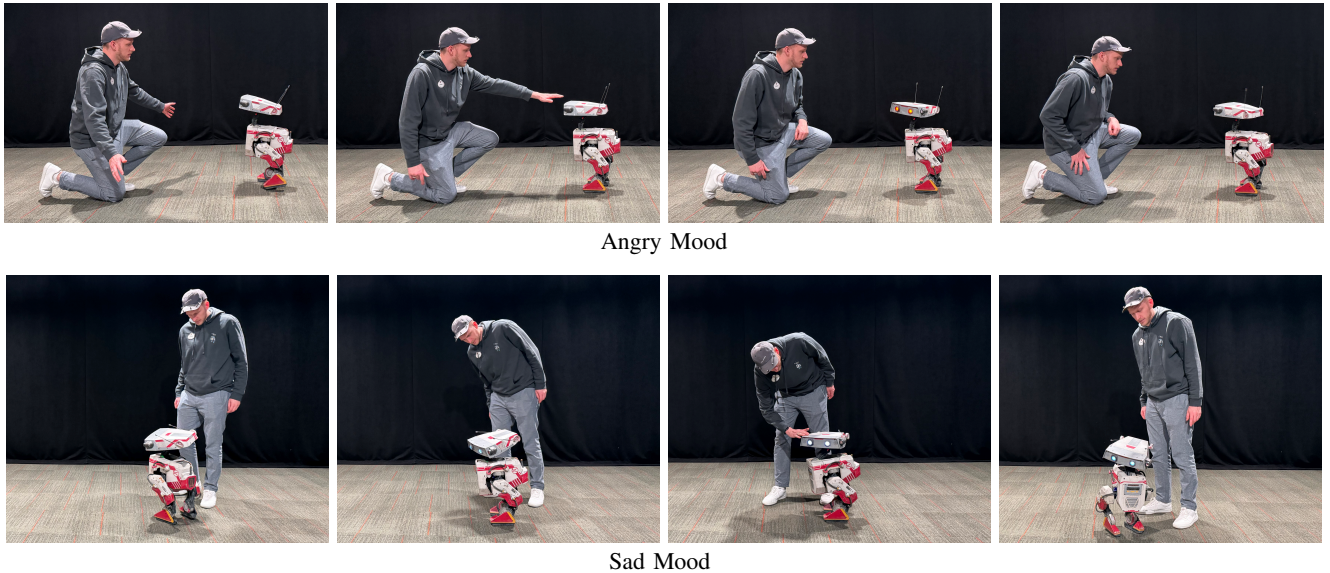


Fig. 4. **Example behavior of different moods.** In the angry mood, the robot refuses interactions and steps away from the human (top). In the sad mood, the robot turns away from the human, has its head tilted towards the ground and occasionally shakes its head (bottom).

the mood), and 2 modes (walking and standing). To improve signal quality, we apply a simple Gaussian smoothing filter to the model’s predictions before sending them to the robot. This helps to reduce the noise introduced by the motion capture conditioning.

All computations, including our diffusion model, motion control policy, state estimation, and the animation engine, run on the robot’s on-board computer. Both human pose and autonomous operator commands are sent to our model at 50 Hz. The predicted operator commands are passed to the animation engine (see Sec. VI-A), which fuses the commands with animation content and sends it to a control policy. This low-level control policy outputs actuator setpoint commands, which are interpolated to 600 Hz and sent to the actuators’ PD controllers. In addition to teleoperation commands, the low-level policy also receives proprioceptive state inputs from the robot.

C. Simulation-Based Evaluation

Quantifying engaging and natural interactions in simulation is challenging, much like defining general heuristics for robot control or rewards for learning behavior. Nevertheless, to measure and compare architectural choices, we define a set of metrics that provide insights into model performance:

Facing Angle Error (FAE): measures the average angle between the robot’s forward direction and the vector from the robot’s root to the human in degrees. We assume that the robot should be facing the human (in the default mode).

Tracking Error (TE): measures the average distance between the human and the robot in the x-y plane. We assume that the robot should follow the human closely.

Mean Squared Derivative (MSD): measures how rapidly the continuous signals change. This helps identify noisy predictions and abrupt transitions between prediction windows. We average over all signals.

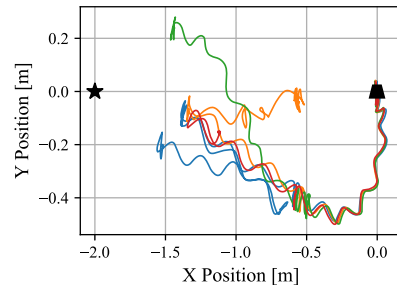


Fig. 5. **Diversity of our framework.** Given the same starting point (black trapezoid) and a fixed human position (black star), we run our model multiple times and plot x-y positions. As can be seen, the model generates different behavior for the same human pose conditions.

We compare our approach to a transformer baseline and ablate design choices of our architecture. For our simulation experiments, we use a 75/25 train-test split for the default mode. During evaluation, we replay human motions from the test set and control the robot using our model, streaming the human-relative pose from simulation directly to our model. Additionally, as shown in Fig. 5, we qualitatively demonstrate that, given a fixed human pose and a single initial robot pose, our model generates diverse outputs.

1) *Baseline Comparison:* Our baseline is a deterministic model that directly maps from human pose as input to operator commands as output. To achieve this, we use the same transformer architecture that is used in our model but remove the diffusion process. To account for the sequential output, we add query tokens as inputs. As can be seen in Tab. II (top), our model achieves better tracking of the human (lower FAE and TE) and leads to a cleaner signal and transitions between prediction windows, as indicated by the lower MSD score. This shows that using a diffusion model is beneficial over using a transformer model without diffusion.

TABLE II
SIMULATION-BASED EVALUATIONS

Variant	FAE [deg.] ↓	TE [m] ↓	MSD ↓
<i>transformer</i>	57.66 ± 19.61	1.49 ± 0.05	4.40 ± 2.29
<i>Ours 25 frames</i>	43.85 ± 2.40	1.47 ± 0.03	2.42 ± 0.40
<i>w/ dropout</i>	39.48 ± 2.19	1.44 ± 0.02	5.76 ± 0.32
<i>w/o human</i>	100.28 ± 4.74	3.20 ± 0.50	2.34 ± 0.26
<i>w/o commands</i>	43.12 ± 2.66	1.44 ± 0.01	14.80 ± 0.69
<i>Ours 75 frames</i>	56.97 ± 1.98	1.54 ± 0.01	4.11 ± 0.39
<i>Ours 50 frames</i>	50.58 ± 12.31	1.48 ± 0.02	2.45 ± 0.33
<i>Ours 25 frames</i>	43.85 ± 2.40	1.47 ± 0.03	2.42 ± 0.40

2) *Ablations*: To ablate the components of our architecture, we train different variants of our model by removing the human pose history \mathbf{c}_p (*w/o human*), command history \mathbf{c}_x (*w/o commands*), dropout (*w/o dropout*). Additionally, we evaluate the influence of different prediction window lengths (25, 50, and 75 frames). We present results in Tab. II. As expected, the variant without information about the human pose leads to high tracking errors. The variant without conditioning on past commands leads to the best tracking of the human, because no constraints are imposed on the coherence of the signal to previous predictions. Hence, the transitions between windows are often abrupt and jerky, as indicated by the high MSD (14.8). Our final variant (*Ours 25 frames*) trades off signal coherence with close tracking of the human. When comparing different prediction windows, we find that 25 frames yield the best tracking performance while staying close to real-time, considering the computational load of the diffusion process. Note that all but the window size ablations are trained with a window of 25 frames.

D. Real-World Evaluation

We perform qualitative and quantitative real-world evaluations. Please see Fig. 4 and our accompanying video for qualitative results. In our empirical evaluations, we aim to answer whether our model can enable 1) interactions with the robot that feel similar to interactions controlled by a trained operator and 2) a robot to express different moods that are recognizable by humans. To analyze these questions, we propose a two-stage in-person user study. We run the user study with 20 participants, aged between 22 and 44 ($M=30$, $SD=4.8$), 12 males and 8 females. The participants were recruited from our organization, but were neither part of the project nor had experience interacting with our robot. After the interactive part of the study, we let participants fill out a small qualitative questionnaire.

1) *Questionnaire*: The anonymous questionnaire comprises three statements. The users are asked how much they agree with each statement on a 5-point scale, with a score of 1 indicating “strongly disagree” and 5 indicating “strongly agree”. The first statement is “I have experience interacting with robots”, which received a mean score of 2.2 ($STD=1.0$), showing that participants are not very experienced interacting with robots. The other statements, “The interaction felt engaging and enjoyable” and “I felt like the robot was reacting to me” received mean scores of 4.9 ($STD=0.3$)

TABLE III
OPERATOR RECOGNITION CONFUSION MATRIX

GT \ User	Operator	Autonomous
Operator	0.55	0.45
Autonomous	0.46	0.54

TABLE IV
MOOD RECOGNITION CONFUSION MATRIX

GT \ User	Happy	Sad	Angry	Shy
Happy	0.74	0.00	0.11	0.16
Sad	0.00	0.74	0.05	0.21
Angry	0.26	0.00	0.74	0.00
Shy	0.05	0.11	0.16	0.68

and 4.7 ($STD=0.5$), respectively. These scores indicate that participants enjoyed interacting with the robot and felt like the robot was reacting to their movements.

2) *Operator Recognition*: In this experiment, we follow a protocol where each participant experiences one of two settings. In the first setting, the operator is controlling the robot to interact with the human in the default mode, following the human without expressing different moods. In the second setting, our model controls the robot to exhibit the same behavior. We demonstrate each setting to the user for 30 seconds, after which the participant has to guess whether it was our model or the operator controlling the robot. In total, we show each setting twice (4 trials in total) in randomized order. To avoid bias, the operator pretends to be actively controlling the robot, even in the autonomous setting. We report the recognition accuracy in a confusion matrix shown in Tab. III, where rows indicate the presented setting (GT) and columns are the user-predicted setting (User). As can be seen, the accuracy is close to 50%, with scores of 55% for the active control by the operator and 54% for our autonomy mode. The autonomous model is classified as operator 46% of the time and the operator is classified as autonomous model 45% of the time. These results indicate that it is difficult for users to tell the two models apart, and that our model in default mode is close to an expert operator’s abilities. Three participants, who were more experienced with robots, figured out tricks to distinguish between the models. For example, these participants realized that if we lose track of the human at the border of the mocap space, the model will only start following once inside the tracked space.

3) *Mood Recognition*: To determine whether the robot can express moods through our approach, we demonstrate the four moods (happy, sad, angry, shy) to the humans (see Tab. I for a detailed description). We randomize the order of moods between participants and let them interact with the robot. An interaction takes one minute, after which humans have to provide a forced-choice answer. We do not explain the moods in detail and simply provide the moods to the humans. To avoid external bias, we use neutral eye coloring during this experiment (as opposed to the colored eyes for different

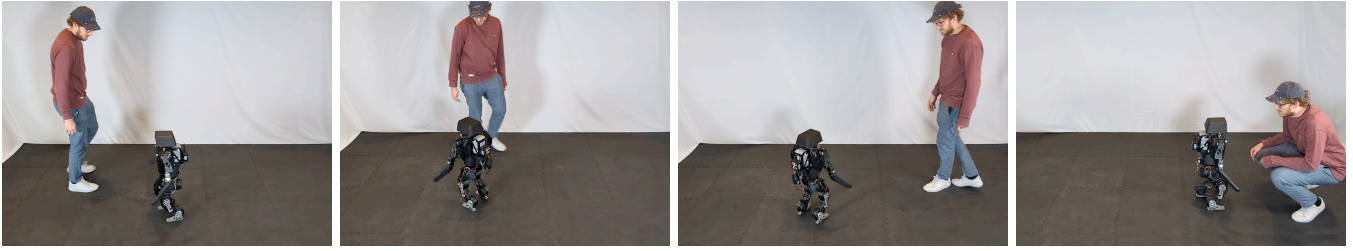


Fig. 6. **Zero-shot transfer to a humanoid robot.** Our method trained with data collected via our non-anthropomorphic bipedal robot can successfully transfer behavior such as tracking the human without retraining.

moods, cf. Fig. 4). We report results in a confusion matrix shown in Tab. IV, where rows indicate the presented setting (GT) and columns are the user-predicted setting. Generally, the users were able to correctly select the presented mood with an accuracy between 68% and 74%. The confusion cases give some interesting insights into the perception of certain moods. The angry mood was considered happy (26%) because certain discrete behaviors, such as body shaking in disagreement, were considered a shake that expresses excitement (cf. video 35 s). In contrast, the happy mood was considered angry by some (11%), because of the sprint towards the human, which was considered as an aggressive “charging” motion (cf. video 1 min45 s). The sad mood was misclassified as shy (21%) due to the robot mostly looking at the ground. We believe these results validate that users can reliably identify the moods, with reasonable explanations for the confused cases. This also shows that there is no clear line between moods, and certain behaviors are interpreted differently between humans.

E. Zero-Shot Cross-Embodiment Transfer

To demonstrate that our method is applicable to different robotic platforms, we perform zero-shot transfer to a humanoid robot. Importantly, the interface mapping is maintained between the two platforms, i.e., the same joystick is used for controlling the walking direction and speed of the two robots. We qualitatively show in Fig. 6 and our supplemental video that the default behavior, such as tracking the human, turning, and walking backwards (cf. Tab. I), can be successfully transferred to a different platform. Notably, the training data used for our model was collected using the robot presented in Sec. VI-A.

VII. DISCUSSION

While our approach takes a step towards autonomous human-robot interactions, it has several limitations. The main limitation is the reliance on a motion capture system to sense the robot-relative human pose. In the future, we hope to integrate perception to enable deployment in real-world environments. Beyond estimating global human pose, future work could predict the full-body human pose or even facial expressions, enabling more nuanced interactions. For example, we manually trigger the different moods in this work. A promising direction is exploring a single model conditioned on detailed human states to autonomously predict mood

transitions. Furthermore, we do not model complex behaviors with long-term dependencies and high-level decision-making, such as combining multiple semantically different interactions into a cohesive interaction. Instead, our focus is on achieving performance comparable to human operators for shorter and simpler interactions, such as expressing moods and reacting to human pose. For more complex behavior, future research could explore the integration of a high-level decision-making module that links multiple interactions together.

Our data was captured with a single operator and two human subjects. Our user study revealed a large diversity in how humans interact with the robot. Expanding the dataset with more participants could improve the model’s robustness to diverse human behavior while increasing the diversity of interactions our model predicts. This would also enhance the adaptability of our model to different human heights, which we currently address through height randomization during training. Future work could also explore interactions that involve physical contact. Finally, current interactions are limited to one robot and one human. An exciting avenue for future research is extending our work to multi-robot and multi-human interactions.

VIII. CONCLUSION

We have proposed a method that enables autonomous human-robot interactions by learning the mapping from human pose to operator commands with a diffusion-based approach, predicting both continuous commands as well as discrete button presses. By relying on operator instead of actuator commands, we train a model with less than 40 minutes of data, drastically reducing data requirements. Our approach also enhances safety, as operator interfaces are typically designed with built-in safety guarantees. Lastly, our experiments demonstrate that users enjoy interacting with the system, can reliably recognize different moods, and find it difficult to distinguish between autonomous or operator-controlled robot behavior for simple interaction.

IX. ACKNOWLEDGEMENTS

We thank Jenny Wang for her exploratory work. We would also like to thank all participants that took part in our user study and provided feedback on our system.

REFERENCES

- [1] M. Hutter, C. Gehring, D. Jud, A. Lauber, C. D. Bellicoso, V. Tsounis, J. Hwangbo, K. Bodie, P. Fankhauser, M. Bloesch, *et al.*, “Anymal—a highly mobile and dynamic quadrupedal robot,” in *International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 38–44.
- [2] Boston Dynamics, “Spot: The agile mobile robot,” n.d. [Online]. Available: <https://www.bostondynamics.com/spot>
- [3] X. Cheng, Y. Ji, J. Chen, R. Yang, G. Yang, and X. Wang, “Expressive whole-body control for humanoid robots,” *arXiv preprint arXiv:2402.16796*, 2024.
- [4] R. Grandia, E. Knoop, M. A. Hopkins, G. Wiedeback, J. Bishop, S. Pickles, D. Müller, and M. Bächer, “Design and control of a bipedal robotic character,” in *Robotics: Science and Systems (RSS)*, 2024.
- [5] A. Thomaz, G. Hoffman, and M. Cakmak, “Computational human-robot interaction,” *Found. Trends Robot.*, vol. 4, no. 2–3, p. 105–223, Dec. 2016.
- [6] T. B. Sheridan, “Human–robot interaction: status and challenges,” *Human factors*, vol. 58, no. 4, pp. 525–532, 2016.
- [7] C. Frith and U. Frith, “Theory of mind,” *Current biology*, vol. 15, no. 17, pp. R644–R645, 2005.
- [8] A. Bauer, B. Gonsior, D. Wollherr, and M. Buss, “Heuristic rules for human-robot interaction based on principles from linguistics—asking for directions,” in *AISB Convention-Symposium on New Frontiers in Human-Robot Interaction*, 2009, pp. 24–30.
- [9] A. E. Block, S. Christen, R. Gassert, O. Hilliges, and K. J. Kuchenbecker, “The six hug commandments: Design and evaluation of a human-sized hugging robot with visual and haptic perception,” in *International Conference on Human-Robot Interaction (HRI)*, 2021, p. 380–388.
- [10] A. H. Qureshi, Y. Nakamura, Y. Yoshikawa, and H. Ishiguro, “Intrinsically motivated reinforcement learning for human–robot interaction in the real-world,” *Neural Networks*, vol. 107, pp. 23–33, 2018.
- [11] J. Thumm, F. Trost, and M. Althoff, “Human-robot gym: Benchmarking reinforcement learning in human-robot collaboration,” in *International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 7405–7411.
- [12] A. H. Qureshi, Y. Nakamura, Y. Yoshikawa, and H. Ishiguro, “Robot gains social intelligence through multimodal deep reinforcement learning,” in *International Conference on Humanoid Robots (Humanoids)*. IEEE, 2016, pp. 745–751.
- [13] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Hausman, B. Ichter, *et al.*, “ π_0 : A vision-language-action flow model for general robot control,” *arXiv preprint arXiv:2410.24164*, 2024.
- [14] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song, “Diffusion policy: Visuomotor policy learning via action diffusion,” *The International Journal of Robotics Research*, p. 02783649241273668, 2023.
- [15] R. Chen, M. Shi, S. Huang, P. Tan, T. Komura, and X. Chen, “Taming diffusion probabilistic models for character control,” in *ACM SIGGRAPH 2024 Conference Papers*, ser. SIGGRAPH ’24. New York, NY, USA: Association for Computing Machinery, 2024.
- [16] G. Tevet, S. Raab, B. Gordon, Y. Shafir, D. Cohen-or, and A. H. Bermano, “Human motion diffusion model,” in *International Conference on Learning Representations (ICLR)*, 2023.
- [17] K. M. Lee, W. Peng, S.-A. Jin, and C. Yan, “Can robots manifest personality?: An empirical test of personality recognition, social responses, and social presence in human–robot interaction,” *Journal of Communication*, vol. 56, no. 4, pp. 754–772, 2006.
- [18] A. Paiva, I. Leite, and T. Ribeiro, “Emotion modeling for social robots,” *The Oxford handbook of affective computing*, p. 296, 2014.
- [19] L. P. Robert Jr, R. Alahmad, C. Esterwood, S. Kim, S. You, Q. Zhang, *et al.*, “A review of personality in human–robot interactions,” *Foundations and Trends® in Information Systems*, vol. 4, no. 2, pp. 107–212, 2020.
- [20] R. Kirby, J. Forlizzi, and R. Simmons, “Affective social robots,” *Robotics and Autonomous Systems*, vol. 58, no. 3, pp. 322–332, 2010, towards Autonomous Robotic Systems 2009: Intelligent, Autonomous Robotics in the UK.
- [21] L. Fussell, K. Bergamin, and D. Holden, “Supertrack: motion tracking for physically simulated characters using supervised learning,” *ACM Trans. Graph.*, vol. 40, no. 6, Dec. 2021.
- [22] Z. Luo, J. Cao, A. W. Winkler, K. Kitani, and W. Xu, “Perpetual humanoid control for real-time simulated avatars,” in *International Conference on Computer Vision (ICCV)*, 2023.
- [23] A. Serifi, R. Grandia, E. Knoop, M. Gross, and M. Bächer, “Vmp: Versatile motion priors for robustly tracking motion on physical characters,” in *Proceedings of the ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, ser. SCA ’24. Goslar, DEU: Eurographics Association, 2024, pp. 1–11.
- [24] Y. Liu, B. Yang, L. Zhong, H. Wang, and L. Yi, “Mimicking-bench: A benchmark for generalizable humanoid-scene interaction learning via human mimicking,” *arXiv preprint arXiv:2412.17730*, 2024.
- [25] J. Braun, S. Christen, M. Kocabas, E. Aksan, and O. Hilliges, “Physically plausible full-body hand-object interaction synthesis,” in *International Conference on 3D Vision (3DV)*, 2024.
- [26] K. Bergamin, S. Clavet, D. Holden, and J. R. Forbes, “Drecon: data-driven responsive control of physics-based characters,” *ACM Trans. Graph.*, vol. 38, no. 6, Nov. 2019.
- [27] T. He, W. Xiao, T. Lin, Z. Luo, Z. Xu, Z. Jiang, C. Liu, G. Shi, X. Wang, L. Fan, and Y. Zhu, “Hover: Versatile neural whole-body controller for humanoid robots,” *International Conference on Robotics and Automation (ICRA)*, 2025.
- [28] C. Tessler, Y. Guo, O. Nabati, G. Chechik, and X. B. Peng, “Masked-mimic: Unified physics-based character control through masked motion inpainting,” *ACM Transactions on Graphics (TOG)*, 2024.
- [29] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” *Conference on Neural Information Processing Systems (NeurIPS)*, vol. 33, pp. 6840–6851, 2020.
- [30] A. Serifi, R. Grandia, E. Knoop, M. Gross, and M. Bächer, “Robot motion diffusion model: Motion generation for robotic characters,” in *SIGGRAPH Asia 2024 Conference Papers*, ser. SA ’24. New York, NY, USA: Association for Computing Machinery, 2024.
- [31] G. Tevet, S. Raab, S. Cohan, D. Reda, Z. Luo, X. B. Peng, A. H. Bermano, and M. van de Panne, “CLOSD: Closing the loop between simulation and diffusion for multi-task character control,” in *International Conference on Learning Representations (ICLR)*, 2025.
- [32] T. Yoneda, L. Sun, G. Yang, B. C. Stadie, and M. R. Walter, “To the noise and back: Diffusion for shared autonomy,” in *Robotics: Science and Systems (RSS)*, 2023.
- [33] E. Ng, Z. Liu, and M. Kennedy, “Diffusion co-policy for synergistic human-robot collaborative tasks,” *Robotics and Automation Letters (RA-L)*, 2023.
- [34] V. Prasad, R. Stock-Homburg, and J. Peters, “Human-robot handshaking: A review,” *International Journal of Social Robotics*, vol. 14, no. 1, pp. 277–293, 2022.
- [35] S. Christen, S. Stevšić, and O. Hilliges, “Demonstration-guided deep reinforcement learning of control policies for dexterous human-robot interaction,” in *International Conference on Robotics and Automation (ICRA)*, 2019.
- [36] S. Christen, W. Yang, C. Pérez-D’Arpino, O. Hilliges, D. Fox, and Y.-W. Chao, “Learning human-to-robot handovers from point clouds,” in *Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [37] S. Christen, L. Feng, W. Yang, Y.-W. Chao, O. Hilliges, and J. Song, “Synh2r: Synthesizing hand-object motions for learning human-to-robot handovers,” in *International Conference on Robotics and Automation (ICRA)*, 2024.
- [38] M. K. I. Eley Ng, Ziang Liu, “It takes two: Learning to plan for human-robot cooperative carrying,” in *International Conference on Robotics and Automation (ICRA)*, 2023.
- [39] S. Nikolaidis, R. Ramakrishnan, K. Gu, and J. Shah, “Efficient model learning from joint-action demonstrations for human-robot collaborative tasks,” in *International Conference on Human-Robot Interaction (HRI)*, 2015, pp. 189–196.
- [40] J. Bütepage, A. Ghadirzadeh, Ö. Öztimur Karadağ, M. Björkman, and D. Kragic, “Imitating by generating: Deep generative models for imitation of interactive tasks,” *Frontiers in Robotics and AI*, vol. 7, p. 47, 2020.
- [41] V. Prasad, D. Koert, R. Stock-Homburg, J. Peters, and G. Chalvatzaki, “Mild: multimodal interactive latent dynamics for learning human-robot interaction,” in *International Conference on Humanoid Robots (Humanoids)*. IEEE, 2022, pp. 472–479.
- [42] E. Valls Mascaró, Y. Yan, and D. Lee, “Robot interaction behavior generation based on social motion forecasting for human-robot interaction,” *International Conference on Robotics and Automation (ICRA)*, 2024.

- [43] C. Wang, C. Pérez-D'Arpino, D. Xu, L. Fei-Fei, K. Liu, and S. Savarese, "Co-gail: Learning diverse strategies for human-robot collaboration," in *Conference on Robot Learning (CoRL)*, vol. 164. PMLR, 08–11 Nov 2022, pp. 1279–1290.
- [44] A. Xie, D. Losey, R. Tolsma, C. Finn, and D. Sadigh, "Learning latent representations to influence multi-agent interaction," in *Conference on Robot Learning (CoRL)*. PMLR, 2021, pp. 575–588.
- [45] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 10 684–10 695.
- [46] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet, "Video diffusion models," *arXiv:2204.03458*, 2022.
- [47] NaturalPoint, Inc., *OptiTrack Motion Capture System*, 2024, available at <https://www.optitrack.com/>.
- [48] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," *International Conference on Learning Representations (ICLR)*, 2020.