# Language Prompt for Autonomous Driving

Dongming Wu[1*‡], Wencheng Han[2*], Tiancai Wang[3], Yingfei Liu[3],
Xiangyu Zhang[3,4], Jianbing Shen[2†]
[1] Beijing Institute of Technology, [2] SKL-IOTSC, CIS, University of Macau,
[3] MEGVII Technology, [4] Beijing Academy of Artificial Intelligence
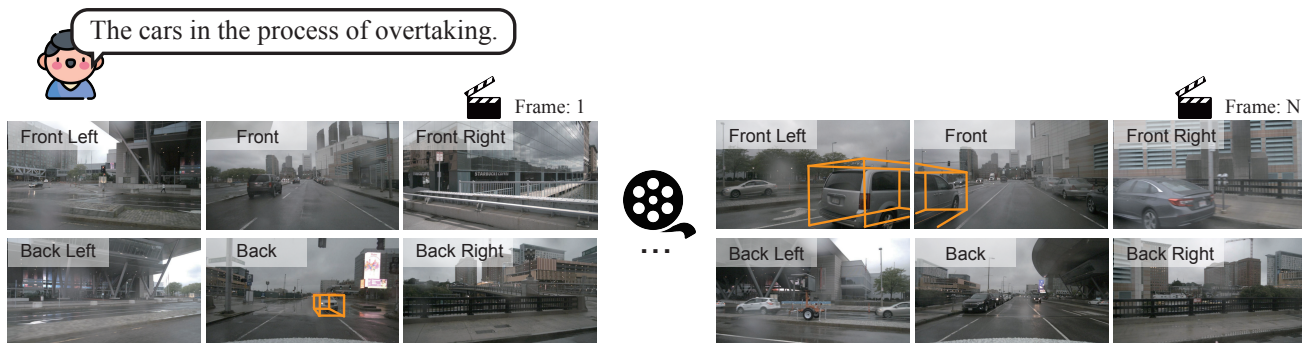{wudongming97, wenchenghan, shenjiangbingcg}@gmail.com, wangtiancai@megvii.com

Figure 1. **A representative example from NuPrompt**. The language prompt "the cars in the process of overtaking" is precisely annotated and matched with the driving objects within the 3D, multi-frame, and multi-view space. NuPrompt contains 35,367 object-prompt paris.

## Abstract

*A new trend in the computer vision community is to capture objects of interest following flexible human command represented by a natural language prompt. However, the progress of using language prompts in driving scenarios is stuck in a bottleneck due to the scarcity of paired prompt-instance data. To address this challenge, we propose the first object-centric language prompt set for driving scenes within 3D, multi-view, and multi-frame space, named NuPrompt. It expands Nuscenes dataset by constructing a total of 35,367 language descriptions, each referring to an average of 5.3 object tracks. Based on the object-text pairs from the new benchmark, we formulate a new prompt-based driving task, i.e., employing a language prompt to predict the described object trajectory across views and frames. Furthermore, we provide a simple end-to-end baseline model based on Transformer, named PromptTrack. Experiments show that our PromptTrack achieves impressive performance on NuPrompt. We hope this work can provide more new insights for the autonomous*

*driving community. Dataset and Code will be made public at https://github.com/wudongming97/Prompt4Driving.*

## 1. Introduction

Leveraging natural language description in visual tasks is one of the recent trends in the vision community [14, 26]. It has garnered significant interest for its potential applications in various downstream tasks, such as embodied intelligence and human-robot interactions [5, 9, 11]. Its core idea is to predict the desired answer by shifting human instruction inputs but not updating model weights, delivering high adaptability in response to varying human demands. A key factor for the progress in 2D scenes is the availability of large-scale image-text pairs [4, 17, 29]. However, this success is hard to replicate in driving scenarios due to the scarcity of 3D instance-text pairs.

Pioneering works like Talk2Car [9], Cityscapes-Ref [31] have started to incorporate natural language into object detection tasks in driving scenes. Unfortunately, these datasets only allow each expression to refer to a single object within an individual image, restricting their usage in scenarios with multiple referred objects or changing temporal states. Furthermore, Refer-KITTI [34] addressed this

---

| Dataset | Basic Task | 3D | #Views | #Videos | #Frames | #Prompts | # Instances per-prompt |
|---|---|---|---|---|---|---|---|
| RefCOCO [37] | Det&Seg | ✗ | 1 | - | 26,711 | 142,209 | 1 |
| Talk2Car [9] | Det | ✓ | 1 | - | 9,217 | 11,959 | 1 |
| Cityscapes-Ref [31] | Det | ✗ | 1 | - | 4,818 | 30,000 | 1 |
| Refer-KITTI [34] | MOT | ✗ | 1 | 18 | 6,650 | 818 | 10.7 |
| Refer-Youtube-VOS [30] | Seg | ✗ | 1 | ∼3.9K | 93,869 | ∼27K | 1 |
| Nuscenes-QA [25] | VQA | ✓ | 6 | 850 | 34,149 | 459,941 | - |
| NuPrompt (Ours) | MOT | ✓ | 6 | 850 | 34,149 | 35,367 | 5.3 |

Table 1. **Comparison of our NuPrompt with existing prompt-based datasets.** '-' means unavailable. NuPrompt provide the nature and complexity of driving scenes, *i.e.*, 3D, multi-view space, and multi-frame domain. Besides, it primarily focuses on object-centric understanding by pairing a language prompt with multiple targets of interest.

issue by extending the KITTI dataset [10] to include expressions that ground multiple video objects. This work mainly focuses on modular images and 2D detection, thereby leaving room for improvement in 3D driving scenes. A recent advancement, namely Nuscenes-QA [25], offers numerous question-answer pairs for 3D multi-view driving scenes, making significant strides in the use of language prompts. However, it primarily contributes to scene-level understanding and overlooks the direct and fine-grained semantic correspondence between 3D instances and text.

To advance the research of prompt learning in driving scenarios, we propose a new large-scale benchmark, named **NuPrompt**. The benchmark is built on the popular multi-view 3D object detection dataset Nuscenes [2]. We assign a language prompt to a collection of objects sharing the same characteristics for grounding them. Essentially, this benchmark provides lots of 3D instance-text pairings with three primary attributes: ❶ *Real-driving descriptions.* Different from existing benchmarks that only represent 2D objects from modular images, the prompts of our dataset describe a variety of driving-related objects from 3D, looking around, and long-temporal space. Fig. 1 illustrates a typical example, *i.e.*, a car surpasses our car from behind towards the front across multiple views. ❷ *Instance-level prompt annotations.* Every prompt indicates a fine-grained and discriminative object-centric description, as well as enabling it to cover an arbitrary number of driving objects. ❸ *Large-scale language prompts.* NuPrompt is comparable to the largest current dataset [7] in terms of the number of prompts, *i.e.*, including 35,367 language prompts.

Along with the benchmark, we formulate a new prompt-based perceiving task, whose primary objective is to predict and track multiple 3D objects in driving environments using a given language prompt. The challenges of this task lie in two aspects: temporal association across frames and cross-modal semantic comprehension. To address the challenges, we propose an end-to-end baseline built on camera-only 3D tracker PF-Track [24], named **PromptTrack**. Note that PF-Track has exhibited excellent spatial-temporal mod-

eling through its past and future reasoning branches. Furthermore, we add one prompt reasoning branch to perform cross-modal fusion and understanding. Specifically, our prompt reasoning involves cross-attention between prompt embedding and query features from past reasoning, further predicting prompt-referred objects.

In summary, our contributions are three-fold:

- We propose a new large-scale language prompt set for driving scenes, named NuPrompt. As far as we know, it is the first dataset specializing in multiple 3D objects of interest from video domain.

- We construct a new prompt-based driving perceiving task, which requires using a language prompt as a semantic cue to predict object trajectories.

- We develop a simple end-to-end baseline model, called PromptTrack, which effectively fuses cross-modal features in a newly built prompt reasoning branch to predict referent objects, showing impressive performance.

## 2. Related Work

**Language Prompt in Driving Scenes.** Language prompt modeling is a broad concept covering various vision tasks, like object detection [12, 37], referring segmentation [35], and text-image generation [27]. In this work, we focus on the application of language prompts within driving scenes. The utilization of human commands within driving scenes allows the system to understand driving systems from the human perspective, thereby facilitating human control over driving procedures. Talk2Car [9], the pioneering benchmark featuring language prompts for autonomous vehicles, is constructed on the base of Nuscenes [2]. However, its annotation only comprises keyframes that catch the eye of annotators. Cityscapes-Ref [31] annotates both language expressions and gaze recordings for each video sequence within the driving dataset Cityscapes [7]. However, the prompts deployed in both Talk2Car and Cityscapes-Ref
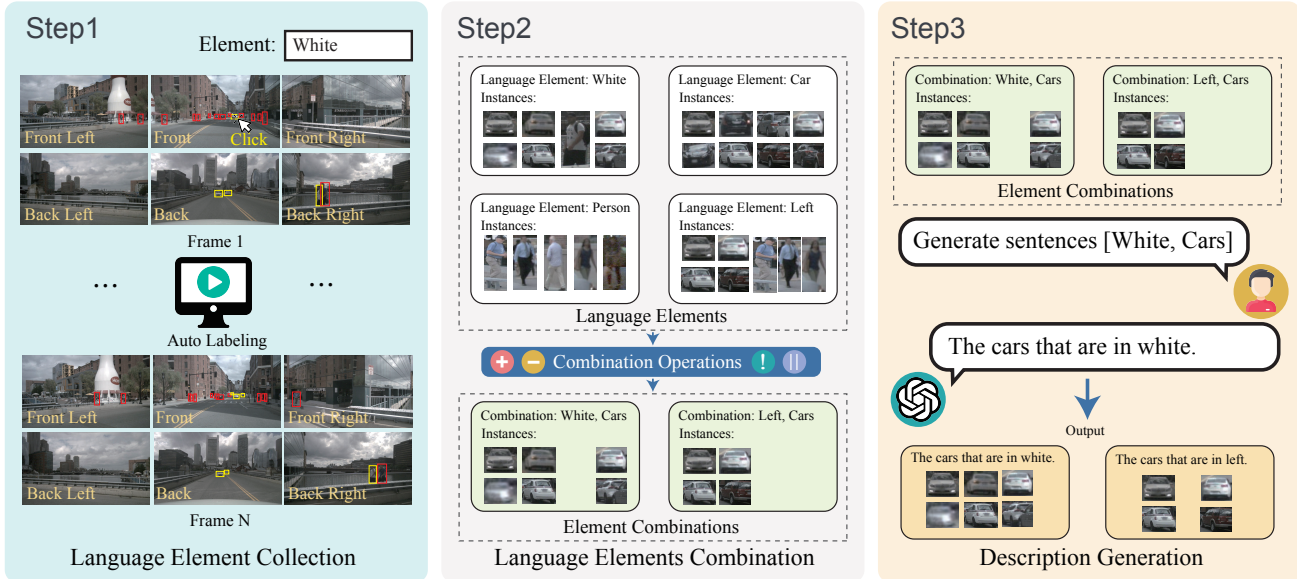
Figure 2. **Pipeline of language prompt annotation procedure**, which includes three steps: language element collection, language element combination, and description generation. Firstly, we pair each language tag with referent objects during the language element collection phase. Following this, certain language elements are selected and combined in the language element combination stage. Finally, with the combinations obtained, we employ Large Language Model (LLM) to create language descriptions in the description generation stage.

tend to represent an individual object. To solve this problem, Refer-KITTI [34] further develops KITTI [10], where each prompt can refer to a collection of referent objects. This makes Refer-KITTI stand out as the first dataset to use a language prompt for specifying an arbitrary number of object predictions. More recently, Nuscenes-QA [25] opens a new avenue, Visual Question Answering (VQA), for understanding scene-level driving scenarios. It contains 459,941 question-answer pairs based on 34,149 visual scenes from Nuscenes [2]. Uniquely, our NuPrompt offers more fine-grained matching for 3D instances and each prompt. A thorough comparison between existing prompt-based driving datasets and ours is summarized in Table 1.

**Referring Expression Understanding.** Given a language prompt, the goal of referring expression understanding is to localize the described objects using boxes or masks, which shares a similar idea with our prompt-based driving benchmark. The initiation of datasets like RefCOCO/+/g [37] has helped stimulate interest in this field. These datasets use succinct yet unambiguous natural language expressions to ground a visual region within an image. Liu *et al*. further improve this dataset by allowing it to support expressions that refer to unlimited target objects [18]. VID-Sentence [6], developed from the ImageNet video object detection dataset [28], annotates sentence descriptions for each video sequence, leading to a video-based referring understanding dataset. Besides, Refer-DAVIS$_{16/17}$ [13] and Refer-Youtube-VOS [30] are another two popular video-based referring expression understanding benchmarks sup-

porting video object segmentation. A recent work in the field called GroOT [22] expands the large-scale multi-object tracking dataset TAO [8] to support referring understanding.

## 3. Dataset Overview

In this section, we will present the details about our NuPrompt dataset. In § 3.1, we will discuss the process of data collection and annotation. In § 3.2, we will go deeper and provide statistics on the newly created dataset. Lastly, in order to evaluate the performance of the methods on the new task, we will introduce some metrics and establish a new benchmark for this dataset in § 3.3.

### 3.1. Data Collection and Annotation

Our NuPrompt is built on one of the most popular datasets for multi-view 3D object detection, Nuscenes [2]. While the original Nuscenes dataset includes visual images and point cloud data, we here focus solely on visual images for NuPrompt. As shown in Fig. 2, the cars collecting the data are equipped with six different cameras: Front, Front Left, Front Right, Back Right, Back Left, and Back. These cameras have some areas of overlap with one another. When combined with their intrinsic and extrinsic matrices, our dataset covers 360° of 3D space for each scene.

To efficiently generate training labels for the new dataset, we designed a three-step semi-automatic labeling pipeline as shown in Fig. 2. The first step aims at identifying language elements and associating them with 3D bounding

3

Figure 3. **Word cloud** of the top 100 words in NuPrompt. It has a large number of words describing driving object appearance, like "black", "white" and "red", as well as covers many motion scenes, such as "walking", "moving", and "crossing'".

boxes. The second step is to combine the language elements using certain rules. In the third step, we base the language element combinations to produce various language prompts using a large language model (LLM). In the next section, we will provide detailed information about these steps.

**Step 1: Language Element Collection.** This paper uses the term "language element" to refer to a basic attribute of objects. Examples of language elements include colors (e.g. red, yellow, and black), actions (e.g. running, stopping, and crossing the road), locations (e.g. left, right, and back), and classes (e.g. car and pedestrian), which cover diverse descriptions of driving scenes. The key problem is how to label the bounding boxes with the corresponding language elements. To solve this, we design a labeling system to manually collect and match language elements with bounding boxes in a video sequence, as demonstrated in Fig. 2. Annotators type the language element texts and click on the corresponding bounding boxes. When the target changes status and no longer belongs to the language elements, annotators click on the target again and remove it from the list. This procedure can efficiently reduce the amount of human labor required. To ensure a variety of expressions, each video is assigned to five independent annotators who manually create descriptive expressions to formulate query sentences. Two other annotators then carefully check the matching between the expressions and the video targets.

**Step 2: Language Elements Combination.** As mentioned earlier, language elements are basic attributes of objects. By combining these attributes, we can create various descriptions for different groups of objects. There are three types of relationships we can use to merge the attributes: AND, OR, and NOT. We use these operations to combine sets of bounding boxes and their language elements, resulting in a new set with the merged attributes. In our dataset, we manually selected some meaningful attribute combinations and also randomly generated a large number of combinations for the objects. To ensure that the combinations are valid, we filtered out those that do not have a certain number of

bounding boxes in the video sequences.

**Step 3: Prompt Generation.** After Step 2, we are able to determine the correspondence between combinations of language elements and a group of bounding boxes. However, getting valid natural language sentences to describe the objects can be expensive using human labor, and there is no guarantee of the desired variety. Large language model (LLM) have recently shown great potential in understanding logistics and producing sentences similar to those generated by humans. Therefore, we determine GPT3.5 [23] as our language model. We prompt it with a request like "Generate a sentence to describe the objects based on the following descriptions: *pedestrians, moving, red, not in the left*," where the italicized words represent the combination of elements. The LLM can respond with a meaningful description of the objects, such as "The objects are red pedestrians, currently in motion, not situated on the left side." To guarantee accuracy, we will ask annotators to filter out any incorrect descriptions generated by the LLM. We also prompt the LLM multiple times to generate multiple descriptions.

### 3.2. Dataset Statistics

Thanks to Nuscenes [2], our language prompts provide a number of comprehensive descriptions for the objects in the 3D, surrounding, and temporal space. Besides, they cover diverse environments comprising pedestrian streets, public roads, and highways in the cities of Boston and Singapore. Furthermore, they encompass different weather conditions (*e.g.*, rain and sun), and illumination (*e.g.*, day and night). To offer deeper insights into NuPrompt, we next present more quantitative statistics.

**Language Prompt.** We manually labeled 13,004 language elements and utilized them in the combination and generation of 22,363 unified descriptions through LLM. In total, the NuPrompt has 35,367 language prompts. On average, each video within the dataset contains 41.6 prompts. As summarized in Table 1, our proposed dataset is comparable to the largest current referring dataset in terms of prompt counts. We also show the word cloud of the top 100 words in Fig. 3. From the word cloud figure, we can observe that NuPrompt dataset has a large number of words that describe motions, like "walking", "moving", and "crossing", and many driving object appearances, like "black", "white" and "red".

**Referent Objects.** In contrast to previous benchmarks that refer to 2D objects in modular images, another feature of NuPrompt is its surrounding 3D space. This indicates that there are lots of objects crossing different views, presenting improved simulation being closer to real driving scenes. More importantly, NuPrompt is designed to involve an arbitrary number of predicted objects. As shown in Fig. 4 (a), most prompts describe 1-10 instances, and the maximum number can be more than 20. From Table 1, each prompt in
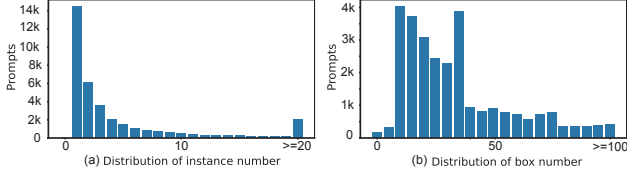
Figure 4. **Statistics of NuPrompt** on (a) distribution of instance number per prompt and (b) distribution of box number per prompt.

our dataset refers to an average of 5.3 instances. In addition, we display the distribution of box per prompt in Fig. 4 (b). As seen, our language prompts mostly have 10∼100 boxes.

### 3.3. Benchmark Protocols

**Evaluation Metrics.** To evaluate the similarity between the predicted tracklet and ground truth tracklet, we primarily use the Average Multiple Object Tracking Precision (AMOTA) metric [1]. However, unlike the original multi-object tracking task that averages AMOTA across different categories, the evaluation on NuPrompt is class-agnostic. Hence, we calculate AMOTA for every test video-prompt pair and take the average of all these AMOTA values. For a more detailed analysis, we also utilize the Average Multi-Object Tracking Precision (AMOTP) and Identity Switches (IDS) metrics.

**Data Split.** The NuPrompt is a large-scale dataset, which contains a total of 850 videos along with 43,223 language prompts. Following the setting of Nuscenes [2], we split NuPrompt into training set and validation set, which contain 700 videos and 150 videos, respectively.

## 4. Method

Given multi-frame multi-view images and a prompt, the goal of our new task is to track the described object. It requires not just temporal association across frames, but also a comprehensive alignment of cross-modal semantics. To accomplish these two objectives, we propose PromptTrack, an end-to-end framework. It modifies the query-based method PF-Track [24] to adapt to the prompt input. Fig. 5 shows the overall pipeline of PromptTrack. Note that PF-Track PF-Track incorporates a past reasoning branch and a future reasoning branch which are based on the decoded queries. These branches aim to refine tack prediction using cached historical queries and improve cross-frame query propagation using motion localization prediction, respectively. In addition to these two branches, we propose a new prompt reasoning branch to predict the prompt-referred tracks. In the following, we will introduce PromptTrack details.

### 4.1. Overall Architecture

Formally, let $F_t$ denote the extracted visual features at timestamp $t$, and $S$ denote the encoded linguistic features.

To enrich the information of visual features, we first incorporate the visual features with the linguistic features in a multiplication way and form enhanced visual feature maps. To capture different views and stereo information, we follow PETR [20] to add position embedding, defined as $F_t'$. Then a set of 3D queries $Q_t$ interact with the position-aware visual features $F_t'$ via a stack of Transformer decoder layers, outputting updated queries $Q_t^D$ and bounding boxes $B_t^D$:

$$Q_t^D, B_t^D = \textbf{Decoder}\left(Q_t, F_t'\right), \tag{1}$$

where each input query $q_t^i \in Q_t$ means an object with a feature vector $f_t^i$ and a 3D localization $c_t^i$, i.e., $q_t^i = \{f_t^i, c_t^i\}$.

To automatically link objects across different frames, the input queries $Q_t$ merge track queries $Q_t^{track}$ from the last frame. The box information from the last frame also provides an excellent spatial position prior to the current frame, benefiting the model for accurately inferring the same object. Besides, to capture new-born objects, a set of fixed 3D queries $Q_t^{fixed}$, also called detection queries, are concatenated with track queries $Q_t^{track}$ to generate $Q_t$. Following the work [24], the number of fixed queries is set to 500. As the first frame has no previous frames, we only utilize the fixed queries to detect objects.

**Past and Future Reasoning.** After the Transformer Decoder, past and future reasoning are sequentially conducted for attending to historical embeddings and predicting future trajectory, respectively. Formally, the past reasoning $\mathcal{F}^p$ integrates two decoded outputs $Q_t^D$ and $B_t^D$ as well as cached historical queries $Q_{t-\tau_h:t-1}$ from past $\tau_h$ frames to produce refined queries $Q_t^R$ and refined bounding boxes $B_t^R$:

$$Q_t^R, B_t^R = \mathcal{F}^p\left(Q_t^D, B_t^D, Q_{t-\tau_h:t-1}\right), \tag{2}$$

where $\mathcal{F}^p$ has a cross-frame attention module for promoting history information integration across $\tau_h$ frames per object. Moreover, it also includes a cross-object attention module to encourage discriminative feature representation for each individual object. The sequential cross-frame attention and cross-object attention modules lead to $Q_t^R$. A multi-layer perceptron (MLP) is used to predict coordinate residuals and adjust the object boxes, leading to $B_t^R$.

Based on the refined results from past reasoning, the future reasoning $\mathcal{F}^f$ uses a cross-frame attention to predict long-term trajectories $M_{t:t+\tau_f}$ for next $\tau_f$ frames:

$$Q_{t+1}^{track}, M_{t:t+\tau_f} = \mathcal{F}^f\left(Q_t^R, Q_{t-\tau_h:t-1}\right), \tag{3}$$

where the position vectors of the refined queries $Q_t^R$ is updated to generate $Q_{t+1}^{track}$ according to the single-step movement $M_{t:t+1}$. The main motivation of future reasoning is that as the ego-car goes forward, the reference position of all objects from the last frame has to be adjusted to align with the new ego-coordinates. In summary, using past and future information can improve the quality of visible object
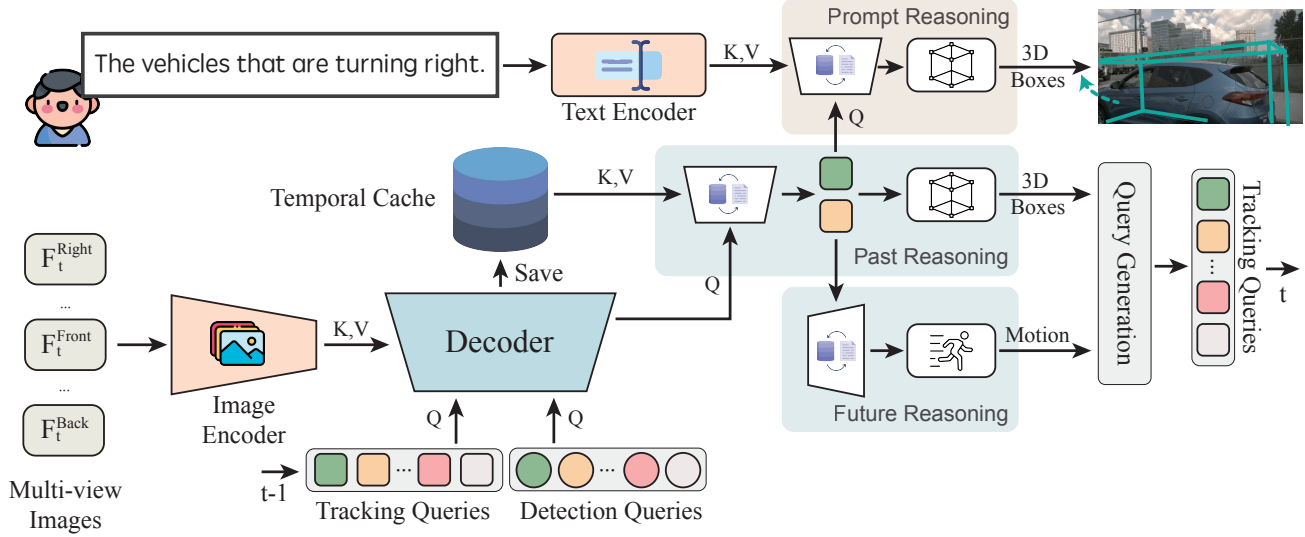
5

Figure 5. **Overall architecture of PromptTrack**. For each frame, the visual features and queries are first fed into the Transformer Decoder to produce decoded queries as like the original DETR framework. Then past reasoning enhances and refines tacks by attending to historical queries, and future reasoning benefits cross-frame query propagation using predicted position. Lastly, the prompt reasoning branch predicts the prompt-referred tracks. The model is evaluated by online mode, while the training can be end-to-end differentiated.

tracking. To further identify the object referred to in the prompt, we outline the new prompt reasoning as follows.

## 4.2. Prompt Reasoning

Prompt reasoning focuses on boosting the coherence between visual and linguistic elements by implementing holistic cross-modal interaction. Since the refined queries $Q_t^R \in \mathbb{R}^{N \times C}$ has integrated the past track information, we use it and the prompt embeddings $S \in \mathbb{R}^{L \times C}$ as inputs, and outputs the prompt-referred object probability $P_t$:

$$P_t = \mathcal{F}^l \left( Q_t^R, S \right) \in \mathbb{R}^{N \times 1}, \qquad (4)$$

where $N$ is the query number, $L$ is the number of prompt words, $C$ indicates the feature dimension.

Specifically, we first use cross-modal attention to encourage the feature fusion between the two kinds of modalities $Q_t^R$ and $S$, generating prompt-aware queries $Q^P \in \mathbb{R}^{N \times C}$:

$$
\begin{aligned}
Q_t^P = \textbf{CrossModalAttn}(\text{Q} &= Q_t^R, \\
\text{K,V} &= S, \qquad (5) \\
\text{PE} &= \text{Pos}(S)),
\end{aligned}
$$

where Pos means the position embedding following [32]. Then we use MLP to output a binary probability that indicates whether the output embedding represents a prompt-referred object:

$$P_t = \textbf{MLP}(Q_t^P). \qquad (6)$$

## 4.3. Instance Matching and Loss

Our method views the detection as a set prediction problem following query-based methods [3], so it requires one-to-one matching between queries and ground-truth before calculating loss. The tracking queries $Q_t^{track}$ and their ground-truth have been matched when propagating queries. As for the fixed queries $Q_t^{fixed}$, we fit them with the new-born objects using a bipartite graph matching. During the process of matching, we only use the queries $Q_t^D$ from decoder outputs, and then implement the correspondence in the overall loss function:

$$
\begin{aligned}
\mathcal{L} = &\lambda_{cls}^D \mathcal{L}_{cls}^D + \lambda_{box}^D \mathcal{L}_{box}^D + \lambda_f \mathcal{L}_f + \\
&\lambda_{cls}^R \mathcal{L}_{cls}^R + \lambda_{box}^R \mathcal{L}_{box}^R + \lambda_p \mathcal{L}_p,
\end{aligned} \qquad (7)
$$

where $\mathcal{L}_{cls}$s are Focal loss [16] for classification, $\mathcal{L}_{box}$s are L1 loss for bounding box regression, $\mathcal{L}_f$ is L1 loss for motion prediction. The settings of their weights $\lambda$s follow [24]. Besides, our prompt reasoning loss $\mathcal{L}_p$ is also Focal loss, weighted by $\lambda_p$.

## 4.4. Inference

During inference, PromptTrack can handle the arbitrary length of videos because it belongs to a frame-by-frame on-line framework. Given the $t^{th}$ frame and language prompt, PromptTrack will predict N instances, each corresponding to true or empty objects. Since the prompt reasoning is based on past reasoning, we first choose all true track objects whose class score exceeds a certain threshold $\gamma^{object}$. Further, the final prompt-referred objects are deter-

6

| Method | Basic Detector | AMOTA ↑ | AMOTP↓ | RECALL↑ | MOTA↑ | IDS↓ |
|---|---|---|---|---|---|---|
| CenterPoint | DETR3D | 0.061 | 1.846 | 17.5% | 0.069 | 527 |
| CenterPoint | PETR | 0.074 | 1.687 | 24.2% | 0.084 | 421 |
| PromptTrack (Ours) | DETR3D | 0.095 | 1.614 | 28.5% | 0.105 | 151 |
| PromptTrack (Ours) | PETR | 0.127 | 1.361 | 43.5% | 0.135 | 146 |

Table 2. **Performance comparison with other methods on NuPrompt**. CenterPoint [36] is a heuristic-based tracking algorithm, while our PromptTrack is an end-to-end framework. ↑ and ↓ represent the direction of better performance with regard to each specific metric.

| Method | AMOTA ↑ | AMOTP↓ | RECALL↑ |
|---|---|---|---|
| PromptTrack | 0.127 | 1.361 | 43.5% |
| *w/o* Prompt Reason | 0.112 | 1.552 | 30.4% |
| *w/o* Past Reason | 0.104 | 1.580 | 29.9% |
| *w/o* Future Reason | 0.105 | 1.453 | 31.5% |

Table 3. **Ablation study of PromptTrack on NuPrompt dataset**. Each reasoning branch is independently removed and evaluated.

| $\gamma^{prompt}$ | AMOTA ↑ | AMOTP↓ | RECALL↑ |
|---|---|---|---|
| 0.1 | 0.115 | 1.371 | 45.1% |
| 0.2 | 0.127 | 1.361 | 43.5% |
| 0.3 | 0.123 | 1.379 | 43.9% |
| 0.4 | 0.112 | 1.382 | 42.8% |

Table 4. **Ablation study of prompt threshold on NuPrompt** in terms of AMOTA, AMOTP and RECALL.

mined from these true objects by filtering another threshold $\gamma^{prompt}$ in terms of prompt scores.

## 5. Experiment

### 5.1. Implementation Details

**Model Details.** We use VoVNetV2 [15] to extract visual features and RoBERTa [19] to embed language prompts. For visual features, the C5 feature is upsampled and fused with the C4 feature, where the fused C4 feature is fed into the Transformer decoder. The feature dimension for visual and linguistic features is $C = 256$. The cross-modal attention in Eq. 4 uses a multi-head technique with 8 heads. For regression and classification of the box and prompt probability, each MLP head contains two fully-connected layers. The classification head for the decoder and past reasoning branch returns 7-class logits, while the prompt head returns a binary value. The length of stored query in past reasoning is $\tau_p = 3$ frames, and the future reasoning predicts the movements of future $\tau_f = 8$ frames.

**Training.** The input images are downsampled into a resolution of $800 \times 320$ for training. For initialization, the parameters in the prompt reasoning branch are randomized, and the rest of the parameters are loaded from the official weights of the PF-Track single-frame detection model, which is pretrained on the Nuscenes [2]. The parameters in the prompt encoder are kept frozen during the training process. We implement training with the AdamW optimizer [21]. Here, the learning rate is initially set at $2.0 \times 10^{-4}$ and scheduled according to cosine annealing. The loss weights are set as $\lambda_{cls}$=2, $\lambda_{reg}$=0.25, $\lambda_f$=0.5, $\lambda_p$=2. We train PromptTrack on three-frame samples for 12 epochs. The overall training is deployed on 8 Nvidia A100 GPUs with batch size of 1.

**Testing.** PromptTrack evaluates each video sequence without any post-processing techniques. The object threshold is $\gamma^{object} = 0.2$ and the prompt threshold is $\gamma^{prompt} = 0.2$.

### 5.2. State-of-the-art Comparison

Since there are no existing methods for the new prompt-based task, we modify one heuristic-based tracker, *i.e.*, CenterPoint [36]. We integrate two state-of-the-art single-frame object detectors, *i.e.*, DETR3D [33] and PETR [20] into CenterPoint. Moreover, all these methods incorporate prompt-based prediction through the implementation of our unique prompt reasoning branch. Besides, we take the place of the basic detector of our PromptTrack using DETR3D. On top of NuPrompt, we test the proposed PromptTrack and these competitors in Table 2. PromptTrack achieves 0.127 on AMOTA and 1.361 on AMOTP. In summary, our method outperforms other counterparts across all metrics.

### 5.3. Ablation Studies

To investigate the effect of each reasoning branch in our model, we conduct ablation studies on NuPrompt, as shown in Table 2. As seen, removing the prompt reasoning leads to performance degradation, resulting in a score of 0.112 on AMOTA and 1.552 on AMOTP. This phenomenon demonstrates the effectiveness of our prompt reasoning. Besides, we validate both past and future reasoning in the last two rows of Table 2. It is obvious that omitting either module contributes to a decline in performance, thereby confirming the benefit of past and future reasoning.

Moreover, we test different settings of prompt threshold $\gamma^{prompt}$ in Table 4. The AMOTA score, which marginally ranges from 0.2 to 0.4, starts to experience a minor decrease as $\gamma^{prompt}$ increases. Empirically, for the best results, we choose $\gamma^{prompt} = 0.2$.

*Language Prompt*: The vehicles that are not moving



*Language Prompt*: A white truck that is stationary in the same direction


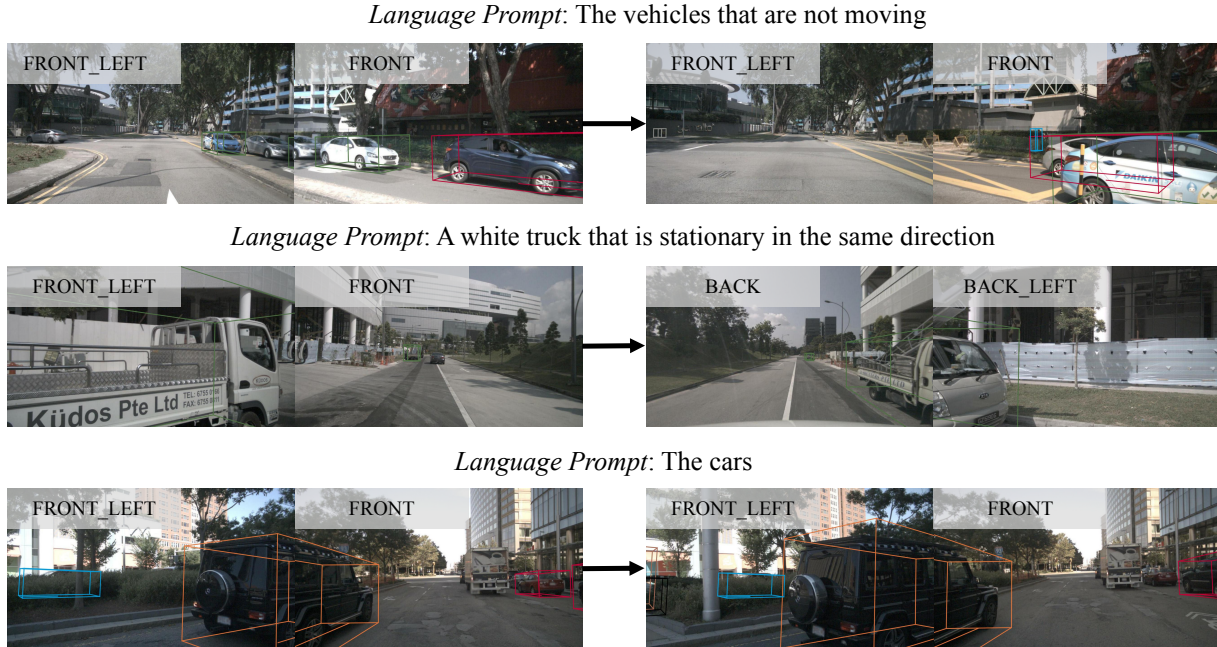
*Language Prompt*: The cars



Figure 6. **Qualitative examples** from PromptTrack on NuPrompt dataset. In terms of the given language prompt, PromptTrack can detect the described objects even if they contain various challenges, like crossing different views and varying object numbers.
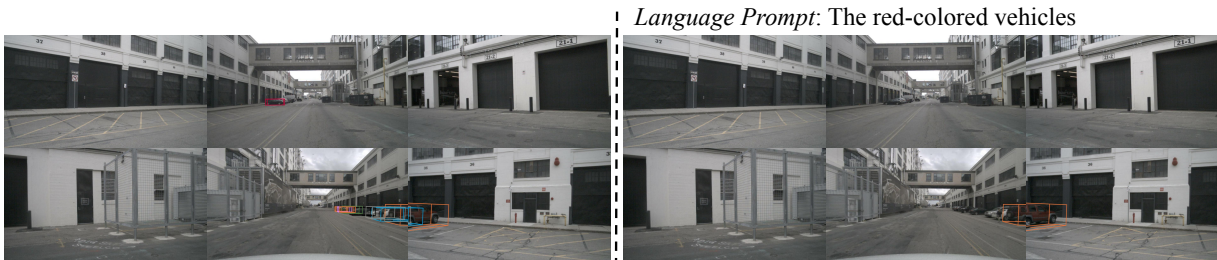
*Language Prompt*: The red-colored vehicles



Figure 7. **Qualitative comparison between all visible objects (left) and the prompt-referred objects (right).** PromptTrack is able to detect all objects as well as highlight the object of interest.

## 5.4. Qualitative Results

We visualize several typical results in Fig. 6. As seen, PromptTrack is able to detect and track prompt-referred targets accurately under various challenging situations, including crossing different views and varying object numbers. Besides, we present a qualitative comparison between all predicted objects and the prompt-referred objects from PormptTrack, as shown in Fig. 7. It is obvious that all visible objects in the driving scene are detected, and the prompt-referred objects, *i.e.*, the red-colored vehicles, are also captured in terms of the language prompt.

## 6. Conclusion and Discussion

In this work, we presented NuPrompt, the first large-scale language prompt set designed specifically for 3D perception in autonomous driving. NuPrompt provides numerous precise 3D object-text pair annotations. Therefore,

we designed a novel tracking task driven by prompts, *i.e.*, grounding objects using these language prompts. To solve this problem, we further proposed an efficient prompt-based tracking model with prompt reasoning modification on PF-Track, called PromptTrack. After conducting a set of experiments on NuPrompt, we verified the effectiveness and promising performance of our algorithm.

**Discussion.** Along with the NuPrompt which provides lots of 3D object-text pairs, there are many interesting research directions that remain to be explored. These include but are not limited to: i) designing more robust algorithms for more comprehensive temporal modeling and reasoning in both visual and linguistic modalities, ii) investigating the potential of text-to-scene generation using our fine-grained language prompts. iii) integrating trajectory prediction and driving planning into a single framework. These problems require more research efforts to promote the development of language prompts for autonomous driving.

# References

[1] Keni Bernardin and Rainer Stiefelhagen. Evaluating multiple object tracking performance: the clear mot metrics. *EURASIP Journal on Image and Video Processing*, 2008. 5

[2] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. *arXiv preprint arXiv:1903.11027*, 2019. 2, 3, 4, 5, 7

[3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 6

[4] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pretraining to recognize long-tail visual concepts. In *CVPR*, 2021. 1

[5] Li Chen, Penghao Wu, Kashyap Chitta, Bernhard Jaeger, Andreas Geiger, and Hongyang Li. End-to-end autonomous driving: Challenges and frontiers. *arXiv preprint arXiv:2306.16927*, 2023. 1

[6] Zhenfang Chen, Lin Ma, Wenhan Luo, and Kwan-Yee K Wong. Weakly-supervised spatio-temporally grounding natural sentence in video. *arXiv preprint arXiv:1906.02549*, 2019. 3

[7] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016. 2

[8] Achal Dave, Tarasha Khurana, Pavel Tokmakov, Cordelia Schmid, and Deva Ramanan. Tao: A large-scale benchmark for tracking any object. In *ECCV*, 2020. 3

[9] Thierry Deruyttere, Simon Vandenhende, Dusan Grujicic, Luc Van Gool, and Marie-Francine Moens. Talk2car: Taking control of your self-driving car. *arXiv preprint arXiv:1909.10838*, 2019. 1, 2

[10] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012. 2, 3

[11] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training. In *CVPR*, 2023. 1

[12] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *ICCV*, 2021. 2

[13] Anna Khoreva, Anna Rohrbach, and Bernt Schiele. Video object segmentation with language referring expressions. In *ACCV*, 2019. 3

[14] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *ICCV*, 2023. 1

[15] Youngwan Lee, Joong-won Hwang, Sangrok Lee, Yuseok Bae, and Jongyoul Park. An energy and gpu-computation efficient backbone network for real-time object detection. In *CVPR workshops*, 2019. 7

[16] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 6

[17] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 1

[18] Chang Liu, Henghui Ding, and Xudong Jiang. Gres: Generalized referring expression segmentation. In *CVPR*, 2023. 3

[19] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019. 7

[20] Yingfei Liu, Tiancai Wang, Xiangyu Zhang, and Jian Sun. Petr: Position embedding transformation for multi-view 3d object detection. In *ECCV*, 2022. 5, 7

[21] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 7

[22] Pha Nguyen, Kha Gia Quach, Kris Kitani, and Khoa Luu. Type-to-track: Retrieve any object via prompt-based tracking. *arXiv preprint arXiv:2305.13495*, 2023. 3

[23] OpenAI. https://chat.openai.com, 2023. 4

[24] Ziqi Pang, Jie Li, Pavel Tokmakov, Dian Chen, Sergey Zagoruyko, and Yu-Xiong Wang. Standing between past and future: Spatio-temporal modeling for multi-camera 3d multi-object tracking. In *CVPR*, 2023. 2, 5, 6

[25] Tianwen Qian, Jingjing Chen, Linhai Zhuo, Yang Jiao, and Yu-Gang Jiang. Nuscenes-qa: A multi-modal visual question answering benchmark for autonomous driving scenario. *arXiv preprint arXiv:2305.14836*, 2023. 2, 3

[26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 1

[27] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, 2023. 2

[28] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 3

[29] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 1

[30] Seonguk Seo, Joon-Young Lee, and Bohyung Han. Urvos: Unified referring video object segmentation network with a large-scale benchmark. In *ECCV*, 2020. 2, 3

[31] Arun Balajee Vasudevan, Dengxin Dai, and Luc Van Gool. Object referring in videos with language and human gaze. In *CVPR*, 2018. 1, 2

[32] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 6

[33] Yue Wang, Vitor Campagnolo Guizilini, Tianyuan Zhang, Yilun Wang, Hang Zhao, and Justin Solomon. Detr3d:

3d object detection from multi-view images via 3d-to-2d queries. In *CoRL*, 2022. 7

[34] Dongming Wu, Wencheng Han, Tiancai Wang, Xingping Dong, Xiangyu Zhang, and Jianbing Shen. Referring multi-object tracking. In *CVPR*, 2023. 1, 2, 3

[35] Linwei Ye, Mrigank Rochan, Zhi Liu, and Yang Wang. Cross-modal self-attention network for referring image segmentation. In *CVPR*, 2019. 2

[36] Tianwei Yin, Xingyi Zhou, and Philipp Krähenbühl. Center-based 3d object detection and tracking. *CVPR*, 2021. 7

[37] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *ECCV*, 2016. 2, 3