

# Reshaping Robot Trajectories Using Natural Language Commands: A Study of Multi-Modal Data Alignment Using Transformers

Arthur Bucker<sup>1</sup>, Luis Figueredo<sup>1</sup>, Sami Haddadin<sup>1</sup>, Ashish Kapoor<sup>2</sup>, Shuang Ma<sup>2</sup>, Rogerio Bonatti<sup>2</sup>

<sup>1</sup>Technische Universität München, <sup>2</sup>Microsoft

**Abstract**—Natural language is the most intuitive medium for us to interact with other people when expressing commands and instructions. However, using language is seldom an easy task when humans need to express their intent towards robots, since most of the current language interfaces require rigid templates with a static set of action targets and commands. In this work, we provide a flexible language-based interface for human-robot collaboration, which allows a user to reshape existing trajectories for an autonomous agent. We take advantage of recent advancements in the field of large language models (BERT and CLIP) to encode the user command, and then combine these features with trajectory information using multi-modal attention transformers. We train the model using imitation learning over a dataset containing robot trajectories modified by language commands, and treat the trajectory generation process as a sequence prediction problem, analogously to how language generation architectures operate. We evaluate the system in multiple simulated trajectory scenarios, and show a significant performance increase of our model over baseline approaches. In addition, our real-world experiments with a robot arm show that users significantly prefer our natural language interface over traditional methods such as kinesthetic teaching or cost-function programming. Our study shows how the field of robotics can take advantage of large pre-trained language models towards creating more intuitive interfaces between robots and machines. Project webpage: [https://arthurfenderbucker.github.io/NL\\_trajectory\\_reshaper/](https://arthurfenderbucker.github.io/NL_trajectory_reshaper/)

## I. INTRODUCTION

Large language models such as BERT [1], GPT3 [2] and Megatron-Turing [3] have radically improved the quality of machine-generated text, along with our ability to solve to natural language processing tasks. Beyond just language, we see a shift in machine learning architectures in multiple domains, as the dominant design paradigm changes from designing task-specific models towards the use of large foundational pre-trained models [4]. Several of these large models already combine multiple data modalities such as text, images, video, depth, and even the temporal dimension [5]–[8]. The use of foundational models is appealing because they are trained on broad datasets over a wide variety of downstream tasks, and therefore provide *general* skills which can be used directly or with minimal fine-tuning to new applications [4].

The field of robotics traditionally uses extremely task and hardware-specific models, which have to be re-trained and even re-designed if there are minor changes in robot dynamics, environment and operational objectives. This inflexible machine learning approach is ripe for innovation with the use of foundational models [4], in particular when it comes



Fig. 1: Our system allows a user to send natural language commands to reshape robot trajectories relative to objects in the environment.

to task specification in ambiguous scenarios (*What should I do?*) and task learning that can generalize across multiple environments (*How should I do it?*). Recent works have just started to explore the use of pre-existing foundational models from language and vision towards robotics [9]–[13], and also the development of robotics-specific foundational models [8, 14].

Our work aims to leverage information contained in existing vision-language foundational models to fill the gap in existing tools for human-robot interaction. Even though natural language is the richest form of communication between humans, modeling human-robot interactions using language is challenging because we often require vast amounts of data [11, 15]–[17], or classically, force the user to operate within a rigid set of instructions [18, 19]. To tackle these challenges, our framework makes use of two key ideas: first, we employ large pre-trained language models to provide rich user intent representations, and second, we align geometrical trajectory data with natural language jointly with the use of a multi-modal attention mechanism.

As seen in Fig 1, we focus our study on robotics applications where a user needs to reshape an existing robot trajectory according to specific operational constraints. This class of use cases arises often in human-robot interaction when autonomous agents that employ traditional motion planners (*e.g.*  $A^*$ , RRT\* [20] or MPC [21] concerned solely about obstacle avoidance and dynamics) need to be corrected by a user according to additional semantic or safety objectives. For instance, our goal is to enable a factory worker to quickly reconfigure a robot arm trajectory further away from fragile objects, or to allow a user to intuitively tell a robot barista to get a little closer to the cup in order to pour a wine bottle.

Then main contribution of this paper is to propose a novel system with a multimodal attention mechanism for semantic trajectory generation. It can effectively align natural language features with geometrical cues jointly, and perform the goal of trajectory reshaping with a predictive trajectory decoder. The use of large pre-trained language models to obtain word embeddings allows us to offer a flexible and intuitive user interface, while lowering the requirements on the number of training examples. We validate the proposed models in a series of experiments in simulation and in real-world tests with a robotic arm. Finally, we show that the proposed trajectory reshaping method is highly preferred by users in comparison with baseline methods both in terms of ease of use and performance.

## II. RELATED WORK

**Robots and language:** As robots become more prevalent in environments outside of laboratories and dedicated manufacturing spaces, it is important to offer non-expert users simple ways of communication with machines. Natural language is an ideal candidate, given that interfaces such as mouse-and-keyboard, touchscreens and programming languages are powerful, but require extensive training for proper usage [22]. Multiple facets of language-based human-robot interaction have been studied in literature, such as instruction understanding [23, 24], motion plan generation [9, 12, 16, 25], human-robot cooperation [26], semantic belief propagation [18, 19], and visual language navigation [11, 27]. Most of the recent works in the field have shifted from representing language in terms of classical grammatical structure towards data-driven techniques, due higher flexibility in knowledge representations [22].

**Multi-modal robotics representations:** Representation learning is a rapidly growing field. The existing visual-language representation approaches primarily rely on BERT-style [1] training objectives to model the cross-modal alignments. Common downstream tasks consist of visual question-answering, grounding, retrieval and captioning etc. [28]–[31]. Learning representations for robotics tasks poses additional challenges, as perception data is conditioned on the motion policy and model dynamics [4]. Visual-language navigation of embodied agents is well-established field with clear benchmarks and simulators [32, 33], and multiple works explore the alignment of vision and language data by combining pre-trained models with fine-tuning [34]–[36]. To better model the visual-language alignment, [37] also proposed a co-grounding attention mechanism. In the manipulation domain we also find the work of [10], which uses CLIP [5] embeddings to combine semantic and spatial information. In this paper we also need to align the semantic information with geometry understanding in order to reshape trajectories according to the desired task specifications.

**Transformers in robotics:** Transformers were originally introduced in the language processing domain [38], but quickly proved to be useful in modeling long-range data dependencies other domains. Within robotics we see the first transformers architectures being used for trajectory

forecasting [39] and reinforcement learning [40, 41]. Our work is the first to present a multi-modal transformer model to align visual-language understanding with robot actions for trajectory reshaping.

## III. APPROACH

### A. Problem Definition

Our overall goal is to provide a flexible language-based interface for human-robot interaction within the context of trajectory reshaping. One typical application for our systems is that of a user re-configuring a robotic arm trajectory that, although already avoids collisions, gets uncomfortably close to a particular fragile obstacles in the environment. We design the trajectory generation system with a sequential waypoint prediction decoder, which takes into account multiple data modalities from geometry and language into a transformer network. The modified trajectory should be as close as possible to the original one throughout its length and respect the original start and goal constraints, while obeying the user’s semantic intent. Fig.2 depicts the expected model behavior in a typical use-case scenario.

Let  $\xi_o : [0, 1] \rightarrow \mathbb{R}^2$  be the original robot trajectory, which is composed by a collection of  $N$  waypoints  $\xi_o = \{(x_1, y_1), \dots, (x_N, y_N)\}$ . We assume that the original trajectory is a reasonable path from the start to the goal positions (*i.e.* avoids collisions) and can be pre-calculated using any desired motion planning algorithm, but falls short of the full task specifications. Let  $L_{in}$  be the user’s natural language input sent to correct the original trajectory, such as  $L_{in} =$  “Stay away from the wine glass”. Let  $\mathcal{O} = \{O_1, \dots, O_M\}$  be a collection of  $M$  objects in the environment, each with a corresponding position  $P(O_i) \in \mathbb{R}^2$  and semantic label, such as  $L(O_i) =$  “glass”. Our goal is to learn a function  $f$  that maps the original trajectory, user command and obstacles towards a modified trajectory  $\xi_{mod}$ , which obeys the user’s semantic objectives:

$$\xi_{mod} = f(\xi_o, L_{in}, \mathcal{O}) \quad (1)$$

### B. Proposed Network Architecture

We approximate function  $f$  from Eq. 1 by a parametrized model  $f_\theta$ , learned directly from data. This mapping is non-trivial since it combines data from multiple distinct modalities, and also ambiguous since there exist multiple

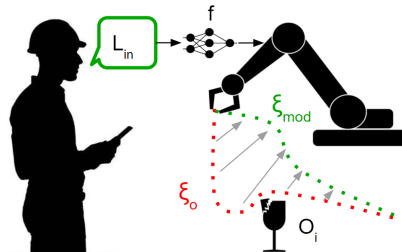


Fig. 2: Typical use case for trajectory reshaping. The user’s natural language command  $L_{in}$  is processed by function  $f$  to reshape the original robot trajectory relative to the target object  $O_i$ .

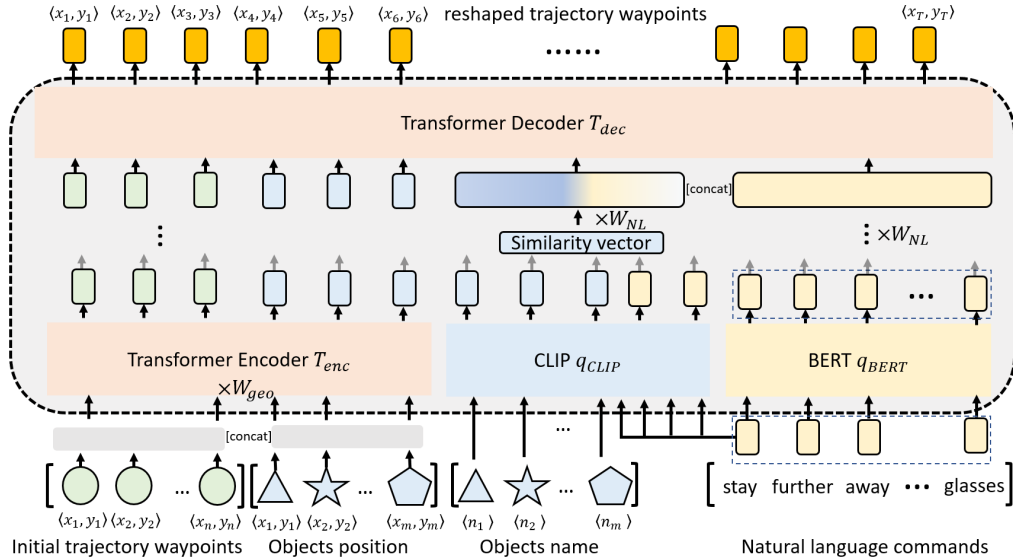


Fig. 3: We proposed a trajectory reshaping system with multimodal attention mechanism. The whole pipeline consists of pretrained CLIP and BERT and trainable transformer encoder and decoders. Given the input of initial trajectory (green), objects with their position and semantic labels (blue), and natural language commands (yellow), the model is trained to generate new trajectories in a predictive manner.

solutions that satisfy the user’s objective. Fig. 3 displays our model architecture, which consists of distinct feature encoders ( $q_{BERT}$ ,  $q_{CLIP}$ ,  $T_{enc}$ ), whose outputs which are fed into a multimodal decoder transformer  $T_{dec}$  for the sequential prediction of the output trajectory  $\xi_{mod}$ . In more detail:

**Language encoding:** We use a pre-trained language model encoder, BERT [1], to produce semantic features  $q_{BERT}(z^{in}|L_{in})$  from the user’s input. The use of a large language model creates more flexibility in the natural language input, allowing the use of synonyms (shown in Section IV-A) and less training data, given that the encoder has already been trained with a massive text corpus. In addition, we use the pre-trained text encoder from CLIP [5] to extract latent embeddings from both the user’s text and the  $M$  object semantic labels ( $q_{CLIP}(z|L)$ ), which enable us compute a similarity vector between the embeddings, and use this information to identify user’s target object. In Section V we discuss how the CLIP model can potentially be used directly with visual data as opposed to textual object labels.

**Geometry encoding:** The original trajectory  $\xi_o$  is composed of low-dimensional tokens  $(x_i, y_i) \in \mathbb{R}^2$ . In order to extract more meaningful information from each waypoint, we apply a linear transform with learnable weights  $W_{geo}$  that projects each waypoint into a higher dimensional features space, following the example of [39]. The poses  $P(O_i)$  of each object are also processed with the same linear transform. We then concatenate both feature vectors and use a transformer-based feature encoder  $T_{enc}$ . The use of a transformer is preferred for sequences because its architecture can attend to multiple time steps simultaneously, as opposed to recurrent networks, which suffer with vanishing gradient issues [39].

**Multi-modal transformer decoder:** Feature embeddings from both language and geometry are combined as input to a multi-modal transformer decoder block  $T_{dec}$ . We generate

the reshaped trajectory  $\xi_{mod}$  sequentially, analogously to common transformer-based approaches in natural language [2, 38]. Section IV-A compares sequential generation with other approaches such as regressing to the entire trajectory at once. We also verify that a fully-connected architecture cannot achieve the same performance as the transformer-based model. We use imitation learning to train the model, using the Huber loss [42] between the predicted and ground-truth waypoint locations.

### C. Synthetic Data Generation

Data collection in the robotics domain is challenging, specially when we require alignment between multiple modalities such as language and trajectories. Different strategies range from large-scale online user studies for language labeling [43] all the way to procedural trajectory-language pairs generation using heuristics [16]. Our work relies on a key hypothesis: the use of large-scale language models for feature encoding ( $q_{BERT}$ ,  $q_{CLIP}$ ) relieves some of the pressure in obtaining a diverse set of vocabulary labels, given that the text encoders are able to find semantic synonyms for different sentence structures. Therefore, we generated a small but meaningful set of examples with semantically-driven trajectory modifications. We employed an  $A^*$  planner to generate reasonable initial trajectories  $\xi_o$  in randomized environments with different object configurations, and based on a set of pre-determined semantic combinations, we used the CHOMP motion planner [44] to compute  $\xi_{mod}$  by modifying weights of different cost functions. Our vocabulary involved different directions relative an object (closer or further away from  $\cdot$ , to the left/right/front/back of  $\cdot$ ), intensity changes (a bit/little, much, very), and a thousand object labels sampled from the ImageNet vocabulary. We generated a total of 10,000 trajectory labels. Fig. 4 displays examples of original and reshaped trajectories.

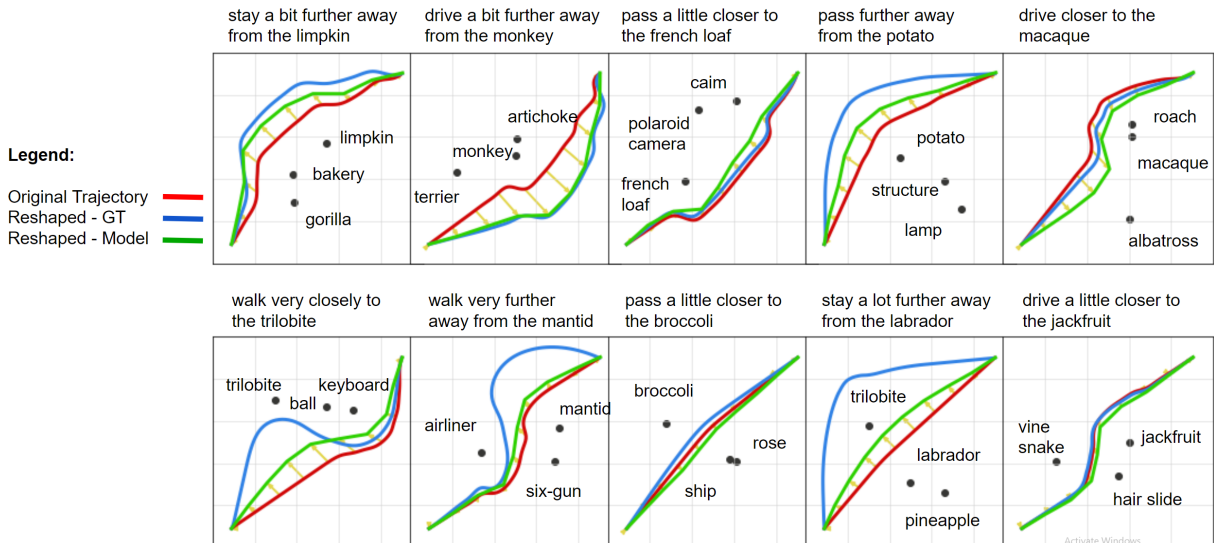


Fig. 4: Randomly picked planning problems extracted from our validation set. Different colors display the original trajectory (calculated using  $A^*$ ), the ground-truth reshaped trajectory (calculated using CHOMP), and the reshaped trajectory outputted by our model.

#### D. Implementation Details

Our language encoders consist of pretrained BERT and CLIP models, with frozen weights. The output sentence embedding  $z^{\text{in}} \in \mathbb{R}^{1 \times 768}$  is concatenated with the similarity vector of size  $1 \times M$  generated from the CLIP embeddings, and scaled to a feature vector of size  $1 \times 256$ . For the geometrical inputs we concatenate an array with  $M$  object poses with a sequence of 100 waypoints from  $\xi_o$ , and upscale the features to a vector of size  $1 \times 256$ .  $T_{\text{enc}}$  is a 2-block transformer encoder, and  $T_{\text{dec}}$  is a 4-block transformer. Each transformer has 3 hidden layers with 512 fully-connected neurons in each. We empirically found it helpful to remove Layer and Batch Normalization [45] steps from both transformers in order to better retain the geometry information. The model is trained in a two-step fashion: first we augment the training data by randomly rotating and re-scaling the planning problem; second we fix the start and goal locations to always lay in the lower left and upper right corner of the map respectively, and fine-tune the network. We use the AdamW [46] optimizer with an initial learning rate  $\gamma = 1e - 4$  and a linear warm-up period of 15 epochs. We use a Nvidia Tesla V100 GPU with batch size of 64, and train the model for 500 epochs.

## IV. EXPERIMENTS

We execute experiments in both simulated and real environments to evaluate our trajectory reshaping model. Our goals are the following: 1) Investigate if the combination of pre-trained large-language models together with multi-modal transformers can create efficient and generalizable human-robot interfaces; 2) Quantitatively and qualitatively compare our approach with other classes of human-robot interfaces from the user’s perspective; 3) Demonstrate our natural language method is applicable to real robots.

#### A. Simulation Experiments

We investigate multiple facets of the model’s capabilities and performance in a series of simulated experiments:

**How language influences trajectory behavior:** Given a fixed environment configuration, we evaluate the model’s ability to follow distinct natural language commands. Fig. 5 displays how a gradient in direction and intensity of language commands correctly modifies the resulting decoded path.

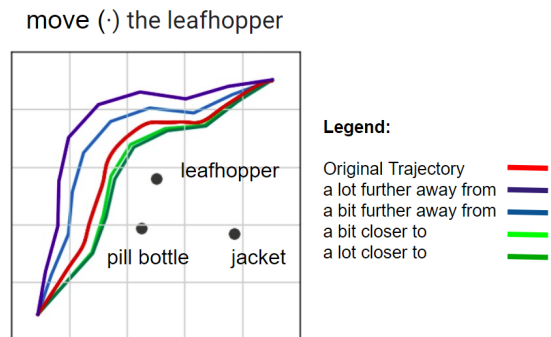


Fig. 5: Influence of different language inputs to the same environment configuration. The network implicitly constrains changes towards the object, as they would cause unsafe behavior.

#### How the model behaves in different planning problems:

To understand the effect of different object configurations and language commands onto the reshaped trajectory, we display in Fig. 4 a set of randomly sampled planning problems from our validation set. We can see from the image that in most cases the reshaped trajectory correctly models the desired user intent, and falls close to the ground-truth reshaped trajectory.

**Vocabulary and object diversity:** One key hypothesis assumed true when designing our model architecture was that the use of pre-trained large language models as feature encoders would make our pipeline amenable to a diverse set of natural language inputs, despite the relatively small amount of training examples. To test this hypothesis we compute



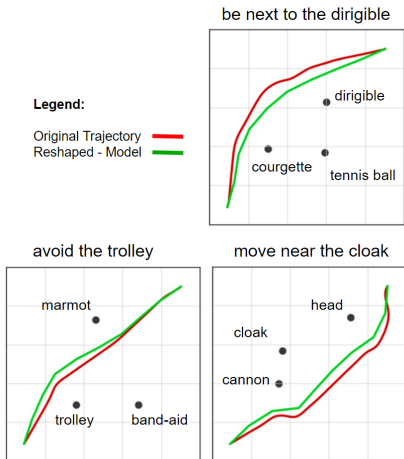


Fig. 6: Trajectory reshaping results using novel vocabulary (not seen in the training data) as the user input. Our model is able to correctly execute the desired semantic commands due to the large capacity of the BERT and CLIP text encoders.

results using with novel user commands, with vocabulary not present in our training language labels. Fig. 6 shows that our model still executes the expected behavior, being able to find the correct semantic meaning despite the new words.

**Baseline architectures:** We perform ablation studies comparing different variations of our proposed multi-modal transformer against baseline architectures. We employ a fully-connected network (FCN) for regression (5 hidden layers, with 512 neurons), that takes as input a single 1D vector composed of the concatenated trajectory, object positions and embedded language features (BERT and similarity vector from CLIP) and outputs a 1D vector with coordinates of all waypoints. The best fully-connected architecture and training procedure was found through a grid-search over the number of hidden layers, neurons, batch size and initial learning rate (64 models in total). Table I summarizes the results, and shows that the best architecture is composed by the multi-modal transformer without layer and batch norm, and using a sequential predictive decoder. The naive predictor referenced in the table shows the loss in the case where we simply copy the original trajectory as the model output, and serves as a baseline loss value. Similarly to previous studies in trajectory forecasting [39], we find that transformers drastically improve model performance, likely due to their unique ability to extract features combining features from distant waypoints.

Model	Features	Test. Loss
Naive predictor		0.0051
FCN w/ regression	2.4M	0.0025
<b>Ours</b>	10.6M	<b>0.002</b>

TABLE I: Baseline architecture comparisons

### B. User Evaluation in Real Robot Experiments

We also evaluate our system with real-world experiments, and compare our method with the use of multiple human-robot interfaces. We use a 7-DOF PANDA Arm robot equipped

with a claw gripper, and execute tasks on a  $1 \times 1\text{m}$  tabletop workspace. A standard desktop computer with an off-the-shelf GPU connected to the robot computes the original trajectories, executes our model, and runs low-level controls for the arm. We operate the model using 2D planar projections of the original robot trajectories, and respect the original waypoint heights when executing the reshaped motion plans. We use object positions given by markers, but we discuss the use of vision-based localization in Section V.

The goal of the study is to have the user control a robotic bartender. A traditional motion planning algorithm calculates an initial trajectory to transport a bottle of wine towards a cocktail shaker and pour the liquid inside (we leave the problem of learning how to make fancy drinks for future iterations of this work). This original trajectory comes dangerously close to toppling over a tower of crystal glasses, and the user needs to interact with the robot to make the end-effector trajectory safer. As seen in Fig. 7 we test 4 different human-robot interfaces: natural language (NL-ours), kinesthetic teaching (KT), trajectory drawing (Draw), and programming obstacle avoidance weights via a keyboard and mouse (Prog). A top-down view of the experimental platform is seen in Fig. 8. All user interactions followed a study protocol approved by the Technical University of Munich’s ethics committee, and we conducted a total of 10 interviews.

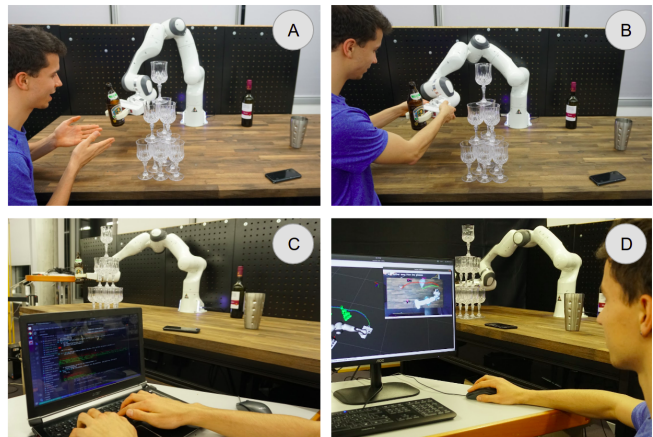


Fig. 7: Human-robot interfaces tested in the user study: a) natural language (NL), b) kinesthetic teaching (KT), c) programming via keyboard, and d) trajectory drawing.

**Quantitative user evaluation:** We measured statistics on the number of iterations, success rate and total time taken for users to modify trajectories using the different interfaces. From Table II we see that the programming interface takes by far the longest for users to master, and requires a large number of iterations. In the meanwhile, NL is the fastest option. We see a large number of failures for kinesthetic teaching and drawing because user inputs are often times kinematically infeasible by the robot joints. The natural language method proved to be the most robust, and we found no failure cases during in the study.

**Qualitative user evaluation:** After the experiments we asked users to rate the trajectories produced by different

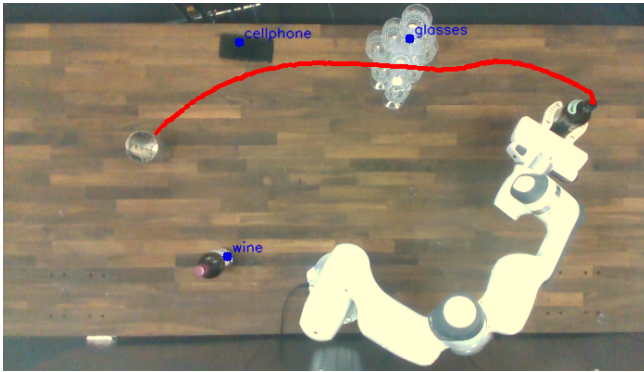


Fig. 8: Experimental platform used for the user study. The tabletop contains three objects (a cellphone, a wine bottle and crystal glasses), and the original robot trajectory (in red) passes dangerously close to the tower of glasses.

Interface	Avg. iterations	Success rate (%)	Avg. Time (s)
NL	<b>1.33</b>	<b>100</b>	<b>81</b>
KT	1.78	56.24	139
Draw	1.89	64.7	120
Prog	4.00	91.66	284

TABLE II: Statistics collected over the user study experiment

interfaces according to different criteria in a psychometric questionnaire:

- 1) How satisfied were you with the final robot motion?
- 2) How easy was refining the robot motion?
- 3) How safe was the final robot trajectory?
- 4) How natural was the human-machine interaction?
- 5) How predictable was the trajectory for you?

Table III summarizes the responses. We can see that most methods present a similar user satisfaction level except for programming, which was rated lower likely due to the difficulty of interaction. NL was rated as the easiest and most natural method, but at the same time was deemed less predictable than KT and drawing because with these two methods users have direct control over the final trajectory.

Interface	Satisfied	Ease of use	Safety	Natural	Predictable
NL	<b>90</b>	<b>92</b>	92	<b>98</b>	72
KT	<b>90</b>	88	88	78	<b>96</b>
Draw	88	74	<b>100</b>	80	88
Prog	62	58	82	62	48

TABLE III: User ratings collected in the user study

**Final experimental remarks:** Overall we find from the experiments that our proposed natural language model stands as a strong alternative to traditional human-robot interfaces. Kinesthetic teaching often not a viable solution for real-world trajectory reshaping depending on the robot’s size, form factor and actuator types. Trajectory drawing is also not a robust solution, as human-defined trajectories often extrapolate acceleration and kinematic constraints. By combining semantic and geometrical information, our method provides a natural and effective interface for trajectory reshaping.

## V. CONCLUSION AND DISCUSSION

In this paper, we present a novel system with multimodal attention mechanism for semantic trajectory reshaping. Given flexible natural language commands and an initial trajectory, it can effectively reshape the trajectory consistent with the language commands. The multimodal attention architecture provides a way for us to jointly align natural language features and the geometrical cues.

We verify through our experiments that by leveraging large pretrained language models like BERT and CLIP, our proposed system creates a flexible and intuitive user interface. Given that these foundational models train on massive corpus of data, we are able to train our robotics system with a smaller dataset, and let the language model find similarities between sentences if novel vocabulary is used.

By evaluating our methods on both simulation and real-world application scenarios, we show that our model outperforms baseline trajectory reshaping approaches in terms of loss values and quality of results. From the user’s perspective, we also perform a user study and show that users significantly prefer our natural language interface in opposition to other methods such as kinesthetic teaching or programming interfaces. Our method is faster to use, and results in a higher success rate.

Even though our study does not address the visual modality, we are confident that our current architecture would also be able to align this additional data through the CLIP encoder. For future iterations of this work we’re interested in using images of the objects as opposed to directly inputting the semantic labels to the model. In addition, in the future we are interested in designing methods that also consider causal relations among objects from the natural language commands in order for the robot to execute more complex task-driven behaviors.

## ACKNOWLEDGMENTS

AB gratefully acknowledges the support from TUM-MIRMI.

## REFERENCES

- [1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [2] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [3] S. Smith, M. Patwary, B. Norick, P. LeGresley, S. Rajbhandari, J. Casper, Z. Liu, S. Prabhunoye, G. Zerveas, V. Korthikanti *et al.*, “Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model,” *arXiv preprint arXiv:2201.11990*, 2022.
- [4] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill *et al.*, “On the opportunities and risks of foundation models,” *arXiv preprint arXiv:2108.07258*, 2021.
- [5] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.
- [6] L. Yuan, D. Chen, Y.-L. Chen, N. Codella, X. Dai, J. Gao, H. Hu, X. Huang, B. Li, C. Li *et al.*, “Florence: A new foundation model for computer vision,” *arXiv preprint arXiv:2111.11432*, 2021.

- [7] H. Alwassel, D. Mahajan, B. Korbar, L. Torresani, B. Ghanem, and D. Tran, "Self-supervised learning by cross-modal audio-video clustering," *Advances in Neural Information Processing Systems*, vol. 33, pp. 9758–9770, 2020.
- [8] S. Ma, S. Vemprala, W. Wang, J. Gupta, Y. Song, D. McDuff, and A. Kapoor, "Compass: Contrastive multimodal pretraining for autonomous systems," February 2022.
- [9] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch, "Language models as zero-shot planners: Extracting actionable knowledge for embodied agents," *arXiv preprint arXiv:2201.07207*, 2022.
- [10] M. Shridhar, L. Manuelli, and D. Fox, "Cliport: What and where pathways for robotic manipulation," in *Conference on Robot Learning*. PMLR, 2022, pp. 894–906.
- [11] Y. Hong, Q. Wu, Y. Qi, C. Rodriguez-Opazo, and S. Gould, "A recurrent vision-and-language bert for navigation. arxiv 2021," *arXiv preprint arXiv:2011.13922*.
- [12] L. Shao, T. Migimatsu, Q. Zhang, K. Yang, and J. Bohg, "Concept2robot: Learning manipulation concepts from instructions and human demonstrations," *The International Journal of Robotics Research*, vol. 40, no. 12-14, pp. 1419–1434, 2021.
- [13] W. Goodwin, S. Vaze, I. Havoutis, and I. Posner, "Semantically grounded object matching for robust robotic scene rearrangement," *arXiv preprint arXiv:2111.07975*, 2021.
- [14] A. S. Chen, S. Nair, and C. Finn, "Learning generalizable robotic reward functions from in-the-wild human videos," *arXiv preprint arXiv:2103.16817*, 2021.
- [15] J. Fu, A. Korattikara, S. Levine, and S. Guadarrama, "From language to goals: Inverse reinforcement learning for vision-based instruction following," *arXiv preprint arXiv:1902.07742*, 2019.
- [16] S. Stepputtis, J. Campbell, M. Phielipp, S. Lee, C. Baral, and H. Ben Amor, "Language-conditioned imitation learning for robot manipulation tasks," *Advances in Neural Information Processing Systems*, vol. 33, pp. 13 139–13 150, 2020.
- [17] P. Goyal, R. J. Mooney, and S. Niekum, "Zero-shot task adaptation using natural language," *arXiv preprint arXiv:2106.02972*, 2021.
- [18] J. Arkin, D. Park, S. Roy, M. R. Walter, N. Roy, T. M. Howard, and R. Paul, "Multimodal estimation and communication of latent semantic knowledge for robust execution of robot instructions," *The International Journal of Robotics Research*, vol. 39, no. 10-11, pp. 1279–1304, 2020.
- [19] M. R. Walter, S. Patki, A. F. Daniele, E. Fahnestock, F. Duvall, S. Hemachandra, J. Oh, A. Stentz, N. Roy, and T. M. Howard, "Language understanding for field and service robots in a priori unknown environments," *arXiv preprint arXiv:2105.10396*, 2021.
- [20] S. M. LaValle, *Planning Algorithms*. Cambridge, U.K.: Cambridge University Press, 2006, available at <http://planning.cs.uiuc.edu/>.
- [21] C. E. Garcia, D. M. Prett, and M. Morari, "Model predictive control: Theory and practice—a survey," *Automatica*, vol. 25, no. 3, pp. 335–348, 1989.
- [22] S. Tellex, N. Gopalan, H. Kress-Gazit, and C. Matuszek, "Robots that use language," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 3, pp. 25–55, 2020.
- [23] M. MacMahon, B. Stankiewicz, and B. Kuipers, "Walk the talk: Connecting language, knowledge, and action in route instructions," *Def*, vol. 2, no. 6, p. 4, 2006.
- [24] N. H. Kirk, D. Nyga, and M. Beetz, "Controlled natural languages for language generation in artificial cognition," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2014, pp. 6667–6672.
- [25] C. Lynch and P. Sermanet, "Language conditioned imitation learning over unstructured data," *arXiv preprint arXiv:2005.07648*, 2020.
- [26] V. Raman, C. Lignos, C. Finucane, K. C. Lee, M. P. Marcus, and H. Kress-Gazit, "Sorry dave, i'm afraid i can't do that: Explaining unachievable robot tasks using natural language," in *Robotics: Science and Systems*, vol. 2, no. 1. Citeseer, 2013, pp. 2–1.
- [27] A. Majumdar, A. Shrivastava, S. Lee, P. Anderson, D. Parikh, and D. Batra, "Improving vision-and-language navigation with image-text pairs from the web," in *European Conference on Computer Vision*. Springer, 2020, pp. 259–274.
- [28] C. Sun, A. Myers, C. Vondrick, K. Murphy, and C. Schmid, "Videobert: A joint model for video and language representation learning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 7464–7473.
- [29] J. Lu, D. Batra, D. Parikh, and S. Lee, "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," *Advances in neural information processing systems*, vol. 32, 2019.
- [30] L. Zhou, H. Palangi, L. Zhang, H. Hu, J. Corso, and J. Gao, "Unified vision-language pre-training for image captioning and vqa," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 13 041–13 049.
- [31] W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai, "Vi-bert: Pre-training of generic visual-linguistic representations," *arXiv preprint arXiv:1908.08530*, 2019.
- [32] A. Szot, A. Clegg, E. Undersander, E. Wijmans, Y. Zhao, J. Turner, N. Maestre, M. Mukadam, D. S. Chaplot, O. Maksymets *et al.*, "Habitat 2.0: Training home assistants to rearrange their habitat," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [33] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. Reid, S. Gould, and A. van den Hengel, "Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3674–3683.
- [34] W. Hao, C. Li, X. Li, L. Carin, and J. Gao, "Towards learning a generic agent for vision-and-language navigation via pre-training," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 13 134–13 143.
- [35] J. Thomason, M. Murray, M. Cakmak, and L. Zettlemoyer, "Vision-and-dialog navigation," in *Conference on Robot Learning*. PMLR, 2020, pp. 394–406.
- [36] K. Nguyen and I. Daumé, "Help, anna! visual navigation with natural multimodal assistance via retrospective curiosity-encouraging imitation learning," 09 2019.
- [37] C.-Y. Ma, Z. Wu, G. AlRegib, C. Xiong, and Z. Kira, "The regretful agent: Heuristic-aided navigation through progress estimation," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 6725–6733.
- [38] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [39] F. Giuliari, I. Hasan, M. Cristani, and F. Galasso, "Transformer networks for trajectory forecasting," in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 10 335–10 342.
- [40] L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch, "Decision transformer: Reinforcement learning via sequence modeling," *Advances in neural information processing systems*, vol. 34, 2021.
- [41] M. Janner, Q. Li, and S. Levine, "Offline reinforcement learning as one big sequence modeling problem," *Advances in neural information processing systems*, vol. 34, 2021.
- [42] P. J. Huber, "Robust estimation of a location parameter," in *Breakthroughs in statistics*. Springer, 1992, pp. 492–518.
- [43] R. Bonatti, A. Buckner, S. Scherer, M. Mukadam, and J. Hodgins, "Batteries, camera, action! learning a semantic control space for expressive robot cinematography," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 7302–7308.
- [44] N. Ratliff, M. Zucker, J. A. Bagnell, and S. Srinivasa, "CHOMP: Gradient optimization techniques for efficient motion planning," in *Robotics and Automation, 2009. ICRA'09. IEEE International Conference on*. IEEE, 2009, pp. 489–494.
- [45] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.
- [46] Y. You, J. Li, S. Reddi, J. Hseu, S. Kumar, S. Bhojanapalli, X. Song, J. Demmel, K. Keutzer, and C.-J. Hsieh, "Large batch optimization for deep learning: Training bert in 76 minutes," *arXiv preprint arXiv:1904.00962*, 2019.