

VELMA: Verbalization Embodiment of LLM Agents for Vision and Language Navigation in Street View

Raphael Schumann¹, Wanrong Zhu², Weixi Feng², Tsu-Jui Fu²,
Stefan Riezler^{1,3}, William Yang Wang²

¹Computational Linguistics, Heidelberg University, Germany

²University of California, Santa Barbara

³IWR, Heidelberg University, Germany

{rschuman, riezler}@cl.uni-heidelberg.de, {wanrongzhu, weixifeng, tsu-juifu, william}@cs.ucsb.edu

Abstract

Incremental decision making in real-world environments is one of the most challenging tasks in embodied artificial intelligence. One particularly demanding scenario is Vision and Language Navigation (VLN) which requires visual and natural language understanding as well as spatial and temporal reasoning capabilities. The embodied agent needs to ground its understanding of navigation instructions in observations of a real-world environment like Street View. Despite the impressive results of LLMs in other research areas, it is an ongoing problem of how to best connect them with an interactive visual environment. In this work, we propose VELMA, an embodied LLM agent that uses a verbalization of the trajectory and of visual environment observations as contextual prompt for the next action. Visual information is verbalized by a pipeline that extracts landmarks from the human written navigation instructions and uses CLIP to determine their visibility in the current panorama view. We show that VELMA is able to successfully follow navigation instructions in Street View with only two in-context examples. We further fine-tune the LLM agent on a few thousand examples and achieve around 25% relative improvement in task completion over the previous state-of-the-art for two datasets.

Introduction

Large language models (LLMs), which have shown impressive reasoning capabilities in traditional natural language processing tasks, are increasingly used as the reasoning engine of embodied agents for, e.g., household robots (Shridhar et al. 2020), video games (Wang et al. 2023) and indoor navigation (Zhou, Hong, and Wu 2023). These tasks are mostly based on simulations that either feature computer-generated images with a fixed set of displayable objects and textures, or are limited in scale and trajectory length. In this paper, we present a verbalization embodiment of an LLM agent (VELMA) for urban vision and language navigation in Street View. The unique challenge of this task is the combination of a large-scale environment derived from an actual road network, real-world panorama images with dense street scenes, and long navigation trajectories. The agent needs to ground its understanding of the navigation instructions in the observable environment and reason about the next action to reach the target location. The navigation instructions

Navigate to the described target location!

Action Space: forward, left, right, turn_around, stop
Navigation Instructions:

"Oriентate yourself such that a blue bench is on your right, go to the end of the block and make a right. Follow the park on your left and make a right at the intersection. Pass the black fire hydrant on your right and stop when you get to a gray door on the brown building."

Action Sequence:

There is a blue bench on your left.

1. **turn_around**

There is a blue bench on your right.

2. **forward**

There is a 3-way intersection.

3. **right**

4. **forward**

There is a park on your left.

5. **forward**

There is a park on your left.

6. **forward**

There is a 4-way intersection.

7. **<next word prediction>**

Figure 1: Prompt sequence used to utilize LLMs for VLN in Street View. Verbalized observations of the visual environment are in green and appended to the prompt at each step. Agent actions (blue) are acquired by LLM next word prediction. Highlighting of text for visual presentation only. Full navigation trajectories are, on average, 40 steps long.

are written by humans and include open-ended landmark references and directional indications intended to guide the agent along the desired path. In order to leverage the reasoning capabilities of LLMs, we use embodiment by verbalization, a workflow where the task, including the agent's trajectory and visual observations of the environment, is verbalized, thus embodying the LLM via natural language. Figure 1 shows the verbalization at step 7 of the current trajectory for a given navigation instance. At each step, the LLM is prompted with the current text sequence in order to predict the next action. Then the predicted action is executed in the environment, and the new observations are verbalized and appended to the prompt. This is repeated until the agent eventually predicts to stop.

The main contributions of our work are as follows:
(i) We introduce VELMA, to our knowledge, the first LLM-

based agent for urban VLN. (ii) We report few-shot results for the urban VLN task and achieve new state-of-the-art performance by finetuning our agent on the training set. (iii) We address and resolve limitations of the commonly used Touchdown environment (Chen et al. 2019), making it amenable for few-shot agents.

Related Work

Outdoor VLN Agent models for the outdoor/urban VLN task (Chen et al. 2019) commonly follow a sequence-to-sequence architecture where encoded text and image representations are fused for each decoder step (Xiang, Wang, and Wang 2020; Mehta et al. 2020; Schumann and Riezler 2022; Sun et al. 2023). Other proposed agents employ pre-trained vision and language transformers that are finetuned on task-specific data (Zhu et al. 2021; Armitage, Impett, and Sennrich 2023). Zhong et al. (2021) represent the visual environment by symbols using semantic segmentation and extreme downsampling of panorama images, but their agent does not improve over previous success rates. Other work uses CLIP to score the presence of extracted landmarks at each panorama node in a graph and uses this information to plan a route for given navigation instructions (Shah et al. 2022). Their non-urban environment has a graph with 300 nodes, and the navigation path is planned *a priori* with full access to all panorama images and landmark scores. In contrast, our agent is embodied and has to plan ad-hoc with access to directly observed information only.

Indoor VLN Indoor agents (Fried et al. 2018; Wang et al. 2019; Tan, Yu, and Bansal 2019; Fu et al. 2020; Zhu et al. 2020; Qi et al. 2020; Hong et al. 2021; Chen et al. 2021; Li, Tan, and Bansal 2022) are used for navigation datasets like R2R (Anderson et al. 2018) and RxR (Ku et al. 2020) or ObjectNav (Ramakrishnan et al. 2021; Zhou et al. 2023). Khan-delwal et al. (2022) showed that using the CLIP encoder for image features improves performance for a range of vision and language tasks. Recently, Zhou, Hong, and Wu (2023) introduced an LLM-based agent for R2R that incorporates image information by transcribing its entire content with an image-to-text model. This is feasible because the navigation trajectories are only six steps on average compared to 40 steps in the urban VLN task considered in our work. Another notable indoor VLN agent uses CLIP to directly predict the next action by scoring the compatibility of the current sub-instruction with available waypoint images (Dorbala et al. 2022).

Urban VLN Environment

We use the Touchdown environment introduced by Chen et al. (2019). The environment is based on Google’s Street View and features 29,641 full-circle panorama images connected by a navigation graph. It covers the dense urban street network spanning lower Manhattan. The navigation graph is a directed graph $G = \langle V, E \rangle$ where each edge $\langle v, v' \rangle \in E$ is associated with $\alpha_{\langle v, v' \rangle}$ which is the heading direction from node v to node v' ranging from 0° to 360° . The agent state $s = (v, \alpha)$ is composed of its current position $v \in V$ and its heading direction α . The agent can move by executing an



Figure 2: The Touchdown environment introduced by Chen et al. (2019) can require action sequences that are semantically inconsistent with the correct navigation instructions. In the depicted subgraph, the action sequence to move from node 1 to node 5 is to move FORWARD four times. The semantically correct sequence of actions would include a right turn in between. We fix the problem by modifying the environment behavior and selecting the desired direction at intersections in relation to all outgoing streets.

action $a \in \{\text{FORWARD}, \text{LEFT}, \text{RIGHT}, \text{STOP}\}$. The state transition function $s_{t+1} = \phi(a_t, s_t)$ defines the behavior of the agent executing an action. In Chen et al. (2019), the agent’s heading α_t at position v is restricted to align with the heading of an outgoing edge $\alpha_{\langle v, v' \rangle}$. In case of the RIGHT action, the new state s_{t+1} is $(v, \alpha_{\langle v, \vec{v} \rangle})$ where \vec{v} is the neighboring node closest to the right of the agent’s current heading. In other words, the agent is rotated in place to the right until it *snaps* to the direction of an outgoing edge. Likewise, for the LEFT action. In the case of the FORWARD action, the agent moves along the edge $\langle v, v' \rangle$ according to its current heading direction $\alpha_t = \alpha_{\langle v, v' \rangle}$. The environment is then forced to automatically rotate the agent’s heading towards an outgoing edge: $\alpha_{t+1} = \alpha_{\langle v', v^* \rangle}$ where v^* is the neighbor node in the direction closest to the previous heading α_t .

Alignment Inconsistencies in Touchdown

As described in Schumann and Riezler (2022), the automatic rotation mentioned above can lead to generalization problems, e.g., when moving towards the flat side of a T-intersection. For example, if the agent is automatically rotated towards the right facing street and subsequently executes the RIGHT action, it rotates towards the direction it came from instead of clearing the intersection in the intended direction. The same problem also occurs at intersections with more than three directions. Figure 2 gives an illustrative example that shows the navigation graph at a 4-way intersection. Because the environment is derived from a real-world street layout, the nodes in the graph are not perfectly arranged as in an artificial grid world. In order to make a right turn at the intersection and to follow the route from v^1 to v^5 , one expects to use the action sequence [FORWARD, FORWARD, RIGHT, FORWARD, FORWARD]. However, when the agent reaches v^3 , it is automatically rotated towards the closest outgoing edge, in this case, $\langle v^3, v^4 \rangle$. This

is because the rotation $20^\circ \rightarrow 50^\circ$ towards v_4 is shorter than the rotation $20^\circ \rightarrow 345^\circ$ towards v_7 . As such, the required sequence of actions to go from v^1 to v^5 in Chen et al. (2019)’s environment is [FORWARD, FORWARD, FORWARD, FORWARD]. This is unpredictable and is not correctly aligned with “*turn right at the intersection*” instructions.¹ To alleviate this problem, Schumann and Riezler (2022) explicitly feed the change of heading at each timestep as additional input to their model. This enables the agent to anticipate the unexpected rotation and to adapt to it. Because adding heading delta values to the text-based interface makes it convoluted and unnecessarily difficult for few-shot learning, we propose a more intuitive way to solve this ambiguity at intersections. We modify the state transition function ϕ such that the agent is not automatically rotated when moving FORWARD. This means the agent’s heading α_t is not automatically aligned with an outgoing edge. Instead, the direction is selected in relation to all outgoing edges. The agent at node v^3 in Figure 2 has the nodes v^6 , v^7 and v^4 in front. The forward direction is selected as the middle one of the three edges, the right direction as the right-most edge, and the left direction as the left-most edge. This means that executing the RIGHT action at position v^3 will now rotate the agent towards node v^4 and allows to use the semantically correct sequence of actions for the depicted route. The proposed modification solves the issue of inconsistent action sequences at intersections and allows to use agents that are not specifically trained in this environment.

Turning Around

We additionally introduce the TURN_AROUND action which lets the agent reverse its direction: $s_{t+1} = (v, \alpha_t - 180^\circ)$. In the unmodified environment, this is achieved using the LEFT or RIGHT action on regular street segments. The new action is better aligned with natural language verbalizations of direction reversal and promotes intuitive communication with the environment.

Navigation Task

The objective of the navigation task is to find the goal location by following the given navigation instructions. A navigation instance is defined by the initial state s_1 , target node \hat{v}_T , gold path $(\hat{v}_1, \hat{v}_2, \dots, \hat{v}_T)$ and navigation instructions text $n = (w_1, w_2, \dots, w_N)$. The agent starts at s_1 and predicts the next action a_1 based on the navigation instructions and current observations. These are the panorama image and number of outgoing edges at the current position. The environment processes the action and puts the agent into a new state: $s_2 = \phi(a_1, s_1)$. This is repeated until the agent predicts STOP at the presumed goal location. If the agent stops within one neighboring node of the target node, the navigation objective is considered accomplished.

Challenges

One main challenge to successfully follow the navigation instructions is to reliably detect landmarks in the panorama

¹In the Appendix we show more examples for 3-way, 4-way and 5-way intersections.

| Egocentric Spatial Reasoning | |
|--------------------------------------|--|
| 1. | ... turn so the orange barrier is on your left ... |
| 2. | ... a red truck in front of you ... |
| 3. | ... a playground on the far right corner ahead ... |
| Allocentric Spatial Reasoning | |
| 4. | ... green metal pole with pink flowers on top ... |
| 5. | ... building with columns around the windows ... |
| 6. | ... stop in between Chase and Dunkin’ Donuts ... |
| Temporal Reasoning | |
| 7. | ... go straight until you see Chipotle and then ... |
| 8. | ... once you passed the underpass ... |
| 9. | ... stop when the park on your right ends ... |
| Other | |
| 10. | ... proceed straight through three intersections ... |
| 11. | ... you should see TD Bank on your left ... |
| 12. | ... if you see Dory Oyster Bar, you are too far ... |

Table 1: Reasoning skills the embodied LLM agent must possess in order to successfully complete the navigation task. Each with three example snippets from the navigation instructions.

images along the route. The landmarks mentioned in the instructions are open-ended and can refer to any object or structure found in street scenes, including vegetation, building features, vehicle types, street signs, construction utilities, company logos and store names. The agent also needs to possess different types of reasoning, most importantly spatial reasoning to follow general directions, locate landmarks and evaluate stopping conditions. The agent also needs to understand the temporal aspect of the task and reason about the sequence of previous observations and actions. See Table 1 for example snippets from the navigation instructions.

Datasets

There are two datasets that provide navigation instructions for the environment described in Section : **Touchdown** (Chen et al. 2019) and **Map2seq** (Schumann and Riezler 2021). Each dataset includes around 10k navigation instances, and we utilize them in the more challenging *unseen* scenario introduced by Schumann and Riezler (2022). This means that generalization is crucial because the training routes are located in an area that is geographically separated from the area of development and test routes. The main difference between the two datasets is that Touchdown instructions were written by annotators who followed the route in Street View, while Map2seq instructions were written by annotators that saw a map of the route. The Map2seq navigation instructions were later validated to also be correct in Street View. Another difference is that the initial state in Map2seq orients the agent towards the correct direction which leads to overall better task completion rates than for Touchdown instances.

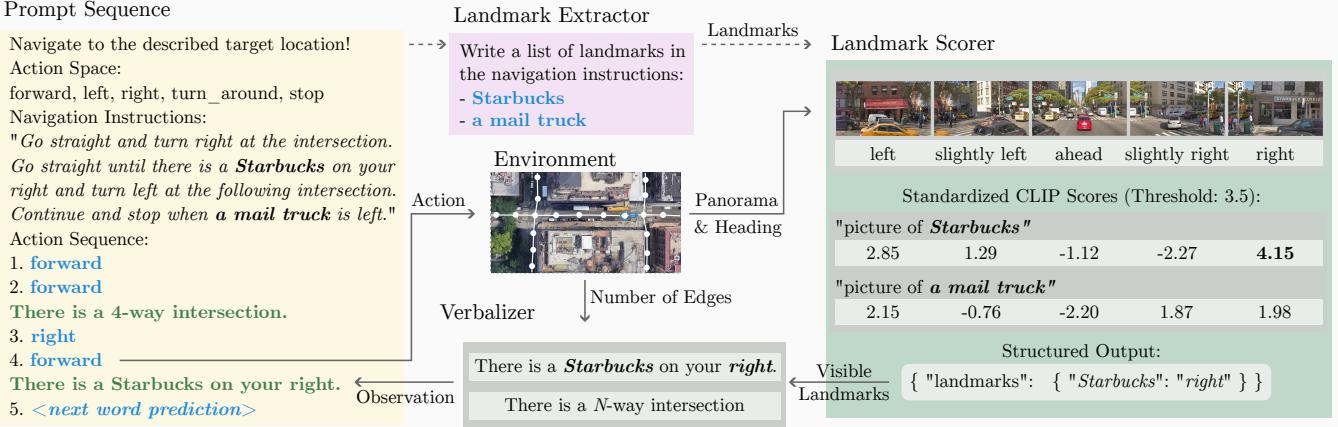


Figure 3: Overview of the proposed agent VELMA navigating in the Street View environment. The prompt sequence includes the task description, navigation instructions, and verbalized navigation trajectory up to the current timestep. The next action is decided by next word prediction utilizing an LLM and subsequently executed in the environment. This puts the agent into a new state, and the landmark scorer determines if an extracted landmark is visible in the current panorama view. The verbalizer takes this landmark information along with the information about a potential intersection and produces the current observations text. This text is then appended to the prompt sequence and again used to predict the next action. This process is repeated until the agent stops and the alleged target location.

LLM Agent

In this section, we propose the urban VLN agent that uses an LLM to reason about the next action. To this end, we verbalize the navigation task, especially the environment observations. The workflow includes the extraction of landmarks that are mentioned in the instructions and determining their visibility in the current panorama image. The verbalizer then integrates the visible landmarks and street intersections into an observation text phrase o_t at each step. The complete text prompt at timestep t is composed as follows:

$$x_t = [d^a, n, d^b, o_1, 1, a_1, o_2, 2, a_2, \dots, o_t, t], \quad (1)$$

where $[]$ denotes string concatenation, d^a and d^b are part of the task description and n is the navigation instructions text. Punctuation and formatting are omitted in the notation for brevity. Figure 3 shows a prompt sequence at $t = 8$ on the left. This formulation of the navigation task enables the agent to predict the next action by next word prediction:

$$a_t = \arg \max_{w \in A} P_{LLM}(w|x_t), \quad (2)$$

where A are the literals of the five defined actions and P_{LLM} is a black-box language model with no vision capabilities.

Landmark Extractor

Each navigation instructions text n mentions multiple landmarks for visual guidance. In order to determine if a mentioned landmark is visible in the current panorama view, we first have to extract them from the instructions text. For this, we create a single prompt that includes five in-context examples of navigation instructions paired with a list of landmarks (shown in the Appendix). It is used by the LLM to automatically generate the list of landmarks (l_1, l_2, \dots, l_L)

mentioned in the given navigation instructions. The landmark extractor is depicted in the top middle of Figure 3 and executed before the navigation starts.

Landmark Scorer

At each step, the agent observes a panorama view p_v^α , defined by its current position v and heading direction α . The view is an 800x460 sized image cut from the panorama with 60° field-of-view. In order to determine if a landmark l_i is visible in the view, we employ a CLIP model (Radford et al. 2021) to embed the image and the caption: "picture of $[l_i]$ ". The similarity score of the two embeddings determines the visibility of the landmark. Because the scores can be biased towards certain types of landmarks, we standardize them using all views p_{train}^* of the ~20k panorama images in the training area. Recall that we operate in the *unseen scenario* where the training area and evaluation area are geographically separated. The standardized score of a landmark is:

$$z(l, p_v^\alpha) = \frac{\text{CLIP}(l, p_v^\alpha) - \mu(C_l)}{\sigma(C_l)} \quad (3)$$

where $C_l = \{\text{CLIP}(l, p_{v'}^{\alpha'}) \mid p_{v'}^{\alpha'} \in p_{train}^*\}$.

If the standardized score is larger than the threshold τ , the landmark is classified as visible in the current view. The process does not require annotations and is completely unsupervised, allowing to score novel landmarks. The threshold is the only tunable parameter in the landmark scorer. Figure 4 shows the distribution of unstandardized CLIP scores and views at different threshold values for two example landmarks. While the views at $\tau = 4.0$ both show the correct landmark, the view at $\tau = 3.0$ for "Bank of America" shows an HSBC branch, and for "yellow truck" it shows a white truck. This suggests that the optimal threshold lies between the two values. As depicted on the right

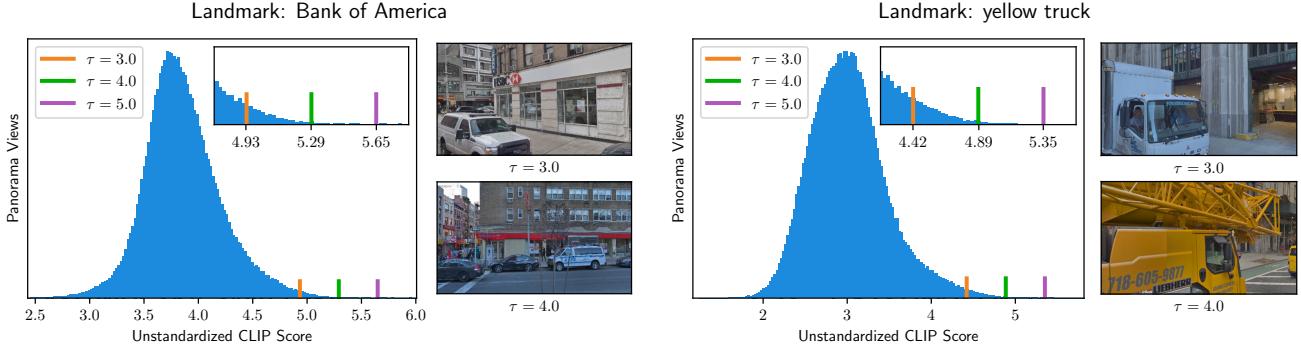


Figure 4: Distribution of CLIP scores between a landmark and panorama images in the training area. The CLIP score represents the semantic similarity of the panorama image and the text caption "picture of [landmark]". The distribution is used to standardize the score of the landmark and a novel panorama. The threshold τ is defined on the standardized score and used to determine the visibility of the landmark in the novel panorama image.

in Figure 3, the agent also evaluates views to the left and right of the current heading. Each panorama view direction ($p_v^{\alpha-90^\circ}, p_v^{\alpha-45^\circ}, p_v^\alpha, p_v^{\alpha+45^\circ}, p_v^{\alpha+90^\circ}$) is associated with a string literal m valued *left*, *slightly left*, *ahead*, *slightly right* or *right*, respectively. A visible landmark l_i and the corresponding direction literal m_i are passed to the verbalizer. A full navigation trajectory includes around 200 image views (40 steps and 5 view directions per step) and each landmark is typically visible in only one or two views.

Verbalizer

The verbalizer is a template-based component that produces environment observations in text form. There are two types of environment observations. First, there are street intersections that are detected based on the number of outgoing edges $N(v)$ at the current node v in the navigation graph. If there are three or more outgoing edges at step t , the verbalizer encodes this information into the observation string o_t^e : "There is a $[N(v)]$ -way intersection". Extracting this information directly from the navigation graph is akin to the junction type embedding used by the ORAR model (Schumann and Riezler 2022) and is motivated by direction arrows displayed in the Street View GUI that human navigators used during data collection. The other type of observations are landmarks visible in the panorama view. The landmark name l_i and direction literal m_i are used to verbalize the observation o_t^l : "There is $[l_i]$ on your $[m_i]$ ". The complete observation is $o_t = [o_t^e, o_t^l]$, where the respective string is empty if no intersection or landmark is detected. The observation is appended to the prompt in Equation 1 and used by the agent to decide the next action.

Experiments

We conducted experiments² to evaluate the navigation performance of the proposed LLM agent in finetuning and in-context learning settings. We used

CLIP-ViT-bigG-14-laion2B-39B-b160k (Schuhmann et al. 2022) as the CLIP model in the landmark scorer. We set the threshold $\tau = 3.5$ for all experiments. The threshold was selected by inspecting the distribution of CLIP scores (as in Figure 4) for a handful of landmarks. On purpose, we did not systematically tune it in order to not violate the premise of few-shot learning.

Landmark Extraction

We ran the landmark extractor once for all instances using GPT-3 (Brown et al. 2020) and used the same extracted landmarks in all experiments. On average, 2.7 landmarks were extracted from a navigation instructions text. Around 58% of the landmarks in the test sets are *novel*, i.e., they are not used in the training instances. In order to estimate the quality of the automatically extracted landmarks, we annotated 50 instances of each development set by hand. For Touchdown we calculated an F1-score of 96.3 (precision: 97.2, recall: 95.4) and the F1-score for Map2seq is 99.6 (precision: 100, recall: 99.3). This shows that GPT-3 reliably extracts landmarks from the instructions text and reusing them for all experiments is minimizing the inaccuracies introduced by this workflow step.

Metrics and Baseline

We use three metrics to measure navigation performance. The task completion (TC) rate is a binary metric that measures whether the agent successfully stopped within one neighboring node of the target location. Shortest-path distance (SPD) calculates the shortest path length between the stopping location and goal location (Chen et al. 2019). Key point accuracy (KPA) measures the ratio of correct decisions at key points. Key points include the initial step, intersections along the gold route, and the target location.

For baselines, we use the current state-of-the-art agent model for urban VLN called ORAR (Schumann and Riezler 2022). The model employs a seq-to-seq architecture where the encoder LSTM reads the navigation instructions text, and the multi-layer decoder LSTM receives image feature

²Project page: <https://velma.schumann.pub/> and code: <https://github.com/raphael-sch/VELMA>

| Model | Development Set | | | | | | Test Set | | | | | |
|--|-----------------|-------------|-----------------|------------|-------------|-----------------|-------------|-------------|-----------------|------------|-------------|-----------------|
| | Touchdown | | | Map2seq | | | Touchdown | | | Map2seq | | |
| | SPD↓ | KPA↑ | TC↑ | SPD↓ | KPA↑ | TC↑ | SPD↓ | KPA↑ | TC↑ | SPD↓ | KPA↑ | TC↑ |
| ORAR-ResNet | 20.0 | - | 15.4±2.2 | 11.9 | - | 27.6±1.8 | 20.8 | - | 14.9±1.2 | 13.0 | - | 30.3±1.8 |
| ORAR [♣] -ResNet | 16.5 | 64.0 | 22.6±0.6 | 10.3 | 74.4 | 29.9±1.7 | 17.4 | 62.3 | 19.1±1.0 | 10.9 | 74.7 | 32.5±1.4 |
| ORAR [♣] -CLIP | 17.5 | 63.7 | 21.5±0.9 | 10.0 | 75.3 | 32.8±1.5 | 17.0 | 63.4 | 20.0±0.1 | 10.5 | 75.1 | 34.0±0.5 |
| 2-Shot In-Context Learning | | | | | | | | | | | | |
| VELMA-Mixtral | 28.4 | 47.2 | 6.5 | 21.1 | 56.8 | 8.0 | - | - | - | - | - | - |
| VELMA-GPT-3 | 22.2 | 49.1 | 6.8 | 19.1 | 58.1 | 9.2 | - | - | - | - | - | - |
| VELMA-GPT-4 | 21.8 | 56.1 | 10.0 | 12.8 | 70.1 | 23.1 | - | - | - | - | - | - |
| LLM Finetuning, full training set | | | | | | | | | | | | |
| VELMA-FT | 18.3 | 62.0 | 23.4±0.2 | 8.7 | 78.7 | 41.3±0.9 | 18.2 | 62.2 | 23.5±0.4 | 9.7 | 78.0 | 40.0±1.0 |
| VELMA-RBL | 15.5 | 63.6 | 26.0±0.6 | 8.3 | 79.5 | 45.3±0.5 | 16.0 | 62.8 | 26.4±1.7 | 8.3 | 79.6 | 47.5±0.7 |

Table 2: Results for the urban VLN task on Touchdown and Map2seq in the *unseen* scenario, meaning the training area is geographically separated from the area where development and test routes are located. ORAR-ResNet (Schumann and Riezler 2022) is the previous best model and follows a seq-to-seq architecture that fuses text and image features during decoding. We re-trained this model in our improved environment (ORAR[♣]-ResNet) and with the same OpenCLIP image embeddings (ORAR[♣]-CLIP) that we use in the landmark scorer. VELMA-GPT-3 and VELMA-GPT-4 models employ our proposed verbalization workflow and are prompted with two in-context examples. Due to cost and data leakage concerns, we evaluate the GPT models on the development sets only. VELMA-FT is LLaMa-7b finetuned on all training text sequences (around 6k for each dataset). The VELMA-RBL finetuning process is described in Section . All experiments are repeated three times with different random seeds (mean/std reported). Bold values are the nominal best results and underlined are best few-shot results.

vectors of the current panorama view as additional input at each action decoding step. The ORAR model is a very strong baseline beating more sophisticated models like the VLN Transformer (Zhu et al. 2021). Because the environment modifications introduced in Section spare the agents from learning specific irregularities, we additionally retrain ORAR in the improved environment for a fair comparison.

Few-Shot Learning Results

The proposed text-only interface allows us to use large language models as reasoners without updating their weights or fusing image representations. The prompt consists of a short task description and two in-context examples (2-shot). The examples are full text sequences along the gold route for randomly selected navigation instances in the training set. The two plots in Figure 5 show that performance scales with parameter count and varies across model families. The FORWARD-Only baseline reveals that OPT (Zhang et al. 2022) can barely compete with a basic heuristic, even at a model size of 65 billion parameters. LLaMa (Touvron et al. 2023a) and especially LLaMa-2 (Touvron et al. 2023b) show promising navigation skills reaching 48.3 and 57.7 key point accuracy (KPA) on Touchdown and Map2seq, respectively. However, this KPA score only translates to task completion (TC) rates of 2.1 and 3.2, revealing that the model is not able to consistently predict correct actions throughout the whole navigation trajectory. Mistral-7b performs on par with a LLaMA-2 model twice its size, but also fails to score task completion rates significantly higher than 3. The only few-shot LLMs that achieve substantial TC rates

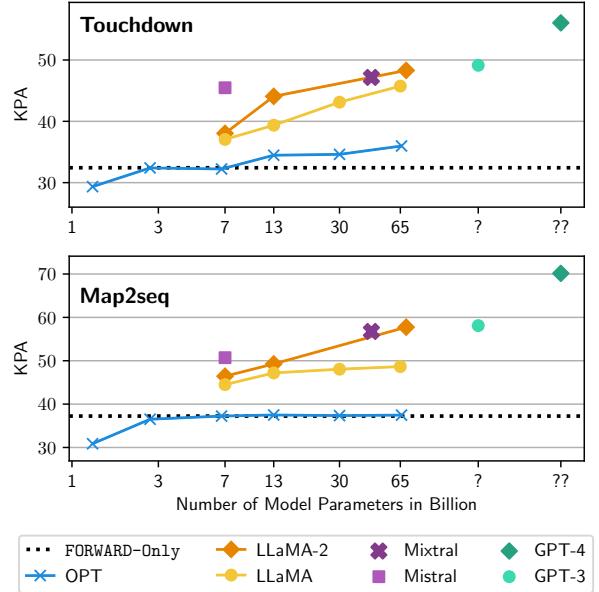


Figure 5: Key point accuracy (KPA) for 2-shot in-context learning of large language models with increasing parameter count. The FORWARD-Only baseline predicts walking forward until the average trajectory length is reached and performs better than predicting random directions.

| Model | Touchdown | | Map2seq | |
|----------|-----------------|-----------------|----------------|-----------------|
| | SPD↓ | TC↑ | SPD↓ | TC↑ |
| no image | 27.4±0.5 | 14.7±0.5 | 9.7±0.2 | 35.2±0.9 |
| CLIP | 21.3±0.5 | 19.5±0.6 | 9.8±0.3 | 37.2±0.5 |
| OpenCLIP | 18.6±0.3 | 22.6±0.4 | 9.8±0.5 | 38.2±0.5 |

Table 3: Vision ablation on the development set. We finetune a separate LLaMa-7b model for each ablation. CLIP refers to `clip-vit-large-patch14` (Radford et al. 2021). The OpenCLIP image model refers to `CLIP-ViT-bigG-14-laion2B-39B-b160k` (Schuhmann et al. 2022).

are GPT-3, GPT-4 (OpenAI 2023) and Mixtral (Mistral AI Team 2023). As listed in Table 2, VELMA-GPT-4 achieves the best results for the 2-shot setting. It reaches 44% and 77% of the TC rate reported for the previous state-of-the-art model ORAR^{*}-ResNet which is a seq-to-seq model that has direct access to image features and was trained on the full training set. In contrast, the LLMs in our work act as a blind agent that solely relies on observation descriptions produced by the verbalizer. This is remarkable because LLMs are not explicitly trained to experience embodiment in a visual environment. This is emergent behavior unearthed by verbalizing the VLN task. We also observe that GPT-4 invokes the TURN_AROUND action in useful ways, e.g. to return a few steps when it notices that it went past the described goal location. This emphasizes the effectiveness of intuitive communication with the environment.

Finetuning Results

To further explore the capabilities of the proposed LLM agent, we finetune LLaMa-7b on all training instances of the respective dataset, denoted by VELMA-FT in Table 2. Each training instance is the full text sequence that is produced by following the gold path. The visibility of landmarks is determined by the landmark scorer during training because gold annotations are not available. There are 6,770 training instances for Touchdown and 5,737 for Map2seq. We finetune for 20 epochs using LoRA (Hu et al. 2022) to adapt query, key and value projections of the attention layer as well as input and output projection of each transformer layer. The best model is selected by task completion on the development set. The resulting agent outperforms the previous state-of-the-art model ORAR^{*} by 10% and 16% relative TC rate. Comparing ORAR^{*} which fuses image features at the vector level to VELMA-FT which finetunes on verbalizations of observations, shows that the text-based environment observations are less prone to overfitting.

Response-Based Learning A navigation task is successfully completed if the agent stops at either the goal location or an adjacent neighboring node. Training the agent with teacher-forcing to exactly follow the gold route penalizes the agent for stopping one step short or one step past the target node, despite accomplishing the navigation objective. Furthermore, the agent can not learn to recover from incor-

Algorithm 1: RBL Optimization of Task Completion

```

Require: mixing ratio  $\lambda$ , training step  $j$ , model weights  $\theta_j$ , gold action sequence  $\hat{\mathbf{a}}$ , prompt  $x_1$ 
if  $\text{random}(0, 1) < \lambda$  then
     $\mathbf{a}_{\theta_j} = \text{StudentForcing}(\theta_j, x_1)$ 
     $\mathbf{a}_j = \arg \max \mathbf{a}_{\theta_j}$ 
    if  $\text{TaskCompletion}(\mathbf{a}_j) = 1$  then
         $\text{loss}_j = \mathcal{L}_{CE}(\mathbf{a}_{\theta_j}, \mathbf{a}_j)$ 
    else
         $\mathbf{a}^* = \text{Oracle}_{\text{stepwise}}(\mathbf{a}_j)$ 
         $\text{loss}_j = \mathcal{L}_{CE}(\mathbf{a}_{\theta_j}, \mathbf{a}^*)$ 
    end if
else
     $\mathbf{a}_{\theta_j} = \text{TeacherForcing}(\theta_j, x_1, \hat{\mathbf{a}})$ 
     $\text{loss}_j = \mathcal{L}_{CE}(\mathbf{a}_{\theta_j}, \hat{\mathbf{a}})$ 
end if

```

rect decisions during inference. We thus train the agent to directly optimize the TC metric while also feeding it its own actions during training, called VELMA-RBL in Table 2. The procedure for VELMA-RBL is inspired by response-based learning (Clarke et al. 2010) and imitation learning (Ross, Gordon, and Bagnell 2011) and is outlined in Algorithm 1. The loss for an instance at training step j is either computed by teacher forcing the gold action sequence $\hat{\mathbf{a}}$, or by student forcing, determined by a mixing parameter λ . In student forcing, the actions decoded by the current model weights θ_j are executed instead of the gold actions. If this trajectory ends within one neighboring node of the target location, the predicted action sequence \mathbf{a}_j is considered correct and used as the reference to train the agent. If the agent stops at the wrong location, an oracle path is computed to provide the optimal counterfactual action at each step in the trajectory. In our case, the oracle’s optimal next action is computed as the shortest path to the goal location. We set $\lambda = 0.5$ to collect training losses in a batch evenly from both training strategies. Manually inspecting trajectories produced by the trained agent, we found improvements of following instructions that have stopping criteria like “Stop a few steps before Y.” or “Stop at X. If you see Y you have gone too far.”. In both cases, the agent learned to walk past the uncertain stopping location and to invoke the TURN_AROUND action in order to walk back once landmark Y appeared. The described training procedure leads to a significant increase of task completion rate by 2.9 and 7.5 for Touchdown and Map2seq, respectively. Overall, our contributions in this work amount to a relative increase of task completion by 77% and 57% over the previously reported state-of-the-art for urban VLN on the Touchdown and Map2seq datasets.

Image Ablation

In this section, we ablate the image model used by the landmark scorer. We finetune a LLaMa-7b model according to Section and use CLIP (Radford et al. 2021), OpenCLIP (Schuhmann et al. 2022) or no image model in the landmark scorer. The latter case means that no landmark observation is passed to the prompt sequence. The results

in Table 3 show that OpenCLIP is better suited for detecting landmarks in our navigation task than the original CLIP model. This is in line with better ImageNet results reported by the OpenCLIP authors and suggests that the agent can directly benefit from further improvements of CLIP models. Appending no landmarks to the prompt sequence further degrades performance, especially on Touchdown.

Conclusion

We introduced VELMA, an agent for urban vision and language navigation, which utilizes a large language model to infer its next action. The LLM is continuously queried with a text prompt that verbalizes the task description, navigation instructions, visual observations, and past trajectory of the agent. In order to include observed landmarks in the prompt, we propose an unsupervised pipeline that extracts landmarks from the instructions and determines their visibility in the current panorama view based on thresholded CLIP scores. We evaluate the embodied LLM agent in a modified version of the commonly used Touchdown environment based on Street View. One proposed modification is fixing a problem at intersections that led to incorrect alignments of action sequences, and another modification adds the TURN_AROUND action which provides a more intuitive way to communicate with the environment. The proposed agent achieves strong few-shot in-context learning results of 10 and 23 task completion rates for Touchdown and Map2seq, respectively, and yields new state-of-the-art results of 26 and 47 task completion rates when finetuned on the full training set. The finetuning results show that verbalization is not an inherent limitation for this task and in-context learning with better base models or improved prompting techniques could outperform our reported few shot results.

Acknowledgments

The research reported in this paper was supported by a Google Focused Research Award.

References

- Anderson, P.; Wu, Q.; Teney, D.; Bruce, J.; Johnson, M.; Sünderhauf, N.; Reid, I.; Gould, S.; and Van Den Hengel, A. 2018. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3674–3683.
- Armitage, J.; Impett, L.; and Sennrich, R. 2023. A Priority Map for Vision-and-Language Navigation with Trajectory Plans and Feature-Location Cues. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 1094–1103.
- Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; Agarwal, S.; Herbert-Voss, A.; Krueger, G.; Henighan, T.; Child, R.; Ramesh, A.; Ziegler, D.; Wu, J.; Winter, C.; Hesse, C.; Chen, M.; Sigler, E.; Litwin, M.; Gray, S.; Chess, B.; Clark, J.; Berner, C.; McCandlish, S.; Radford, A.; Sutskever, I.; and Amodei, D. 2020. Language Models are Few-Shot Learners. In Larochelle, H.; Ranzato, M.; Hadsell, R.; Balcan, M.; and Lin, H., eds., *Advances in Neural Information Processing Systems*, volume 33, 1877–1901. Curran Associates, Inc.
- Chen, H.; Suhr, A.; Misra, D.; Snavely, N.; and Artzi, Y. 2019. TOUCHDOWN: Natural Language Navigation and Spatial Reasoning in Visual Street Environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Long Beach, California.
- Chen, S.; Guhur, P.-L.; Schmid, C.; and Laptev, I. 2021. History aware multimodal transformer for vision-and-language navigation. *Advances in neural information processing systems*, 34: 5834–5847.
- Clarke, J.; Goldwasser, D.; Chang, M.-W.; and Roth, D. 2010. Driving Semantic Parsing from the World’s Response. In *Proceedings of the Fourteenth Conference on Computational Natural Language Learning*, 18–27. Uppsala, Sweden: Association for Computational Linguistics.
- Dorbala, V. S.; Sigurdsson, G.; Piramuthu, R.; Thomason, J.; and Sukhatme, G. S. 2022. Clip-nav: Using clip for zero-shot vision-and-language navigation. In *CoRL 2022 Workshop on Language and Robot Learning*.
- Fried, D.; Hu, R.; Cirik, V.; Rohrbach, A.; Andreas, J.; Morency, L.-P.; Berg-Kirkpatrick, T.; Saenko, K.; Klein, D.; and Darrell, T. 2018. Speaker-Follower Models for Vision-and-Language Navigation. In *Neural Information Processing Systems (NeurIPS)*.
- Fu, T.-J.; Wang, X. E.; Peterson, M. F.; Grafton, S. T.; Eckstein, M. P.; and Wang, W. Y. 2020. Counterfactual vision-and-language navigation via adversarial path sampler. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16*, 71–86. Springer.
- Hong, Y.; Wu, Q.; Qi, Y.; Rodriguez-Opazo, C.; and Gould, S. 2021. Vln bert: A recurrent vision-and-language bert for navigation. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, 1643–1653.
- Hu, E. J.; yelong shen; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; and Chen, W. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on Learning Representations*.
- Khandelwal, A.; Weihs, L.; Mottaghi, R.; and Kembhavi, A. 2022. Simple but Effective: CLIP Embeddings for Embodied AI. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 14829–14838.
- Ku, A.; Anderson, P.; Patel, R.; Ie, E.; and Baldridge, J. 2020. Room-Across-Room: Multilingual Vision-and-Language Navigation with Dense Spatiotemporal Grounding. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 4392–4412.
- Li, J.; Tan, H.; and Bansal, M. 2022. Envedit: Environment editing for vision-and-language navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 15407–15417.
- Mehta, H.; Artzi, Y.; Baldridge, J.; Ie, E.; and Mirowski, P. 2020. Retouchdown: Releasing Touchdown on StreetLearn

- as a Public Resource for Language Grounding Tasks in Street View. In *Proceedings of the Third International Workshop on Spatial Language Understanding (SpLU)*. Online.
- Mistral AI Team. 2023. Mixtral of Experts: A High Quality Sparse Mixture-of-Experts. *Mistral AI Blog*. Accessed: December 18, 2023.
- OpenAI. 2023. GPT-4 Technical Report. *ArXiv*, abs/2303.08774.
- Qi, Y.; Wu, Q.; Anderson, P.; Wang, X.; Wang, W. Y.; Shen, C.; and Hengel, A. v. d. 2020. Reverie: Remote embodied visual referring expression in real indoor environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9982–9991.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; Krueger, G.; and Sutskever, I. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *ICML*.
- Ramakrishnan, S. K.; Gokaslan, A.; Wijmans, E.; Maksymets, O.; Clegg, A.; Turner, J. M.; Undersander, E.; Galuba, W.; Westbury, A.; Chang, A. X.; Savva, M.; Zhao, Y.; and Batra, D. 2021. Habitat-Matterport 3D Dataset (HM3D): 1000 Large-scale 3D Environments for Embodied AI. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Ross, S.; Gordon, G. J.; and Bagnell, J. A. 2011. A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning. In *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS)*. Fort Lauderdale, FL, USA.
- Schuhmann, C.; Beaumont, R.; Vencu, R.; Gordon, C. W.; Wightman, R.; Cherti, M.; Coombes, T.; Katta, A.; Mullis, C.; Wortsman, M.; Schramowski, P.; Kundurthy, S. R.; Crowson, K.; Schmidt, L.; Kaczmarczyk, R.; and Jitsev, J. 2022. LAION-5B: An open large-scale dataset for training next generation image-text models. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Schumann, R.; and Riezler, S. 2021. Generating Landmark Navigation Instructions from Maps as a Graph-to-Text Problem. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 489–502. Online: Association for Computational Linguistics.
- Schumann, R.; and Riezler, S. 2022. Analyzing Generalization of Vision and Language Navigation to Unseen Outdoor Areas. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 7519–7532. Dublin, Ireland: Association for Computational Linguistics.
- Shah, D.; Osinski, B.; Ichter, B.; and Levine, S. 2022. LM-Nav: Robotic Navigation with Large Pre-Trained Models of Language, Vision, and Action. *arXiv:2207.04429*.
- Shridhar, M.; Thomason, J.; Gordon, D.; Bisk, Y.; Han, W.; Mottaghi, R.; Zettlemoyer, L.; and Fox, D. 2020. ALFRED: A Benchmark for Interpreting Grounded Instructions for Everyday Tasks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Sun, Y.; Qiu, Y.; Aoki, Y.; and Kataoka, H. 2023. Outdoor Vision-and-Language Navigation Needs Object-Level Alignment. *Sensors*, 23(13).
- Tan, H.; Yu, L.; and Bansal, M. 2019. Learning to Navigate Unseen Environments: Back Translation with Environmental Dropout. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2610–2621. Minneapolis, Minnesota: Association for Computational Linguistics.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; Rodriguez, A.; Joulin, A.; Grave, E.; and Lample, G. 2023a. LLaMA: Open and Efficient Foundation Language Models. *arXiv:2302.13971*.
- Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; Bikel, D.; Blecher, L.; Ferrer, C. C.; Chen, M.; Cucurull, G.; Esiobu, D.; Fernandes, J.; Fu, J.; Fu, W.; Fuller, B.; Gao, C.; Goswami, V.; Goyal, N.; Hartshorn, A.; Hosseini, S.; Hou, R.; Inan, H.; Kardas, M.; Kerkez, V.; Khabsa, M.; Kloumann, I.; Korenev, A.; Koura, P. S.; Lachaux, M.-A.; Lavril, T.; Lee, J.; Liskovich, D.; Lu, Y.; Mao, Y.; Martinet, X.; Mihaylov, T.; Mishra, P.; Molybog, I.; Nie, Y.; Poulton, A.; Reizenstein, J.; Rungta, R.; Saladi, K.; Schelten, A.; Silva, R.; Smith, E. M.; Subramanian, R.; Tan, X. E.; Tang, B.; Taylor, R.; Williams, A.; Kuan, J. X.; Xu, P.; Yan, Z.; Zarov, I.; Zhang, Y.; Fan, A.; Kambadur, M.; Narang, S.; Rodriguez, A.; Stojnic, R.; Edunov, S.; and Scialom, T. 2023b. Llama 2: Open Foundation and Fine-Tuned Chat Models. *arXiv:2307.09288*.
- Wang, G.; Xie, Y.; Jiang, Y.; Mandlekar, A.; Xiao, C.; Zhu, Y.; Fan, L.; and Anandkumar, A. 2023. Voyager: An Open-Ended Embodied Agent with Large Language Models. *arXiv preprint arXiv: Arxiv-2305.16291*.
- Wang, X.; Huang, Q.; Celikyilmaz, A.; Gao, J.; Shen, D.; Wang, Y.-F.; Wang, W. Y.; and Zhang, L. 2019. Reinforced Cross-Modal Matching and Self-Supervised Imitation Learning for Vision-Language Navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Xiang, J.; Wang, X.; and Wang, W. Y. 2020. Learning to Stop: A Simple yet Effective Approach to Urban Vision-Language Navigation. In *Findings of the Association for Computational Linguistics (ACL Findings)*. Online.
- Zhang, S.; Roller, S.; Goyal, N.; Artetxe, M.; Chen, M.; Chen, S.; Dewan, C.; Diab, M.; Li, X.; Lin, X. V.; et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.
- Zhong, V.; Hanjie, A. W.; Wang, S.; Narasimhan, K.; and Zettlemoyer, L. 2021. SILG: The Multi-domain Symbolic Interactive Language Grounding Benchmark. In Ranzato, M.; Beygelzimer, A.; Dauphin, Y.; Liang, P.; and Vaughan,

- J. W., eds., *Advances in Neural Information Processing Systems*, volume 34, 21505–21519. Curran Associates, Inc.
- Zhou, G.; Hong, Y.; and Wu, Q. 2023. NavGPT: Explicit Reasoning in Vision-and-Language Navigation with Large Language Models. *arXiv:2305.16986*.
- Zhou, K.; Zheng, K.; Pryor, C.; Shen, Y.; Jin, H.; Getoor, L.; and Wang, X. E. 2023. ESC: Exploration with Soft Commonsense Constraints for Zero-shot Object Navigation. *arXiv preprint arXiv:2301.13166*.
- Zhu, W.; Hu, H.; Chen, J.; Deng, Z.; Jain, V.; Ie, E.; and Sha, F. 2020. BabyWalk: Going Farther in Vision-and-Language Navigation by Taking Baby Steps. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2539–2556.
- Zhu, W.; Wang, X.; Fu, T.-J.; Yan, A.; Narayana, P.; Sone, K.; Basu, S.; and Wang, W. Y. 2021. Multimodal Text Style Transfer for Outdoor Vision-and-Language Navigation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, 1207–1221. Online: Association for Computational Linguistics.