VLN UBERT: A Unified Object and Scene Aware Transformer for Vision-and-Language Navigation

Anonymous Author(s) Affiliation Address email

Abstract

1	Natural language instructions for visual navigation use <i>object references</i> and <i>scene</i>
2	descriptions to provide a breadcrumb trail to a goal location. This work presents a
3	multimodal transformer for Vision-and-Language Navigation (VLN) that processes
4	visual observations using both an object detector and a scene classification network,
5	which produce features that mirror these two distinct types of visual cues. In
6	our approach, scene features provide high-level contextual information to support
7	object-level processing. With this design, our model is able to use vision-and-
8	language pretraining – i.e., learning the alignment between images and text from
9	large-scale web data - to substantially improve performance over a strong baseline
10	on the Room-to-Room (R2R) [1] and Room-Across-Room (RxR) [2] benchmarks.

Dummy Teaser Label

Figure 1: Dummy Teaser caption

11 **1 Introduction**

The vision-and-language navigation (VLN) task [1] requires an agent to find a goal location within 12 an environment by following natural language navigation instructions. A central component of this 13 task is associating the instruction with visual landmarks in the environment. Figure 1 provides an 14 illustrative example from the Room-to-Room (R2R) dataset [1] in which the agent needs to follow 15 the indoor navigation instructions: "Exit the bedroom and turn left. Continue down the hall and 16 into the room straight ahead and stop before the desk with two green chairs." The visual landmarks 17 in this instruction range from scene-level descriptions (e.g., "bedroom" and "hall") to specific 18 19 object references (e.g., "desk" and "two green chairs"). Accordingly, VLN agents should be able to 20 recognize this diverse range of visual cues.

Submitted to 35th Conference on Neural Information Processing Systems (NeurIPS 2021). Do not distribute.

In most VLN methods [1, 3-12], visual observations are first encoded with a convolutional network 21 that was trained to solve an image classification task – either using ImageNet [13] or the Places [14] 22 scene recognition dataset. While ImageNet features may identify objects mentioned in the instructions 23 and Places features might match the scene descriptions, neither solution was explicitly trained to 24 recognize both types of visual cues. Another limitation of standard VLN methods is that the alignment 25 (or grounding) between instructions and visual features is typically learned using only the limited 26 supervision provided by VLN datasets (e.g., the 14k or 42k path-instruction pairs from the R2R [1] 27 and RxR [2] datasets, respectively). Recent work has overcome this challenge by transferring visual 28 grounding from large-scale web corpora to VLN models operating over either scene [15] or object [16] 29 features. We advance this line of research by developing a VLN model that transfers visual grounding 30 to make effective use of both scene and object features. 31

Our proposed approach extends the VLN CBERT framework [15] by incorporating object features 32 as an input to the model and for action prediction. We find that naïvely using both scene and object 33 features does not significantly improve VLN performance over a VLN OBERT baseline. Thus, we 34 propose architectural changes that allow the model to take better advantage of these two distinct types 35 of visual information. Specifically, our model treats scene features as contextual information that is 36 used for object-level processing. This design is similar to the object-level processing paradigm that 37 has been adopted by several recent vision-and-language pretraining (VLP) methods (e.g., [17, 18]), 38 which learn visual grounding by pretraining on large image-text datasets collected from the web. 39 Accordingly, unlike VLN BERT, our model consistently uses object features during vision-and-40 language pretraining and VLN task-specific fine-tuning. 41

We experiment with our proposed approach on the Room-to-Room (R2R) [1] and Room-Across-Room 42 (RxR) [2] datasets. Empirically we find that our model substantially improves VLN metrics over a 43 strong baseline on R2R and outperforms state-of-the-art methods on English language instructions 44 in RxR. Concretely, our proposed approach improves Success Rate (SR) on the unseen validation 45 split in the R2R dataset by 1.4 absolute percentage points. On RxR, which is a more challenging 46 dataset due to indirect paths and larger variations in path length, we see even larger improvements 47 in SR of 3.7 absolute percentage points, alongside a gain of 1.9 absolute percentage points on the 48 Normalized Dynamic Time Warping (NDTW) metric. Through ablation experiments we discover 49 that vision-and-language pretraining is vital to our approach, suggesting that strong visual grounding 50 is key for using object features in VLN. 51

- ⁵² To summarize, we make the following contributions:
- We propose a unified multi-modal transformer model for vision-and-language navigation that
 leverages vision-and-language pretraining for jointly processing scene and object features.
- We show that our approach outperforms strong baselines on the R2R and RxR datasets. [SL]
 Would be nice to make these a little punchier and include key results.
- We provide analysis that demonstrates the effectiveness of our design decisions and highlights the importance of vision-and-language pretraining to the success of our model.

59 2 Related Work

Vision-and-Language Navigation. The Room-to-Room (R2R) [1] and Room-Across-Room 60 (RxR) [2] datasets both situate the VLN task within Matterport3D [19] indoor environments. Since 61 the release of R2R there has been steady improvement in VLN task performance [3–12, 15]. Some of 62 the key innovations include using instruction-generation via "speaker" models for data augmenta-63 tion [3, 8], combining imitation and reinforcement learning [4], using auxiliary losses [5, 11], and 64 different pretraining strategies [10, 12, 15]. All of these methods have one thing in common - they 65 process visual observations with a convolutional network pretrained to solve an image classification 66 task (using either the ImageNet [13] or Places [14] datasets). In contrast, this work explores using a 67 combination of features from visual encoders pretrained for scene classification and object detection. 68

Object Features for VLN. Intuitively, object detections should naturally match the object cues
 mentioned in VLN instructions. Indeed, several recent studies [16, 20–22] have demonstrated the
 utility of using object detectors for VLN. In [20, 21] detections are converted into object "*tags*" (i.e.,
 object classification labels), which are encoded using a GLoVe [23] embedding. Similarly, [22]

convert detections into a feature vector using the classification label, object area, and detector
 confidence. Unlike these methods, our approach directly uses object features produced by a detector
 and takes advantage of vision-and-language pretraining to transfer visual grounding. In [16], object
 features are used in a model that solves a path selection task in VLN, which requires pre-exploring an
 environment before executing the navigation task. By comparison, this work focuses on combining

⁷⁸ scene and object features to solve instruction-guided visual navigation without pre-exploration.

79 **3** Preliminaries

80 [AM] TODO: Update this paragraph.

In this section, we review the Vision-and-Language Navigation (VLN) task and the recently proposed VLNOBERT [15] model for this task. We describe our approach that builds on top of the VLNOBERT architecture in Section 4.

84 3.1 Vision-and-Language Navigation

In VLN, an agent is placed in a photo-realistic 3D environment and must navigate to a goal location that is specified through natural language navigation instructions I (as illustrated in Figure 1). At each timestep t, the agent receives a panoramic observation $O_t = \{o_{t,i}\}_{i=1}^{36}$ (introduced by Fried et al. [3]), which is composed of RGB images from 36 viewing angles (12 headings \times 3 elevations). In this work we focus on the "*nav-graph*" setting, in which the agent has access to a navigation graph that specifies a discrete set of navigable locations for each viewpoint in the environment. Using this information the agent selects an action from the set $A_t = \{a_{t,i}\}_{i=0}^{N_t}$ consisting of N_t navigable locations and the stop action. The agent is successful if it calls stop within 3m of the goal location.

93 3.2 Visual Encoders for VLN

⁹⁴ The goal of this work is to develop a VLN model that can process both scene and object features

95 (as discussed in Section 1). This is accomplished by processing each RGB image $o_{t,i}$ with two

⁹⁶ different visual encoders that are adopted from prior work. The feature extraction methods used by

97 VLNOBERT and our model are summarized here; additional details are provided in the Appendix.

Scene Features. Following [1], scene features are produced using a ResNet-152 [24] CNN that was trained on the Places [14] scene recognition dataset. In VLN, understanding the relative angle to each image $o_{t,i}$ is important for following instructions with directional cues (e.g., "to your left"). Accordingly, the relative heading $\theta_{t,i}$ and elevation $\phi_{t,i}$ to each RGB image is encoded into a 4dimensional vector $[\sin\theta_{t,i}, \cos\theta_{t,i}, \sin\phi_{t,i}, \cos\phi_{t,i}]$ and combined with the CNN features to construct the final scene feature vector $f_{t,i}$.

Object Features. To generate object features we use the approach taken in [16], in which a Faster R-CNN [25] detector trained on Visual Genome [26] using the training procedure from [27] is used to produce a set of M region features for each RGB image. Heading and elevation angles to each image region are encoded (as above) and combined with the region features, resulting in a set of object features { $g_{t,i,j}$ }^M_{j=1} for each RGB image.

109 3.3 Multimodal Transformers

Here we provide a brief overview of multimodal transformers (e.g., OSCAR [17]), which provide a 110 basis for VLNOBERT and our approach. Multimodal transformers are an extension of transformer-111 based language models such as BERT [28] that process paired image and text data. Like their 112 language counterparts, the inputs to these models are a set of tokens, which are processed by a series 113 of transformer encoder layers [29] to output representations for each token. Commonly, language 114 tokens are constructed using the approach taken in BERT [28] and "visual tokens" are generated 115 using region features from an object detector (such as the specific model described above). Additional 116 details can be found in [17], which describes OSCAR – the multimodal transformer used in this work. 117

[[]AM] TODO: add a sentence about how multimodal transformers are used in our work

119 **3.4 Multimodal Transformers for VLN**

VLN requires sequentially following navigation instructions (e.g., "Exit the bedroom ... " then "Con-120 tinue down the hall..." then "stop before the desk with two green chairs."). Accordingly, maintaining 121 a history of the agent's state is helpful for understanding which sub-instruction to follow at each 122 timestep. Traditional VLN agents use recurrent neural networks (e.g., LSTM [30]) to model state 123 history. In contrast, VLN BERT [15] introduces a generic recurrence mechanism that can, in 124 principle, be added to any multimodal transformer model to refashion it for the VLN task. In this 125 work we extend the OSCAR [17] instantiation of VLN OBERT, which allows us to study the effect 126 of vision-and-language pretraining on a VLN agent that processes both scene and object features. 127

VLN \bigcirc BERT [15] uses a multimodal transformer to process a panoramic observation O_t alongside the navigation instructions I to selection an action from A_t . The recurrence mechanism in VLN \bigcirc BERT is operationalized using a state token s_t , which is passed as an input to the multimodal transformer and updated at each timestep using self-attention (as discussed below).

Initialization. The state token is initialized by passing the word tokens from the navigation instructions I along with special [CLS] and [SEP] tokens through the multimodal transformer. The output corresponding to the [CLS] token is used to set s_0 and the outputs for the word tokens (denoted as X) are used as the language representation during navigation.

Visual Tokens. For each panoramic observation $O_t = \{o_{t,i}\}_{i=1}^{36}$, VLN \bigcirc BERT runs a Places [14] CNN on each of the 36 RGB images to extract scene features $\{f_{t,i}\}_{i=1}^{36}$ (see Section 3.2). Since the agent can only navigate to one of N_t navigable viewpoints, only features for images pointing towards these locations are used as visual tokens. To facilitate the stop action an all zeros feature vector is added to the set of visual tokens (denoted as F_t).

Navigation. At each timestep (t > 0), the input to the multimodal transformer is composed of the 141 previous state s_{t-1} , the encoded instruction X, and the visual tokens F_t . To reduce computational 142 complexity, the instruction tokens are not updated during navigation and only serve as keys within the 143 attention-based processing. Actions are selected using attention scores for the visual tokens from the 144 last layer of the multimodal transformer model. Specifically, attention scores for the visual tokens are 145 calculated with respects to the state token (i.e., the state token serves as the query and visual tokens 146 serve as keys). [AM] Should we add an equation describing this? The scores are passed through a 147 softmax function and treated as the probability of moving to that navigable viewpoint or stopping. 148 Finally, the state history is maintained by using the output representation for the state token s_{t-1} as 149 the state token input at the next timestep (after a "refinement" procedure described in [15]). 150

151 4 Approach

One key limitation of VLN OBERT is that it only uses high-level scene features, which might fail to encode specific object references such as "*desk*" or "*two green chairs*" (see Figure 1). Second, there is a mismatch between the pretraining and fine-tuning paradigm while training the model. OSCAR [17] is pre-trained on Conceptual Captions [31] using image region features from an object detector as input. However, during fine-tuning on the VLN dataset, VLN OBERT simply feeds high-level scene features as input and hopes that the model can adapt to the new feature distribution.

To overcome these two limitations, we present a unified multi-modal transformer model which uses both scene and object features as shown in Figure 2. Interestingly, we find that simply adding objects as additional inputs in VLN OBERT does not help. To address this challenge, we propose a selective attention mechanism that only refines a subset of the input tokens. This in turn shifts the processing in our model closer to the base model (i.e., OSCAR [17]) used for vision-and-language pretraining.

163 4.1 Input Tokens

Our approach expands the set of input tokens used in VLN \bigcirc BERT to include visual tokens for object features.¹ Specifically, we use object features from the Faster R-CNN detector described in Section 3.2. Similar to the filtering done on scene features, only object features corresponding with navigable viewpoints are retained to form the set of object tokens G_t . Unlike scene tokens F_t , each

¹[AM] Add a note about how VLN OBERT using object features.

Dummy Approach Label

Figure 2: Dummy Approach caption

navigable viewpoint may include zero or more object tokens (based on the outputs from the object detector). Thus, we modify the VLN \bigcirc BERT approach for calculating action probabilities (discussed in Section 4.3). At each timestep, inputs to the multimodal transformer consists of four sets of tokens: the previous state s_{t-1} , encoded instruction X, scene tokens F_t , and object tokens G_t .

172 4.2 Selective Object Attention

In initial experiments (detailed in Section 5) we found that simply adding object tokens to the set of inputs used by VLNOBERT did not improve performance. We hypothesize that scene tokens dominate the VLN training, which does not allow the model to properly learn how to use object features. To mitigate this issue, we modify the attention pattern used within the multimodal transformer as illustrated in Figure 2. With our attention pattern, both the language I and scene F_t tokens only serve as keys (and values) during the attention-based processing, while the state token s_t and object tokens G_t operate as queries, keys and values (which is standard in transformer models).

Intuitively, this attention pattern focuses processing on the object tokens, while allowing the scene tokens to providing high-level contextual information to support this object-level processing. Additionally, the scene tokens are still used for action prediction (see section 4.3), which provides the model with flexibility in using both sources of information. We hypothesize that this design leads to more effective transfer learning (demonstrated in Section 5), because the object-level processing more closely matches the pretraining setup, in which the model only receives language and object tokens as input.

187 4.3 Action Prediction

Recall that in VLN \bigcirc BERT actions are selected using attention scores over the scene tokens F_t (see Section 3.4), which correspond to each of the navigable viewpoints and the stop action. These attention scores are also used to update the state token (as described in [15]), so all of the values (not just the maximum) are used by the model. In our approach attention scores are also available for object tokens G_t that correspond to image regions pointing towards each of the navigable viewpoints. However, it is not immediately evident how the scene and object attention scores should be combined.

Through initial experiments, we discovered that aggregating attention scores for each viewpoint by selecting the maximum over the scene and object scores yielded the best results. Intuitively, this approach allows the model to either select a relevant object (which may be mentioned in the instructions such as "*green chairs*") or the full scene (which might match scene descriptions such as "*hallway*") to represent each navigable viewpoint. Accordingly, the model is able to select actions using features that match both types of visual cues.

200 [AM] We tried many alternatives here, should we provide some examples?

			2R		RxR				
	Methods	TL	$\text{NE}\downarrow$	SR \uparrow	SPL ↑	NE↓	SR \uparrow	SDTW \uparrow	NDTW ↑
1 2	Random Human	9.77 -	9.23	16 -	-	9.5 1.32	5.1 90.4	3.8 74.3	27.6 77.7
3 4 5 6 7	RxR baseline[2] EnvDrop [8] PREVALENT [12] VLN OBERT [15] (init. OSCAR) VLN OBERT [15] (init. PREVALENT)	- 10.70 10.19 11.86 12.01	- 5.22 4.71 4.29 3.93	37 52 58 59 63	32 48 53 53 57	10.1 - - -	25.6 - - -	20.3	41.3
8 9	VLNĊBERT [15] † Ours	12.16 12.15	4.40 4.28	58 59	51 53	7.31 6.72	40.5 44.2	50.8 54.6	65.9 67.8

Table 1: R2R and RxR val-unseen results. † indicates reproduced results.

201 5 Experiments

202 5.1 Training and Evaluation

Datasets. We initialize our model with a pre-trained OSCAR model [17] for training. We evaluate our 203 method on R2R [1] and RxR [2] datasets. R2R dataset contains 21,567 instruction-path pairs which 204 are divided into four splits: training (14,025), val-seen (1,020), val-unseen (2,349) and test-unseen 205 (4,173). Val-seen split uses environments from the training split but the instruction-path pairs are 206 novel. Val-unseen and test-unseen tests the generalization capacity of agents on unseen environments 207 with novel instruction-path pairs. For R2R, we also use augmented dataset generated from speaker 208 model as done in prior works [8, 12, 15]. RxR [2] is a recently introduced multi-lingual VLN 209 dataset. It contains total 126K instruction-path pairs in 3 languages (Hindi, English and Telugu) on 210 the same 90 Matterport scenes which were used in R2R. Moreover, it contains English instructions 211 from two regions: India (en-IN) and US (en-US). Since our model is pre-trained on Conceptual 212 Captions dataset which is in English language, we just use English instructions (combining en-IN 213 and en-US instructions) from RxR dataset for training and evaluation. Specifically, we use 26,464 214 English instructions from RxR train split for training and test on 4,551 English instructions from RxR 215 val-unseen split. 216

Evaluation. We follow standard evaluation protocol for R2R and RxR datasets. On R2R, we report (\uparrow indicates higher is better and \downarrow indicates lower is better): Trajectory Length (TL), Navigation Error (NE \downarrow), Success Rate (SR \uparrow) and Success weighted by Path Length (SPL \uparrow). On RxR dataset, in addition to NE and SR metrics, we also report Normalized Dynamic Time Warping (NDTW \uparrow) and Success weighted by normalized Dynamic Time Warping (SDTW \uparrow) metrics which explicitly measure path adherence. We refer the reader to [1, 32, 33] for a detailed explanation of these metrics.

223 5.2 Implementation Details

Our model has been implemented in PyTorch [34] on a single Nvidia TitanX GPU. For both R2R and 224 RxR, we fine-tune our model with behaviour cloning and reinforcement learning objectives. Each 225 minibatch consists of 50% rollouts from behaviour cloning and 50% from reinforcement learning 226 (policy gradient). We train using a constant learning rate of 1e-5 with AdamW optimizer and batch 227 size 16 for 300k iterations. We extract visual features for scene tokens with a pre-trained ResNet-152 228 on Places [14] dataset. For object tokens, we extract object proposals from a pre-trained bottom-up 229 attention model [27] and adopt filtering procedure from [16] to discard redundant proposals. We use 230 the exact same hyperparameters for VLN OBERT baseline and our approach in all the experiments. 231

232 5.3 Results and Analysis

How does our method compare to state-of-the-art methods? We compare our approach with the recently published state-of-the-art approaches in Table 1 on R2R [35] and RxR [2] val-unseen splits. EnvDrop [8] trains an encoder-decoder model [3] with mixture of imitation and reinforcement learning on augmented data (in addition to R2R training split) with back-translated instructions and "dropped out" environments in order to generalize well to unseen environments. PREVALENT [12]

			R2R Test Unseen						
	Methods	TL	NE \downarrow	SR ↑	SPL ↑				
1	Random	9.89	9.79	13	12				
2	Human	11.85	1.61	86	76				
3	EnvDrop [8]	11.66	5.23	51	47				
4	PREVALENT [12]	10.51	5.30	54	51				
5	VLN OBERT [15] (init. OSCAR)	12.34	4.59	57	53				
6	VLN OBERT [15] (init. PREVALENT)	12.35	4.09	63	57				
7	VLNĊBERT [15] †	12.78	4.55	58	52				
8	Ours	12.26	4.49	58	53				

Table 2: R2R test results. † indicates reproduced results.

Table 3: R2R Ablations.

		Object Score		Selective	Val Seen				Val Unseen			
	Models	Features	Aggregation	Attention	TL	$\text{NE}\downarrow$	SR \uparrow	SPL ↑	TL	$NE\downarrow$	SR ↑	SPL ↑
1	Baseline [15]				11.02	3.11	69.93	65.65	12.16	4.40	57.90	51.43
2		\checkmark			11.05	2.89	71.69	67.22	12.25	4.35	57.26	50.96
3		\checkmark	\checkmark		11.12	3.24	70.03	66.59	12.05	4.34	57.85	51.73
4	Ours	\checkmark	\checkmark	\checkmark	11.81	3.63	62.78	58.01	12.15	4.28	58.71	53.24

²³⁸ builds on top of EnvDrop method by pre-training a multi-modal transformer model on back-translated
²³⁹ augmented VLN data. After pre-training, PREVALENT method feeds contextual word embeddings
²⁴⁰ from this pre-trained transformer model into EnvDrop [8] encoder-decoder model for fine-tuning on
²⁴¹ R2R dataset. VLN BERT extends PREVALENT by directly fine-tuning the pre-trained multi-modal
²⁴² transformer model for navigation instead of using an ensemble of encoder-decoder model and a
²⁴³ multi-modal transformer.

RxR is more challenging dataset than R2R since it contains much longer instructions and paths 244 (which are often not shortest) with dense object and scene references. On RxR, our approach (row 6) 245 outperforms the previous state-of-the-art RxR method (row 1) by 26.33% absolute on NDTW and 246 19.03% absolute on Success Rate (SR) metrics. We also train and evaluate VLN OBERT baseline on 247 248 RxR dataset with the same hyperparameters as ours. Our unified approach VLNUBERT which uses both object and scene features for navigation also beats (only scene features) VLNOBERT baseline 249 (row 5) by 3.7% on SR, 3.8% on SDTW and 1.9% on NDTW metrics. Moreover, the gap between 250 our approach (row 6) and VLN OBERT baseline (row 5) is significantly increased in RxR than R2R 251 which highlights that object features significantly help as you move to a larger dataset with dense 252 instructions containing more entity references. On R2R, our approach (row 6) shows 1.81% gain in 253 SPL over the reproduced VLN OBERT baseline (row 5) and is also competitive with the published 254 VLN () BERT result. 255

256 Table 2 reports results on R2R test-unseen split which again demonstrates that our method is competitive with the previous state-of-the-art approaches across all the metrics. For completeness, 257 we also report VLN () BERT with PREVALENT initialization in Table 1 (row 5) and Table 1 (row 258 4) since it achieves state-of-the-art results on R2R. However, comparison with it is unfair since it is 259 pre-trained and fine-tuned with high-level scene features on augmented VLN data which has much 260 less domain gap between pre-training and fine-tuning than VLN OBERT with OSCAR initialization 261 (row 4) which is pre-trained on Conceptual Captions and fine-tuned on VLN data. However, the focus 262 of our work is to show how we can effectively leverage pre-training on large-scale V&L data to boost 263 performance on VLN task without any additional pre-training on VLN data. Pre-training on VLN 264 data requires large number of instruction-path pairs which is very expensive collection process [1, 2]. 265 PREVALENT resort to using noisy augmented data generated by speaker model for pre-training 266 which leads to improvements over OSCAR pre-training. We posit that this gap can be overcome by 267 pre-training OSCAR [17] model on larger V&L data whose collection process is relatively much 268 cheaper and scalable than VLN. 269

Dummy Qualitative Figure

Figure 3: Dummy caption

Learning curves. Figure 4 plots the learning curves for VLN OBERT baseline and our method
 on RxR val-unseen split. [AB] Compare with human baseline on NDTW metric. Our model gives
 NDTW score 69 and human baseline is at 77.

Do object features help just as additional input? In Table 3 (row 2), we show the results of an 273 ablation experiment on val-seen and val-unseen split of R2R dataset in which we used object features 274 just as additional input to VLN OBERT baseline. In this ablation, all visual tokens (object and scene 275 tokens) get refined via self-attention through multiple layers in the transformer and we get action 276 probabilities and next hidden state in the similar fashion as VLN () BERT baseline [15]. We find out 277 that adding object tokens as additonal input does not help in generalization - we see improvements 278 by $\sim 2\%$ in SPL and SR on val-seen but it is slightly worse than baseline on val-unseen by 0.64% 279 on SR and 0.47% on SPL. Since additional detected object proposals are noisy, this experiment 280 demonstrates that updating high-scene tokens from these object tokens introduces additional noise 281 and causes overfitting. Moreover, during navigation, state token [CLS] gets updated from these 282 updated scene tokens in experiment 2 which introduces more noise. In 3 (row 3), we fix this noise in 283 [CLS] token by aggregating scores from objects and scenes for getting action probabilities and using 284 corresponding token embeddings for updating state token [CLS]. Incorporating score aggregation 285 (row 3) improves object features alation (row 2) SPL by 0.77% on val-unseen and reduces overfitting 286 since it has lower gap between val-seen and val-unseen metrics. 287

Do object features help without any pre-training? Table 4 reports results of VLN OBERT baseline and our approach without V&L pre-training on R2R val-useen split. As evident from the Table 4, both the approaches perform similar across all the metrics. This further demonstrates that object grounding is crucial for our approach to work; our approach improves performance on RxR and R2R datasets compared to VLN OBERT scene features baseline predominantly because it effectively utilizes the large-scale V&L pre-training.

Table 4: Results on R2R without V&L pre-training.

		R2R Val Unseen								
	Models	TL	$\text{NE}\downarrow$	SR ↑	SPL↑					
1	VLN OBERT [15]	10.31	5.21	49.21	45.53					
2	Ours	12.00	4.99	50.36	44.69					

293

²⁹⁴ How often does the model pick up objects rather than scenes for navigation?



Figure 4: NDTW learning curve comparing Recurrent VLN-BERT, our method and human baseline on RxR val-unseen split.

295 6 Conclusion

296 **References**

- [1] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen
 Gould, and Anton van den Hengel. Vision-and-language navigation: Interpreting visually-grounded
 navigation instructions in real environments. In *CVPR*, 2018.
- [2] Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and Jason Baldridge. Room-Across-Room:
 Multilingual vision-and-language navigation with dense spatiotemporal grounding. In *Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2020.
- [3] 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*, pages 3314–3325, 2018.
- [4] Xin Wang, Qiuyuan Huang, Asli Celikyilmaz, 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6629–6638, 2019.
- [5] Chih-Yao Ma, Jiasen Lu, Zuxuan Wu, Ghassan AlRegib, Zsolt Kira, Richard Socher, and Caiming Xiong.
 Self-monitoring navigation agent via auxiliary progress estimation. *arXiv preprint arXiv:1901.03035*, 2019.
- [6] Liyiming Ke, Xiujun Li, Yonatan Bisk, Ari Holtzman, Zhe Gan, Jingjing Liu, Jianfeng Gao, Yejin Choi, and
 Siddhartha Srinivasa. Tactical rewind: Self-correction via backtracking in vision-and-language navigation.
 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6741–6749,
 2019.
- [7] Chih-Yao Ma, Zuxuan Wu, Ghassan AlRegib, Caiming Xiong, and Zsolt Kira. The regretful agent:
 Heuristic-aided navigation through progress estimation. In *CVPR*, pages 6732–6740, 2019.
- [8] Hao Tan, Licheng Yu, and Mohit Bansal. Learning to navigate unseen environments: Back translation with environmental dropout. *arXiv preprint arXiv:1904.04195*, 2019.
- [9] Peter Anderson, Ayush Shrivastava, Devi Parikh, Dhruv Batra, and Stefan Lee. Chasing ghosts: Instruction following as bayesian state tracking. *arXiv*, 2019.
- [10] Xiujun Li, Chunyuan Li, Qiaolin Xia, Yonatan Bisk, Asli Celikyilmaz, Jianfeng Gao, Noah Smith,
 and Yejin Choi. Robust navigation with language pretraining and stochastic sampling. *arXiv preprint arXiv:1909.02244*, 2019.
- [11] Fengda Zhu, Yi Zhu, Xiaojun Chang, and Xiaodan Liang. Vision-language navigation with self-supervised
 auxiliary reasoning tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10012–10022, 2020.
- [12] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13137–13146, 2020.

- [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image
 Database. In *CVPR*, 2009.
- [14] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million
 image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
 2017.
- [15] Yicong Hong, Qi Wu, Yuankai Qi, Cristian Rodriguez-Opazo, and Stephen Gould. A recurrent vision-andlanguage bert for navigation, 2020.
- [16] 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, 2020.
- [17] Xiujun Li, Xi Yin, C. Li, X. Hu, Pengchuan Zhang, Lei Zhang, Longguang Wang, H. Hu, Li Dong, Furu
 Wei, Yejin Choi, and Jianfeng Gao. Oscar: Object-semantics aligned pre-training for vision-language tasks.
 In *ECCV*, 2020.
- [18] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic
 representations for vision-and-language tasks, 2019.
- [19] 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.
 In *International Conference on 3D Vision (3DV)*, 2017.
- [20] Ronghang Hu, Daniel Fried, Anna Rohrbach, Dan Klein, Trevor Darrell, and Kate Saenko. Are you
 looking? grounding to multiple modalities in vision-and-language navigation, 2019.
- [21] Yicong Hong, Cristian Rodriguez-Opazo, Yuankai Qi, Qi Wu, and Stephen Gould. Language and visual
 entity relationship graph for agent navigation, 2020.
- Yubo Zhang, Hao Tan, and Mohit Bansal. Diagnosing the environment bias in vision-and-language
 navigation. *arXiv preprint arXiv:2005.03086*, 2020.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. GloVe: Global Vectors for Word
 Representation. In *EMNLP*, 2014.
- [24] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. 2016
 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 770–778, 2016.
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015.
- [26] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yan nis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. Visual genome: Connecting
 language and vision using crowdsourced dense image annotations. *arXiv preprint arXiv:1602.07332*, 2016.
- Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei
 Zhang. Bottom-up and top-down attention for image captioning and visual question answering. *arXiv preprint arXiv:1707.07998*, 2017.
- J. Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional
 transformers for language understanding. In *NAACL-HLT*, 2019.
- [29] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [30] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. Neural Computation, 1997.
- [31] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned,
 hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565,
 2018.
- [32] 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, 2018.
- [33] Gabriel Ilharco, Vihan Jain, Alexander Ku, Eugene Ie, and Jason Baldridge. General evaluation for
 instruction conditioned navigation using dynamic time warping, 2019.

[34] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, 382 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, 383 Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit 384 Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-385 performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-386 Buc, E. Fox, and R. Garnett, editors, Advances in Neural Information Processing Systems 32, 387 pages 8024-8035. Curran Associates, Inc., 2019. URL http://papers.neurips.cc/paper/ 388 9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf. 389

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

393 Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or [N/A]. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? [Yes] See Section ??.
- Did you include the license to the code and datasets? [No] The code and the data are proprietary.
- Did you include the license to the code and datasets? [N/A]

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

- 1. For all authors...
- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's 406 contributions and scope? [TODO] 407 (b) Did you describe the limitations of your work? **[TODO]** 408 (c) Did you discuss any potential negative societal impacts of your work? [TODO] 409 (d) Have you read the ethics review guidelines and ensured that your paper conforms to 410 them? [TODO] 411 2. If you are including theoretical results... 412 (a) Did you state the full set of assumptions of all theoretical results? [N/A]413 (b) Did you include complete proofs of all theoretical results? [N/A] 414 3. If you ran experiments... 415 (a) Did you include the code, data, and instructions needed to reproduce the main experi-416 mental results (either in the supplemental material or as a URL)? [TODO] 417 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they 418 were chosen)? [TODO] 419 (c) Did you report error bars (e.g., with respect to the random seed after running experi-420 ments multiple times)? [TODO] 421 (d) Did you include the total amount of compute and the type of resources used (e.g., type 422 of GPUs, internal cluster, or cloud provider)? [TODO] 423 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets... 424 (a) If your work uses existing assets, did you cite the creators? **[TODO]** 425 (b) Did you mention the license of the assets? **[TODO]** 426 (c) Did you include any new assets either in the supplemental material or as a URL? 427 [TODO] 428 (d) Did you discuss whether and how consent was obtained from people whose data you're 429 using/curating? [TODO] 430

431 432	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [TODO]
433	5. If you used crowdsourcing or conducted research with human subjects
434 435	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
436 437	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
438 439	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

440 A Appendix

Table 5: Notation							
Symbol	Description						
Ι	Natural language navigation instruction						
$m{O}_t = \{m{o}_{t,i}\}_{i=1}^{N_t}$	Panoramic observation at time t for N_t navigable viewpoints						
$A_t = \{a_{t,k}\}_{i=0}^{N_t}$	Actions available at time t consisting of N_t navigable viewpoints and stop						
$ heta_{t,i}$	Heading angle for view <i>i</i> at time <i>t</i>						
$\phi_{t,i}$	Elevation angle for view <i>i</i> at time <i>t</i>						
$oldsymbol{f}_{t,i}$	CNN scene features for view i at time t						
$oldsymbol{g}_{t,i,j}$	Faster R-CNN features for region j in view i at time t						
$p(a_{t,j})$	Probability of taking action $a_{t,k}$ at time t						
$oldsymbol{s}_t$	Agent state at time t						

441

Role of attention pattern on performance. [AB] TODO: Update this paragraph when all the experiments finish.

 Table 6: R2R attention ablations. [AB] TODO: Update numbers when experiments finish.

 Val Seen
 Val Unseen

					vai Seen				var Oliseen			
	Models	Input Tokens	Query Tokens	TL	$\text{NE}\downarrow$	SR \uparrow	SPL ↑	TL	NE \downarrow	SR ↑	SPL ↑	
1	Baseline [15]	$\langle \pmb{s}, \pmb{X}, \pmb{V} angle$	$\langle \pmb{s}, \pmb{V} angle$	11.02	3.11	69.93	65.65	12.16	4.40	57.90	51.43	
2	Scene attention	$\langle \pmb{s}, \pmb{X}, \pmb{V}, \pmb{O} angle$	$\langle m{s},m{V} angle$	11.12	3.24	70.03	66.59	12.05	4.34	57.85	51.73	
3	All attention	$\langle m{s}, m{X}, m{V}, m{O} angle$	$\langle \pmb{s}, \pmb{V}, \pmb{O} angle$	-	-	-	-	-	-	-	-	
4	Object attention	$\langle {oldsymbol s}, {oldsymbol X}, {oldsymbol O} angle$	$\langle m{s},m{O} angle$	-	-	-	-	-	-	-	-	
5	Ours	$\langle \pmb{s}, \pmb{X}, \pmb{V}, \pmb{O} angle$	$\langle m{s},m{O} angle$	11.81	3.63	62.78	58.01	12.15	4.28	58.71	53.24	