

---

# The Argoverse Trajectory Retrieval Benchmark

---

**Eric Zhan**                      **Jagjeet Singh**                      **Yisong Yue**                      **Andrew Hartnett**  
Caltech                              Argo AI                              Caltech                              Argo AI  
ezhan@caltech.edu      jsingh@argo.ai      yyue@caltech.edu      ahartnett@argo.ai

## Abstract

1                      As tracking data becomes more readily available in many domains such as sports,  
2                      animal tracking, and autonomous vehicles, so does the need for effective informa-  
3                      tion access and retrieval of those growing datasets. To that end, we develop the  
4                      Argoverse Trajectory Retrieval Benchmark for contextual trajectory retrieval of  
5                      driving scenarios. The goal of this task is to find similar trajectories from within a  
6                      large dataset given a query trajectory. This task is challenging because there are  
7                      many dimensions of variation in which two trajectories can be similar, such as  
8                      vehicle kinematics, social causality, and road configurations. To our knowledge,  
9                      this is the first standardized benchmark for trajectory retrieval of driving scenarios.  
10                     We also provide an evaluation of baseline approaches based on representation  
11                     learning and relevance feedback, and highlight several areas for improvement for  
12                     which machine learning can play a large role in future work.

## 13 **1 Introduction**

14 Behavioral tracking data is growing rapidly in many domains, including sports analytics [9, 44, 41],  
15 pedestrian crowds [26, 33, 36], traffic scenes [11, 15, 8], and animal behavior [6, 17, 45]. As  
16 behavioral track datasets grow, it becomes increasingly important to develop retrieval systems to  
17 organize and access information from the data. In this paper, we focus on traffic scenes collected in  
18 contexts involving autonomous vehicles (AVs). AV fleets have gathered millions of miles of such log  
19 data [11, 15, 8]; having effective information retrieval systems is important for extracting value from  
20 the data and accelerating the development of AV technologies.

21 An effective retrieval system can impact numerous applications, similar to the ubiquity of use cases  
22 that exist for web search systems [7, 29]. One use case that motivates our work is dataset curation.  
23 For instance, suppose we found an example of an unusual and rare driving maneuver, such as the  
24 one depicted in Figure 1a. We can then use our trajectory retrieval system to obtain similar scenes  
25 for many possible downstream tasks that can: 1) reveal a better understanding of how often such  
26 maneuvers arise; 2) refine our taxonomy of driving behaviors; 3) construct training data to train  
27 forecasting models that can more accurately capture such behavior; or 4) create simulations that  
28 include such scenarios for safety-critical testing of rare events. Similar retrieval needs arise in related  
29 fields such as sports analytics [41, 14, 51].

30 When designing retrieval systems, especially when using machine learning, it is important to establish  
31 standardized benchmarks. To that end, we present the **Argoverse Trajectory Retrieval Benchmark**,  
32 which is, to our knowledge, the first standardized retrieval dataset and task for traffic scenes. Our  
33 dataset consists of 2,795 scenarios from the Argoverse Motion Forecasting 1.1 validation set [11].  
34 Each scenario has been augmented with relevance labels for 13 complex retrieval intents pertaining to  
35 the focal agent and both its social and map contexts. The remaining 240,000+ training and validation  
36 scenarios are available for unsupervised methods. We believe that this dataset and retrieval task will  
37 stimulate research in designing retrieval systems for behavioral tracking data.

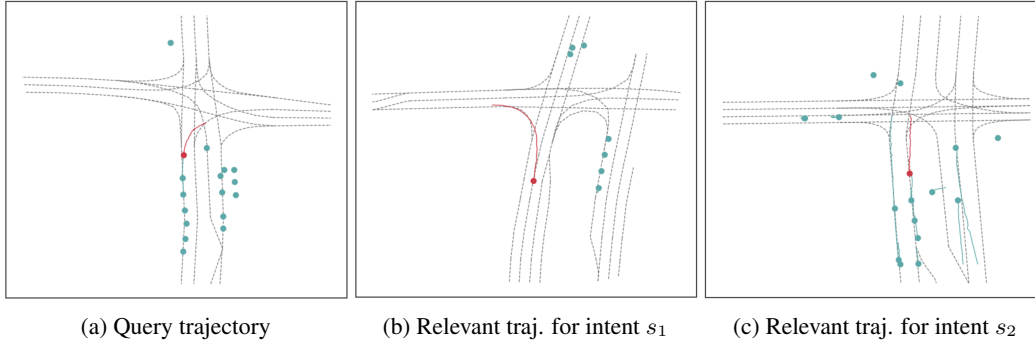


Figure 1: (a) Query trajectory with the focal agent in red. The task is to retrieve trajectories similar to this one. (b) Relevant trajectory for the intent of “turn then change lanes”. (c) Relevant trajectory for the intent of “decelerate for moving lead vehicle”. Both retrievals are valid for the query, depending on the underlying intent. One challenge with this task is inferring the underlying (hidden) intent.

38 Benchmarking for traffic scene retrieval poses new design decisions compared to more conventional  
 39 retrieval domains such as text. The first is defining a suitable similarity measure for retrieval. While  
 40 such query/item similarity measures are commonly used in retrieval [29], defining such a measure for  
 41 traffic scenes is challenging. For multi-agent systems like autonomous vehicles, there can be many  
 42 reasons why two scenes are similar, such as the kinematics of the ego agent (e.g. a significant juke),  
 43 atypical maneuvers (e.g. k-point turn), an agent’s interactions with nearby actors (e.g. yielding to a  
 44 jaywalker), or the road configuration itself. Importantly, the notion of similarity can differ depending  
 45 on who is using the retrieval system. For instance, an engineer focused on behavioral understanding  
 46 and motion prediction for other actors may define similarity by quantifying social influence between  
 47 actors. Conversely, an engineer focused on motion planner for the AV would likely key in on the  
 48 maneuver executed by the ego vehicle.

49 A second challenge is how to model and evaluate longer “information gathering” retrieval sessions  
 50 that can optionally include relevance feedback. Many of the use cases we have in mind do not fall into  
 51 the categories for short retrieval sessions such as navigational queries (e.g., the home page of Argo  
 52 AI) or very specific informational queries (e.g., the number of bridges in Pittsburgh) [7]. For instance  
 53 in Figure 1, the query trajectory (left) by itself is not enough to distinguish between two possible  
 54 intents (center, right) without any input from the user. We find in our experiments that catering to  
 55 multiple definitions of similarity simultaneously without any user feedback is very difficult for a  
 56 retrieval system, and in fact user feedback may be crucial for developing a practical retrieval system.

57 To summarize, our contributions are:

- 58 1. We present the Argoverse Trajectory Retrieval Benchmark, which enables studying trajectory  
 59 retrieval for multi-agent systems in a standardized way. We discuss design decisions and  
 60 provide a benchmark for public use at our dataset page.<sup>1</sup>
- 61 2. We establish a suite of baselines, including both hand-crafted feature-based approaches and  
 62 learned embeddings with state-of-the-art model architectures for AV trajectories.
- 63 3. We propose an initial retrieval system that leverages learned embeddings and conduct an  
 64 evaluation on both the standard and interactive (with relevance feedback) retrieval settings.
- 65 4. We conclude with a thorough discussion of our findings and directions for future work.

## 66 2 Related Work

67 **Information Retrieval.** Broadly speaking, information retrieval is the study of how to access specific  
 68 pieces of information within a data repository [29]. The canonical setting is: given a query, retrieve  
 69 a ranked list of (relevant) results. To date, information retrieval has been studied in many contexts,  
 70 including web search [7, 29], media retrieval (e.g., music or images) [31, 13, 49], and recently  
 71 in sports analytics [41, 14, 51]. Sports play retrieval is perhaps the most related to traffic scene

<sup>1</sup><https://github.com/ezhan94/argoverse-trajectory-retrieval-benchmark>

72 retrieval, although sports settings tend to be much more structured (fixed number of players, two  
73 teams, well-specified objectives, etc.). Furthermore, while some sports trajectory data is publicly  
74 available [44] there are currently no standardized retrieval datasets or benchmarks.

75 Our task is reminiscent of classic retrieval tasks that involve multiple intents or subtopics [32, 39, 55].  
76 In such tasks, two new considerations arose. The first is to be able to “cover” all the different intents  
77 or subtopics in order to have some minimal coverage over all intents in a single static ranking [55].  
78 The second is to study interactive ranking settings where users provide so-called relevance feedback  
79 [38, 57], after which the retrieval system responds by returning a modified ranking.

80 **Learning to Rank: Benchmarks & Methods.** Existing benchmarks for information retrieval largely  
81 fall under the category of “learning to rank”, where there is a set of supervised labels of the form  
82 (query, item, relevance level), in addition to a large repository of items [12, 35, 50]. Some datasets  
83 may also include information about global query types or genres [3], or query-specific information  
84 like intent and subtopics [32]. A related set of benchmarks is based on collaborative filtering, where  
85 one is also provided user information [2, 20].

86 This availability in data has led to significant interest in developing learning algorithms for retrieval  
87 (see [27] a broad overview). For retrieval over multiple or ambiguous intents, prior work includes  
88 learning for static rankings [54, 42] as well as for dynamic rankings that utilize relevance feedback  
89 [5, 52]. These prior work largely use engineered features based on text or metadata (e.g., URL), which  
90 can be hard to translate well to our setting. More recent methods that study continuous tracking data  
91 typically utilize learned embeddings [47, 21, 51], which we will also use to establish our baselines.

92 **Trajectory Datasets & Benchmarks.** The rapid growth of AV research opportunities has led to the  
93 release of many high-quality large-scale trajectory datasets. These datasets are typically focused on  
94 perception issues (detection, segmentation, and tracking) or on issues related to motion forecasting  
95 and have widely used benchmarks focused on these tasks. Significant examples include nuScenes  
96 [8], the Waymo Open Dataset [48, 15], Lyft Level 5 Dataset [23], and Argoverse [11]. We chose  
97 to build our retrieval benchmark on top of the Argoverse Motion Forecasting dataset because it has  
98 been widely used by the research community as evidenced by the active leaderboard with more than  
99 225 unique teams as of June 7, 2021. Trajectory benchmarks in other domains include behavior  
100 recognition, such as for laboratory animals [6, 17, 45] and human poses [37, 43]. Recent work by  
101 Segal et al. [40] is closely related to our proposed benchmark. Segal et al. proposed a method for  
102 learning spatio-temporal tags for driving scenes that could then be used for search, and present results  
103 on an internal dataset (SDVScenes).

104 **Trajectory Representation Learning & Modeling.** Modern research on trajectory modeling via  
105 representation learning has concentrated on forecasting of future behaviors (e.g., sequential generative  
106 modeling) [25, 56, 10, 18, 34], detection of pre-specified behavior categories (e.g., classification)  
107 [22, 1, 17], and open-ended knowledge extraction (e.g., unsupervised learning such as clustering)  
108 [4, 30, 16]. The study of methods for information access and retrieval of tracking data has received  
109 comparatively much less attention, with some exceptions for pose retrieval [47].

## 110 3 Contextual Multi-Intent Trajectory Retrieval

### 111 3.1 Problem Description

112 Let  $\tau$  denote a traffic scene trajectory, which can track multiple agents as well as contain contextual  
113 information (see Section 4.1). Let  $\mathcal{S}$  denote the set of possible intents, i.e. notions of similarity. A  
114 *query* is a trajectory-intent pair  $(\tau, s)$ ,  $s \in \mathcal{S}$ . Our retrieval task is to find and rank trajectories in a  
115 retrieval set  $\mathcal{R}$  that are similar to  $\tau$  with respect to intent  $s$ . We will denote  $\mathcal{Q}$  as the set of queries.  
116 The key challenge with our task is that the relevance of a retrieval depends on the intent  $s$ , but  $s$  is  
117 hidden from the retrieval system (see Figure 1 for an example). Furthermore, the set of intents  $\mathcal{S}$  is  
118 also not known ahead of time and can be extended to include new intents in the future.

### 119 3.2 Quantitative Evaluation

120 Retrieval systems will be evaluated on how well they rank the trajectories in  $\mathcal{R}$  for queries in  $\mathcal{Q}$ . Let  
121  $rel(\tau, s, \tau_q)$  be a scoring function that rates how relevant trajectory  $\tau$  is to query  $(\tau_q, s)$ , with higher  
122 scores being more relevant. The ranking metric we use to evaluate a ranked retrieval  $\{\tau_1, \dots, \tau_n\}$  is



Figure 2: High-level summary of the Argoverse trajectory data format. See [11] for complete details.

123 the normalized discounted cumulative gain (NDCG):

$$NDCG = \frac{DGC}{iDCG}, \quad DCG = \sum_{i=1}^n \frac{rel(\tau_i, s, \tau_q)}{\log_2(i+1)}, \quad (1)$$

124 where iDCG is computed with respect to the ideal/optimal ranking of  $n$  trajectories in  $\mathcal{R}$ . NDCG is  
 125 bounded between 0 and 1 is larger for retrievals that rank more relevant trajectories higher. We will  
 126 compute NDCG and average them over all queries in  $\mathcal{Q}$ .

### 127 3.3 Relevance Feedback

128 Achieving a high NDCG score can be difficult without knowing the hidden intent, as intents can have  
 129 very different meanings and correspond to different types of trajectories (Figure 1). To address this  
 130 challenge, we introduce one round of relevance feedback in our benchmark to allow retrieval systems  
 131 to infer the hidden intent, outlined below:

- 132 1. Query  $(\tau_q, s)$ , retrieval system receives  $\tau_q$ ,  $s$  is hidden.
- 133 2. Retrieval system returns initial set  $\{\tau_1, \dots, \tau_m\}$ .
- 134 3. Relevance feedback given to retrieval system  $\{rel(\tau_1, s, \tau_q), \dots, rel(\tau_m, s, \tau_q)\}$ .
- 135 4. Retrieval system returns new set  $\{\tau_1, \dots, \tau_n\}$ , which can have overlap with the initial set  
 136  $\{\tau_1, \dots, \tau_m\}$ , and is then scored with NDCG.

137 These steps simulate a user providing feedback to the retrieval system to allow it to hone in on the  
 138 hidden intent. In principle, multiple rounds of feedback are possible, but our benchmark will only  
 139 include one. Instructions for this step will be provided on our dataset page (see appendix).

## 140 4 The Argoverse Trajectory Retrieval Benchmark

141 We design our benchmark with the following goals in mind:

- 142 1. Multi-intent trajectory retrieval is challenging in domains where data is plentiful, as there  
 143 can be many dimensions in which two trajectories are similar. To this end, we derive our  
 144 dataset from the Argoverse Motion Forecasting 1.1 dataset [11], a real-world dataset for  
 145 trajectory forecasting that contains rich map information with each trajectory (see Figure 2).  
 146 We describe this process in Section 4.1.
- 147 2. Our retrieval task is already very challenging even for simple notions of similarity, so we  
 148 consider simple intents with scoring functions  $rel(\tau, s, \tau_q) = rel(\tau, s)$  to focus on whether  
 149 or not we’re retrieving trajectories for the right intent (the original query trajectory will not  
 150 affect the score). We describe the labeling process for our intents in Section 4.2. Future  
 151 iterations of our dataset can consider more complex intents.
- 152 3. Lastly, we highlight that the set of intents  $\mathcal{S}$  is not fixed. As more data is obtained and  
 153 annotated (e.g. maps for drivable areas, maps for ground height, etc.), new intents will  
 154 ultimately be introduced. Ideally, retrieval systems should adapt and be somewhat robust to  
 155 new intents. To simulate this scenario, we select a subset our intents to only appear in the  
 156 test query set, described in Section 4.3.

157 The Argoverse Trajectory Retrieval Benchmark dataset will consist of train/test query sets  $\mathcal{Q}_{\text{train}}/\mathcal{Q}_{\text{test}}$ ,  
 158 train/test retrieval sets  $\mathcal{R}_{\text{train}}/\mathcal{R}_{\text{test}}$ , the intent set  $\mathcal{S}$ , and relevance labels  $rel(\tau, s)$ . We summarize key  
 159 information about our dataset in Table 1 and Figure 3.

Intent in $\mathcal{S}$	All	$\mathcal{Q}_{\text{train}}$	$\mathcal{Q}_{\text{test}}$	$\mathcal{R}_{\text{train}}$	$\mathcal{R}_{\text{train}}$
Turn then change lanes	393	18	18	178	179
Straight then turn	354	18	19	164	153
Decelerate then turn	176	10	7	85	74
Turn then decelerate	133	39	10	38	46
Decelerate for stationary LV	251	25	10	115	101
Decelerate for moving LV	425	65	8	173	179
Decelerate to a stop	610	82	15	260	253
Decelerate after intersection	237	76	14	75	72
Test intent #1	228	0	12	(104)	112
Test intent #2	526	0	12	(258)	256
Test intent #3	815	0	21	(393)	401
Test intent #4	305	0	20	(136)	149
Test intent #5	245	0	17	(113)	115
Total # trajectories	2,795	100	50	1,323	1,322
Avg. # intents/trajectory	1.68	3.33	3.66	1.58	1.58
# trajectory-intent queries	n/a	321	170	n/a	n/a

Table 1: # of trajectories with each intent (counted if  $rel(\tau, s) > 0$ ) for all query and retrieval sets. LV = leading vehicle.  $\mathcal{Q}_{\text{train}}$  contains no trajectories with test intents while  $\mathcal{R}_{\text{train}}$  does, but the labels are not provided. Summary statistics are included in the last 3 rows. We only consider a trajectory-intent pair as a query if  $rel(\tau, s) = 2$ .

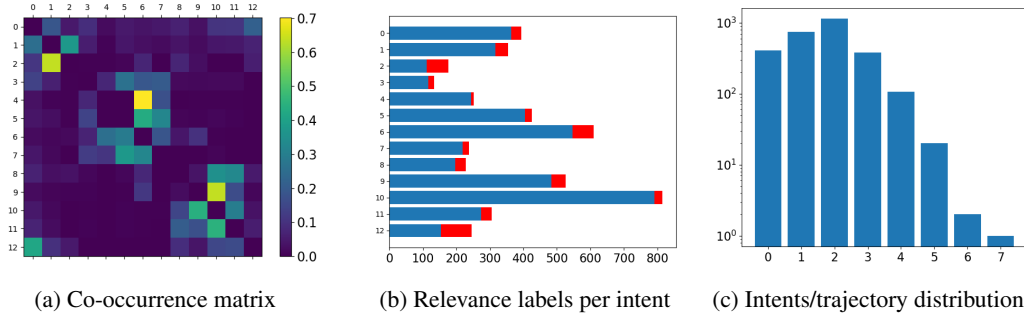


Figure 3: (a) Co-occurrence matrix of all intents of 2,795 trajectories. The matrix is not symmetric because it is row-normalized, i.e. cell  $o_{ij}$  is the percentage of trajectories with intent  $s_i$  that also have intent  $s_j$ . The order of intents is same as in Table 1. (b) Counts of relevance 2 (blue) and relevance 1 (red) labels for each intent. The order of intents is same as in Table 1. (c) Distribution of the # of intents per trajectory in log-scale. Trajectories have at most 7 intents in our dataset.

#### 160 4.1 Constructing the Dataset

161 Our dataset is derived from the Argoverse Motion Forecasting 1.1 dataset [11], which extracts  
162 planar trajectories and centerlines from sequences of LiDAR and camera images (see Figure 2).  
163 Each trajectory is 5 seconds long and tracks  $K > 1$  agents at 10 Hz ( $T = 50$ ). We let  $\mathbf{x}_t^k \in \mathbb{R}^2$   
164 denote the  $k$ -th agent’s planar  $(x, y)$  coordinates at time  $t$ . Similarly, denote  $\mathbf{X}_t := \{\mathbf{x}_t^1, \dots, \mathbf{x}_t^K\}$ ,  
165  $\mathbf{X}^k := \{\mathbf{x}_1^k, \dots, \mathbf{x}_T^k\}$ , and  $\mathbf{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_T\} = \{\mathbf{X}_1, \dots, \mathbf{X}_K\}$ .  $\mathbf{X}^1$  will always denote the focal  
166 agent and is visualized in red, while all other agents in teal (see Figure 1). Each trajectory also  
167 contains contextual information  $\mathbf{C} = \{\mathbf{C}_1, \dots, \mathbf{C}_T\}$ , some of which may change over time (e.g.  
168 nearest centerline to focal agent) while others remain static (e.g. lane connectivity graph). We refer  
169 to the original Argoverse paper [11] for the complete details. In summary, our trajectories  $\tau$  consist  
170 of tracking information  $\mathbf{X}$  and contextual information  $\mathbf{C}$ :  $\tau = (\mathbf{X}, \mathbf{C})$ .

171 We filter the Argoverse Motion Forecasting validation set that initially contains 39,472 trajectories  
172 using automatic labeling functions to find "interesting" trajectories that contain more complex ma-  
173 neuvers and/or social interactions. We filter for features such as large acceleration, large deceleration,  
174 leading vehicles, traffic control, etc., and refine the validation set down to 2,975 trajectories that we  
175 label with intents in Section 4.2. This filtering step will be included with our code release.

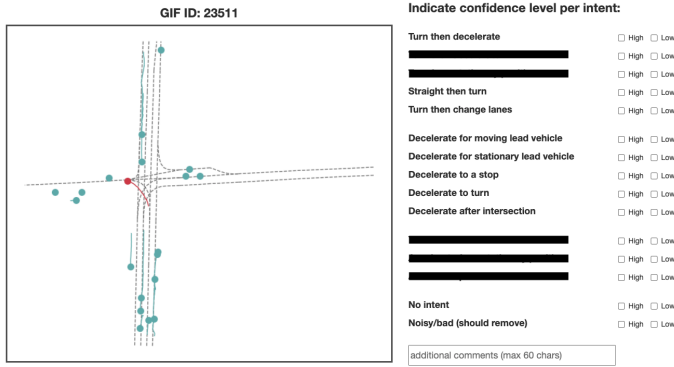


Figure 4: Interface for labeling trajectories. Annotators are shown a video of a trajectory on the left and asked to label the relevance of each intent on the right, including test intents (redacted in this figure). Options are high relevance (2), low relevance (1), or no selection (0). There is an additional option at the bottom to remove a trajectory.

## 176 4.2 Labeling Trajectories with Intents from $\mathcal{S}$

177 We focus on intents that describe or lead to complex behaviors, such as intents with an "A THEN  
 178 B" structure (e.g. turn THEN change lanes, then THEN decelerate) and intents that capture social  
 179 interaction (e.g. decelerate for leading vehicle). In total,  $\mathcal{S}$  contains 13 intents in  $\mathcal{S}$ , listed in Table 1.  
 180 We annotate all trajectory-intent pairs by labeling  $rel(\tau, s)$  as one of 3 degrees of relevance:  $\{0, 1, 2\}$   
 181 for  $\{\text{not, somewhat, highly}\}$  relevant. The labeling details are as follows:

- 182 1. Initially, 2 domain expert each labeled roughly half of the 2,975 filtered trajectories using  
 183 the interface depicted in Figure 4. An option to remove a trajectory was available to address  
 184 issues like trajectory jitter or over-segmentation.
- 185 2. Trajectories with more than 2 labeled intents were then labeled again by the other expert.
- 186 3. For trajectories that were labeled twice, we considered labels to be in agreement if they were  
 187 the same, were 0 and 1 (defaulted to 0) or were 1 and 2 (defaulted to 2). Roughly 80% of  
 188 double-labels were in agreement.
- 189 4. For labels that were in disagreement (0 and 2), the domain experts resolved them together.

190 Label statistics are summarized in Table 1 and visualized in Figure 3. 175 of the initial 2,970  
 191 trajectories were ultimately removed, bringing the final total to 2,795 trajectories. 1,144 trajectories  
 192 were labeled by both experts while 1,651 have a single set of labels. Labeling took 30hrs combined.

## 193 4.3 Query Selection and Train/Test Split

194 We first split the 2,795 labeled trajectories into query and retrieval sets. For queries, we manually  
 195 select 150 trajectories that have multiple intents (at least 2 intents with a high relevance label of  
 196 2) and such that the intents have good coverage over the intent set  $\mathcal{S}$ . The remaining trajectories  
 197 comprise the retrieval set.

198 Next, we split  $\mathcal{Q}$  and  $\mathcal{R}$  into train and test sets such that only 8 of the 13 intents appear in  $\mathcal{Q}_{\text{train}}$ , while  
 199 all 13 are represented in  $\mathcal{Q}_{\text{test}}$ . This simulates the real-world scenario of having to adapt to newly  
 200 encountered intents. Then we split the retrieval set into  $\mathcal{R}_{\text{train}}$  and  $\mathcal{R}_{\text{test}}$  such that they have similar  
 201 distributions over intents. Refer to Table 1 for full details about our query and retrieval sets.

202 Our dataset will provide all relevance labels  $rel(\tau, s)$  for the 8 train intents for trajectories in  $\mathcal{Q}_{\text{train}}$   
 203 and  $\mathcal{R}_{\text{train}}$ . Queries in  $\mathcal{Q}_{\text{test}}$  will be provided with masked intents, and retrieval systems will be  
 204 evaluated on how well they rank the trajectories in  $\mathcal{R}_{\text{test}}$  via NDCG score.

## 205 5 Baseline Experiments

206 Defining a similarity measure between two trajectories can be challenging due to dealing with many  
 207 modalities (maps, trajectories, logged metadata, etc.). Traditional methods that rely on feature  
 208 engineering and feature matching may have trouble scaling as more data is collected. Recent work  
 209 has instead focused on learning embedding functions that encode input data into a lower-dimensional

210 vector (or embedding) space. The advantage is that similarity can be more intuitively understood  
 211 as distance in embedding space, but we lose the ability to interpret what information is retained in  
 212 the embeddings. Nevertheless, learning trajectory embeddings have been shown to be effective for  
 213 many downstream tasks [24, 18, 46], and so we establish our retrieval baselines in this way. Our main  
 214 evaluation results are described Table 2, which we will discuss throughout this section.

## 215 5.1 Retrieval via Nearest Neighbors in Embedding Space

216 We define an embedding function as  
 217  $f_\theta$  parameterized by  $\theta$  that encodes a  
 218 trajectory  $\tau$  into a lower-dimensional  
 219 embedding vector  $\mathbf{z} = f_\theta(\tau)$ . The in-  
 220 formation retained in the embedding  
 221 ultimately depends on the auxiliary  
 222 task used to train the model (e.g. fore-  
 223 casting vs. autoencoding).  $f_\theta$  itself  
 224 can take on many forms and contain  
 225 multiple components, such as hand-crafted features, sliding window operations [53], recurrent neural  
 226 networks [24], and graph attention networks for capturing social interactions [18].

---

### Algorithm 1 Nearest-Neighbor( $(\tau_q, s), \mathcal{R}, n, \mathbf{d}(\cdot), f_\theta$ )

---

- 1: **Inputs:** query  $(\tau_q, s)$ , retrieval set  $\mathcal{R}$ , top- $n$  trajectories
  - 2: **Inputs:** distance function  $\mathbf{d}(\cdot)$ , embedding function  $f_\theta$
  - 3: Compute query embedding  $\mathbf{z}_q = f_\theta(\tau_q)$ .
  - 4: Compute embeddings  $\mathbf{z}_i = f_\theta(\tau_i)$  for  $\tau_i \in \mathcal{R}$ .
  - 5: Rank  $\tau_i \in \mathcal{R}$  in increasing order of  $\mathbf{d}(\mathbf{z}_i, \mathbf{z}_q)$ .
  - 6: **Output:** top- $n$  closest trajectories to query
- 

227 Given an embedding function  $f_\theta$ , we can design a ranked retrieval system that returns the top- $n$   
 228 trajectories in the retrieval set with embeddings closest to the query embedding with respect to  
 229 some distance function (a common choice is Euclidean distance:  $\mathbf{d}_{\text{Euclidean}}(\mathbf{z}_i, \mathbf{z}_j) = \|\mathbf{z}_i - \mathbf{z}_j\|_2$ ,  
 230 [41]). This algorithm is outlined in Algorithm 1 and has time complexity  $O(|\mathcal{R}| \log n)$  per query if  
 231 implemented with a heap. Note that the algorithm does not take into account a hidden intent and  
 232 always returns the same retrievals for each query trajectory.

233 We consider 4 embeddings functions  $f_\theta$ , described below:

- 234 1. **FEAT** - a naive embedding function that simply computes 15 domain-specific features, such  
 235 as average speed of the focal agent, the curvature of its trajectory, and its distances to other  
 236 agents. There are no parameters to be learned for this embedding function.
- 237 2. **AE** - a simple autoencoder for the focal agent trajectory  $\mathbf{X}^1$  implemented with a recur-  
 238 rent neural network for both the encoder and decoder. Other agents  $\{\mathbf{X}^2, \dots, \mathbf{X}^K\}$  and  
 239 contextual information  $\mathbf{C}$  are ignored.
- 240 3. **WIMP** [24] - a state-of-the-art model for trajectory forecasting that encodes all agents as  
 241 well as the nearest centerline to the focal agent. We use the same model architecture but  
 242 train it to reconstruct the focal agent trajectory  $\mathbf{X}^1$ .
- 243 4. **VNET** [18] - VectorNet, another state-of-the-art model trained for trajectory forecasting  
 244 (forecast next 3sec given a 2sec history) that includes a graph attention network for encoding  
 245 all contextual map information and also a node reconstruction task in its objective. We use  
 246 the node embedding for the focal agent full 5sec trajectory.

247 The embedding functions are developed using the Argoverse Motion Forecasting 1.1 training set.  
 248 We evaluate nearest-neighbor retrieval using NDCG and the query and retrieval sets constructed in  
 249 Section 4.3 and report our results in Table 2 (rows with  $m = 0$ , “standard” columns). We observe  
 250 that NDCG decreases as the number of retrievals  $n$  increases because retrieving a larger optimal  
 251 set is more difficult. Out of all the embeddings, WIMP performs the best. We hypothesize that  
 252 this is the case because WIMP is trained to reconstruct the focal agent trajectory and all queries  
 253 pertain to said focal agent. On the other hand, VectorNet is trained for trajectory forecasting and  
 254 performs the worst. We reason that this occurs because VectorNet embeddings must retain some  
 255 information about possible futures (and also information for node completion), which can be irrelevant  
 256 for comparing embeddings of trajectory histories. We note that the hand-crafted FEAT embedding  
 257 performs reasonably well, although noticeably worse than the best learned embedding. We conclude  
 258 that embeddings trained for trajectory reconstruction are better suited for our retrieval task.

## 259 5.2 Triplet Loss Fine-tuning with $\mathcal{Q}_{\text{train}}, \mathcal{R}_{\text{train}}$

260 In our next set of experiments, we use the relevance labels given in  $\mathcal{Q}_{\text{train}}$  and  $\mathcal{R}_{\text{train}}$  to fine-tune  
 261 embeddings with a triplet loss. Our motivation is that having trajectories with the same intent labels

Query Set	NDCG		standard (Section 5.1)				triplet fine-tuning (Section 5.2)			
	$n$	$m$	FEAT	AE	WIMP	VNET	FEAT	AE	WIMP	VNET
$\mathcal{Q}_{\text{train}}$	10	0	.379	.349	.391	.230	.385	.370	.385	.243
	30	0	.352	.334	.374	.231	.367	.359	.376	.238
	50	0	.345	.330	.368	.231	.361	.358	.371	.236
$\mathcal{Q}_{\text{train}}$	10	5	.411	.425	.410	.236	.399	.429	.414	.256
	30	5	.365	.367	.366	.209	.352	.377	.375	.219
	50	5	.343	.346	.348	.203	.336	.361	.360	.203
$\mathcal{Q}_{\text{test}}$	10	0	.337	.355	.371	.273	.334	.331	.355	.273
	30	0	.310	.324	.343	.261	.318	.318	.336	.257
	50	0	.305	.310	.328	.254	.311	.313	.331	.250
$\mathcal{Q}_{\text{test}}$	10	5	.409	.429	.436	.236	.378	.394	.396	.272
	30	5	.353	.373	.390	.208	.339	.359	.365	.246
	50	5	.332	.343	.367	.202	.324	.343	.351	.238

Table 2: NDCG scores for queries in  $\mathcal{Q}_{\text{train}}$ ,  $\mathcal{Q}_{\text{test}}$  and retrievals from  $\mathcal{R}_{\text{train}}$ ,  $\mathcal{R}_{\text{test}}$  respectively.  $n$  is the # of retrievals,  $m$  is the # of trajectories for relevance feedback. 1) NDCG decreases as  $n$  increases, as retrieving a larger optimal set is more difficult. 2) Utilizing relevance feedback leads to clear improvement for all embeddings except VNET. 3) Triplet fine-tuning does *not* lead a clear improvement. 4) There is generally not a big difference in performance between train and test queries, but the difference is larger for fine-tuned embeddings, possibly because of overfitting to training intents. 5) Overall, WIMP embeddings without fine-tuning appear to be the best for our retrieval task.

262 closer together in embedding space will improve nearest neighbor retrieval.<sup>2</sup> In particular, we train  
263 an autoencoder ( $\mathbf{g}_{\text{enc}}, \mathbf{g}_{\text{dec}}$ ) that minimizes the following objective:

$$\max(\underbrace{\|\mathbf{g}_{\text{enc}}(\mathbf{z}) - \mathbf{g}_{\text{enc}}(\mathbf{z}_{\text{pos}})\|_2 - \|\mathbf{g}_{\text{enc}}(\mathbf{z}) - \mathbf{g}_{\text{enc}}(\mathbf{z}_{\text{neg}})\|_2}_{\text{triplet loss}} + \alpha, 0) + \underbrace{\|\mathbf{z} - \mathbf{g}_{\text{dec}}(\mathbf{g}_{\text{enc}}(\mathbf{z}))\|_2}_{\text{reconstruction loss}}. \quad (2)$$

264  $(\mathbf{z}, \mathbf{z}_{\text{pos}}, \mathbf{z}_{\text{neg}})$  is a triplet of embeddings where  $\mathbf{z}, \mathbf{z}_{\text{pos}}$  share the same label while  $\mathbf{z}, \mathbf{z}_{\text{neg}}$  do not. The  
265 triplet loss in (2) encourages embeddings with the same label to be closer together than embeddings  
266 with different labels, up to some margin  $\alpha$ . At the same time, we aim to retain the same information  
267 encoded in the original embeddings by including the standard autoencoder reconstruction loss in (2).

268 We construct triplets  $(\mathbf{z}, \mathbf{z}_{\text{pos}}, \mathbf{z}_{\text{neg}})$  by considering every trajectory-intent pair  $(\tau, s)$  with  $rel(\tau, s) = 2$   
269 in  $\mathcal{Q}_{\text{train}} \cup \mathcal{R}_{\text{train}}$ . For each pair, we sample a positive trajectory  $\tau_{\text{pos}}$  from those in  $\mathcal{Q}_{\text{train}} \cup \mathcal{R}_{\text{train}}$  that  
270 share the same label ( $rel(\tau_{\text{pos}}, s) = 2$ ), and similarly we sample a negative trajectory ( $rel(\tau_{\text{neg}}, s) = 0$ ).  
271 Triplets are re-sampled at the beginning of every epoch (e.g. offline triplet mining).

272 The new embedding function we use for our retrieval system is then  $\mathbf{z} = \mathbf{g}_{\text{enc}}(\mathbf{f}_{\theta}(\tau))$  and we report  
273 our results in Table 2 (“triplet fine-tuning” columns). We observe that the results are inconsistent:  
274 NDCG can both increase/decrease compared to the “standard” embedding columns. Furthermore, we  
275 see that there is generally a drop in performance on the query test set, which likely occurs because  
276 the query test set contains test intents that were not fine-tuned with our triplet loss. We note that our  
277 fine-tuning step is applied after training the initial embeddings so there might be some information  
278 loss (that we tried to mitigate with the autoencoding loss in (2)). Future work should consider jointly  
279 training embeddings with the triplet loss.

### 280 5.3 Retrieval with Relevance Feedback

281 Our previous two experiments ignore a main challenge of our problem setting by disregarding that  
282 there is a hidden intent and will always return the same set of trajectories for each query. In our final  
283 experiment, we design a retrieval system that utilizes the relevance feedback procedure described in  
284 Section 3.3 to address this challenge. We consider a version of nearest neighbor retrieval in Algorithm  
285 1 that uses an updated distance function given the relevance feedback, as described in Algorithm 2.

286 Let  $(\tau_q, s)$  be our initial query and  $\mathcal{M} = \{\tau_1, \dots, \tau_m\}$  be our initial set of  $m$  retrievals for which  
287 we receive relevance feedback  $\{rel(\tau_1, s), \dots, rel(\tau_m, s)\}$ . We construct two sets: relevant set

<sup>2</sup>Indeed, triplet or contrast loss has been used in other related retrieval settings, such as for human poses [47].



288  $\mathcal{A} = \{\tau | rel(\tau, s) > 0, \tau \in \mathcal{M}\} \cup \{\tau_q\}$  and non-relevant set  $\mathcal{B} = \{\tau | rel(\tau, s) = 0, \tau \in \mathcal{M}\}$ . We  
 289 consider an updated distance function that prioritizes trajectories with embeddings close to the  
 290 relevant set and far from the non-relevant set:

$$\mathbf{d}_{\mathcal{AB}}(\mathbf{z}, \mathbf{z}_q) = \frac{1}{|\mathcal{A}|} \sum_{\tau \in \mathcal{A}} \mathbf{d}(\mathbf{z}, \mathbf{f}_\theta(\tau)) - \frac{1}{|\mathcal{B}|} \sum_{\tau \in \mathcal{B}} \mathbf{d}(\mathbf{z}, \mathbf{f}_\theta(\tau)). \quad (3)$$

291 (3) is reminiscent of the Rocchio algorithm [28] except we compute the distances to the relevant set  
 292 rather than update the query embedding directly. Algorithm 2 summarizes our approach incorporating  
 293 relevance feedback and has time complexity  $O(|\mathcal{R}|(\log n + \log m))$  per query.

294 We observe in Table 2 (rows with  $m = 5$ ) that leveraging relevance feedback improves NDCG for all  
 295 embeddings (except VectorNet). This matches our intuition because our approach in Algorithm 2 uses  
 296 feedback given by the simulated user to refine the retrieval for the hidden intent. These results suggest  
 297 user feedback, even a limited amount, can be crucial for efficient multi-intent trajectory retrieval.

---

**Algorithm 2** Nearest-Neighbor-with-Relevance-Feedback( $(\tau_q, s), \mathcal{R}, n, m, \mathbf{d}(\cdot), \mathbf{f}_\theta$ )

---

- 1: **Inputs:** query  $(\tau_q, s)$ , retrieval set  $\mathcal{R}$ , top- $n$  trajectories,  $m$  feedback
  - 2: **Inputs:** distance function  $\mathbf{d}(\cdot)$ , embedding function  $\mathbf{f}_\theta$
  - 3:  $\mathcal{M} = \text{Nearest-Neighbor}((\tau_q, s), \mathcal{R}, m, \mathbf{d}(\cdot), \mathbf{f}_\theta)$  using Algorithm 1.
  - 4: Receive relevance feedback for trajectories in  $\mathcal{M}$ .
  - 5: Construct sets  $\mathcal{A}, \mathcal{B}$ , and update distance function  $\mathbf{d}_{\mathcal{AB}}$  in (3).
  - 6: **Output:** Nearest-Neighbor( $(\tau_q, s), \mathcal{R}, n, \mathbf{d}_{\mathcal{AB}}(\cdot), \mathbf{f}_\theta$ ) using Algorithm 1.
- 

## 298 6 Discussion and Future Work

299 We have introduced the Argoverse Trajectory Retrieval Benchmark for standardizing the challenging  
 300 task of multi-intent retrieval in the domain of AV trajectories. We explore initial baseline retrieval  
 301 algorithms that use trajectory embeddings and summarize our findings:

- 302 1. Embeddings trained to reconstruct rather than forecast the focal agent trajectory are better-  
 303 suited for queries that pertain to the focal agent (Section 5.1).
- 304 2. Triplet loss fine-tuning with relevance labels does not appear to be effective, but a joint  
 305 training approach has yet to be explored (Section 5.2).
- 306 3. Incorporating relevance feedback may be key for this retrieval setting (Section 5.3).

307 Our benchmark is the first iteration of what we expect to be a promising research area. There are  
 308 many directions for future work and many more challenges to overcome as we continue to scale  
 309 up. For instance, the trajectory data provided in Argoverse is only a small subset of the data that’s  
 310 available, such as richer map information like ground height, agent type (vehicle vs. pedestrian), and  
 311 the state of traffic control. As more data is incorporated, intents will grow in number and complexity  
 312 and retrieval systems may fail to scale accordingly. Another direction for future work is to consider  
 313 more diverse queries beyond those that pertain to the focal agent, as embeddings trained to reconstruct  
 314 the focal agent trajectory is unlikely to be the best solution for all query types. Potential solutions may  
 315 use multiple embeddings trained with different auxiliary tasks within their retrieval systems. A third  
 316 direction is explore other forms of relevance feedback, such as pairwise comparisons or ranking an  
 317 initial retrieval set. It is unclear what form of relevance feedback is the most informative for retrieval  
 318 systems and also easy for users to provide.

319 Ultimately, further progress in this research direction will come from scaling up our benchmark.  
 320 For instance, approaches may overfit to our set of intents that all pertain to the focal agent. We try  
 321 to prevent this by having held-out test intents, and we also expect future versions of our dataset to  
 322 include more diverse queries. Lastly, it’s important to understand that the usefulness of retrieval  
 323 systems is tied to the underlying data and can be subject to biases of the data. Thus, some scenarios  
 324 may intrinsically be harder to retrieve than others. Diagnosing biases in retrieval systems could be  
 325 another interesting direction for future work.

## References

- 326  
327 [1] David J Anderson and Pietro Perona. Toward a science of computational ethology. *Neuron*, 84(1):18–31,  
328 2014.
- 329 [2] Robert M Bell and Yehuda Koren. Lessons from the netflix prize challenge. *Acm Sigkdd Explorations*  
330 *Newsletter*, 9(2):75–79, 2007.
- 331 [3] Thierry Bertin-Mahieux, Daniel P.W. Ellis, Brian Whitman, and Paul Lamere. The million song dataset. In  
332 *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011)*, 2011.
- 333 [4] Alina Bialkowski, Patrick Lucey, Peter Carr, Yisong Yue, Sridha Sridharan, and Iain Matthews. Large-scale  
334 analysis of soccer matches using spatiotemporal tracking data. In *2014 IEEE international conference on*  
335 *data mining*, pages 725–730. IEEE, 2014.
- 336 [5] Christina Brandt, Thorsten Joachims, Yisong Yue, and Jacob Bank. Dynamic ranked retrieval. In  
337 *Proceedings of the fourth ACM international conference on Web search and data mining*, pages 247–256,  
338 2011.
- 339 [6] Kristin Branson, Alice A Robie, John Bender, Pietro Perona, and Michael H Dickinson. High-throughput  
340 ethomics in large groups of drosophila. *Nature methods*, 6(6):451–457, 2009.
- 341 [7] Andrei Broder. A taxonomy of web search. In *ACM Sigir forum*, volume 36, pages 3–10. ACM New York,  
342 NY, USA, 2002.
- 343 [8] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan,  
344 Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In  
345 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631,  
346 2020.
- 347 [9] Julen Castellano, David Alvarez-Pastor, and Paul S Bradley. Evaluation of research using computerised  
348 tracking systems (amisco® and prozone®) to analyse physical performance in elite soccer: A systematic  
349 review. *Sports medicine*, 44(5):701–712, 2014.
- 350 [10] Rohan Chandra, Tianrui Guan, Srujan Panuganti, Trisha Mittal, Uttaran Bhattacharya, Aniket Bera,  
351 and Dinesh Manocha. Forecasting trajectory and behavior of road-agents using spectral clustering in  
352 graph-lstms. *IEEE Robotics and Automation Letters*, 5(3):4882–4890, 2020.
- 353 [11] Ming-Fang Chang, John Lambert, Patsorn Sangkloy, Jagjeet Singh, Slawomir Bak, Andrew Hartnett,  
354 De Wang, Peter Carr, Simon Lucey, Deva Ramanan, et al. Argoverse: 3d tracking and forecasting with  
355 rich maps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
356 pages 8748–8757, 2019.
- 357 [12] Olivier Chapelle and Yi Chang. Yahoo! learning to rank challenge overview. In *Proceedings of the learning*  
358 *to rank challenge*, pages 1–24. PMLR, 2011.
- 359 [13] Ritendra Datta, Dhiraj Joshi, Jia Li, and James Z Wang. Image retrieval: Ideas, influences, and trends of  
360 the new age. *ACM Computing Surveys (Csur)*, 40(2):1–60, 2008.
- 361 [14] Mingyang Di, Diego Klabjan, Long Sha, and Patrick Lucey. Large-scale adversarial sports play retrieval  
362 with learning to rank. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 12(6):1–18, 2018.
- 363 [15] Scott Ettinger, Shuyang Cheng, Benjamin Caine, Chenxi Liu, Hang Zhao, Sabeek Pradhan, Yuning Chai,  
364 Ben Sapp, Charles Qi, Yin Zhou, et al. Large scale interactive motion forecasting for autonomous driving:  
365 The waymo open motion dataset. *arXiv preprint arXiv:2104.10133*, 2021.
- 366 [16] Eyrun Eyjolfssdottir, Kristin Branson, Yisong Yue, and Pietro Perona. Learning recurrent representations  
367 for hierarchical behavior modeling. In *International Conference on Learning Representations*, 2017.
- 368 [17] Eyrun Eyjolfssdottir, Steve Branson, Xavier P Burgos-Artizzu, Eric D Hoopfer, Jonathan Schor, David J  
369 Anderson, and Pietro Perona. Detecting social actions of fruit flies. In *European Conference on Computer*  
370 *Vision*, pages 772–787. Springer, 2014.
- 371 [18] Jiyang Gao, Chen Sun, Hang Zhao, Yi Shen, Dragomir Anguelov, Congcong Li, and Cordelia Schmid.  
372 Vectornet: Encoding hd maps and agent dynamics from vectorized representation. In *Proceedings of the*  
373 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11525–11533, 2020.
- 374 [19] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna M. Wallach,  
375 Hal Daumé III, and Kate Crawford. Datasheets for datasets. abs/1803.09010, 2018.

- 376 [20] F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *Acm transactions*  
377 *on interactive intelligent systems (tiis)*, 5(4):1–19, 2015.
- 378 [21] Chih-Hui Ho, Pedro Morgado, Amir Persekian, and Nuno Vasconcelos. Pies: Pose invariant embeddings.  
379 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12377–  
380 12386, 2019.
- 381 [22] Mayank Kabra, Alice A Robie, Marta Rivera-Alba, Steven Branson, and Kristin Branson. Jaaba: interactive  
382 machine learning for automatic annotation of animal behavior. *Nature methods*, 10(1):64, 2013.
- 383 [23] R Kesten, M Usman, J Houston, T Pandya, K Nadhamuni, A Ferreira, M Yuan, B Low, A Jain, P Ondruska,  
384 et al. Lyft level 5 perception dataset 2020, 2019.
- 385 [24] Siddhesh Khandelwal, William Qi, Jagjeet Singh, Andrew Hartnett, and Deva Ramanan. What-if motion  
386 prediction for autonomous driving. arXiv preprint arXiv:2008.10587, 2020.
- 387 [25] Hoang M Le, Peter Carr, Yisong Yue, and Patrick Lucey. Data-driven ghosting using deep imitation  
388 learning. In *MIT Sloan Sports Analytics Conference*, 2017.
- 389 [26] Alon Lerner, Yiorgos Chrysanthou, and Dani Lischinski. Crowds by example. In *Computer graphics*  
390 *forum*, volume 26, pages 655–664. Wiley Online Library, 2007.
- 391 [27] Tie-Yan Liu. Learning to rank for information retrieval. 2011.
- 392 [28] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. An introduction to information  
393 retrieval. pages 163–167, 2009.
- 394 [29] Christopher D Manning, Hinrich Schütze, and Prabhakar Raghavan. *Introduction to information retrieval*.  
395 Cambridge university press, 2008.
- 396 [30] Andrew Miller, Luke Bornn, Ryan Adams, and Kirk Goldsberry. Factorized point process intensities:  
397 A spatial analysis of professional basketball. In *International conference on machine learning*, pages  
398 235–243. PMLR, 2014.
- 399 [31] Nicola Orio. Music retrieval: A tutorial and review. 2006.
- 400 [32] Paul Over. The trec interactive track: an annotated bibliography. *Information Processing & Management*,  
401 37(3):369–381, 2001.
- 402 [33] Stefano Pellegrini, Andreas Ess, and Luc Van Gool. Improving data association by joint modeling  
403 of pedestrian trajectories and groupings. In *European conference on computer vision*, pages 452–465.  
404 Springer, 2010.
- 405 [34] Tung Phan-Minh, Elena Corina Grigore, Freddy A Boulton, Oscar Beijbom, and Eric M Wolff. Covernet:  
406 Multimodal behavior prediction using trajectory sets. In *Proceedings of the IEEE/CVF Conference on*  
407 *Computer Vision and Pattern Recognition*, pages 14074–14083, 2020.
- 408 [35] Tao Qin, Tie-Yan Liu, Jun Xu, and Hang Li. Letor: A benchmark collection for research on learning to  
409 rank for information retrieval. *Information Retrieval*, 13(4):346–374, 2010.
- 410 [36] Alexandre Robicquet, Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. Learning social etiquette:  
411 Human trajectory understanding in crowded scenes. In *European conference on computer vision*, pages  
412 549–565. Springer, 2016.
- 413 [37] Matteo Ruggero Ronchi, Joon Sik Kim, and Yisong Yue. A rotation invariant latent factor model for  
414 moveme discovery from static poses. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*,  
415 pages 1179–1184. IEEE, 2016.
- 416 [38] Gerard Salton and Chris Buckley. Improving retrieval performance by relevance feedback. *Journal of the*  
417 *American society for information science*, 41(4):288–297, 1990.
- 418 [39] Rodrygo LT Santos, Craig Macdonald, and Iadh Ounis. Intent-aware search result diversification. In *ACM*  
419 *SIGIR conference on Research and development in Information Retrieval*, pages 595–604, 2011.
- 420 [40] Sean Segal, Eric Kee, Wenjie Luo, Abbas Sadat, Ersin Yumer, and Raquel Urtasun. Universal embeddings  
421 for spatio-temporal tagging of self-driving logs. arXiv preprint arXiv:2011.06165, 2020.
- 422 [41] Long Sha, Patrick Lucey, Yisong Yue, Peter Carr, Charlie Rohlf, and Iain Matthews. Chalkboarding: A  
423 new spatiotemporal query paradigm for sports play retrieval. In *Proceedings of the 21st International*  
424 *Conference on Intelligent User Interfaces*, pages 336–347, 2016.

- 425 [42] Ruben Sipos, Pannaga Shivaswamy, and Thorsten Joachims. Large-margin learning of submodular  
426 summarization models. In *Proceedings of the 13th Conference of the European Chapter of the Association  
427 for Computational Linguistics*, pages 224–233, 2012.
- 428 [43] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions  
429 classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- 430 [44] LLC Stats. Stats sportvu basketball player tracking. *SportVU website, available at: <https://www.stats.com/sportvu-basketball/>, last accessed: Feb, 12, 2019.*
- 432 [45] Jennifer J Sun, Tomomi Karigo, Dipam Chakraborty, Sharada P Mohanty, David J Anderson, Pietro Perona,  
433 Yisong Yue, and Ann Kennedy. The multi-agent behavior dataset: Mouse dyadic social interactions. *arXiv  
434 preprint arXiv:2104.02710*, 2021.
- 435 [46] Jennifer J Sun, Ann Kennedy, Eric Zhan, David J Anderson, Yisong Yue, and Pietro Perona. Task program-  
436 ming: Learning data efficient behavior representations. In *Proceedings of the IEEE/CVF Conference on  
437 Computer Vision and Pattern Recognition*, 2021.
- 438 [47] Jennifer J Sun, Jiaping Zhao, Liang-Chieh Chen, Florian Schroff, Hartwig Adam, and Ting Liu. View-  
439 invariant probabilistic embedding for human pose. In *European Conference on Computer Vision*, pages  
440 53–70. Springer, 2020.
- 441 [48] Pei Sun, Henrik Kretschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James  
442 Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving:  
443 Waymo open dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
444 Recognition*, pages 2446–2454, 2020.
- 445 [49] Douglas Turnbull, Luke Barrington, David Torres, and Gert Lanckriet. Semantic annotation and retrieval of  
446 music and sound effects. *IEEE Transactions on Audio, Speech, and Language Processing*, 16(2):467–476,  
447 2008.
- 448 [50] Ellen M Voorhees, Donna K Harman, et al. *TREC: Experiment and evaluation in information retrieval*,  
449 volume 63. MIT press Cambridge, MA, 2005.
- 450 [51] Zheng Wang, Cheng Long, Gao Cong, and Ce Ju. Effective and efficient sports play retrieval with deep  
451 representation learning. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge  
452 Discovery & Data Mining*, pages 499–509, 2019.
- 453 [52] Zuobing Xu, Ram Akella, and Yi Zhang. Incorporating diversity and density in active learning for relevance  
454 feedback. In *European Conference on Information Retrieval*, pages 246–257. Springer, 2007.
- 455 [53] Di Yao, Chao Zhang, Zhihua Zhu, Jianhui Huang, and Jingping Bi. Trajectory clustering via deep  
456 representation learning. In *2017 international joint conference on neural networks (IJCNN)*, pages  
457 3880–3887. IEEE, 2017.
- 458 [54] Yisong Yue and Thorsten Joachims. Predicting diverse subsets using structural svms. In *Proceedings of  
459 the 25th international conference on Machine learning*, pages 1224–1231, 2008.
- 460 [55] ChengXiang Zhai, William W Cohen, and John Lafferty. Beyond independent relevance: methods and  
461 evaluation metrics for subtopic retrieval. In *ACM SIGIR Forum*, volume 49, pages 2–9. ACM New York,  
462 NY, USA, 2015.
- 463 [56] Eric Zhan, Stephan Zheng, Yisong Yue, Long Sha, and Patrick Lucey. Generating multi-agent trajectories  
464 using programmatic weak supervision. In *International Conference on Learning Representations*, 2019.
- 465 [57] Xiang Sean Zhou and Thomas S Huang. Relevance feedback in image retrieval: A comprehensive review.  
466 *Multimedia systems*, 8(6):536–544, 2003.

467 **Checklist**

- 468 1. For all authors...
- 469 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
470 contributions and scope? [Yes]
- 471 (b) Did you describe the limitations of your work? [Yes] See Section 6.
- 472 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See  
473 Section 6.
- 474 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
475 them? [Yes]
- 476 2. If you are including theoretical results...
- 477 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 478 (b) Did you include complete proofs of all theoretical results? [N/A]
- 479 3. If you ran experiments (e.g. for benchmarks)...
- 480 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
481 mental results (either in the supplemental material or as a URL)? [Yes] All code and  
482 data will be provided at our dataset page once released.
- 483 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
484 were chosen)? [Yes] All training details will be specified in our dataset page.
- 485 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
486 ments multiple times)? [No] We present an initial set of experiments. We will include  
487 any additional results in our dataset page.
- 488 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
489 of GPUs, internal cluster, or cloud provider)? [Yes] These details will be specified in  
490 our dataset page.
- 491 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 492 (a) If your work uses existing assets, did you cite the creators? [Yes] Our dataset is derived  
493 from the Argoverse Motion Forecasting 1.1 dataset [11].
- 494 (b) Did you mention the license of the assets? [Yes] See our dataset page.
- 495 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
496 These will be available at our dataset page once released.
- 497 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
498 using/curating? [Yes]
- 499 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
500 information or offensive content? [Yes]
- 501 5. If you used crowdsourcing or conducted research with human subjects...
- 502 (a) Did you include the full text of instructions given to participants and screenshots, if  
503 applicable? [N/A]
- 504 (b) Did you describe any potential participant risks, with links to Institutional Review  
505 Board (IRB) approvals, if applicable? [N/A]
- 506 (c) Did you include the estimated hourly wage paid to participants and the total amount  
507 spent on participant compensation? [N/A]

## 508 A Key Information

509 **Dataset page:** <https://github.com/ezhan94/argoverse-trajectory-retrieval-benchmark>.  
510 All relevant information can be found at our dataset page linked above (dataset download, code,  
511 license, instructions for submitting a benchmark, additional supplementary materials, etc.)

512 **Dataset documentation and intended uses:** we use the datasheets for datasets framework [19] in  
513 Appendix B.

514 **Author statement:** We bear all responsibility in case of violation of rights, etc., and confirmation of  
515 the data license.

516 **Hosting, licensing, and maintenance plan:** This information will be provided on our dataset page.

## 517 B Datasheets for Datasets [19]

### 518 B.1 Motivation

- 519 • **For what purpose was the dataset created?** The task of finding “similar” scenes or  
520 trajectories within a large corpus of log data has proven challenging. Existing “learning to  
521 rank” systems do not readily port to this trajectory domain. This dataset was created to enable  
522 and encourage further research on trajectory retrieval in the setting of AV development.
- 523 • **Who created the dataset (e.g., which team, research group) and on behalf of which**  
524 **entity (e.g., company, institution, organization)?** The prediction team at Argo AI in  
525 collaboration with Caltech.
- 526 • **Who funded the creation of the dataset?** Argo AI.
- 527 • **Any other comments?** None.

### 528 B.2 Composition

- 529 • **What do the instances that comprise the dataset represent (e.g., documents, photos,**  
530 **people, countries)?** In this work we add relevance labels to a subset scenarios of the  
531 Argoverse Motion Forecasting validation set. The underlying scenarios represent the planar  
532 centroid positions of actors in a traffic scene. Each 5s (10Hz) scenario has been derived  
533 from a AV log and contains at least one agent that is present for the entire 5s and performs a  
534 significant action.
- 535 • **How many instances are there in total (of each type, if appropriate)?** We add 13  
536 relevance labels to 2,795 scenarios. Label statistics are shown in Table 1 and Figure 3.
- 537 • **Does the dataset contain all possible instances or is it a sample (not necessarily ran-**  
538 **dom) of instances from a larger set?** We label 2,795 of the 39,472 scenarios comprising  
539 the Argoverse 1.1 Motion Forecasting validation set. 2,970 scenarios were selected using  
540 automatic labeling functions to find “interesting” trajectories that contain more complex  
541 maneuvers or social interactions. We detect features like the presence of acceleration, decel-  
542 eration, leading vehicles, traffic control, etc. 175 scenarios were removed during labeling.  
543 These scenarios eliminated for tracking errors such as id-swaps or over-segmentation of the  
544 focal track.
- 545 • **What data does each instance consist of?** Each Argoverse scenario consists of planar  
546 centroid positions for actors in a traffic scene. These centroids are sampled at 10Hz and  
547 the full duration of the scene is 5s. A lane graph and underlying lane centerlines are also  
548 provided. Here we add relevance labels  $\in \{0, 1, 2\}$  for each of 13 intents to each selected  
549 scenario.
- 550 • **Is there a label or target associated with each instance?** For each of 2,795 there are 13  
551 relevance labels associated with the underlying scenario.
- 552 • **Is any information missing from individual instances?** Relevance labels corresponding  
553 to 5 of the 13 intents are hidden for all training examples. All test set labels are also hidden.
- 554 • **Are relationships between individual instances made explicit (e.g., users’ movie rat-**  
555 **ings, social network links)?** Not applicable.

- 556 • **Are there recommended data splits (e.g., training, development/validation, testing)?**  
557 Yes. Items are split into test queries, test retrievals
- 558 • **Are there any errors, sources of noise, or redundancies in the dataset?** The underlying  
559 Argoverse scenarios represent a real urban driving dataset; there is an expected degree of  
560 tracking noise and segmentation errors. Relevance labels were provided by domain experts  
561 but nevertheless may contain noise due to human error or subjective judgement.
- 562 • **Is the dataset self-contained, or does it link to or otherwise rely on external resources**  
563 **(e.g., websites, tweets, other datasets)?** The retrieval benchmark is built upon another  
564 existing dataset. However, the retrieval benchmark labels will be hosted with the requisite  
565 forecasting scenarios.
- 566 • **Does the dataset contain data that might be considered confidential (e.g., data that is**  
567 **protected by legal privilege or by doctor/patient confidentiality, data that includes the**  
568 **content of individuals' non-public communications)?** No.
- 569 • **Does the dataset contain data that, if viewed directly, might be offensive, insulting,**  
570 **threatening, or might otherwise cause anxiety?** No.
- 571 • **Does the dataset relate to people?** No.

### 572 B.3 Collection Process

- 573 • **How was the data associated with each instance acquired?** The underlying Argoverse  
574 Motion Forecasting scenarios were captured by an AV (part of the Argo AI fleet). Each AV  
575 is equipped with multiple cameras, lidar, and radar. Raw sensor data is processed to produce  
576 tracks localized on a pre-constructed map. Full details are available in [11]. The relevance  
577 labels provided in this work were provided by two domain expert labelers using the labeling  
578 tool depicted in Figure 4.
- 579 • **What mechanisms or procedures were used to collect the data (e.g., hardware appa-**  
580 **ratus or sensor, manual human curation, software program, software API)?** Scenarios  
581 were selected through hand crafted labeling functions. Relevance scores were added using  
582 the web-app labeling tool shown in Figure 4.
- 583 • **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., de-**  
584 **terministic, probabilistic with specific sampling probabilities)?** Scenarios for labeling  
585 were chosen using a set of labeling functions designed to identify complex and interesting  
586 scenarios.
- 587 • **Who was involved in the data collection process (e.g., students, crowdworkers, con-**  
588 **tractors) and how were they compensated (e.g., how much were crowdworkers paid)?**  
589 Data was collected by employees of Argo AI.
- 590 • **Over what timeframe was the data collected?** Source logs Argoverse scenarios were  
591 collected over several months in 2019.
- 592 • **Were any ethical review processes conducted (e.g., by an institutional review board)?**  
593 Not applicable to the relevance labels outlined in this work.
- 594 • **Does the dataset relate to people?** No.

### 595 B.4 Preprocessing/cleaning/labeling

- 596 • **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or buck-**  
597 **eting, tokenization, part-of-speech tagging, SIFT feature extraction, removal of in-**  
598 **stances, processing of missing values)?** 175 of the programmatically selected scenarios  
599 were excluded at the discretion of the labelers. Additionally, automated procedures were  
600 used to resolve a significant set of slightly disparate results across labelers. See section 4.2  
601 for details.
- 602 • **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g.,**  
603 **to support unanticipated future uses)?** All raw labels were saved but are not part of the  
604 publicly released benchmark.
- 605 • **Is the software used to preprocess/clean/label the instances available?** No. These are  
606 simple heuristics outlined in section 4.2.
- 607 • **Any other comments?** None.

## 608 B.5 Uses

- 609 • **Has the dataset been used for any tasks already?** The underlying scenarios from Ar-  
610 goverse 1.1 have been used extensively for tracking and motion forecasting competitions.  
611 The new relevance labels for the retrieval task have not been used outside of the presented  
612 baselines.
- 613 • **Is there a repository that links to any or all papers or systems that use the dataset?**  
614 Not applicable.
- 615 • **What (other) tasks could the dataset be used for?** Our intent labels can also be used as  
616 the first step towards establishing a taxonomy of driving behaviors.
- 617 • **Is there anything about the composition of the dataset or the way it was collected  
618 and preprocessed/cleaned/labeled that might impact future uses?** Our dataset does not  
619 include full contextual information (e.g. camera images and 3D shapes), which impacts  
620 what conclusions we can draw about this dataset.
- 621 • **Are there tasks for which the dataset should not be used? No. Any other comments?**  
622 None.

## 623 B.6 Distribution

- 624 • **Will the dataset be distributed to third parties outside of the entity (e.g., company,  
625 institution organization) on behalf of which the dataset was created?** The benchmark  
626 will be publicly available under a non-commercial license.
- 627 • **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?** Data  
628 will be availble for tarball download through the existing Argoverse website. Test labels are  
629 hidden and test performance can only be obtained via API calls to an evaluation server.
- 630 • **When will the dataset be distributed?** Our current plan is to publicly release the dataset  
631 by July 1, 2021 on our dataset page.
- 632 • **Will the dataset be distributed under a copyright or other intellectual property (IP)  
633 license, and/or under applicable terms of use (ToU)?** We intend to release the data  
634 under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International  
635 Public License (“CC BY-NC-SA 4.0”). Terms of use for all Argoverse data are posted at  
636 <https://www.argoverse.org/about.html#terms-of-use>
- 637 • **Have any third parties imposed IP-based or other restrictions on the data associated  
638 with the instances?** No.
- 639 • **Do any export controls or other regulatory restrictions apply to the dataset or to indi-  
640 vidual instances?** No.
- 641 • **Any other comments?** None.

## 642 B.7 Maintenance

- 643 • **Who is supporting/hosting/maintaining the dataset?** Data will be supported and hosted  
644 as part of the Argoverse project by Argo AI.
- 645 • **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**  
646 Dataset owners can be contacted via email ([ahartnett@argo.ai](mailto:ahartnett@argo.ai) or [ezhan@caltech.edu](mailto:ezhan@caltech.edu)) or via  
647 github issues at <https://github.com/argoai/argoverse-api/issues>
- 648 • **Is there an erratum?** Not currently, though one can be added if errors are discovered.
- 649 • **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete  
650 instances)?** Yes and the version number will be incremented.
- 651 • **If the dataset relates to people, are there applicable limits on the retention of the data  
652 associated with the instances (e.g., were individuals in question told that their data  
653 would be retained for a fixed period of time and then deleted)?** Not applicable.
- 654 • **Will older versions of the dataset continue to be supported/hosted/maintained?** Dep-  
655 recated version of the dataset will be hosted but labeled as deprecated.



656  
657  
658  
659

- **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?** Please reach out via <https://github.com/argoai/argoverse-api/issues>.
- **Any other comments?** None.