

# A New Path: Scaling Vision-and-Language Navigation with Synthetic Instructions and Imitation Learning

Aishwarya Kamath<sup>\*†1</sup> Peter Anderson<sup>\*2</sup> Su Wang<sup>2</sup> Jing Yu Koh<sup>†3</sup> Alexander Ku<sup>2</sup>  
 Austin Waters<sup>2</sup> Yinfei Yang<sup>†4</sup> Jason Baldridge<sup>2</sup> Zarana Parekh<sup>2</sup>  
<sup>1</sup>New York University <sup>2</sup>Google Research <sup>3</sup>Carnegie Mellon University <sup>4</sup>Apple

## Abstract

Recent studies in Vision-and-Language Navigation (VLN) train RL agents to execute natural-language navigation instructions in photorealistic environments, as a step towards robots that can follow human instructions. However, given the scarcity of human instruction data and limited diversity in the training environments, these agents still struggle with complex language grounding and spatial language understanding. Pretraining on large text and image-text datasets from the web has been extensively explored but the improvements are limited. We investigate large-scale augmentation with synthetic instructions. We take 500+ indoor environments captured in densely-sampled 360° panoramas, construct navigation trajectories through these panoramas, and generate a visually-grounded instruction for each trajectory using Marky [63], a high-quality multilingual navigation instruction generator. We also synthesize image observations from novel viewpoints using an image-to-image GAN [27]. The resulting dataset of 4.2M instruction-trajectory pairs is two orders of magnitude larger than existing human-annotated datasets, and contains a wider variety of environments and viewpoints. To efficiently leverage data at this scale, we train a simple transformer agent with imitation learning. On the challenging RxR dataset, our approach outperforms all existing RL agents, improving the state-of-the-art NDTW from 71.1 to 79.1 in seen environments, and from 64.6 to 66.8 in unseen test environments. Our work points to a new path to improving instruction-following agents, emphasizing large-scale training on near-human quality synthetic instructions.

## 1. Introduction

Developing intelligent agents that follow human instructions is a long-term, formidable challenge in AI [66]. A recent focus addressing this problem space is Vision-and-Language Navigation (VLN) [3, 9]. Navigation is an ideal

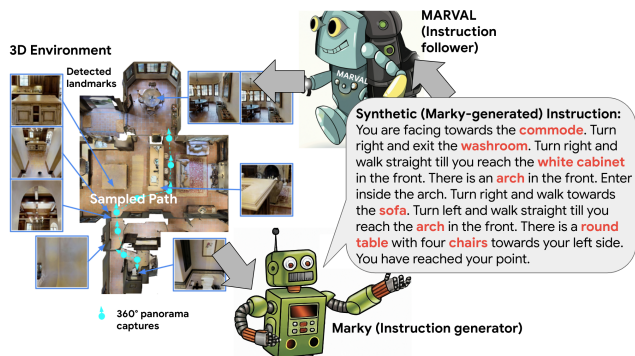


Figure 1. Simpler agents with more data: We investigate large-scale augmentation using 500+ environments annotated with synthetic instructions that approach human quality.

test bed for studying instruction-following, since the task can be simulated photo-realistically at scale and evaluation is straightforward. However, datasets that capture the linguistic diversity and idiosyncrasies of real human instructors are small and expensive to collect.

Shortages of human-annotated training data for other vision-and-language tasks have been partially addressed by pretraining transformers on up to billions of image-text pairs. This has underpinned dramatic improvements in image captioning [65, 70], visual question answering [59], phrase grounding [26, 35], text-to-video retrieval, video question answering [32] and text-to-image synthesis [49, 69]. However, these are all static image or video tasks, whereas VLN agents *interact* with 3D environments. In VLN, pretraining on large image-text and text-only datasets has been thoroughly explored [21, 22, 40, 45], but improvements are more limited. Arguably, progress in VLN has plateaued while still leaving a large gap between machine and human performance [73]. We hypothesize that static image-text and text-only datasets – despite their size – lack the spatially grounded and action-oriented language needed for effective VLN pretraining. Consider instructions from the Room-across-Room (RxR) dataset [30], which illustrate that wayfinding requires an understanding of allocen-

<sup>\*</sup>Equal Contribution. <sup>†</sup>Work done while at Google Research.

tric and egocentric spatial expressions (*near a grey console table behind you*), verbs (*climb the stairs*), imperatives and negations (*do not enter the room in front*) and temporal conditions (*walk until you see an entrance on your left*). Such expressions are rarely found in image-text datasets. Though similar expressions are found in text-only corpora, their meaning as it relates to the physical world is hard to infer from text alone (without sensorimotor context) [6].

To address this problem, we investigate large-scale augmentation with synthetic in-domain data, i.e., model-generated navigation instructions for trajectories in realistic 3D environments using previously developed components [27, 63]. We construct a large dataset using Marky [63], which generates VLN instructions that approach the quality of human instructors. [63] released the 1M Marky instruction-trajectory pairs situated in 61 Matterport3D [7] environments. To increase the diversity of the environments (and thus the scenes and objects available in them), we automatically annotate an additional 491 environments from the Gibson dataset [67]. Gibson environments have been underutilized in prior VLN work due to the lack of navigation graphs indicating navigable trajectories through its densely-sampled 360° panoramas. We train a model that classifies navigable directions for Matterport3D and use it to construct the missing navigation graphs. We sample 3.2M trajectories from these graphs and annotate them with Marky. To further increase the variability of trajectories, we synthesize image observations from novel viewpoints using an image-to-image GAN [27]. The resulting dataset is two orders of magnitude larger than existing human-annotated ones, and contains a wider variety of scenes and viewpoints. We have released our Gibson navigation graphs and the Marky-Gibson dataset.<sup>2</sup>

With orders of magnitude more training examples and environments, we explore VLN agent performance with imitation learning (IL), i.e., behavioral cloning and DAGGER [53] IL can take advantage of high-throughput transformer frameworks such as T5 [48] and thus efficiently train on 4.2M instructions (accumulating over 700M steps of experience). This is a departure from most prior VLN work in low-data settings, e.g. [10] report that pure IL underperforms by 8.5% success rate compared to agents trained with both IL and online reinforcement learning (RL) algorithms such as A3C [44]. However, IL outperforms RL in related tasks with sufficient training data [50]. Online RL also requires interacting with the environment at each step; this precludes efficient data prefetching and parallel preprocessing and thus imposes unavoidable overhead compared to IL. Empirically, we confirm that training existing models such as HAMT [10] on 4.2M instructions is infeasible without ground-up re-engineering, though we do find incorporating 10K additional synthetic instructions into

HAMT training modestly improves performance. Training with IL aligns with the trend towards large-scale multi-task vision-and-language models trained with supervised learning; these have unified tasks as diverse as visual question answering, image captioning, object detection, image classification, OCR and text reasoning [12] – and could include VLN in future.

Experimentally, in detailed ablation studies we show that adding Gibson environments, synthesizing additional image observations from novel viewpoints, increasing the capacity of the transformer, and finetuning with DAGGER all improve agent performance. On the challenging RxR dataset – which contains multilingual instructions with a median trajectory length of 15m – our best agent *using only imitation learning* outperforms all prior RL agents. Evaluating on novel instruction-trajectories in seen environments (Val-Seen), we improve over the state-of-the-art by 8%, reaching 79.1 NDTW. In new, unseen environments (Test), we improve by 2%, achieving 66.8 NDTW. We also show that that self-training with synthetic instructions in new environments (still without human annotations) improves performance by an additional 2% to 68.6 NDTW. Overall, our RxR results point to a new path to improving instruction-following agents, emphasizing large-scale training on near-human quality synthetic instructions. Perhaps surprisingly, on the English-only R2R dataset [3], our IL agent achieves strong but not state-of-the-art results. Marky was trained on RxR, so we attribute this to domain differences between R2R and RxR, underscoring the domain dependence of synthetic instructions.

## 2. Related Work

**Vision-and-Language Navigation** Agents that follow instructions by navigating to a prescribed location were initially studied in simple settings requiring limited or no perception, using instructions that were often procedurally generated [4, 5, 8, 42, 43]. More recent work has explored photorealistic 3D settings and natural language instructions [3, 9], using environments such as Matterport3D [7] and Streetview [41]. This instantiation of the problem, known as Vision-and-Language Navigation (VLN), raises the prospect of sim-to-real transfer to physical robots [2], and encouraged further datasets exploring dialog [17, 62], object search [46] and multilingual instructions [30].

**Pretraining and Transfer Learning** The use of realistic imagery and language in VLN, combined with the cost of collecting human instruction annotations, leads to a natural focus on pretraining and transfer learning to improve performance. Majumdar et al. [40] formulate VLN as an instruction-trajectory alignment problem, and initialize a transformer model using pretrained BERT weights [13] then perform additional pretraining on image-text pairs from Conceptual Captions [57]. Li et al. [36] also use a

<sup>2</sup>[//github.com/google-research-datasets/RxR/tree/main/marky-mT5](https://github.com/google-research-datasets/RxR/tree/main/marky-mT5)

BERT model, although more recent approaches have favored learned text encoders from scratch by pretraining with Masked Language Modeling (MLM) and related objectives on instruction-trajectory data [10, 18]. In terms of image representations, early work [15] used ResNet features [19] pretrained on ImageNet [54], although pretrained object detectors have also been explored [22, 33, 40] (typically a Faster-RCNN [51]). More recently, Chen et al. [10] use a vision transformer (ViT) [14] and current-state-of-the-art agents [34, 58] use CLIP [47], obtaining improvements over similarly sized encoders pretrained on ImageNet. However, although pretraining and transfer learning from large text and image-text datasets has been thoroughly explored, a significant gap to human performance remains.

**Data Augmentation** Fried et al. [15] were the first to demonstrate that performance following human instructions could be improved by augmenting training with synthetic (model-generated) instructions. A variety of other data augmentation approaches have been investigated, including modifying existing environments before generating new instructions [34, 60], training on a synthetic dataset of path-instruction pairs generated using online rental listings [16], and training with a generative model that infills and outpaints spatially perturbed panos of indoor environments to generate new observations [27, 28] (we also use this in Section 5). Notwithstanding these contributions, prior work incorporating synthetic instructions has been severely limited by instruction quality and scale. In human wayfinding evaluations, the instructions used [15, 60] were shown to be surprisingly weak, being poorly grounded and mostly unfollowable by people [72]. The recently proposed Marky model [63] (an instruction generator trained with text-aligned visual landmark correspondences) addresses this limitation, achieving near-human quality on R2R-style paths in unseen environments. We address the second limitation (scale) by developing an automated pipeline for scaling navigation graphs to 500+ new environments which we annotate with 3.2M instructions, and training agents on two orders of magnitude more data than before. Using this approach, we achieve state-of-the-art results in a VLN setting using a purely imitation learning agent. In contrast, most recent VLN work focuses on RL agents. An exception is DUET [11], which uses imitation learning in conjunction with a global action space based on a topological map.

### 3. Approach

**Problem set up** The agent is instantiated in an environment and must follow a natural language instruction  $\mathcal{W}$ . At time step  $t$ , the agent receives observation  $o_t$  and chooses action  $a_t$  that transitions it from state  $s_t$  to new state  $s_{t+1}$ . Following prior work, each observation is a photorealistic panoramic image (hereafter, pano) encoded as 36 image feature vectors  $o_t = \{I_{t,1}^o, I_{t,2}^o, \dots, I_{t,K}^o\}$ . These features

are extracted from perspective projections at 36 view angles (12 headings  $\times$  3 elevations at  $30^\circ$  intervals). The agent moves by choosing an action  $a_t$  from a set of candidates  $\mathcal{A}_t = \{I_{t,1}^a, I_{t,2}^a, \dots, I_{t,J}^a\}$  given by the environment. Action candidates are determined by the adjacent panos in a predefined navigation graph; each is represented by the image feature vector extracted from the perspective projection looking towards the adjacent pano. Selecting an action teleports the agent a few meters to the new pano. Alternatively, the agent can choose ‘STOP’ to end the episode. On average agents have 5 actions available at each step, including ‘STOP’. See [3] for more details.

**Agent architecture** Our imitation-learning agent is a transformer encoder which predicts the next action  $a_{t+1}$  by jointly combining all four input modalities: the instruction text  $\mathcal{W}$ , the history of observations  $o_{1:t-1}$  and actions  $a_{1:t-1}$ , the current observation  $o_t$ , and the action candidates  $\mathcal{A}_t$  (see Figure 2). At each step, all input features are concatenated into a single multimodal sequence with no attention masking, allowing every input to attend to every other input. For biasing interactions between different input modalities we include learned attention biases for each pair of input types, e.g. the instruction and the observation/action history. Like HAMT [10], our approach is not autoregressive: every forward pass predicts a single action using the full history. Given our emphasis on data augmentation, we name our agent MARVAL for Maximum Augmentation Regime for Vision And Language navigation. Our implementation is based on mT5 [68], a multilingual variant of the T5 transformer architecture [48].

**Image features** As noted above, pano observations  $o_t$  and action candidates  $\mathcal{A}_t$  are represented with sets of image features. We use precomputed, fixed 640-d features from MURAL-large [24], an EfficientNet-B7 [61] backbone trained on 1.8B multilingual image-text pairs and 6B translation pairs. MURAL’s image encoder’s representational power is similar to CLIP [47], which is used in previous work [10, 58] and is trained on 400M English image-text pairs with a ViT [14] backbone. To provide orientation information, each feature is combined with two learned embeddings: an *absolute direction embedding* capturing the feature’s orientation in the environment’s fixed coordinate system, and a *relative direction embedding* based on orientation relative to the agent’s heading. The agent’s initial heading at  $t=0$  is given by the dataset, and is typically random. We also augment the action candidates  $\mathcal{A}_t$  with a ‘STOP’ action. This is convenient for modeling action classification over the candidates (refer ‘Action classification’, below) and is represented by a zero image vector with unique direction embeddings. We use 37 absolute and relative direction embeddings, and snap features to the closest.

**Instruction encoding** The instruction  $\mathcal{W}$  is encoded as a sequence of WordPiece [56] tokens using the mT5 vocab-

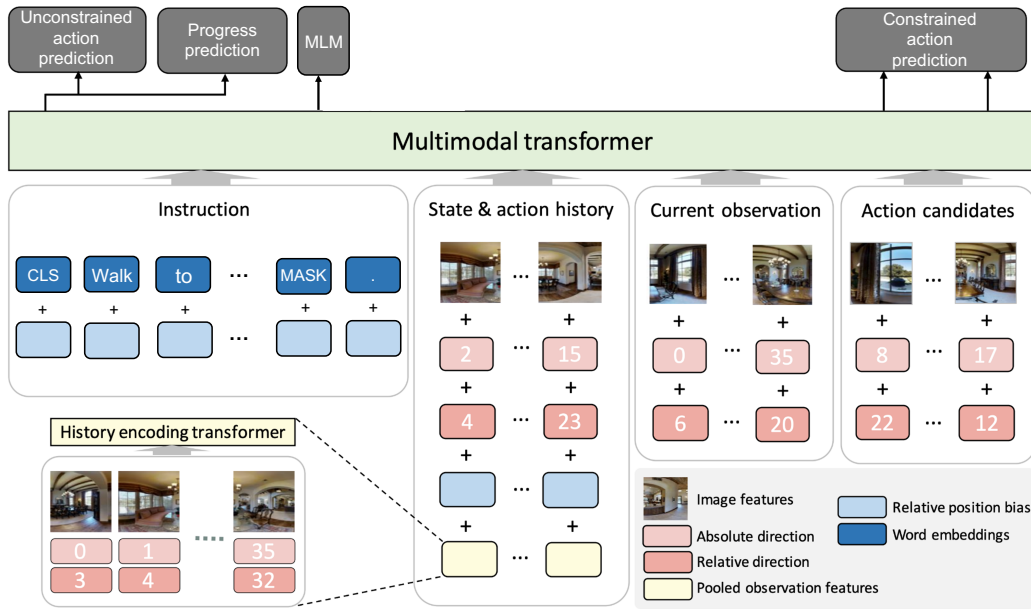


Figure 2. Agent architecture. At each time step, we combine the instruction, state and action history, current observation and action candidates into a multimodal transformer to predict the next action. Since observations consist of 36 image feature vectors (representing different views from a 360° camera) we compress each previous observation into a single vector, similar to [10].

ulary which supports up to 101 languages via a SentencePiece [31] model trained on mC4. Following T5, position information within the instruction is derived from relative position biases applied to the transformer’s attention logits.

**History encoding** The history of agent observations  $o_{1:t-1}$  and actions  $a_{1:t-1}$  is computationally expensive to process, since each pano observation  $o_t$  is comprised of 36 feature vectors. Similar to [10] we embed the 36 features from each previous pano observation into a single vector, based on the mean-pooled output of a separate transformer applied to the image features and their direction embeddings. This is added to the action candidate selected at each previous step. Position information for the state and action history is provided by relative position biases.

**Pretraining** We train the agent in two stages. We first pretrain on a large dataset of instruction-trajectory pairs, including both model-generated instructions and trajectories containing synthesized image observations from novel viewpoints (refer Section 4). We then finetune on a single dataset of human-annotated instruction-trajectory pairs to maximize performance on that dataset. Unlike [40] and [37], our transformer weights are initialized from scratch – we do not use any image-caption datasets or text corpora to train the transformer. Since the model is not autoregressive, each training trajectory is broken down into  $T$  training examples, where  $T$  is the number of time steps in the trajectory. Each training example requires the model to predict the next action for a single step in a trajectory, given the full instruction  $\mathcal{W}$ , the action history  $a_{1:t-1}$ , the observa-

tion history  $o_{1:t-1}$ , the current observation  $o_t$  and the set of action candidates  $\mathcal{A}_t$ . To increase the amount of supervision, during pretraining we combine four tasks:

- **Masked language modeling (MLM)** [13]: 15% of instruction tokens are masked, with all consecutive spans of masked tokens replaced by a single MASK token. Similar to [10], the model predicts the masked tokens using the surrounding text and visual clues from the observation/action history and the current observation.
- **Progress prediction**: A small MLP is added to the output representation of the CLS token (a special symbol capturing the fused representation of the entire sequence) to predict the proportion of the trajectory that is completed (based on 20 discretized classes). Progress monitoring has been shown to improve instruction grounding [39].
- **Constrained action prediction**: A classification task to predict the correct action from the constrained set of available action candidates  $\mathcal{A}_t$ . Since action candidates are inputs to the encoder (refer Figure 2), we compute the logit for each action as a learned projection of its output representation and normalize with softmax (a simplification of [10]).
- **Unconstrained action prediction**: A second small MLP is added to the CLS output to directly predict the next action from all 36 discretized agent-relative directions or ‘STOP’. Hence, these predictions are not constrained to  $\mathcal{A}_t$ , similar to the approach in [18]. The constrained and unconstrained action prediction tasks

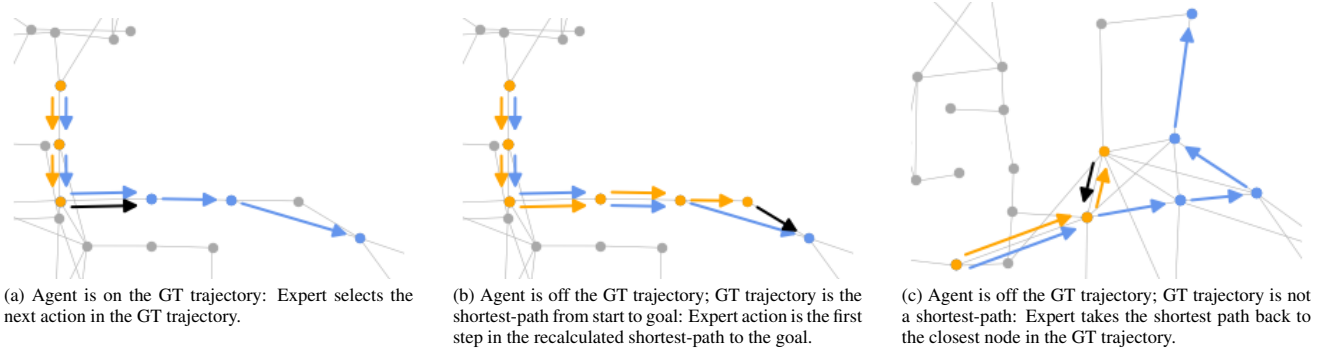


Figure 3. Calculation of the DAGGER expert action (black) given the ground-truth (GT) trajectory (blue) and an agent trajectory (yellow).

are highly related but complimentary; in early experiments we found that equally weighting the logits from both improves accuracy by 1-2%, so we adopt this approach in all experiments.

**Finetuning (behavioral cloning)** During finetuning we adapt our pretrained agent for best performance on a smaller human-annotated dataset. We update only the WordPiece embeddings in the agent and keep all other transformer weights frozen, as this makes finetuning more stable and less prone to overfitting (especially on the smaller R2R dataset). We consider two finetuning strategies. The first is *behavioral cloning*. In this setting, we simply drop the instruction text masking and the MLM objective, retaining the progress prediction and constrained and unconstrained action prediction losses used in pretraining. We then finetune the agent to predict the next action at each step along ground-truth trajectories, treating imitation learning as supervised learning.

**DAGGER training** The main weakness of behavioral cloning is that the state distribution seen in training differs from the state distribution induced by the agent during inference [52]. Previous works [10, 60] report substantial improvements by combining behavioral cloning with on-line reinforcement learning algorithms such as A3C [44]. We use *DAGGER* [11, 52] to help train the agent to better recover from errors, since it is simple to implement and requires no environment interaction during training. In DAGGER, during each iteration of finetuning the dataset is augmented with trajectories of states visited by the current agent policy and actions given by an expert. Figure 3 explains the calculation of expert actions. We find that most of the gains are captured in a single DAGGER iteration.

**Pre-Exploration** While most of the focus in VLN is on instruction-following in new, unseen environments, in reality environments persist over time providing opportunities for pre-exploration. Similar to [60, 64] we consider a *pre-exploration* setting in which the agent may explore unseen environments with self-supervision before evaluation. Our synthetic-instruction approach is readily applicable to this scenario; we simply sample paths from the Val-Unseen or

Test environments, annotate them with Marky instructions, and include them in the training data.

## 4. Datasets and Augmentation

The datasets used for training and evaluation are described below and summarized in Table 1.

- **Room-to-Room (R2R)** [3] consists of 22K human-annotated English language navigation instructions, each describing a trajectory that traverses multiple rooms in Matterport3D [7]. This was the first dataset to use a photo-realistic environment for the instruction guided navigation task. R2R trajectories average 10m, and the trajectories are always the shortest path between the start point and the goal.
- **Room-across-Room (RxR)** [30] is a larger human-annotated dataset containing 126K instructions in English, Hindi and Telugu. To mitigate goal seeking behaviour and to ensure that agents are faithful to the instruction, RxR includes Matterport3D trajectories that are diverse in terms of length (average is 15m) and the landmarks that are referred to, and it also includes trajectories that do not go directly to the goal.
- **Speaker-Matterport (S-MP)** [15] is a set of 178K sampled trajectories in Matterport3D environments, annotated with synthetic instructions generated with an LSTM [20] Speaker model trained on R2R.
- **Marky-Matterport (M-MP)** Marky [63] is a landmark-aware multilingual instruction generator trained on RxR, used to generate 1M instructions in English, Hindi and Telugu for 330K sampled Matterport3D trajectories. In human wayfinding evaluations in unseen environments Marky achieves close to human performance on shortest-path trajectories (e.g., R2R’s paths). On the more challenging RxR paths a gap remains: human wayfinders obtain a 62% success rate with Marky vs. 78% with human instructions.
- **Marky-Gibson (M-Gib)** The Gibson [67] dataset consists of 572 indoor 3D environments. Despite its large size compared to Matterport3D, prior work has under-

Dataset	Instruction Count	Avg steps	Avg words	Environment	Model-generated	Language
Room-to-Room (R2R) [3]	14K	5.0	26	Matterport	✗	en
Room-across-Room (RxR) [30]	79K	8.1	78	Matterport	✗	en/hi/te
Speaker-Matterport [15]	178K	5.1	21	Matterport	✓	en
Marky-Matterport [63]	1.0M	9.5	87	Matterport	✓	en/hi/te
Marky-Gibson (ours)	3.2M	7.1	71	Gibson	✓	en/hi/te

Table 1. Training data. Existing datasets are situated in the 61 train environments from Matterport3D [7]. We introduce the multilingual Marky-Gibson dataset containing 3.2M model-generated navigation instructions in 491 Gibson [67] environments.

utilized Gibson data for training VLN agents. This is primarily due to a lack of navigation trajectories and instruction annotations. To alleviate this, and unlock the Gibson dataset for VLN training, we propose an automated process to label these environments with high quality navigation graphs (see below). We then sample 1.3M trajectories and annotate them with 3.2M Marky instructions in English, Hindi and Telugu.

**Gibson navigation graphs** In the standard VLN setting, agents are trained and evaluated using panos as observations (refer Section 3). Movement in the environment requires a graph with panos as nodes and edges indicating navigability. Navigation graphs for Matterport3D were generated by [3], using a semi-automated process combined with human visual inspection. However, there are no navigation graphs for Gibson environments and the size of the dataset precludes human inspection. We therefore train a model on panos and navigation graphs from the Matterport3D train split to classify whether a patch of pixels in a pano constitutes a navigable direction. The model is based on RedNet [25], an RGB-D encoder-decoder first proposed for image segmentation, using a ResNet-50 [19] backbone. The output space is discretized into  $8 \times 16 \times 5$  pitch, heading and distance buckets. During training each bucket is assigned a positive value if the corresponding location corresponds to a navigable node, and 0 otherwise.

To compute Gibson navigation graphs, we combine model edge predictions with obstacle information from the dataset’s 3D meshes. We add an edge between pano nodes  $i$  and  $j$  if the following boolean expression evaluates to true:

$$e(i, j) = (\lambda_d \frac{g_{i,j}}{s_{i,j}} - \lambda_p p_{i,j} \leq 1) \wedge (s_{i,j} \leq 3.5) \wedge (|z_i - z_j| \leq 3)$$

where  $g_{i,j}$  is the geodesic distance (accounting for obstacles) between nodes  $i$  and  $j$  calculated using the Habitat Simulator [55],  $s_{i,j}$  is the straight-line Euclidean distance between nodes  $i$  and  $j$ ,  $p_{i,j}$  is the model probability of an edge connecting nodes  $i$  and  $j$ ,  $z_i$  is the vertical coordinate of pano  $i$ , and  $\lambda_d$  and  $\lambda_p$  are weighting parameters. The first term captures model predictions and encourages edges between panos that have few intervening obstacles. The second term ensures that nodes are within 3.5m, and the third

term ensures that nodes are within 3m in the vertical axis (these values are chosen based on [3]). Finally, to ensure that the navigation graph for each environment is fully connected, we compute the minimum spanning tree (MST) [29] of the graph with the edge weights given by the first term in the equation for  $e(i, j)$ , and apply a logical ‘OR’ operation over  $e(i, j)$  and the MST.

To set the weighting parameters  $\lambda_d$  and  $\lambda_p$ , we perform grid search to maximize the  $F_1$  score when predicting edges in Matterport3D val environments.  $F_1$  uses the standard calculation where the population includes all pairs of nodes in each environment, and the true condition is positive if an edge exists in the hand-crafted navigation graphs from [3]. Our approach achieves an  $F_1$  score of 0.70, precision of 0.695, and recall of 0.713. The average edge length in the generated Gibson graphs is 3.02m (median of 2.06m), and the average node degree is 4.15 (median of 4).

**Trajectory sampling and instruction generation** Using the generated navigation graphs, we sample trajectories from 491 Gibson train and val environments (we do not use test environments). Unlike Matterport3D, Gibson lacks room annotations, which precludes us from using the two-step sampling approach from RxR. Instead, we use a simpler approach: we randomly sample 3 panos, and use a TSP solver to find the shortest path that visits all 3 panos. Trajectories longer than 40m or 16 steps are discarded, and no more than 3K paths are sampled per environment. This procedure generates 1.06M paths, with an average of 7.1 steps and length of 19.3m. Using Marky we annotate each trajectory with English, Hindi and Telugu instructions to create the Marky-Gibson dataset.

**Synthesizing image observations with SE3DS** One weakness of training VLN agents on pano images is that training trajectories are constrained to the locations of the captured images. VLN agents tend to overfit to these trajectories [71], contributing to a performance drop in unseen environments. [27, 28] showed that a strong generative model is capable of successfully rendering high resolution panos from novel viewpoints, and that training VLN agents with spatially-perturbed panos could improve the success rate of the agent on R2R Val-Unseen by 1.5%. To assess if this approach is complimentary to instruction augmentation, we

	Pretraining Data						Iterations	RxR VAL-UNSEEN				R2R VAL-UNSEEN			
	Size	R2R	RxR	S-MP	M-MP	M-Gib		SE3DS	NE	SR	NDTW	SDTW	NE	SR	SPL
1	Base	✓	✓					630K	11.16	23.2	40.1	19.0	7.58	33.5	31.9
2	Base	✓	✓	✓				1.68M	11.17	25.0	40.3	21.0	6.50	41.0	39.0
3	Base	✓	✓	✓	✓			2.94M	7.53	45.0	56.9	39.2	7.07	40.3	37.2
4	Base	✓	✓		✓			2.00M	6.90	50.2	60.3	43.9	6.07	50.0	46.5
5	Base	✓	✓	✓	✓		✓	1.94M	6.68	50.9	61.6	45.0	5.34	52.2	49.0
6	Base	✓	✓	✓	✓	✓		3.00M	5.90	55.6	64.4	49.0	5.24	52.5	47.9
7	Base	✓	✓	✓	✓	✓	✓	2.80M	5.75	55.9	65.1	49.3	5.08	54.2	50.1
8	Large	✓	✓	✓	✓	✓	✓	4.80M	5.72	58.3	66.1	51.7	4.87	55.6	52.8
9	Large	✓	✓		✓	✓	✓	5.14M	<b>5.56</b>	<b>59.4</b>	<b>67.0</b>	<b>52.7</b>	<b>4.84</b>	<b>57.4</b>	<b>54.6</b>

Table 2. Comparison of pretraining settings. Best results (without finetuning) are obtained by combining the **R2R** and **RxR** datasets *with* Marky-generated instructions in both Matterport3D (**M-MP**) and Gibson (**M-Gib**) environments, *without* Speaker-generated instructions (**S-MP**), and *with* synthesis of observations from novel viewpoints using **SE3DS** (row 9).

use the proposed SE3DS (Simple and Effective 3D Synthesis) model to augment panoramas from the Matterport environments. Following [27], we create 200 variations of each environment which are randomly sampled during training. In each environment variation, with 50% probability a pano will be spatially-perturbed by up to 1.5m and re-rendered at the new location using SE3DS.

## 5. Experiments

**Pretraining settings** In Table 2 we explore pretraining using varying amounts and types of augmented data. During pretraining, we monitor one-step action prediction accuracy on ground-truth trajectories using held-out instructions from RxR and R2R Val-Unseen. Each setting is trained until convergence, requiring more iterations (**Its**) for larger models and datasets. We select the best checkpoint based on one-step prediction accuracy then perform a full evaluation using standard VLN path-fidelity metrics [1, 23]: Navigation Error (**NE** ↓, the average distance in meters between the agent’s final position and the goal), Success Rate (**SR** ↑, the proportion of trajectories with  $NE < 3m$ ), Success rate weighted by normalized inverse Path Length (**SPL** ↑), normalized dynamic time warping (**NDTW** ↑), and success weighted DTW (**SDTW** ↑).

**Speaker vs. Marky** Consistent with previous work, we find that data augmentation with synthetic instructions from the [15] Speaker model improves performance (row 2 vs. 1, +2% SR on RxR and +8% on R2R), but instructions from Marky [63] are far more effective (row 4 vs. 1, +27% SR on RxR and +17% on R2R). This is consistent with human evaluations of instruction quality, confirming that improvements in instruction-generation flow through to instruction-following performance. Interestingly, we find that combining the Speaker model with Marky leads to worse performance on both RxR and R2R (row 3 vs. 4, and also row 8 vs. 9), which we attribute to the introduction of noise from the lower-quality Speaker instructions.

**Gibson, SE3DS and model size** Augmentation with Marky instructions in Gibson environments (row 6 v. 3) provides a substantial boost (+11% SR on RxR and +12% on R2R), suggesting that the returns from scaling synthetic instructions to more environments are not exhausted. Using SE3DS to synthesize image observations from novel viewpoints improves +6% SR on RxR and +12% on R2R (row 5 vs. 3), but this benefit is substantially reduced (+0% SR on RxR and +2% on R2R, row 7 vs. 6) if Gibson is included, presumably because new environments also increase viewpoint variety. Most experiments use the mT5-base [68] model; switching to mT5-large provides a further performance boost (+2% SR on RxR and +1% on R2R, row 8 vs. 7). Our best pretraining results on both RxR and R2R are achieved using an mT5-large model with all the previously mentioned data, but leaving out the Speaker instructions (row 9). We use this checkpoint in all finetuning experiments. This agent pretrains for 5.14M iterations, which, using a batch size of 128, represents over 650M steps of experience (over 700M including finetuning).

**Finetuning** In Tables 3 and 4 we compare results for our MARVAL agent after finetuning to previous work on the RxR and R2R datasets. On both datasets, finetuning with behavioral cloning on just human-annotated data (Finetuned-BC) substantially improves the pretrained model. The improvement from using DAGGER over behavioral cloning is small but consistent. On the RxR dataset, MARVAL outperforms all prior work. Evaluating on novel instruction-trajectories in seen environments (Val-Seen), we improve over the state-of-the-art by 8%, reaching 79.1 NDTW. In new, unseen environments (Test), we improve by 2%, achieving 66.8 NDTW. Self-training with Marky synthetic instructions in the Test environments (a form of privileged access, but still without human annotations) improves performance by an additional 2% to 68.6 NDTW.

**RxR vs. R2R** On the English-only R2R dataset (Table 4), MARVAL achieves strong performance but not state-of-

Agent	VAL-SEEN				VAL-UNSEEN				TEST (UNSEEN)			
	NE	SR	NDTW	SDTW	NE	SR	NDTW	SDTW	NE	SR	NDTW	SDTW
LSTM [30]	10.7	25.2	42.2	20.7	10.9	22.8	38.9	18.2	12.0	21.0	36.8	16.9
EnvDrop+ [58]	-	-	-	-	-	42.6	55.7	-	-	38.3	51.1	32.4
CLEAR-C [33]	-	-	-	-	-	-	-	-	-	40.3	53.7	34.9
HAMT [10]	-	59.4	65.3	50.9	-	56.5	63.1	48.3	6.2	53.1	59.9	45.2
EnvEdit* [34]	-	67.2	71.1	58.5	-	62.8	68.5	54.6	<b>5.1</b>	60.4	64.6	51.8
MARVAL (Pretrained)	3.62	72.7	77.0	65.9	5.56	59.4	67.0	52.7	-	-	-	-
MARVAL (Finetuned-BC)	3.25	75.4	79.0	68.7	4.80	63.7	70.6	56.9	-	-	-	-
MARVAL (DAGGER)	<b>3.01</b>	<b>75.9</b>	<b>79.1</b>	<b>68.8</b>	<b>4.49</b>	<b>64.8</b>	<b>70.8</b>	<b>57.5</b>	5.5	<b>60.7</b>	<b>66.8</b>	<b>53.5</b>
MARVAL (Pre-Explore)†	3.33	73.7	77.6	66.6	4.19	66.5	72.2	59.1	5.2	61.8	68.6	54.8
Human [30]	-	-	-	-	-	-	-	-	0.9	93.9	79.5	76.9

\*Results from an ensemble of three agents.

Table 3. Results on RxR. Our MARVAL agent trained with imitation learning – behavioral cloning (BC) or DAGGER – outperforms all existing RL agents. Pre-Exploration in the eval environments († a form of privileged access, but still without human annotations) can provide a further boost.

Agent	VAL-SEEN				VAL-UNSEEN				TEST (UNSEEN)			
	TL	NE	SR	SPL	TL	NE	SR	SPL	TL	NE	SR	SPL
PREVALENT [18]	10.32	3.67	69	65	10.19	4.71	58	53	10.51	5.30	54	51
RecBERT [22]	11.13	2.90	72	68	12.01	3.93	63	57	12.35	4.09	63	57
EnvDrop+ [58]	-	-	-	-	-	-	59.2	52.9	-	-	59	53
AirBERT [16]	11.09	2.68	75	70	11.78	4.01	62	56	12.41	4.13	62	57
HAMT [10]	11.15	2.51	76	72	11.46	2.29	66	61	12.27	3.93	65	60
REM [38]	10.88	2.48	75.4	71.8	12.44	3.89	63.6	57.9	13.11	3.87	65	59
EnvEdit* [34]	11.18	<b>2.32</b>	<b>76.9</b>	<b>73.9</b>	11.13	<b>3.24</b>	<b>68.9</b>	<b>64.4</b>	11.90	<b>3.59</b>	<b>68</b>	<b>64</b>
MARVAL (Pretrained)	10.32	3.73	68.2	64.9	9.71	4.84	57.4	54.6	-	-	-	-
MARVAL (Finetuned-BC)	10.43	3.11	72.3	68.9	9.72	4.20	63.0	60.0	-	-	-	-
MARVAL (DAGGER)	10.60	2.99	73.0	69.1	10.15	4.06	64.8	60.7	10.22	4.18	62	58
Human [3]	-	-	-	-	-	-	-	-	11.90	1.61	86	76

\*Results from an ensemble of three agents.

Table 4. Results on R2R. MARVAL achieves strong performance but not state-of-the-art, which we attribute to domain differences between the R2R and RxR (which was used to train Marky).

the-art. Surprisingly, the Val-Unseen success rate (SR) of 64.8% is the same for both RxR and R2R, whereas typically RxR performance is lower since the trajectories are longer and more varied. Noting that Marky was trained on RxR, we attribute lower relative performance on R2R to domain differences between R2R and RxR. While the average length of instructions in R2R is 26 words, RxR has an average of 87 words — 3 times more. This is partly because RxR instructions are more verbose, often describing objects in more detail and including state verification. Further, cultural differences arising from the data collection process (annotators from USA or from India) may also contribute to the distribution shift due to subtle differences in the vocabulary and structure of language used to form the instructions. We note however, that while our augmentation approach focuses on scaling up in terms of high quality instructions, EnvEdit [34] focuses on generalization through

augmentation of visual features. These two approaches may ultimately prove to be complementary.

## 6. Conclusion

We build a purely imitation learning agent that achieves state-of-the-art results on the RxR benchmark. This result paves a new path towards improving instruction-following agents, emphasizing large-scale imitation learning with generic architectures, along with a focus on developing synthetic instruction generation capabilities – which are shown to directly improve instruction-following performance. We find that aligning synthetic instructions to the target domain is essential, as seen through the gap in performance on R2R. On RxR, the performance improvement over the state-of-the-art is much larger in seen environments (+8%) than unseen test environments (+2%). Scaling to even more indoor environments might improve generalization further.



## References

- [1] Peter Anderson, Angel Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, and Amir R. Zamir. On Evaluation of Embodied Navigation Agents. *ArXiv preprint*, 2018. 7
- [2] Peter Anderson, Ayush Shrivastava, Joanne Truong, Arjun Majumdar, Devi Parikh, Dhruv Batra, and Stefan Lee. Sim-to-real transfer for vision-and-language navigation. In *CoRL*, 2021. 2
- [3] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian D. Reid, Stephen Gould, and Anton van den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *CVPR*, 2018. 1, 2, 3, 5, 6, 8
- [4] Jacob Andreas and Dan Klein. Alignment-based compositional semantics for instruction following. In *EMNLP*, 2015. 2
- [5] Yoav Artzi and Luke Zettlemoyer. Weakly supervised learning of semantic parsers for mapping instructions to actions. *TACL*, 2013. 2
- [6] Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, Nicolas Pinto, and Joseph Turian. Experience grounds language. In *EMNLP*, 2020. 2
- [7] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3D: Learning from RGB-D data in indoor environments. *3DV*, 2017. 2, 5, 6
- [8] David L. Chen and Raymond J. Mooney. Learning to interpret natural language navigation instructions from observations. In *AAAI*, 2011. 2
- [9] Howard Chen, Alane Suhr, Dipendra Misra, Noah Snaveley, and Yoav Artzi. TOUCHDOWN: natural language navigation and spatial reasoning in visual street environments. In *ICCV*, 2019. 1, 2
- [10] Shizhe Chen, Pierre-Louis Guhur, Cordelia Schmid, and Ivan Laptev. History aware multimodal transformer for vision-and-language navigation. In *NeurIPS*, 2021. 2, 3, 4, 5, 8
- [11] Shizhe Chen, Pierre-Louis Guhur, Makarand Tapaswi, Cordelia Schmid, and Ivan Laptev. Think global, act local: Dual-scale graph transformer for vision-and-language navigation. In *CVPR*, 2022. 3, 5
- [12] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, Alexander Kolesnikov, Joan Puigcerver, Nan Ding, Keran Rong, Hassan Akbari, Gaurav Mishra, Linting Xue, Ashish Thapliyal, James Bradbury, Weicheng Kuo, Mojtaba Seyedhosseini, Chao Jia, Burcu Karagol Ayan, Carlos Riquelme, Andreas Steiner, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. PaLI: A jointly-scaled multilingual language-image model. *ArXiv preprint*, 2022. 2
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019. 2, 4
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 3
- [15] Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe Morency, Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. Speaker-follower models for vision-and-language navigation. In *NeurIPS*, 2018. 3, 5, 6, 7
- [16] Pierre-Louis Guhur, Makarand Tapaswi, Shizhe Chen, Ivan Laptev, and Cordelia Schmid. Airbert: In-domain pretraining for vision-and-language navigation. *ICCV*, 2021. 3, 8
- [17] Meera Hahn, Jacob Krantz, Dhruv Batra, Devi Parikh, James Rehg, Stefan Lee, and Peter Anderson. Where are you? localization from embodied dialog. In *EMNLP*, 2020. 2
- [18] Weituo Hao, Chunyuan Li, Xiujun Li, Lawrence Carin, and Jianfeng Gao. Towards learning a generic agent for vision-and-language navigation via pre-training. In *CVPR*, 2020. 3, 4, 8
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3, 6
- [20] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 1997. 5
- [21] Yicong Hong, Cristian Rodriguez Opazo, Yuankai Qi, Qi Wu, and Stephen Gould. Language and visual entity relationship graph for agent navigation. In *NeurIPS*, 2020. 1
- [22] Yicong Hong, Qi Wu, Yuankai Qi, Cristian Rodriguez-Opazo, and Stephen Gould. VLN-BERT: A recurrent vision-and-language bert for navigation. *CVPR*, 2021. 1, 3, 8
- [23] Gabriel Ilharco, Vihan Jain, Alexander Ku, Eugene Ie, and Jason Baldridge. Effective and general evaluation for instruction conditioned navigation using dynamic time warping. *NeurIPS Visually Grounded Interaction and Language Workshop*, 2019. 7
- [24] Aashi Jain, Mandy Guo, Krishna Srinivasan, Ting Chen, Sneha Kudugunta, Chao Jia, Yinfei Yang, and Jason Baldridge. MURAL: Multimodal, multitask representations across languages. In *EMNLP Findings*, 2021. 3
- [25] Jindong Jiang, Lunan Zheng, Fei Luo, and Zhijun Zhang. Rednet: Residual encoder-decoder network for indoor rgb-d semantic segmentation. *ArXiv preprint*, 2018. 6
- [26] Aishwarya Kamath, Mannat Singh, Yann LeCun, Ishan Misra, Gabriel Synnaeve, and Nicolas Carion. Mdetr - modulated detection for end-to-end multi-modal understanding. *ICCV*, 2021. 1
- [27] Jing Yu Koh, Harsh Agrawal, Dhruv Batra, Richard Tucker, Austin Waters, Honglak Lee, Yinfei Yang, Jason Baldridge, and Peter Anderson. Simple and effective synthesis of indoor 3d scenes. In *AAAI*, 2023. 1, 2, 3, 6, 7
- [28] Jing Yu Koh, Honglak Lee, Yinfei Yang, Jason Baldridge, and Peter Anderson. Pathdreamer: A world model for indoor navigation. In *ICCV*, 2021. 3, 6

- [29] Joseph B Kruskal. On the shortest spanning subtree of a graph and the traveling salesman problem. *Proceedings of the American Mathematical Society*, 1956. 6
- [30] Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and Jason Baldridge. Room-Across-Room: Multilingual vision-and-language navigation with dense spatiotemporal grounding. In *EMNLP*, 2020. 1, 2, 5, 6, 8
- [31] Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *EMNLP*, 2018. 4
- [32] Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L Berg, Mohit Bansal, and Jingjing Liu. Less is more: Clipbert for video-and-language learning via sparse sampling. In *CVPR*, 2021. 1
- [33] Jialu Li, Hao Tan, and Mohit Bansal. CLEAR: Improving vision-language navigation with cross-lingual, environment-agnostic representations. In *Findings of the Association for Computational Linguistics: NAACL*, 2022. 3, 8
- [34] Jialu Li, Hao Tan, and Mohit Bansal. EnvEdit: Environment editing for vision-and-language navigation. In *CVPR*, 2022. 3, 8, 13
- [35] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, Kai-Wei Chang, and Jianfeng Gao. Grounded language-image pre-training. In *CVPR*, 2022. 1
- [36] Xiujun Li, Chunyuan Li, Qiaolin Xia, Yonatan Bisk, Asli Celikyilmaz, Jianfeng Gao, Noah A. Smith, and Yejin Choi. Robust navigation with language pretraining and stochastic sampling. In *EMNLP*, 2019. 2
- [37] Xiujun Li, Chunyuan Li, Qiaolin Xia, Yonatan Bisk, Asli Celikyilmaz, Jianfeng Gao, Noah A. Smith, and Yejin Choi. Robust navigation with language pretraining and stochastic sampling. In *EMNLP*, 2019. 4
- [38] Chong Liu, Fengda Zhu, Xiaojun Chang, Xiaodan Liang, and Yi-Dong Shen. Vision-language navigation with random environmental mixup. *ICCV*, 2021. 8
- [39] Chih-Yao Ma, Jiasen Lu, Zuxuan Wu, Ghassan AlRegib, Zolt Kira, Richard Socher, and Caiming Xiong. Self-monitoring navigation agent via auxiliary progress estimation. In *ICLR*, 2019. 4
- [40] Arjun Majumdar, Ayush Shrivastava, Stefan Lee, Peter Anderson, Devi Parikh, and Dhruv Batra. Improving vision-and-language navigation with image-text pairs from the web. In *ECCV*, 2020. 1, 2, 3, 4
- [41] Harsh Mehta, Yoav Artzi, Jason Baldridge, Eugene Ie, and Piotr Mirowski. Retouchdown: Adding touchdown to streetlearn as a shareable resource for language grounding tasks in street view. *ArXiv preprint*, 2020. 2
- [42] Hongyuan Mei, Mohit Bansal, and Matthew R. Walter. Listen, attend, and walk: Neural mapping of navigational instructions to action sequences. In *AAAI*, 2016. 2
- [43] Dipendra Misra, John Langford, and Yoav Artzi. Mapping instructions and visual observations to actions with reinforcement learning. In *EMNLP*, 2017. 2
- [44] Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *ICML*, 2016. 2, 5
- [45] Abhinav Moudgil, Arjun Majumdar, Harsh Agrawal, Stefan Lee, and Dhruv Batra. Soat: A scene-and object-aware transformer for vision-and-language navigation. *NeurIPS*, 2021. 1
- [46] Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang, William Yang Wang, Chunhua Shen, and Anton van den Hengel. REVERIE: remote embodied visual referring expression in real indoor environments. In *CVPR*, 2020. 2
- [47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 3
- [48] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*, 2020. 2, 3
- [49] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, 2021. 1
- [50] Ram Ramrakhya, Eric Undersander, Dhruv Batra, and Abhishek Das. Habitat-web: Learning embodied object-search strategies from human demonstrations at scale. In *CVPR*, 2022. 2
- [51] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 3
- [52] Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. In *AISTATS*, 2010. 5
- [53] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *AISTATS*, 2011. 2
- [54] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 3
- [55] Manolis Savva, Jitendra Malik, Devi Parikh, Dhruv Batra, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, and Vladlen Koltun. Habitat: A platform for embodied AI research. In *ICCV*, 2019. 6
- [56] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *ACL*, 2016. 3
- [57] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypemymed, image alt-text dataset for automatic image captioning. In *ACL*, 2018. 2
- [58] Sheng Shen, Liunian Harold Li, Hao Tan, Mohit Bansal, Anna Rohrbach, Kai-Wei Chang, Zhewei Yao, and Kurt Keutzer. How much can CLIP benefit vision-and-language tasks? *ArXiv*, 2021. 3, 8

- [59] Hao Tan and Mohit Bansal. LXMERT: Learning cross-modality encoder representations from transformers. In *EMNLP*, 2019. 1
- [60] Hao Tan, Licheng Yu, and Mohit Bansal. Learning to navigate unseen environments: Back translation with environmental dropout. In *NAACL-HLT*, 2019. 3, 5
- [61] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, Proceedings of Machine Learning Research, 2019. 3
- [62] Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer. Vision-and-dialog navigation. In *CoRL*, 2019. 2
- [63] Su Wang, Ceslee Montgomery, Jordi Orbay, Vighnesh Birodkar, Aleksandra Faust, Izzeddin Gur, Natasha Jaques, Austin Waters, Jason Baldrige, and Peter Anderson. Less is more: Generating grounded navigation instructions from landmarks. In *CVPR*, 2022. 1, 2, 3, 5, 6, 7
- [64] Xin Wang, Qiuyuan Huang, Asli Çelikyilmaz, Jianfeng Gao, Dinghan Shen, Yuan-Fang Wang, William Yang Wang, and Lei Zhang. Reinforced cross-modal matching and self-supervised imitation learning for vision-language navigation. In *CVPR*, 2019. 5
- [65] Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. SimVLM: Simple visual language model pretraining with weak supervision. *ArXiv*, 2021. 1
- [66] Terry Winograd. Procedures as a representation for data in a computer program for understanding natural language. Technical report, Massachusetts Institute of Technology, 1971. 1
- [67] Fei Xia, Amir Roshan Zamir, Zhi-Yang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gibson env: Real-world perception for embodied agents. In *CVPR*, 2018. 2, 5, 6
- [68] Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mT5: A massively multilingual pre-trained text-to-text transformer. In *NAACL*, 2021. 3, 7
- [69] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldrige, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation. *ArXiv preprint*, 2022. 1
- [70] Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. VinVL: Revisiting visual representations in vision-language models. *CVPR*, 2021. 1
- [71] Yubo Zhang, Hao Tan, and Mohit Bansal. Diagnosing the environment bias in vision-and-language navigation. In *IJCAI*, 2020. 6
- [72] Ming Zhao, Peter Anderson, Vihan Jain, Su Wang, Alexander Ku, Jason Baldrige, and Eugene Ie. On the evaluation of vision-and-language navigation instructions. In *EACL*, 2021. 3
- [73] Wanrong Zhu, Yuankai Qi, Pradyumna Narayana, Kazuo Sone, Sugato Basu, Xin Eric Wang, Qi Wu, Miguel Eckstein, and William Yang Wang. Diagnosing Vision-and-Language Navigation: What Really Matters. In *NAACL-HLT*, 2022. 1