# Room-Across-Room: Multilingual Vision-and-Language Navigation with Dense Spatiotemporal Grounding

**Alexander Ku**[1][*]  **Peter Anderson**[1][*]  **Roma Patel**[2]  **Eugene Ie**[1]  **Jason Baldridge**[1]

[1]Google Research    [2]Brown University

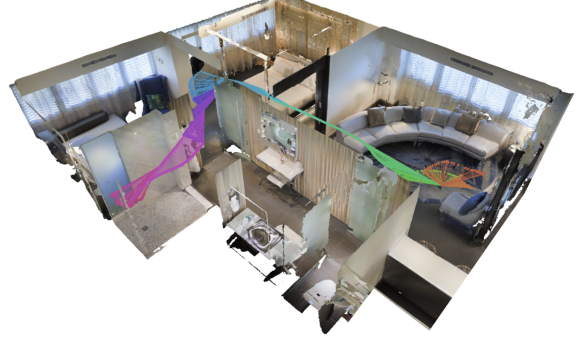{alexku, pjand, eugeneie, jridge}@google.com romapatel@brown.edu

## Abstract

We introduce Room-Across-Room (RxR), a new Vision-and-Language Navigation (VLN) dataset. RxR is multilingual (English, Hindi, and Telugu) and larger (more paths and instructions) than other VLN datasets. It emphasizes the role of language in VLN by addressing known biases in paths and eliciting more references to visible entities. Furthermore, each word in an instruction is time-aligned to the virtual poses of instruction creators and validators. We establish baseline scores for monolingual and multilingual settings and multitask learning when including Room-to-Room annotations (Anderson et al., 2018b). We also provide results for a model that learns from synchronized pose traces by focusing only on portions of the panorama attended to in human demonstrations. The size, scope and detail of RxR dramatically expands the frontier for research on embodied language agents in simulated, photo-realistic environments.

Our starting point is in a living room, we're facing towards a long beige sofa, and in front of the sofa there are three glass coffee tables, turn around and exit through the doorway that's in front of you, walk pass the bed that's on your right and then turn left, we're now facing towards another living room, and on the left there's an open door, walk towards that open door enter the bathroom that's in front of you, turn towards the right into the shower area, and that's your destination.

Figure 1: RxR's instructions are densely grounded to the visual scene by aligning the annotator's virtual pose to their spoken instructions for navigating a path.

## 1  Introduction

Vision-and-Language Navigation (VLN) tasks require computational agents to mediate the relationship between language, visual scenes and movement. Datasets have been collected for both indoor (Anderson et al., 2018b; Thomason et al., 2019b; Qi et al., 2020) and outdoor (Chen et al., 2019; Mehta et al., 2020) environments; success in these is based on clearly-defined, objective task completion rather than language or vision specific annotations. These VLN tasks fall in the Goldilocks zone: they can be tackled – but not solved – with current methods, and progress on them makes headway on real world grounded language understanding.

We introduce Room-across-Room (RxR), a VLN dataset that addresses gaps in existing ones by (1) including more paths that (2) counter known biases in existing datasets, and (3) collecting an order of magnitude more instructions for (4) three languages (English, Hindi and Telugu) while (5) capturing annotators' 3D pose sequences. As such, RxR includes dense spatiotemporal grounding for every instruction, as illustrated in Figure 1.

We provide monolingual and multilingual baseline experiments using a variant of the Reinforced Cross-Modal Matching agent (Wang et al., 2019). Performance generally improves by using monolingual learning, and by using RxR's follower paths as well as its guide paths. We also concatenate R2R and RxR annotations as a simple multitask strategy (Wang et al., 2020): the agent trained on both datasets obtains across the board improvements.

RxR contains 126K instructions covering 16.5K sampled guide paths and 126K human follower

---
[*]First two authors contributed equally.

demonstration paths. The dataset is available.[1] We plan to release a test evaluation server, our annotation tool, and code for all experiments.

## 2 Motivation

A number of VLN datasets situated in photo-realistic 3D reconstructions of real locations contain human instructions or dialogue: R2R (Anderson et al., 2018b), Touchdown (Chen et al., 2019; Mehta et al., 2020), CVDN (Thomason et al., 2019b) and REVERIE (Qi et al., 2020). RxR addresses shortcomings of these datasets—in particular, multilinguality, scale, fine-grained word grounding, and human follower demonstrations (Table 1). It also addresses path biases in R2R. More broadly, our work is also related to instruction-guided household task benchmarks such as ALFRED (Shridhar et al., 2020) and CHAI (Misra et al., 2018). These synthetic environments provide interactivity but are generally less diverse, less visually realistic and less faithful to real world structures than the 3D reconstructions used in VLN.

**Multilinguality.** The dominance of high resource languages is a pervasive problem as it is unclear that research findings generalize to other languages (Bender, 2009). The issue is particularly severe for VLN. Chen and Mooney (2011) translated(~1K) English navigation instructions into Chinese for a game-like simulated 3D environment. Otherwise, all publicly available VLN datasets we are aware of have English instructions.

To enable multilingual progress on VLN, RxR includes instructions for three typologically diverse languages: English (en), Hindi (hi), and Telugu (te). The English portion includes instructions by speakers in the USA (en-US) and India (en-IN). Unlike Chen and Mooney (2011) and like the TyDi-QA multilingual question answering dataset (Clark et al., 2020), RxR's instructions are not translations: all instructions are created from scratch by native speakers. This especially matters for VLN, as different languages encode spatial and temporal information in idiosyncratic ways–e.g., how contact/support relationships are expressed (Munnich et al., 2001), frame of reference (Haun et al., 2011), and how temporal accounts are expressed (Bender and Beller, 2014).

**Scale.** Embodied language tasks suffer from a relative paucity of training data; for VLN, this has

| | Number of: | | | | Includes: | | |
|---|---|---|---|---|---|---|---|
| | Lang | Instruct | Words | Paths | Text | Ground | Demos |
| CVDN | 1 | 2K[†] | 167K | 7K | ✓ | | |
| R2R | 1 | 22K | 625K | 7K | ✓ | | |
| Touchdown | 1 | 9K | 1.0M | 9K | ✓ | ✓[‡] | |
| REVERIE | 1 | 22K | 388K | 7K | ✓ | ✓[‡] | |
| RxR | 3 | 126K | 9.8M | 16.5K | ✓ | ✓ | ✓ |

[†] The number of dialogues. [‡] Grounding limited to one object per instruction.

Table 1: VLN dataset comparison. RxR is larger, multilingual, and includes dense spatiotemporal groundings (Ground) and follower demonstrations (Demos).

led to a focus on data augmentation (Fried et al., 2018; Tan et al., 2019), pre-training (Wang et al., 2019; Huang et al., 2019; Li et al., 2019), multi-task learning (Wang et al., 2020) and better generalization through piece-wise curriculum design (Zhu et al., 2020). To address this shortage, for each language RxR contains 14K paths with 3 instructions per path, for a total of 126K instructions and 10M words (based on whitespace tokenization). As illustrated in Table 1, this is *an order of magnitude* larger than previous datasets.

**Fine-Grained Grounding.** Like R2R, RxR's instructions are collected by immersing *Guide* annotators in a simulated first-person environment backed by the Matterport3D dataset (Chang et al., 2017) and asking them to describe predefined paths. RxR also enhances each instruction with dense spatiotemporal groundings. Guides speak as they move and later transcribe their audio; our annotation tool records their 3D poses and time-aligns the entire *pose trace* with words in the transcription. Instructions and pose traces can thus be aligned with any Matterport data including surface reconstructions (Figure 1), RGB-D panoramas (Figure 4), and 2D and 3D semantic segmentations.

**Follower Demonstrations.** Annotators also act as *Followers* who listen to a Guide's instructions and attempt to follow the path. In addition to verifying instruction quality, this allows us to collect a play-by-play account of how a human interpreted the instructions, represented as a pose trace. Guide and Follower pose traces provide dense spatiotemporal alignments between instructions, visual percepts and actions – and both perspectives are useful for agent training.

**Path Desiderata.** R2R paths span 4–6 edges and are the shortest paths from start to goal. Thomason et al. (2019a) showed that agents can exploit effective priors over R2R paths, and Jain et al. (2019) showed that R2R paths encourage goal seeking
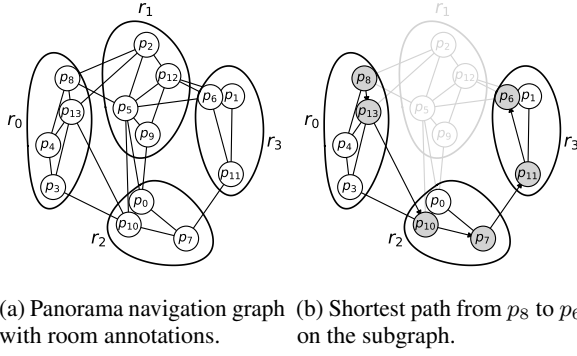
(a) Panorama navigation graph with room annotations. (b) Shortest path from $p_8$ to $p_6$ on the subgraph.

Figure 2: Given the panorama navigation graph $P$ with room graph $R$ in Figure 2a, we sample a simple room path $(r_0, r_2, r_3)$ inducing the subgraph in Figure 2b. The generated panorama path is the shortest path in the subgraph linking sampled panoramas $r_8$ and $r_6$.

over path adherence. These matter both for generalization to new environments and fidelity to the descriptions given in the instruction—otherwise, strong performance might be achieved by agents that mostly ignore the language. RxR addresses these biases by satisfying four *path desiderata*:

1. High variance in path length, such that agents cannot simply exploit a strong length prior.
2. Paths may approach their goal indirectly, so agents cannot simply go straight to the goal.
3. Naturalness: paths should not enter cycles or make continual direction changes that would be difficult for people to describe and follow.
4. Uniform coverage of environment viewpoints, to maximize the diversity of references to visual landmarks and objects over all paths.

This increases RxR's utility for testing agents' ability to ground language. It also makes RxR a more challenging VLN dataset—but one for which human followers still achieve a 93.9% success rate.

## 3 Two-Level Path Sampling

We satisfy desiderata 1-3 using a two-level procedure. At a high-level, each path visits a sequence of rooms; these are *simple paths* with no repeated (room) vertices. Such paths are *not necessarily* shortest paths. The low-level sequence is then the shortest panorama path, constrained by the room sequence. Given the set of all such paths across all houses, the fourth desiderata is satisfied by iteratively selecting the path that most improves coverage while maintaining a bias against shortest paths.

**Preliminaries** Movement in the simulator is based on a navigation graph. Vertices correspond

to 360-degree panoramic images, captured at approximately 2.2m intervals throughout 90 indoor environments. Edges are navigable links between panoramas. Chang et al. (2017) also partition panoramas via human-defined *room* annotations.

Let $P$ be an undirected graph of interconnected panoramas, with vertices $p_i \in \mathcal{V}(P)$ and edges $(p_i, p_j) \in \mathcal{E}(P)$. Let $A_R$ be a set of disjoint room annotations; each room $r_i \in A_R$ is a non-overlapping subset of panoramas $r_i \subseteq \mathcal{V}(P)$, as shown in Figure 2a. We abbreviate $(p_1, \cdots, p_m)$ as $p_{1:m}$.

We create $R$, an undirected room graph with vertices $\mathcal{V}(R) = \{\bigcup \mathcal{C}(P[r_i]) \mid r_i \in A_R\}$. $P[r_i]$ is the subgraph of $P$ induced by room annotation $r_i$ and $\mathcal{C}$ returns a graph's connected components. Simply put, each vertex in $R$ encompasses a subgraph of $P$. An edge $(r_i, r_j) \in \mathcal{E}(R)$ exists if the subgraph of $P$ induced by $\mathcal{V}(r_i) \cup \mathcal{V}(r_j)$ is connected.

**Path Generation** We generate the set of all simple paths in $R$ that traverse at most 5 rooms and two building levels. Let $r_{p_i} \in \mathcal{V}(R)$ be the room containing panorama $p_i$. As shown in Figure 2b, for each room path $r_{1:n}$, we construct a directed graph $P[r_{1:n}]$ in which an edge $(p_i, p_j)$ exists if $r_{p_i} = r_{p_j}$ ($p_i$ and $p_j$ are in the same room) or $(r_{p_i}, r_{p_j})$ is an edge in the room path. Given $P[r_{1:n}]$, we sample the start $p_1$ and goal $p_m$ uniformly from $r_1$ and $r_n$, respectively. The full panorama path $p_{1:m}$ is then the shortest path between $p_1$ and $p_m$ in $P[r_{1:n}]$.

Room size varies greatly, so this approach produces high path length variance. It also satisfies naturalness because people tend to ground instructions at the room level (e.g., *Exit through the carved wooden door on the other side of the room*). We find such paths easy to describe even with as many as 20 edges. Finally, these paths can approach their goal indirectly, as exemplified in Figure 2b.

**Greedy Selection for Coverage** The final path dataset $D$ is constructed by repeatedly selecting a panorama path $p_{1:m}$ from all sampled paths (without replacement) until a desired size is reached. After selecting $k$ paths, let $\mathcal{O}(p_i, D_k)$ be the number of occurrences of panorama $p_i$ in the paths in $D_k$. At step $k+1$, we select the path with the minimum value for $\frac{d(p_1, p_m)}{L(p_{1:m})} + \frac{1}{m} \sum_{p_i \in p_{1:m}} \mathcal{O}(p_i, D_k)$, where $L$ is path length in $P$ and $d(p_1, p_m)$ is the shortest path distance between $p_1$ and $p_m$ in $P$. The first term prefers non-shortest paths while the second encourages selection of paths that cover panoramas with low coverage in $D_k$. This selection step is also
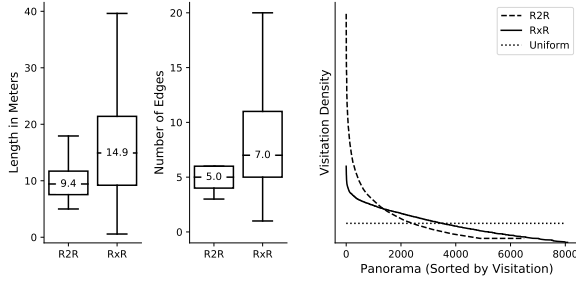
Figure 3: RxR's paths are longer on average than R2R's, exhibiting far greater variation in length (measured in both meters and edges) while achieving more uniform panorama coverage. Comparisons shown are for train and val since the R2R test set is sequestered.

subject to a maximum path length of 40m, and a maximum of 500 paths per building environment.

**Path Statistics** In total, we sample 16522 paths, which are split: 11089 train, 1232 val-seen (train environments), 1517 val-unseen (val environments), and 2684 test, following the same environment splits as Matterport3D and R2R. Compared to R2R, RxR paths are longer, spanning 8 edges and 14.9m on average, vs. 5 edges and 9.4m in R2R. More importantly, as shown in Figure 3, RxR paths exhibit much greater variation in length while also achieving more uniform coverage of the panoramas (and edges). Furthermore unlike R2R, 44.5% of RxR paths are *not* the shortest path from the start to the goal location. RxR paths are on average 27.4% longer than the shortest path.

## 4 Data Collection and Metrics

We immerse annotators in our own web-based version of the Matterport3D simulator using the panoramic images and the navigation graph. Compared to Anderson et al. (2018b), our annotation tool has additional capabilities including speech collection, virtual pose tracking, and time-alignment between transcript and pose. Figure 4 gives an example instruction with accompanying Guide and Follower pose traces. Here, we describe our collection process, analysis of the data, path evaluation metrics and simple baselines.

**Guide Task** Like R2R, our simulator has camera controls allowing continuous heading and elevation changes and movement between panoramas. Guides look around and move to explore a provided path and attempt to create an instruction others can follow. R2R's Guides create written instructions.

**Guide Alignment**          **Follower Alignment**



Now you are standing in-front of a closed door, turn to your left, you can see two wooden steps, climb the steps and walk forward by crossing a...

Now you are standing in-front of a closed door, turn to your left, you can see two wooden steps, climb the steps and walk forward by...

...crossing a wall painting which is to your right side, you can see open door enter...

...by crossing a wall painting which is to your right..

...enter into it. This is a gym room, move forward, walk...

...right side, you can see open door enter into it. This is a gym room, move forward, walk...

...walk till the end of the room, you can see a grey...

...walk till the end of the room, you can...

...grey colored ball to the corner of the room, stand there, that's...

...can see a grey colored ball to the corner...

...that's your end point.

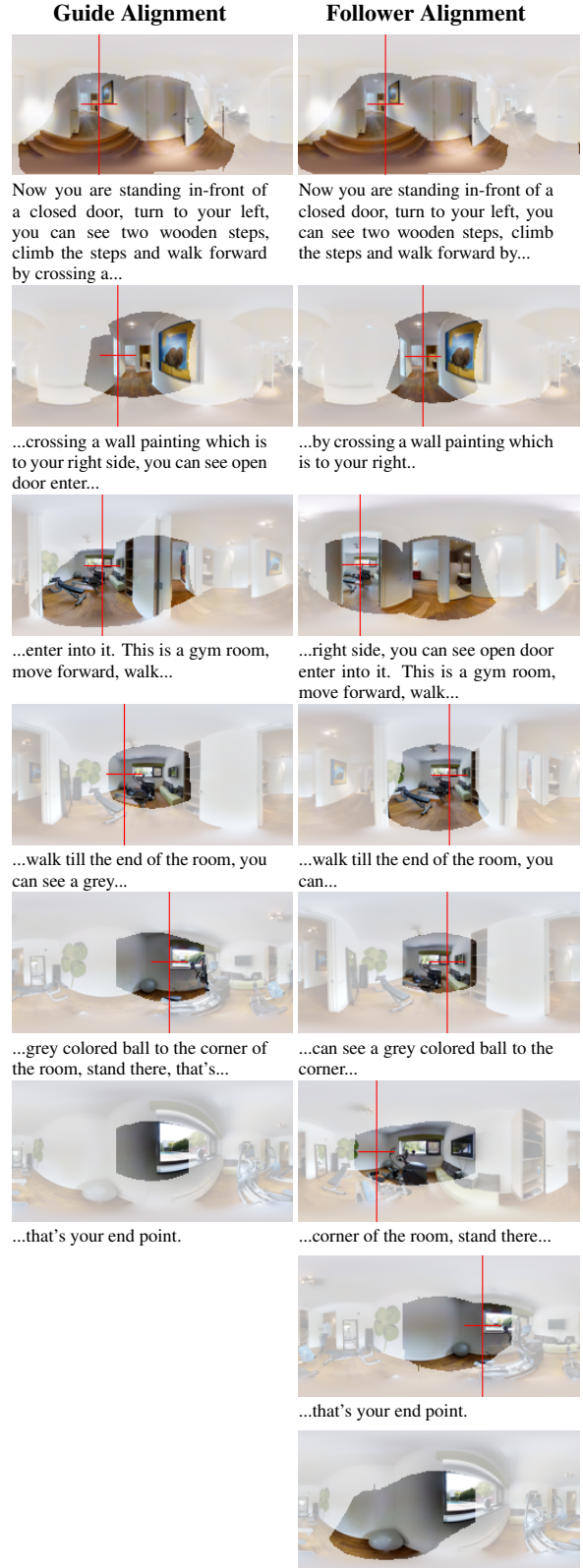...corner of the room, stand there...

...that's your end point.

Figure 4: Example spatiotemporal alignment of textual instructions, visual percepts and actions for an en-US Guide and the corresponding Follower. The next selected action is indicated in red and unseen pixels in the equirectangular panoramic images are faded. The Follower takes a slightly longer path but produces similar visual-textual alignments. Best viewed enlarged.

| | en-IN | en-US | en | hi | te | Total |
|---|---|---|---|---|---|---|
| **Counts** | | | | | | |
| Instructions | 28010 | 13992 | 42002 | 42068 | 41999 | 126069 |
| Paths | 14005 | 13992 | 14005 | 14026 | 14003 | 16522 |
| **Averages** | | | | | | |
| Words | 87 | 129 | 101 | 76 | 56 | 78 |
| WordPieces | 104 | 159 | 123 | 143 | 184 | 150 |
| Characters | 457 | 659 | 524 | 355 | 395 | 425 |
| Audio (s) | 64 | 80 | 69 | 53 | 58 | 60 |
| Guide (s) | 431 | 509 | 457 | 451 | 465 | 458 |
| Follower (s) | 134 | 202 | 156 | 110 | 132 | 132 |

Table 2: RxR summary statistics. Times in seconds (s).

In contrast, RxR's Guides *speak* and the tool logs their entire virtual camera pose sequence. We use a 640 × 480 pixel viewing canvas and a camera vertical field of view of 75 degrees. This process is inspired by Localized Narratives (Pont-Tuset et al., 2020), an image captioning dataset for which annotators move mouse pointers around images while talking about them.

As with Localized Narratives, RxR Guides transcribe their own recordings; this produces high quality text versions of the instructions. To align text and pose traces, we generate a time-stamped transcription using automatic speech recognition.[2] The transcription and ASR output are aligned using dynamic time warping. The output of the Guide task is an audio file, a tokenized, timestamped, manually-transcribed instruction, and a *pose trace* (a series of timestamped 6-DOF camera poses). On average, Guide task annotations (including both steps, performed back-to-back) take 458 seconds.

For each language (English, Hindi and Telugu) we annotate 14K paths with three instructions each. In the English dataset, each path gets one US English instruction and two Indian English instructions. Of the 14K paths per language, 12.8K paths are common across all three languages, and 1.2K paths in each language are unique (equaling 16.5K paths in total). The fact that most paths are annotated 9 times (3 per language) creates interesting opportunities to study aligned instructions across languages. Unique paths add variety and coverage.

**Follower Task** As Followers, annotators begin at the start of an unknown path and try to follow the Guide's instruction. They observe the environment and navigate in the simulator as the Guide's audio plays. They can pause, rewind and skip forward in the instruction. If they believe they have reached the the end of the path, or give up, they indicate they

---

[2]https://cloud.google.com/speech-to-text

are done and rate the instruction's clarity and their confidence in their own navigation. On average, Follower tasks take 132 seconds.

The Follower tasks objectively validate the quality of Guide instructions based on whether the Follower can succeed (i.e., reaching within 3m of the last panorama in the path). If the Follower doesn't succeed, the Guide instruction is paired with a second Follower. If the second Follower succeeds, the first Follower annotation is discarded and replaced. If the second Follower also fails, then the path is re-enqueued to generate another Guide and Follower annotation. The most successful of the three resulting Guide-Follower pairs is selected for inclusion in RxR and the others are discarded.

In addition to validating data quality, the Follower task also trains annotators to be better Guides—following bad instructions often helps one see how to produce better instructions. Most importantly, we collect the pose trace of the Follower as they execute the instruction. This provides an alternative path with dense grounding that we can compare to the Guide's pose trace and use as an additional training signal.

**Dataset Analysis** Table 2 provides summary statistics for RxR. The average words per instruction (using whitespace tokenization) is 78 vs R2R's 29. US English instructions are the longest on average. We attribute this to conventions developed by each annotator pool rather than language specific properties. On average Guide tasks take much longer than Follower tasks (458 vs. 132 seconds). Most of the Guide's time is spent transcribing audio (Guide audio recordings average 60 seconds).

Following a similar analysis as Chen et al. (2019), Table 3 gives examples and statistics for linguistic phenomena, based on manual analysis of instructions for 25 paths. All RxR subsets produce a higher rate of *entity references* compared to R2R. This is consistent with the extra challenge of RxR's paths and our annotation guidance that instructions should help followers stay on the path as well as reach the goal. Doing so requires more extensive use of objects in the environment. RxR's higher rate of both *coreference* and *sequencing* indicates that its instructions have greater discourse coherence and connection than R2R's. RxR also includes a far higher proportion of *allocentric relations* and *state verification* compared to R2R, and matches Touchdown (navigation instructions). Hindi contains less coreference, sequencing, and temporal

| Phenomenon | R2R | | RxR | | | | | | | | | RxR Example (en-US) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | en | | hi | | te | | en-IN | | en-US | | | |
| | $p$ | $\mu$ | $p$ | $\mu$ | $p$ | $\mu$ | $p$ | $\mu$ | $p$ | $\mu$ | | |
| Reference | 100 | 3.7 | 100 | 5.8 | 100 | 6.6 | 100 | 6.4 | 100 | 8.3 | | ...there is **a white chair** and **a table stand**... |
| Coreference | 32 | 0.5 | 40 | 0.4 | 76 | 2.9 | 76 | 6.4 | 64 | 5.3 | | ...hallway with black curtains, towards **that**... |
| Comparison | 4 | 0.0 | 0 | 0.0 | 4 | 0.1 | 4 | 0.0 | 8 | 0.0 | | ...the large archway with the **smaller** archway in... |
| Sequencing | 16 | 0.2 | 24 | 0.2 | 44 | 0.6 | 44 | 0.5 | 52 | 0.9 | | ...the **next** room... turn to see the **next** door... |
| Allocentric Relation | 20 | 0.2 | 68 | 2.1 | 76 | 3.2 | 92 | 3.4 | 76 | 2.4 | | ...a window with a black folding table **under** that... |
| Egocentric Relation | 80 | 1.2 | 96 | 2.9 | 80 | 2.3 | 64 | 2.8 | 60 | 2.3 | | ...chairs on **your right**, closet doors on **your left**. |
| Imperative | 100 | 4.0 | 100 | 5.6 | 100 | 6.5 | 100 | 8.4 | 100 | 6.3 | | **Do not** go down the stairs. Instead, **look** further... |
| Direction | 100 | 2.8 | 96 | 5.8 | 96 | 4.9 | 100 | 7.0 | 96 | 6.3 | | ...**veer to the left** of the fireplace and you will... |
| Temporal Condition | 28 | 0.4 | 32 | 0.4 | 36 | 0.7 | 44 | 1.0 | 52 | 0.8 | | Move around the island **until** you come to the... |
| State Verification | 8 | 0.1 | 72 | 1.7 | 68 | 1.6 | 80 | 2.3 | 84 | 3.1 | | ...**you are in** the balcony area facing towards... |

Table 3: Linguistic phenomena in a manually annotated random sample of 25 paths from RxR and R2R. $p$ is the % of sentences that contain the phenomena while $\mu$ is the average number of times they occur within each sentence.
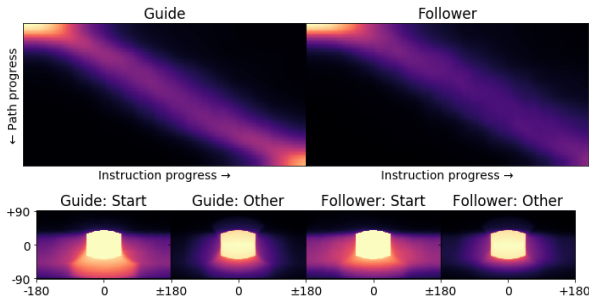


Figure 5: Top: Instruction and path progress alignment for Guides and Followers. Bottom: Equirectangular heatmap of Guide and Follower camera poses, centered on their initial perspective at each viewpoint.

| | PL | NE↓ | SR↑ | SPL↑ | SDTW↑ | NDTW↑ |
|---|---|---|---|---|---|---|
| 1. Random walk | | | | | | |
| R2R | 10.4 | 9.5 | 5.1 | 3.6 | 3.8 | 27.6 |
| RxR | 16.8 | 12.4 | 8.8 | 2.5 | 3.8 | 18.2 |
| 2. Random heading then go straight | | | | | | |
| R2R | 9.7 | 9.9 | 8.2 | 7.2 | 6.6 | 28.3 |
| RxR | 15.1 | 13.5 | 8.0 | 3.4 | 3.9 | 16.3 |
| 3. Given correct first step then go straight | | | | | | |
| R2R | 9.5 | 6.2 | 27.2 | 25.7 | 23.6 | 52.6 |
| RxR | 15.3 | 11.4 | 13.7 | 7.5 | 8.3 | 25.9 |

Table 4: Simple baselines on val-unseen paths. RxR proves more difficult than R2R overall, and less amenable to agents that tend to go straight (baselines 2 and 3). Note: Baseline 3 partly exploits the gold path.

conditions than the other languages. That said, it is not clear how much the differences *within* RxR exhibited in Table 3 can be attributed to language, dialect, annotator pools, or other factors.

Figure 5 (top) illustrates the close alignment between instruction progress (measured in words) and path progress (measured in steps). Figure 5 (bottom) indicates that both Guide and Tourist annotators orient themselves by looking around at the first panoramic viewpoint, after which they maintain a narrower focus. On average, Guides / Tourists observe 43% / 44% of the available spherical visual signal at the first viewpoint, and 27% / 28% at subsequent viewpoints. These findings stand in contrast to standard VLN agents that routinely consume the entire panoramic image and attend over the entire instruction sequence at each step. Inputs that the Guide / Tourist have not observed cannot influence their utterances / actions, so pose traces offer rich opportunities for agent supervision.

**Evaluation**   We use the following standard evaluation metrics (with arrows indicating improvement): *Path Length* (**PL**), *Navigation Error* (**NE** ↓) *Success Rate* (**SR** ↑), *Success weighted by inverse Path Length* (**SPL** ↑), *Normalized Dynamic Time Warping* (**NDTW** ↑), and *Success weighted by normalized Dynamic Time Warping* (**SDTW** ↑). See Anderson et al. (2018a) and Ilharco et al. (2019) for discussion of VLN metrics. Since RxR was designed to include paths that approach their goal indirectly, we focus primarily on **NDTW** and **SDTW** which explicitly capture path adherence. See Table 4 for a comparison of the performance of several simple baselines on R2R and RxR. Each simple baseline requires a stopping criteria; we choose to stop after $N$ steps where $N$ is the average number of steps in the train set paths (5 in R2R and 8 in RxR). Consistent with our motivation to reduce biases in paths, these simple baselines show that going straight is far less effective in RxR than R2R.

## 5 Experiments

**Agent** We use a model architecture similar to that of the Reinforced Cross-Modal Matching (RCM) agent (Wang et al., 2019), consisting of an instruction encoder and a sequential LSTM (Hochreiter and Schmidhuber, 1997) decoder that computes a distribution over actions at each step. However, since RxR instructions are much longer than R2R, we replace the bidirectional LSTM instruction encoder with a more parallelizable CNN encoder. In preliminary experiments on R2R we find that encoding word embeddings via successive 1D convolutions with rectified linear (ReLU) activations and residual connections (He et al., 2016) is equally effective and more time and space efficient. We denote the output of the instruction encoder by $x \in \mathbb{R}^{l \times d}$ where $l$ is the instruction length and $d$ is the feature dimension. In both monolingual and multilingual experiments we use features extracted from a pre-trained multilingual BERT model (Devlin et al., 2019) for the word embeddings.

At each time step $t$, the agent receives a panoptic encoding of its viewpoint $v_t \in \mathbb{R}^{k \times d}$ (where $k = 36$ is the number of $30°$ intervals that span the panorama) along with a visual encoding of navigable directions $a_t \in \mathbb{R}^{n \times d}$ (where $n$ is the number of navigable directions). Each feature of dimension $d$ is a pre-trained CNN feature concatenated with an angle encoding (Fried et al., 2018). The LSTM decoder computes an updated hidden state $h_t$ by conditioning on the previous selected action in $a_{t-1}$ and attending over the panoptic encoding $v_t$ and the instruction $x$ using dot-product attention (Luong et al., 2014). The distribution over next actions is computed via a similarity ranking $h_t \cdot a_{t,i}$ between hidden state $h_t$ and each direction encoding in $a_t$.

For the image features we use an EfficientNet-B4 CNN (Tan and Le, 2019). Following Parekh et al. (2020), we pretrain the CNN in an image-text dual encoder setting using the Conceptual Captions dataset (Sharma et al., 2018). In preliminary experiments, we found that pretraining the CNN in this way gave noticeable improvements over the same CNN pretrained for image classification on ImageNet (Russakovsky et al., 2015).

**Grounding Supervision** To incorporate spatiotemporal groundings into agent training, for each Guide path (G-path) and Follower path (F-path) we convert the corresponding pose trace into: (1) a sequence of text masks $b_t \in \{0, 1\}^l$ indicating which words in instruction $x$ the Guide spoke / Follower heard *at or prior to* step $t$, and (2) a sequence of visual masks $M_t \in \{0, 1\}^{h \times w}$ indicating which pixels were observed in the panoramic image at $t$ (like Figure 5 bottom). We then project and max-pool $M_t$ to a vector mask $m_t \in \{0, 1\}^k$ aligning to the agent's visual input features $v_t$. Zeros in $b_t$ and $m_t$ indicate irrelevant textual and visual inputs that were not observed by the annotators, and are therefore not related to their utterances and actions.

To help prevent the agent from overfitting to superficial correlations in the training data, we use $b_t$ and $m_t$ to supervise the normalized textual and visual attention weights in the model. Specifically, during training whenever the agent is on the gold path we apply a cross-entropy loss to the visual attention weights given by $\mathcal{L}(z, m_t) = \log \sum_{i=1}^{k} \exp(z_i) - \log \sum_{i=1}^{k} m_{t,i} \exp(z_i)$, where $z$ is the vector of unnormalized logits determining attention weights via a softmax. This loss forces the attention weights on irrelevant input features towards zero. The textual version is analogous.

**Implementation Details** Agents are implemented in VALAN (Lansing et al., 2019), a distributed reinforcement learning framework designed for VLN. We use a mix of supervised learning and policy gradients. Each minibatch is constructed from 50% behavioural cloning roll-outs (following the gold paths while minimizing cross-entropy loss), and 50% policy gradient rollouts with reward (following paths sampled from the agent's policy). As in Ilharco et al. (2019), the reward at each step is the incremental difference in NDTW, plus a linear function of navigation error after stopping. All agents are trained with Adam (Kingma and Ba, 2014) to convergence (100K iterations with batch size of 32 and initial learning rate of 1e-4).

**Monolingual Results** Table 5 provides results on the val-unseen split for several training settings, as well as human performance from Follower annotations. We report en-US and en-IN results together as en. Experiments 1–3 compare agents trained (1) only on G-paths, (2) only on F-paths, and (3) on both. In contrast to algorithmically generated G-paths, each F-path reflects a grounded human interpretation of an instruction, which may deviate from the G-path because multiple correct interpretations are possible (e.g., Figure 4). For training, we do not differentiate F-paths from G-paths, and each

| Exp. | Method | Setting | | | Training | NE ↓ | | | SR ↑ | | | SDTW ↑ | | | NDTW ↑ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | G | F | X | Pairs (K) | en | hi | te | en | hi | te | en | hi | te | en | hi | te |
| (1) | Mono | ✓ | | | 42 | 10.1 | 9.7 | 9.4 | 25.6 | 24.8 | 28.0 | 20.3 | 19.7 | 22.7 | 41.3 | 38.8 | 43.7 |
| (2) | Mono | | ✓ | | 42 | 10.3 | **9.2** | 9.5 | 23.9 | 28.0 | 27.0 | 18.5 | 22.7 | 22.0 | 37.0 | **45.9** | 43.9 |
| (3) | Mono | ✓ | ✓ | | 84 | **9.8** | **9.2** | **9.1** | **26.1** | **29.6** | **29.8** | **21.0** | **24.0** | **24.2** | **42.4** | 45.5 | **45.6** |
| (4) | Multi | ✓ | ✓ | | 252 | **11.0** | 10.9 | 11.0 | **22.2** | **23.0** | 23.1 | **17.8** | **18.3** | **18.4** | **38.6** | 39.2 | 38.8 |
| (5) | Multi | ✓ | ✓ | ✓ | 504 | 11.5 | 11.4 | 11.4 | 20.0 | 18.7 | 20.3 | 15.9 | 14.9 | 16.1 | 36.3 | 36.0 | 36.7 |
| (6) | Multi* | ✓ | ✓ | | 252 | **11.0** | **10.7** | **10.7** | 21.9 | 22.6 | **23.2** | 17.5 | 18.1 | **18.4** | **38.6** | **39.9** | **39.7** |
| (H) | Human | | | | - | 1.32 | 0.59 | 0.79 | 90.4 | 96.8 | 94.7 | 74.3 | 80.6 | 76.5 | 77.7 | 82.2 | 79.2 |

Settings – G: instruction paired with Guide paths, F: instructions paired with Follower paths, X: cross-translated instructions.

Table 5: RxR val-unseen: Monolingual vs. multilingual results. Training with both Guide and Follower paths benefits all languages (exp. 3 vs. 1 and 2), monolingual outperforms multilingual (exp. 3 vs. 4), training with cross-translations hurts performance (exp. 5 vs. 4), and training with visual attention supervision gives mixed results (Multi* in exp. 6 vs 4).

| Exp. | Train Data | | SR ↑ | | | | SPL ↑ | | | | SDTW ↑ | | | | NDTW ↑ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | R2R | RxR | R2R | en | hi | te | R2R | en | hi | te | R2R | en | hi | te | R2R | en | hi | te |
| (7) | ✓ | | 36.5 | 14.5 | 9.6 | 9.7 | 31.7 | 11.2 | 7.5 | 7.4 | 29.5 | 9.8 | 6.3 | 6.1 | 48.1 | 29.0 | 25.4 | 25.2 |
| (4) | | ✓ | 19.2 | 22.2 | 23.0 | **23.1** | 17.7 | 19.8 | 20.7 | **20.7** | 16.0 | 17.8 | 18.3 | **18.4** | 43.2 | 38.6 | 39.2 | **38.8** |
| (8) | ✓ | ✓ | **37.8** | 22.5 | **23.6** | **23.1** | **34.3** | 20.1 | 21.0 | 20.5 | **32.0** | 18.3 | **19.2** | **18.4** | **52.3** | 38.8 | **39.4** | 38.4 |

Table 6: Multitask and transfer learning results on RxR and R2R val-unseen. A multitask model (exp. 8) performs best on both datasets, but domain differences thwart simple transfer learning (i.e., train on X, evaluate on Y).

instruction-path pair is treated as an independent example. Experiment (3) shows that including both G- and F-paths in training benefits every metric. Given the overall positive impact of F-paths, we use both path types in our further experiments.

**Multilinguality** For experiment (4) in Table 5, we train a single multilingual agent on all three languages simultaneously. While the multilingual agent sees substantially more instructions than each monolingual agent, performance is worse across all metrics. This is consistent with results in multilingual machine translation (MT) and automatic speech recognition (ASR) where adding more languages can also lead to degradation for high-resource languages (Aharoni et al., 2019; Pratap et al., 2020). Experiment (5) takes this one step further by obtaining translations from every instruction into the two other languages (e.g., en → hi, te) using a MT service.[3] Including these translations hurts performance for all languages. The fact that most G-paths are shared across languages may limit the value of automatic cross-translations. Notwithstanding the higher performance of the monolingual approaches, in the remaining experiments we focus on multilingual agents for greater scalability.

---
[3]https://cloud.google.com/translate
These translations are included in the RxR data release.

**Spatiotemporal Grounding Supervision** Table 5 experiment (6) incorporates a loss for spatiotemporal grounding over visual attention which gives mixed results on val-unseen (better on NDTW, NE and worse on success-based metrics) compared to (4). Applying the same approach to textual attention did not improve performance. However, we stress that this is only a preliminary investigation. Using human demonstrations to supervise visual groundings is an active area of research (Wu and Mooney, 2019; Selvaraju et al., 2019). As one of the first large-scale spatially-temporally aligned language datasets, RxR offers new opportunities to extend this work from images to environments.

**Multitask and Transfer Learning** Table 6 reports the performance of the multilingual agent under multitask and transfer learning settings. For simplicity, the R2R model (exp. 7) is trained without data augmentation from model-generated instructions (Fried et al., 2018; Tan et al., 2019) and with hyperparameters tuned for RxR. Under these settings, the multitask model (exp. 8) performs best on both datasets. However, transfer learning performance (RxR → R2R and vice-versa) is much weaker than the in-domain results. Although RxR and R2R share the same underlying environments, we note that RxR → R2R cannot exploit R2R's

| Exp. | Method | Input Modalities | | NE ↓ | | | SR ↑ | | | SDTW ↑ | | | NDTW ↑ | | |
| | | Vision | Language | en | hi | te | en | hi | te | en | hi | te | en | hi | te |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| (4) | Multi | ✓ | ✓ | **11.0** | **10.9** | **11.0** | **22.2** | **23.0** | **23.1** | **17.8** | **18.3** | **18.4** | **38.6** | **39.2** | **38.8** |
| (9) | Multi | | ✓ | 12.3 | 11.9 | 12.0 | 16.0 | 18.0 | 16.9 | 12.3 | 14.2 | 13.3 | 30.9 | 33.1 | 32.8 |
| (10) | Multi | ✓ | | 15.7 | 15.7 | 15.7 | 7.8 | 7.8 | 7.8 | 4.3 | 4.3 | 4.3 | 16.5 | 16.5 | 16.5 |

Table 7: Language-only and vision-only model ablations on RxR val-unseen. The language-only agent is much better than random, but both modalities are required for best performance.

| Split | Method | NE ↓ | | | | SR ↑ | | | | SDTW ↑ | | | | NDTW ↑ | | | |
| | | en | hi | te | avg | en | hi | te | avg | en | hi | te | avg | en | hi | te | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Val-Seen | Mono | 9.5 | 9.2 | 9.3 | 9.3 | 28.6 | 29.5 | 28.3 | 28.8 | 23.2 | 24.6 | 23.7 | 23.8 | 45.4 | 47.9 | 47.1 | 46.8 |
| | Multi | 11.0 | 10.4 | 10.6 | 10.7 | 23.9 | 26.7 | 25.1 | 25.2 | 19.6 | 21.9 | 20.5 | 20.7 | 41.2 | 43.4 | 42.0 | 42.2 |
| Val-Unseen | Mono | 9.8 | 9.2 | 9.1 | 9.4 | 26.1 | 29.6 | 29.8 | 28.5 | 21.0 | 24.0 | 24.2 | 23.1 | 42.4 | 45.5 | 45.6 | 44.5 |
| | Multi | 11.0 | 10.9 | 11.0 | 10.9 | 22.2 | 23.0 | 23.1 | 22.8 | 17.8 | 18.3 | 18.4 | 18.2 | 38.6 | 39.2 | 38.8 | 38.9 |
| Test-Std | Mono | 11.0 | 10.5 | 10.5 | 10.6 | 25.3 | 26.1 | 26.2 | 25.9 | 20.5 | 21.0 | 21.5 | 21.0 | 40.3 | 41.9 | 42.4 | 41.5 |
| | Multi | 12.0 | 11.8 | 11.8 | 11.9 | 20.8 | 21.4 | 21.6 | 21.3 | 16.8 | 17.3 | 17.3 | 17.1 | 36.7 | 37.6 | 37.4 | 37.2 |
| | Random | 14.1 | 14.1 | 14.1 | 14.1 | 7.5 | 7.5 | 7.5 | 7.5 | 3.1 | 3.1 | 3.1 | 3.1 | 15.4 | 15.4 | 15.4 | 15.4 |
| | Human | 1.4 | 0.6 | 0.7 | 0.9 | 90.2 | 96.7 | 94.9 | 93.9 | 73.6 | 80.5 | 76.6 | 76.9 | 77.2 | 82.0 | 79.2 | 79.5 |

Table 8: RxR test set results, based on the monolingual agents (3) and the multilingual agent (4).

path bias, and for R2R → RxR, the much longer paths and richer language are out-of-domain.

**Unimodal Ablations**    Table 7 reports the performance of the multilingual agent under settings in which we ablate either the vision or the language inputs during both training and evaluation, as advocated by Thomason et al. (2019a). The multimodal agent (4) outperforms both the language-only agent (9) and the vision-only agent (10), indicating that both modalities contribute to performance. The language-only agent performs better than the vision-only agent. This is likely because even without vision, parts of the instructions such as 'turn left' and 'go upstairs' still have meaning in the context of the navigation graph. In contrast, the vision-only model has no access to the instructions, without which the paths are highly random.

**Test Set**    RxR includes a heldout test set, which we divide into two splits: test-standard and test-challenge. These splits will remain sequestered to support a public leaderboard and a challenge so the community can track progress and evaluate agents fairly. Table 8 provides test-standard performance of the mono and multilingual agents using Guide and Follower paths, along with random and human Follower scores. While the learned agent is clearly much better than a random agent, there is a great deal of headroom to reach human performance.

## 6    Conclusion

RxR represents a significant evolution in the scale, scope and possibilities for research on embodied language agents in simulated, photo-realistic 3D environments. RxR's paths better ensure that language itself will play a fundamental role in better agents. Evaluating on three typologically diverse languages will help the community avoid overfitting to a particular language and dataset.

We have only begun to explore the possibilities opened up by pose traces. Whereas others have retro-actively refined R2R's annotations to get alignments between sub-instructions and panorama sequences (Hong et al., 2020), RxR provides *word-level* alignments to *specific pixels* in panoramas. This is obtained as a by-product of significant work on the annotation tooling itself and designing the process to be more natural for Guides. Finally, every instruction is accompanied by a Follower demonstration, including a perspective camera pose trace that shows a play-by-play account of how a human interpreted the instructions given their position and progress through the path. We have shown that these can help with agent training, but they also open up new possibilities for studying grounded language pragmatics in the VLN setting, and for training VLN agents with perspective cameras – either in the graph-based simulator or by lifting RxR into a continuous simulator (Krantz et al., 2020).

## Acknowledgments

## References

Roee Aharoni, Melvin Johnson, and Orhan Firat. 2019. Massively multilingual neural machine translation. In *NAACL-HLT*.

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. 2018a. On evaluation of embodied navigation agents. *arXiv:1807.06757*.

Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton van den Hengel. 2018b. Vision-and-Language Navigation: Interpreting visually-grounded navigation instructions in real environments. In *CVPR*.

Andrea Bender and Sieghard Beller. 2014. Mapping spatial frames of reference onto time: A review of theoretical accounts and empirical findings. *Cognition*, 132(3):342–382.

Emily M. Bender. 2009. Linguistically naïve != language independent: Why NLP needs linguistic typology. In *EACL Workshop on the Interaction between Linguistics and Computational Linguistics: Virtuous, Vicious or Vacuous?*

Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. 2017. Matterport3D: Learning from RGB-D data in indoor environments. *3DV*.

David L Chen and Raymond J Mooney. 2011. Learning to interpret natural language navigation instructions from observations. In *AAAI*.

Howard Chen, Alane Suhr, Dipendra Misra, Noah Snavely, and Yoav Artzi. 2019. Touchdown: Natural language navigation and spatial reasoning in visual street environments. In *CVPR*.

Jonathan H. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. TyDi QA: A benchmark for information-seeking question answering in typologically diverse languages. *TACL*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*.

Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe Morency, Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. 2018. Speaker-follower models for vision-and-language navigation. In *NeurIPS*.

Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumeé III, and Kate Crawford. 2020. Datasheets for Datasets. *arXiv:1803.09010*.

Daniel BM Haun, Christian J Rapold, Gabriele Janzen, and Stephen C Levinson. 2011. Plasticity of human spatial cognition: Spatial language and cognition covary across cultures. *Cognition*, 119(1):70–80.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Computation*.

Yicong Hong, Cristian Rodriguez-Opazo, Qi Wu, and Stephen Gould. 2020. Sub-instruction aware vision-and-language navigation. In *EMNLP*.

Haoshuo Huang, Vihan Jain, Harsh Mehta, Alexander Ku, Gabriel Magalhaes, Jason Baldridge, and Eugene Ie. 2019. Transferable representation learning in vision-and-language navigation. In *ICCV*.

Gabriel Ilharco, Vihan Jain, Alexander Ku, Eugene Ie, and Jason Baldridge. 2019. Effective and general evaluation for instruction conditioned navigation using dynamic time warping. *NeurIPS Visually Grounded Interaction and Language Workshop (ViGIL)*.

Vihan Jain, Gabriel Magalhaes, Alexander Ku, Ashish Vaswani, Eugene Ie, and Jason Baldridge. 2019. Stay on the path: Instruction fidelity in vision-and-language navigation. In *ACL*.

Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv:1412.6980*.

Jacob Krantz, Erik Wijmans, Arjun Majundar, Dhruv Batra, and Stefan Lee. 2020. Beyond the nav-graph: Vision and language navigation in continuous environments. In *ECCV*.

Larry Lansing, Vihan Jain, Harsh Mehta, Haoshuo Huang, and Eugene Ie. 2019. VALAN: Vision and language agent navigation. *arXiv:1912.03241*.

Xiujun Li, Chunyuan Li, Qiaolin Xia, Yonatan Bisk, Asli Celikyilmaz, Jianfeng Gao, Noah Smith, and Yejin Choi. 2019. Robust navigation with language pretraining and stochastic sampling. In *EMNLP*.

Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2014. Effective approaches to attention-based neural machine translation. In *EMNLP*.

Harsh Mehta, Yoav Artzi, Jason Baldridge, Eugene Ie, and Piotr Mirowski. 2020. Retouchdown: Adding touchdown to streetlearn as a shareable resource for language grounding tasks in street view. *EMNLP Workshop on Spatial Language Understanding (SpLU)*.

Dipendra Kumar Misra, Andrew Bennett, Valts Blukis, Eyvind Niklasson, Max Shatkhin, and Yoav Artzi. 2018. Mapping instructions to actions in 3d environments with visual goal prediction. In *EMNLP*.

Edward Munnich, Barbara Landau, and Barbara Anne Dosher. 2001. Spatial language and spatial representation: A cross-linguistic comparison. *Cognition*, 81(3):171–208.

Zarana Parekh, Jason Baldridge, Daniel Cer, Austin Waters, and Yinfei Yang. 2020. Crisscrossed captions: Extended intramodal and intermodal semantic similarity judgments for MS-COCO. *arXiv:2004.15020*.

Jordi Pont-Tuset, Jasper Uijlings, Soravit Changpinyo, Radu Soricut, and Vittorio Ferrari. 2020. Connecting vision and language with localized narratives. In *ECCV*.

Vineel Pratap, Anuroop Sriram, Paden Tomasello, Awni Hannun, Vitaliy Liptchinsky, Gabriel Synnaeve, and Ronan Collobert. 2020. Massively multilingual asr: 50 languages, 1 model, 1 billion parameters. *arXiv:2007.03001*.

Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang, William Yang Wang, Chunhua Shen, and Anton van den Hengel. 2020. REVERIE: Remote embodied visual referring expression in real indoor environments. In *CVPR*.

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. 2015. Imagenet large scale visual recognition challenge. *IJCV*.

Ramprasaath R Selvaraju, Stefan Lee, Yilin Shen, Hongxia Jin, Shalini Ghosh, Larry Heck, Dhruv Batra, and Devi Parikh. 2019. Taking a hint: Leveraging explanations to make vision and language models more grounded. In *ICCV*.

Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*.

Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020. ALFRED: A Benchmark for Interpreting Grounded Instructions for Everyday Tasks. In *CVPR*.

Hao Tan, Licheng Yu, and Mohit Bansal. 2019. Learning to navigate unseen environments: Back translation with environmental dropout. In *NAACL*.

Mingxing Tan and Quoc Le. 2019. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*.

Jesse Thomason, Daniel Gordon, and Yonatan Bisk. 2019a. Shifting the baseline: Single modality performance on visual navigation & QA. In *NAACL*.

Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer. 2019b. Vision-and-dialog navigation. In *CoRL*.

Xin Wang, Qiuyuan Huang, Asli Celikyilmaz, Jianfeng Gao, Dinghan Shen, Yuan-Fang Wang, William Yang Wang, and Lei Zhang. 2019. Reinforced cross-modal matching and self-supervised imitation learning for vision-language navigation. In *CVPR*.

Xin Wang, Vihan Jain, Eugene Ie, William Yang Wang, Zornitsa Kozareva, and Sujith Ravi. 2020. Environment-agnostic multitask learning for natural language grounded navigation. In *ECCV*.

Jialin Wu and Raymond Mooney. 2019. Self-critical reasoning for robust visual question answering. In *NeurIPS*.

Wang Zhu, Hexiang Hu, Jiacheng Chen, Zhiwei Deng, Vihan Jain, Eugene Ie, and Fei Sha. 2020. BabyWalk: Going farther in vision-and-language navigation by taking baby steps. In *ACL*.

## A  Supplementary Material

**Annotators**  In total, 247 annotators contributed to RxR, with 97 based in the USA and the remainder based in India and contributing to the Indian English, Hindi and Telugu annotations. The annotators were paid hourly wages that are competitive for their locale. They have standard rights as contractors. They were fluent in the language they were tasked with.

We ensure that a Guide does not annotate the same path twice. As Followers, annotators do not follow their own Guide instructions. Furthermore, we have provided annotators multiple forms of feedback as they complete tasks. After a round of pilot instructions were collected, we provided detailed analysis of common patterns that produced poor instructions and clear guidelines for producing better instructions. Annotators provided UI suggestions and interesting corner cases to us that allowed us to refine the simulator and annotation process before kicking off the full annotation process. Throughout the process, annotators have had access to a dashboard that shows them their success rate as both Guide and Follower. We indicated that their success as Guide and a Follower should be above 80%. Any annotator whose success is lower is either given further training or is taken off the task.

Unfortunately, we cannot release the audio instructions yet due to the impact of COVID-19: our annotators had to complete the tasks from home, so we need to review all recordings for safety and privacy. We hope to include the audio in a future release.

**Instructions**

A friend needs your help to follow an invisible path only you can see, in a house full of dangerous traps! To ensure that our friend is safe, you must record instructions for them to follow. **Before starting, please read ALL INSTRUCTIONS below and watch the training video** . You may be asked to speak in a language other than English.

**Path:** The path is shown to you as a trail of marbles.

- The marbles transition in color from **blue** to **red** as you approach the end of the path.
- Our friend's starting orientation at the first marble is shown with a **green arrow**. The other arrows indicate the direction to the next marble.
- The last marble in the path is **outlined in green** for extra visibility.
- These marbles are invisible to your friend, and should not be mentioned in your instructions.

**Recording:** Please be aware that after you press the start button, your voice will be recorded.

- As you move through the building following the path, speak clearly and audibly to record your instructions for our friend to follow.
- **The recording must only include your voice.** If other voices intrude during the recording, please stop the path and start over using the Restart button.
- **Try to mention objects and landmarks, to make your instructions as clear and unambiguous as possible.** To ensure our friend is safe, they must follow the exact same path.
  - Good example:
    "Okay look around until you can see the open door and go outside. You'll see an open balcony with a city behind it and a bed in front of you. Turn around so you can see the stairs. Go up the stairs. And stop on top of the carpet right next to a large plant that's resting against the wall."
  - Bad example:
    "Turn left. Go straight. Turn slight left again. Go a little bit forward. Exit the door."

**Movement:** Look around the room by dragging the screen, and move around by double-clicking. We recommend using a mouse instead of a touchpad.

- **Record your instructions as you move, while looking at the path and objects you are describing.** IMPORTANT: DO NOT write down your instructions and DO NOT speak all at once at the end. DO NOT traverse the path multiple times to check your work.
- The **green square** indicates where you will move to next if you double-click.
- When you reach the end of the path, press the Done button. This will load a window that allows you to transcribe your voice recording.
- Your instructions will be evaluated by whether our friend can later follow the path and reach the end without encountering any traps.

**Now please watch the training video.**

[ Start ]

(a) Guide Worker Instructions (49.1 seconds)

**Instructions**

Oh no, you're in house full of dangerous traps! Fortunately, a friend has recorded spoken instructions to enable you to follow a safe path in the house. **Before starting, please read ALL INSTRUCTIONS below and watch the training video.**

**Start:** After you press the start button, your friend's spoken audio instructions will begin to play.

- **Listen to the instructions as you move, while trying to follow the path as it is described.**
- To control the audio progress, click the audio waveform or press the ESC key to play/pause.
- Please follow the instructions as closely as you can. DO NOT explore the building unnecessarily. DO NOT traverse the path multiple times unless you are lost.
- You should never hear instructions that you personally recorded. Please tell somebody if this happens.

**Movement:** You can look around the room by dragging the screen, and move around by double-clicking. We recommend using a mouse instead of a touchpad.

- The **green square** indicates where you will move to next if you double-click.
- When you believe you've reached the end of the path, press the Done button.
- Both you and your friend will be evaluated by whether you were able to follow the path and reach the end correctly.

**Now please watch the training video.**

[ Start ]

(b) Follower Worker Instructions (38.5 seconds)

**Voice Recording: Please speak in English**

[ Restart ] [ Info ] [ Done ]

(c) Guide Annotation (64.4 seconds)

**Annotation (English audio)**

[ Info ] [ Done ]

Playback speed: [ Normal ]

[ Play/Pause ]

(d) Follower Verification (89.8 seconds)

**Transcribe in English text**

**Type literally what you just said.**

Playback speed: [ Normal ]

[ Play/Pause ]

Include disfluencies and filler words if you said them (e.g. "I think", "alright", "turn left... I mean right") but not filler sounds (e.g. "um", "uh", "er"). Feel free to separate the text into multiple sentences and add punctuation. Press the submit button when you are done.

**Audio key bindings:** ESC will play and pause the audio. ALT plus the left or right arrow key will skip the audio backwards or forwards.

How confident are you that our friend can follow the path given your instructions?

Not confident                                              Very confident

[ Submit ]

(e) Guide Transcription (358.3 seconds)

**Submission**

In the audio, did you hear any voices other than the main annotator?
○ Yes
○ No

How confident are you that you stopped at the correct goal location?

Not confident                                              Very confident

How confident are you that you stayed exactly on the path described in the instructions?

Not confident                                              Very confident

[ Submit ]

(f) Follower Survey (15.2 seconds)

Figure 6: Screenshots of the Guide (a, c, e) and Follower (b, d, f) views in our annotation tool, and the average duration for each phase during collection of the first 33K instructions.

## Guide Alignment
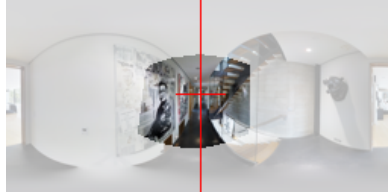### ordered left-to-right →

You're starting in a closet, facing an abstract painting on your right. Just slightly to your left will be an open, wooden door next to an amp. Walk through that wooden door.
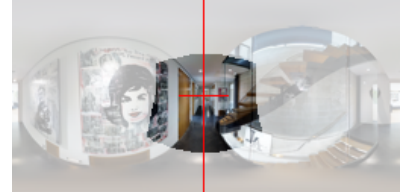
This will take you to a hallway with stairs going up on the right hand side. Just go straight down the hallway...

...about five steps...

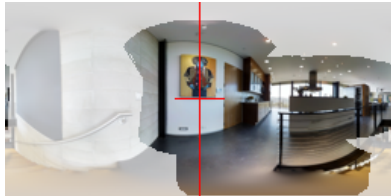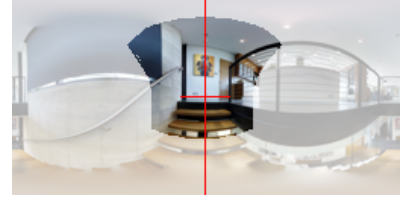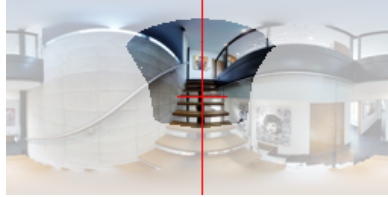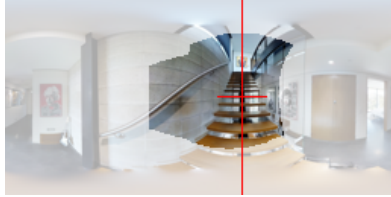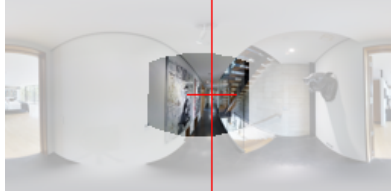...steps. You're going to pass the...

...the stairs.

Go one...

...one step past the stairs. You'll just pass the Albert Einstein painting on your right, and an open doorway on your left.

There will be a guitar on the floor. At this point, turn around and go up the stairs.

Once you get to the Jimi Hendrix painting, turn...

...turn to your right and walk between the stair railing and the white kitchen cabinet toward...

...toward the refrigerator. Take a step in front of the refrigerator.

Take another step toward the windows overlooking the trees.

Then take a right at the end of the refrigerator. You'll take three steps...

...toward the fireplace.

Once you get...

...get to the fireplace, it will be on your right hand side. This is where you stop.

Figure 7: Spatiotemporal alignment of textual instructions, visual percepts and actions for a long (19-step) en-US **Guide** path. The next action is indicated in red and unseen pixels in the panoramic images are faded.

**Follower Alignment**          ordered left-to-right →

You're starting in a closet, facing an abstract painting on your right. Just slightly to your left will be an open, wooden door next to an amp. Walk through that wooden door.
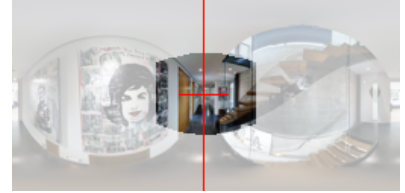
This will take you to a hallway with stairs going up on the right hand side. Just...

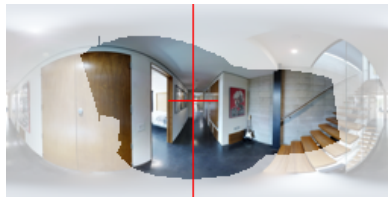Just go straight down the hallway...

...about five steps.

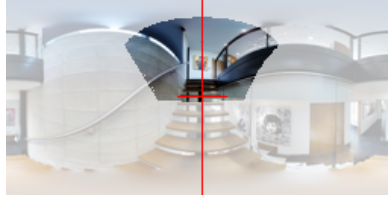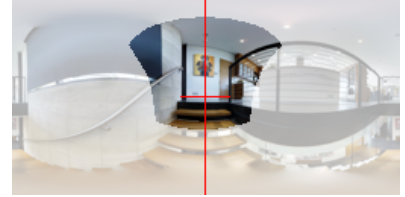You're going to pass the stairs.

...stairs.

Go one step past the stairs. You'll just pass the Albert Einstein painting on your right, and an open doorway...

...doorway on your left. There will be a guitar on the floor. At this point, turn around and go up the stairs.

Once you get to the...

Jimi Hendrix painting, turn to your right and...

...and walk between the stair railing and the white kitchen cabinet toward the refrigerator.

Take a step in front of the refrigerator.

Take another step toward the windows...

...windows overlooking the trees.

Then take a right at the end of the refrigerator. You'll take three steps toward the fireplace.

...fireplace. Once you get to the fireplace, it will be on your right hand side.

...side. This is where you stop.

Figure 8: Spatiotemporal alignment of textual instructions, visual percepts and actions for a long en-US **Follower** path. The next action is indicated in red and unseen pixels in the panoramic images are faded.

You begin in a large wooden room with a dining table, an immense fireplace, and a lovely carpet. turn to your right, and move along the edge of that carpet you're nearest to, towards the wooden doorway into another interior room. You should see a large circular table with an urn in the center of it when you enter that room. Skirt the edge of that table to the left, moving towards the staircase. Don't go to the staircase, but instead proceed to the left of it, down the large rectangular rug. Continue through the open glass door,, and the second glass door across the small hallway from it. Step inside this small... Dining area? If you are just inside the room with the circular table in the middle of it, a couch on the left hand wall, two armchairs across from the entrance, and one armchair, striped, just next to you, you're in the right place, and you are done.

Starting facing a large ornate vase with gold leaf on it as well as a curtained window, we are going to turn towards the dinning room table we are face. We are going to hop around it and come to just beside the painting in the background. Once we're behind the head of the table chair, we are going to face forward and notice that there is a marble staircase before us. Head towards that, but don't head up the stairs and don't exit the room, instead we're going to turn to the left and you should see a kitchen before you. Let's go ahead and enter the kitchen through the archway, and here walk to the right of the China cabinet, and towards the island with the dark cabinets and the granite countertop. Once we've turned the corner, and we're beside the large gas range and the stainless steel hood, we're going to walk between the stove and the kitchen island, towards the refrigerator, and you should see an open doorway before you, to the right of the fridge. Go ahead and walk towards this open door and through it. Walk all the way down and turn to the right, passed the closed door, until you're faced with another flight of stairs. Let's go ahead and move up them. When you've ascended the stairs, turn and face your right, and walk towards the music room that we can see in the distance with its grande piano. We are going to come to a stop right at the base of another small flight of stairs, and looking into the sitting room with a grande piano and marble mantle over a fireplace.

You are beside the bed in your bed room, turn towards your left and keep moving forward. Go near the stair case support and turn towards your right, keep moving forwards and you can find a long corridor on your left. Go through the corridor and the opposite end you can find a gaming room. Go through that gaming room and opposite end of a room, towards your slight right, you can find a air hockey table. Go and stand near that table and you reached your destination.

You are facing towards the white door. Turn left and walk towards the swimming pool. Turn left and walk towards the gym equipment. Turn right. Walk a few steps ahead and stand beside the swimming pool. There is a window towards your left side. You have reached your point.

now you are on a stair case facing the stairs, climb up the stair case, now you will enter a big hall, now walk to the other end of the hall and now you will see two doors which are wide opened, exit through the doors and take a right turn and walk on the corridor, to the send window from the right is your destination.

Right now you're facing towards a curtain. Now turn behind and move towards the wall which is in front of you. Now turn left and exit the room, there are portraits to your left. Now turn right and move forward in the walkway. You can see an open door to your left, move towards the door and turn left. Now enter in to the room, there are two washing machines to your right and you can see shelves in front of you, move towards the shelves and stand in front of it and it is your end point.

You are in a living area, facing towards the corner of a door. Turn towards your slight right and keep moving forward. In front, towards your slight right, you find an other section. Go near that section and turn towards your left. You find a brown door, go pass through the door and move forward. You enter into your bedroom. In front, you find a bed, walk towards the bed and stand near it. You reached your destination.

Right now you're facing towards a bed. Now slightly turn right, there is an open door in front of you, move towards the door and exit the room. There is a walkway in front of you and some portraits on the wall to your right and a staircase to your left, move forward in the walkway, continue moving forward in the walkway, until you reach an open door in front of you, there is an open door to your right, move towards the door and turn right. Now enter in to the room, there is a portrait in between two windows in front of you. Now slightly turn left, there is a sliding door in front of you, which guides to the balcony, move towards the door and enter in to the balcony. Now turn left, you can see a sliding door in front of you, move towards the sliding door and enter in to the room and this is your end point.

Figure 9: Randomly selected English navigation instructions from RxR train. The first two examples are US English and the others are Indian English.

हलका सा बाए मुड़े और सीधे आगे बड़े सामने आपको मेज़ और कुर्सी नज़र आ रही होगी हलका सा दाहिने ले और और आखरी वाली कुर्सी के पीछे तक जाए फिर हलका सा बाए ले और सीधे आगे बड़े दीवार के कोने तक फिर हलका सा दाहिने ले सामने आपको काले रंग का सोफा नज़र आ रहा होगा सीधे सोफे के पास जाए और रुक जाए।

आप अभी सीढ़ी के बाज़ू में हैं,आप को सीधे आगे बढ़ना है,सामने मेज़ के बाज़ू से बाएं मुड़ना है,और सामने कमान में प्रवेश करना है,बिलकुल सामने भूरे रंग का बंद द्वार है,आप को वहीं पे जाकर रुक जाना है.

आप अभी फ्रिज के पास खड़े हुवे हैं आप को हल्का सा बाएं मुड़ना है,और सामने आगे बढ़ना है,आगे बढ़ने के बाद छोटी सी अलमारी है जिसपर किताबें रखी हुई हैं, उसके बाज़ू में दीवार पर एक चित्र लटका हुवा है,और उसके नीचे एक बैठक भी है,बैठक तक जाने के बाद आप को पीछे मुड़ना है,पीछे आप को किताबों की एल्मा ॠ दिखाई देगी उस के बाज़ू में एक द्वार है आप को उस द्वार तक जाना है,फिर उस द्वार से वापस आप को दाहिने की तरफ बाहर निकलना है,जहाँ पर गोल मेज़ है उस पर एक गुलदस्ता है,वहां पर दो कुर्सियाँ भी हैं, आप को वहां पर आ करके रुक जाना है.

आप भूरे रंग के दिवार के ओर मुड़कर खड़े हुए है। वहा से हल्का सा दाए ओर मुड़कर एक कदम आगे बढे। और वहा से बाए ओर मुड़े और सामने आपको भूरे रंग का खुला दरवाजा नज़र आएगा। उस से कमरे के बाहर जाकर दाए ओर मुड़े और सीधा आगे बढे। सामने आपको कांच का बंद दरवाजा नज़र आएगा।आप वहा से दाए ओर मुड़कर भूरे रंग के खुले दरवाजे में प्रवेश करके सामने आपको एक कुर्सी नज़र आएगी। उस कुर्सी के बाज़ू रुक जाए। आपके बाए ओर भूरे रंग का कांच का खुला दरवाजा होगा। और आपके सामने सफेद रंग का पलंग होगा।

पीछे मुड़े और द्वार के अंदर प्रवेश करे हल्का बाए मुड़े और गलियारे मे आगे बढे आपके दोनों तरफ भूरे रंग की अलमारी है आगे बढे और रुक जाना है आपकी बाए तरफ अलमारी में कपडे है.

आप बाहर खड़े हो। पीछे मुड़कर खुले दरवाज़े से अंदर जाए। अन्दर जाते ही आपके बाए तरफ भूरे रंग का खुला दरवाज़ा हैं। आपको उस दरवाज़े से अंदर जाना हैं। अंदर जाते ही आपके सामने बिस्तर के बाजु भूरे रंग के मेज़ पे दीपक हैं। आपको उस दीपक के पास जाना हैं। वह जाते ही आपके बाए तरफ भूरे रंग का खुला दरवाज़ा हैं। आपको उस दरवाज़े से अंदर जाना हैं। वह जाते ही आपको दाहिने मुड़ना हैं। दाहिने मुड़ते ही आपके सामने गलियारा हैं। आपको गलियारे मे सीधा चलना हैं। सीधा चलते हुवे आपके अंत मे भूरे रंग का खुला दरवाज़ा दिखाए देगा। आपको उस दरवाज़े से अंदर जाना हैं। अंदर जाते ही आपके बाए तरफ बटी हैं। आपको बटी के सामने जाना हैं। वह जते ही पको दाहिने मुड़ना हैं। दाहिने मुड़ते ही आपके सामने गोल आकर के मेज़ के साथ दो सफेद रंग के सोफे हैं। आपको बाए वाले सोफे के पास जाना हैं। वह जाते ही आपको दाहिने मुड़कर सफ़ेद रंग के फूलदान के पास जाकर रुक जाना हैं।

आप सामने वाले छोटे भूरे रंग के अलमारी के ओर मुड़कर खड़े हुए है। वहा से दाए ओर मुड़े और सीधा आगे बढे। आपके सामने सफ़ेद रंग का कुर्सी होगा। उस तक जाकर वहा से बाए ओर मुड़े और सामने वाले दिवार जिसपर एक चित्र है,उस तक जाए। वहा से बाए ओर मुड़े और सामने वाले खुले सफ़ेद रंग के दरवाजे में प्रवेश करके फिरसे उल्टा मुड़कर बाए ओर वाले खुले सफ़ेद रंग के दरवाजे में प्रवेश करके फिरसे बाए ओर मुड़े और कमोड के सामने रुक जाए। आपके सामने स्नान घर होगा।

पीछे मुड़े, सीढ़ियों से नीचे उतरे, नीचे जाए, दाए लेकर मेज़ के दाए तरफ जाए, यहाँ से सामने की ओर जाए, दाए लेकर आगे जाए, सामने की ओर जाए, बाए ले और आगे जाए, आपको यहाँ पर मेज़ के पास बाए लेना है, आप इस वक़्त एक भूरे रंग की कुर्सी के सामने खड़े है, आपके पीछे किताबें है, आपको यही पर कोने में इस भूरे रंग की कुर्सी के सामने जिसके पीछे किताबें है वही रुक जाना है, आप देख सकते है आपको दाए तरफ भी एक कुर्सी दिख रही है आपको यहाँ पर बाए तरफ़ कुर्सी के सामने रुकना है।

आपके सामने स्नान कक्ष है। पीछे मुड़कर सफ़ेद खुले दरवाज़े की चौखट पर जाए। सामने आपको सीढ़िया नज़र आएँगी। सीढ़ियों के दाए और से आगे बड़े। सामने आपको तीन सीढ़िया नज़र आएँगी। नीचे उतरे। गलियारे में सीधा आगे बड़े। दाए तरफ के दो काले रंग के कुर्सियों को पार करते हुए आगे बड़े। हलके दाए मुड़कर दाए तरफ के सफ़ेद खुले दरवाज़े की चौखट पर जाए। आगे बढ़कर कमरे में प्रवेश करे। बाए तरफ के पलंग को पार करते हुए कांच की खिड़की के सामने जाकर रुक जाए। बाए तरफ आपको एक गोल मेज़ नज़र आएगा। जिस पर संगणक है। दाए तरफ हरे रंग का सोफा है। जिस पर दो तकिये है।

आप एक सफ़ेद गोल मेज़ के पास खड़े हैं। दाएं मुड़िये। सामने गलीचे के अंत तक जाइये। फिर हल्का बाएं मुड़िये और सामने दिख रहे सीढ़ी से ३ कदम निचे उतर कर रुक जाइये।

आपके आगे दिवार होगी वह से बाए मुड़े। आपके आगे सोफा होगा।थोड़ा सा आगे प्रवेश करे।आपके आगे दूरदर्शन यंत्र होगा। वह से बाए मुड़े आगे प्रवेश करे। आपके बाए में सीढ़ियां होगी।आप उस सीढ़ियां से पूरा ऊपर आकर खड़े होजाइये।आपके आगे चित्र होगा।

पीछे पलट कर भूरे रंग के दरवाज़े से बाहर निकले फिर दाए मुड़े एक कदम आगे जाए फिर दाए मुड़े भूरे रंग के दरवाज़े से अंदर जाए और रुक जाए।

Figure 10: Randomly selected Hindi navigation instructions from RxR train.

మీరు ఉన్న దిక్కు నుండి కొంచెం ఎడమవైపు తిరిగి, నేరుగా ద్వారం బయటకి వెళ్ళి, వెంటనే ఎడమవైపు తిరిగి నేరుగా నడుస్తూ వెళ్ళి, వెంటనే ఎడమవైపు తిరిగి నేరుగా ముందుకు నడిచి ఆగండి. మీకు ఎదురుగా కొంచెం దూరంలో మూసి ఉన్న ద్వారం ఉంటుంది.

ఇప్పుడు మనం ఉన్న దగ్గర నుండి కొంచెం ముందుకు వెళ్ళండి, అక్కడ కుడిపక్కకు తిరిగి కొంచెం ముందుకు వెళ్తే ఒక గది నుండి బయటకి వెళ్ళాము, అక్కడ ఎడమపక్కకు తిరిగి కొంచెం మునుకు వెళ్ళండి, అక్కడ నుండి ఇంకా కొంచెం ముందుకు వెళ్ళండి, అక్కడ నుండి ఇంకా కొంచెం ముందుకు వెళ్ళండి, అక్కడ నుండి ఇంకా కొంచెం ముందుకు వెళ్ళండి, అక్కడ నుండి ఇంకా కొంచెం ముందుకు వెళ్తే ఒక గది నుండి బయటికి వెళ్ళాము, అక్కడ నుండి కుడిపాకకు తిరిగి కొంచెం ముందుకు వెళ్ళి ఒక మూడు మెట్లు ఎక్కండి, అక్కడ నుండి కుడిపక్కకు తిరిగి ఒక ఐదు మెట్లు ఎక్కండి, అక్కడ నుండి ఒక రెండు మెట్లు ఎక్కి, కుడిపక్కకు తిరిగి ఒక నాలుగు మెట్లు ఎక్కండి, అక్కడ నుండి ఇంకొక రెండు మెట్లు ఎక్కి, కొంచెం ముందుకు వెళ్ళండి, అక్కడ ఎడమపక్కకు తిరిగి కొంచెం ముందుకు వెళ్ళండి, అక్కడ కొంచెం కుడిపక్కకుగా తిరిగి కొంచెం ముందుకు వెళ్ళండి, ఇంకా కొంచెం ముందుకు వెళ్ళి ఒక గది ముందు ఆగండి.

కుడివైపుకు తిరిగి, ఎదురుగా ఉన్న తివాచీ మీద నడుచుకుంటూ నేరుగా ముందుకు వెళ్ళండి. ఇక్కడ ఎదురుగా ఉన్న ద్వారము గుండా నేరుగా ముందుకు వెళ్ళండి. ఎదురుగా ఉన్న మరొక ద్వారము గుండా బయటకు వెళ్ళండి. ఇక్కడ ఎదురుగా మెట్లు ఉన్నాయి. వాటిని ఎక్కుతూ పైకి వెళ్ళండి. ఎదురుగా ఉన్న ద్వారము గుండా గది లోపలికి వెళ్ళిన వెంటనే, ఎడమవైపుకు తిరిగి, నేరుగా ముందుకు వెళ్ళండి. ఎదురుగా ఉన్న చొక్కాలు అల్మారా ముందు కుడివైపుకు తిరిగి, కొంచెము ముందుకు వెళ్ళి, కుడివైపు ఉన్న చొక్కాలు అల్మారా ముందు ఆగండి.

ఇప్పుడు మీరు ద్వరంకి ఎదురుగా ఉన్నారు. కాస్త కుడివైపుకి జరిగి, ముందుకు వెళ్ళి ద్వారం నుండి లోపలికి వెళ్ళండి. కాస్త లోపలికి వెళ్ళి భోజన బల్ల మరియు మీ ఎడమవైపు ద్వారం మధ్యలో ఆగండి.

మీరు నిల్చున్న చోటు నుండి కుడివైపు తిరిగి, మంచం ముందుకి వెళ్ళి, కుడివైపు తెరిచి ఉన్న ద్వారం దగ్గర ఆగండి.

మీరు గాజు ద్వారంపైపు ఉన్నారు. వెనుకకి తిరిగి, తెరిచి ఉన్న ద్వారంలోంచి బయటికి వెళ్ళండి. మీ ఎదురుగా చిత్రపటం ఉంటుంది. కుడివైపు ఉన్న ద్వారంలోంచి బయటికి వెళ్ళి, ఎదురుగా ఉన్న తివాచీపైకి వెళ్ళండి. కాస్త ఎడమవైపుగా ఉన్న మెట్లను పూర్తిగా ఎక్కండి. కాస్త కుడివైపుకి తిరిగి, ముందుకు వెళ్ళండి. ఎడమవైపు ఉన్న ద్వారంలోంచి లోపలికి వెళ్తే, మీ ముందు మంచం ఉంటుంది. కాస్త కుడివైపుకి తిరిగి, మంచం ముందుకు వెళ్ళి ఆగండి. మీ ఎడమవైపు మంచం, కుడివైపు సోఫా ఉంటాయి.

మీరు ఉన్న దిక్కు నుండి నేరుగా ఉన్న గదిలోనికి వెళ్ళి, తివాచీని దాటుకొని వెళ్ళి, ఎదురుగా తెరిచి ఉన్న ద్వారం ముందు ఆగండి.

మీరు దూరదర్శిని వైపు ముఖం చేసి ఉన్నారు. అక్కడ నుండి కుడివైపుకు తిరిగి నేరుగా ద్వారం నుండి బయటకి వెళ్ళండి. ఇప్పుడు మీకు ఎడమవైపున చిత్రపట్టం కుడివైపున ఒక్క ద్వారం ఉంటుంది. ఆ ద్వారం లోపటికి వెళ్ళి ఆగండి. అదే మీ గమ్యం.

ఉన్న చోట నుంచి కుడివైపుకు తిరిగి, మరల ఎడమవైపు తిరిగి మెట్లు దిగండి. అలా మెట్లు దిగాక, మీకు ఎడమ పక్కకు కనిపిస్తున్న దారిగుండా వెళ్ళండి . ఇప్పుడు మీ ఎడమ పక్క ఒక చెట్టు ఉంటుంది. ఆ చెట్టు దగ్గరకు వెళ్ళి కుడివైపుకు తిరగండి. నేరుగా ముందుకు నడవండి. ఎడమ వైపు కనిపిస్తున్నా ఉన్న ద్వారం లోపలి కి వెళ్ళి మరల నేరుగా నడవండి. ఇప్పుడు మీకు చెట్టు కనిపిస్తుంది. ఆ చెట్టు దగ్గరికి వెళ్ళి కుర్చీ వెనక ఆగండి. ఇదే మీ గమ్య స్థానం

మనకు ఎడమవైపున ఒక బల్ల ఉంది. కొంచెం కుడివైపుకు మళ్ళి , కొంచెం ముందుకు నడిచి , అక్కడ మనకు ఒక కుర్చీ మరియు చెత్త వేసే బుట్ట మరియు ఒక బల్ల కనపడుతుంది. అక్కడ నుంచి కొంచెం ముందుకు నడిచి , మనకి ఎడమవైపున ఒక ఆట ఆడే బల్ల కనపడుతుంది . దాని వెనుక వైపు నుంచి కొంచెం దూరం నడిచి , ఎడమవైపు మళ్ళాలి. అక్కడ మనకు రెండు రకాల బల్లలు కనపడతాయి . కొంచెం ముందుకు నడిచి , ఎడమకు మళ్ళి, అక్కడ మనకు కుర్చీలతో కూడిన ఒక బల్ల కనపడుతుంది. ఆ బల్ల వైపున పెనుక నుంచి ముందుకు నడిచి , ఇంకా కొంచెం ముందుకు నడిచి అక్కడ మనకు కుడివైపున ఒక తెల్లటి కుర్చీ మరియు ఒక బల్ల కనపడుతుంది . కుడివైపు మెట్లు కూడా ఉంటాయి. మనము ఆ బల్లదగ్గర ఆగండి .

ఇప్పుడు మీరు ఉన్న దేగ్గర నుండి కొద్దిగా కుడి పక్కకి తిరిగి ఎదురుగా తెరిచి ఉన్న గాజు ద్వారము నుండి ట్రైటికి రావాలి. ఇప్పుడు మీరు ఉన్న దేగ్గర నుండి చూస్తే మీ కుడి వైపు మళ్ళి తెరిచి ఉన్న గాజు ద్వారము కనిపిస్తుంది ఇదే మీ గమ్యం.

మీకు ఎదురుగా ఉన్న ద్వారంలో నుంచి నేరుగా వెళ్ళి. అక్కడ వెనుక తిరిగి, కుడివైపుగా నేరుగా వెళ్ళి. అక్కడి నుంచి నేరుగా వెళ్ళి, ఎడమవైపుగా తిరిగి నేరుగా వెళ్తే. మీరు ఒక్క ద్వారాన్ని చేరుకుంటారు. అక్కడి నుంచి నేరుగా వెళ్ళి, అక్కడ స్నానాల తొట్టిలో ఆగండి.

నిల్చున్న చోటు నుండి ఎడమవైపు తిరిగి, కుడివైపు తెరిచి ఉన్న ద్వారం దగ్గరికి వెళ్ళి, కుడివైపు తిరిగి, ఎడమవైపు ఉన్న సోఫా పక్కకు వెళ్ళి, ఎదురుగా ఉన్న మరో ద్వారం దగ్గరికి వెళ్ళి, ఎదురుగా మూసి ఉన్న గాజు ద్వారం ముందర ఆగండి.

Figure 11: Randomly selected Telugu navigation instructions from RxR train.

# B DATASHEET: ROOM-ACROSS-ROOM (RxR)

This document is based on *Datasheets for Datasets* (Gebru et al., 2020). Please see the most updated version here.

## MOTIVATION

### For what purpose was the dataset created?

RxR was created to advance progress on vision-and-language navigation (VLN) in multiple languages (English, Hindi, Telugu). It addresses gaps in existing datasets by including more paths that counter known biases and an order of magnitude more navigation instructions for three languages plus annotators' 3D virtual pose sequences.

### Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?

This dataset was created by Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, Jason Baldridge and the Google Data Compute team on behalf of Google Research.

### What support was needed to make this dataset?

Funding was provided by Google Research.

## COMPOSITION

### What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?

The instances in RxR are natural language navigation instructions paired with trajectories in reconstructed 3D buildings. Each navigation instruction has been recorded as speech and transcribed by the speaker. The dataset includes the text transcriptions, but not the audio files, although they may be released in future. The trajectories are provided as paths, consisting of sequences of viewpoint ids corresponding to navigation graphs from Anderson et al. (2018b), and pose traces, consisting of sequences of virtual camera poses situated in the underlying building reconstructions which are from the Matterport3D dataset (Chang et al., 2017). Pose traces and text transcriptions are timestamped and aligned. Pose traces are provided for both the instruction

annotator (the Guide), and a second annotator charged with following the Guide's instructions (the Follower).

### How many instances are there in total (of each type, if appropriate)?

RxR contains 126K Guide instructions covering 16.5K sampled paths and 126K human Follower demonstration paths. Annotations are split equally across the three languages in the dataset. Refer to Table 1 for a comparison of the number of instances to previous datasets and Table 2 for summary statistics.

### Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?

Refer to Section 3 for a detailed description of the sampling procedure used to select the paths for annotation.

### What data does each instance consist of?

Each instance consists of a trajectory through a building from the Matterport3D dataset (Chang et al., 2017) paired with a natural language navigation instruction. A trajectory can be visualized as a sequence of 360-degree panoramic images, or as path traversing a 3D reconstruction of the building represented as a textured mesh. Refer to Table 3 for an analysis of linguistic phenomena in the instructions and Figures 9, 10 and 11 for instruction examples in English, Hindi and Telugu respectively.

### Is there a label or target associated with each instance?

When training wayfinding agents to navigate from natural language instructions, the trajectory is the target. Instructions and paths are annotated with unique identifiers.

### Is any information missing from individual instances?

We do not provide the Guide audio recordings, for reasons outlined in Appendix A.

### Are relationships between individual instances made explicit (e.g., users' movie ratings, social network links)?

Trajectories may belong to the same building or different buildings; each instance is annotated with

a scan (building) identifier.

**Are there recommended data splits (e.g., training, development/validation, testing)?**
Yes. We follow the same building splits as Matterport3D and R2R. Refer to Section 3 for details regarding the RxR train/validation/test instance splits.

**Are there any errors, sources of noise, or redundancies in the dataset?**
The process we followed to validate instruction quality using Follower annotations is described in Section 4.

**Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)?**
This dataset is based on building reconstructions from the Matterport3D dataset (Chang et al., 2017) and viewpoint navigation graphs from the R2R dataset (Anderson et al., 2018b). Apart from these dependencies, RxR is self-contained, i.e., it does not rely on web resources.

**Does the dataset contain data that might be considered confidential (e.g., data that is protected by legal privilege or by doctor-patient confidentiality, data that includes the content of individuals' non-public communications)?**
No.

**Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?**
No.

**Does the dataset identify any subpopulations (e.g., by age, gender)?**
No.

**Is it possible to identify individuals (i.e., one or more natural persons), either directly or indirectly (i.e., in combination with other data) from the dataset?**
No.

**Does the dataset contain data that might be considered sensitive in any way (e.g., data that reveals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or union memberships, or locations; financial or health data; biometric or genetic data; forms of government identification, such as social security numbers; criminal history)?**
Each natural language instruction in the sample is either in English, Hindi or Telugu, thus potentially revealing linguistic origin. However, no other annotator data is included in the dataset.

## COLLECTION

**How was the data associated with each instance acquired?**
Refer to Section 4 for details of the annotation procedure, as well as measures undertaken to validate the data.

**Over what timeframe was the data collected?**
The dataset was collected between March 2020 and September 2020.

**What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual human curation, software program, software API)?**
We developed a web-based annotation tool to collect the data. It is described further in Section 4 and screenshots are included in Figure 6.

**If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)?**
Please see Section 3 and Figure 2 for details of the strategy for selecting paths for annotation.

**Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)?**
Refer to Appendix A.

**Does the dataset relate to people?**
Yes.

**Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (e.g., websites)?**
Directly from the individuals.

**Were the individuals in question notified about the data collection?**
Yes.

**Did the individuals in question consent to the collection and use of their data?**
Yes.

## PREPROCESSING / CLEANING / LABELING

**Was any preprocessing/cleaning/labeling of the data done(e.g.,discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)?**
Please see Section 4 describing the process used to remove and re-annotate instances in which the Follower was not able to correctly follow the path described in the Guide's instruction.

**Was the "raw" data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)?**
Yes.

**Is the software used to preprocess/clean/label the instances available?**

We plan to publicly release our web-based annotation tool.

## USES

**Has the dataset been used for any tasks already?**
We have used RxR to train vision-and-language navigation (VLN) agents as described in the paper.

**Is there a repository that links to any or all papers or systems that use the dataset?**
No, although we plan to release a test server and leaderboard to support the research community using the dataset.

**What (other) tasks could the dataset be used for?**
Training models to generate natural language navigation instructions, visual referring expression grounding and comprehension, grounded dialog tasks, pre-training for various other vision-and-language tasks, multilingual learning and so on.

**Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?**
No.

## DISTRIBUTION

**Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?**
Yes, this dataset is open to use by the research community.

**How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?**
Via GitHub and Google Cloud Storage.

**When will the dataset be distributed?**
This dataset has been distributed on publication.

**Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)?**
RxR is released under a CC-BY license.

**Have any third parties imposed IP-based or other restrictions on the data associated with the instances?**
Yes, the Matterport3D dataset is governed by the Matterport3D Terms of Use.

## MAINTENANCE

**How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**
Email contact: `rxrvln@google.com`.

**Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?**
No.