

Robust Scenario Mining Assisted by Multimodal Semantics

Anonymous ICCV submission

Paper ID *****

Abstract

001 Scenario mining from large-scale autonomous driving
 002 datasets, such as Argoverse 2, is crucial for the develop-
 003 ment and validation of autonomous driving systems. The
 004 RefAV framework represents a promising approach by em-
 005 ploying Large Language Models (LLMs) to translate natural-
 006 language queries into executable code for identifying rele-
 007 vant scenarios. However, the performance of this method
 008 is constrained by its reliance on the quality of upstream 3D
 009 multi-object tracking data, the absence of a direct linkage
 010 between natural-language descriptions and RGB images,
 011 runtime errors stemming from LLM-generated code, and
 012 inaccuracies in interpreting parameters for functions that
 013 describe complex multi-object spatial relationships. To ad-
 014 dress these issues, we introduce a method that utilizes a
 015 CLIP encoder for multimodal semantic similarity filtering,
 016 first performing a coarse-grained selection by comparing
 017 raw images against the natural-language description, fol-
 018 lowed by fine-grained mining using an LLM-generated script
 019 composed of atomic functions. Additionally, a fault-tolerant
 020 iterative code generation mechanism is introduced, which
 021 refines code by reprompting the LLM with error feedback,
 022 along with specialized prompt engineering to enhance the
 023 LLM’s comprehension and correct application of spatial-
 024 relationship functions. Experiments on Argoverse 2 with
 025 various LLMs show that our method achieves consistent
 026 improvements across multiple metrics. These results under-
 027 score the efficacy of the proposed techniques for reliable,
 028 high-precision scenario mining.

029 1. Introduction

030 The deployment of Autonomous Vehicles (AVs) necessitates
 031 rigorous testing and validation, for which the identification
 032 of interesting, rare, or safety-critical scenarios from vast op-
 033 erational data is paramount. This process is vital not only
 034 for evaluating ego-behavior and safety testing but also for
 035 enabling active learning at scale [12]. Traditional methods
 036 relying on manual inspection or predefined heuristics are of-
 037 ten prohibitively time-consuming and prone to errors when

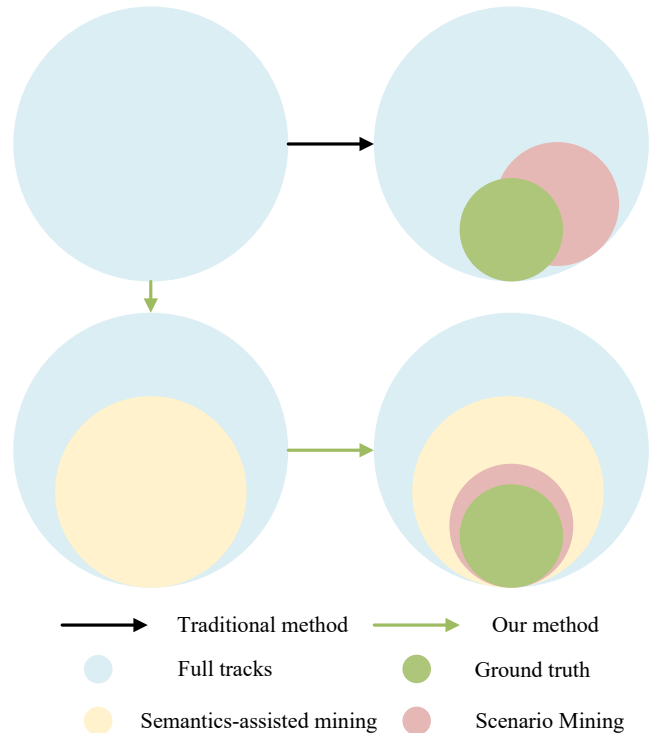


Figure 1. In contrast to traditional scenario-mining pipelines that interrogate the entire collection of 3D tracks—risking substantial drift from the intended query when trajectories share confounding similarities in colour, object class, or event labels—our multimodal, semantics-assisted method first subjects the raw RGB imagery to a semantic filter, isolates a candidate subset of 3D tracks, and only then executes the natural-language query within this reduced search space, which markedly enhances retrieval precision and curtails computational overhead.

038 faced with the terabytes of multimodal data collected by AV
 039 fleets [29]. Previous methods that used database queries
 040 for scenario mining lacked flexibility compared to methods
 041 based on LLMs [10, 13, 16]. The sheer volume and com-
 042 plexity of this data pose a major challenge, making efficient
 043 and accurate scenario mining a major ongoing challenge.

044 The RefAV [6, 30] framework is a method for retriev- 097
 045 ing specific scenarios from sensor data via natural language 098
 046 queries, leveraging the powerful zero-shot capabilities of 099
 047 Large Language Models. RefAV translates natural-language 100
 048 descriptions of scenarios into composable function calls, 101
 049 which are then executed to identify relevant events within 102
 050 driving logs. This approach offers flexibility and expressive- 103
 051 ness beyond structured query languages. 104

052 Despite the promise of LLM-based scenario mining, prac- 105
 053 tical implementations like RefAV encounter specific limita- 106
 054 tions. First, the method hand-crafts 28 atomic functions 107
 055 that detect trajectory states, articulate relations between a 108
 056 focal object and its surrounding entities, and implement ba- 109
 057 sic Boolean logic; an LLM then composes these atoms into 110
 058 scripts that operate directly on the 3D tracks. As a result, min- 111
 059 ing accuracy is tied to the quality of the upstream data, since 112
 060 the preceding 3D object-detection and tracking modules 113
 061 determine the object labels—and therefore whether the tra- 114
 062 jectories themselves are computed correctly. Consequently, 115
 063 poor performance in the upstream 3D multi-object track- 116
 064 ing directly degrades the performance of scenario mining. 117
 065 Furthermore, this approach contravenes the conventional 118
 066 intuition of video retrieval by neglecting the association be- 119
 067 tween the raw image and the natural-language description. 120
 068 Secondly, code generated directly by LLMs can frequently 121
 069 contain syntactic or logical errors, leading to runtime fail- 122
 070 ures. These failures disrupt the mining pipeline and result 123
 071 in incomplete scenario discovery. LLMs may also struggle 124
 072 with the nuanced semantics of functions describing relative 125
 073 spatial relationships between multiple objects. For instance, 126
 074 functions such as *has objects in relative direction()* or *facing*
 075 *toward()* require precise parameter assignment to reflect the 127
 076 intended meaning (e.g., distinguishing "a car in front of a 128
 077 pedestrian" from "a pedestrian in front of a car"). Misinter-
 078 pretation of these parameters leads to semantic inaccuracies 129
 079 in the retrieved scenarios, even if the code executes with-
 080 out error. This issue is a manifestation of a known failure 130
 081 mode in LLMs, often termed ‘factual hallucination’ or a 131
 082 breakdown in understanding relational knowledge [17, 23]. 132
 083 These represent fundamental hurdles in reliably converting 133
 084 complex human language into precise and correct machine-
 085 executable instructions. 134

086 In light of the aforementioned limitations, we propose 135
 087 a robust, multimodally-aware scenario mining methodol- 136
 088 ogy that enhances the RefAV framework. Our approach 137
 089 introduces a dual-branch architecture comprising an image- 138
 090 semantic branch and a text-semantic branch. The image- 139
 091 semantic branch employs the YOLOv8 model for object 140
 092 detection on raw RGB frames. Subsequently, a pre-trained 141
 093 CLIP image encoder [22] is utilized to extract offline fea- 142
 094 ture embeddings for each detected object. Concurrently, the 143
 095 text-semantic branch processes the input natural-language 144
 096 query using the spaCy [14] NLP toolkit to perform keyword 145

097 extraction. This process isolates critical terms, including 098
 099 colors, nouns, and spatial prepositions, which are then en- 100
 101 coded into offline feature embeddings using a CLIP text 102
 103 encoder. During inference, a coarse-grained filtering stage 104
 105 is executed by computing the cosine similarity between the 106
 107 keyword features and the object features across all frames 108
 109 within a complete log. Based on a Top-K selection, the 110
 111 tracklets corresponding to frames with the highest similar- 112
 113 ity scores are shortlisted. Simultaneously, leveraging the 114
 115 Fault-Tolerant Iterative Code Generation mechanism and 116
 117 spatially-aware prompting, the LLM generates an executable 118
 119 script from the natural-language description. This script then 120
 121 performs a fine-grained search directly on the shortlisted 122
 123 tracklets—subsets of the overall 3D tracks—to precisely 124
 125 identify the target scenario. 126
 127

128 In summary, our main contributions are the following: 129

- We propose a multimodal semantics enhancement to the 130
 RefAV methodology. This method addresses a critical 131
 deficiency in the original pipeline by establishing a direct 132
 association between RGB images and natural language 133
 descriptions. 134
- We introduce the Fault-Tolerant Iterative Code Genera- 135
 tion (FT-ICG) mechanism, specifically designed for the 136
 paradigm of using Large Language Models to compose 137
 atomic functions. This contribution significantly enhances 138
 the robustness of the method. 139
- We propose the integration of enhanced prompting for 140
 spatial relationship functions. This technique mitigates 141
 the propensity of the LLM to misinterpret parameters for 142
 atomic functions that describe complex spatial relation- 143
 ships. 144

2. Related Works 129

2.1. Scenario mining 130

131 The safety and reliability of autonomous driving (AD) are 132
 133 of paramount importance, necessitating rigorous testing and 134
 135 validation protocols before deployment. While real-world 136
 137 road testing is indispensable, it is prohibitively expensive, 138
 139 time-consuming, and fails to provide sufficient coverage 140
 141 of rare but critical "edge cases." Consequently, simulation- 142
 143 based testing has emerged as an essential component of 144
 145 the verification and validation pipeline. A core challenge 146
 147 in this paradigm is the generation of a comprehensive and 148
 149 challenging suite of test scenarios. This has given rise to the 149
 150 field of Scenario Mining, which focuses on systematically 151
 152 creating diverse, critical, and realistic driving scenarios to 153
 154 test AD in simulation efficiently. [9, 29] The primary goal is 154
 155 to pinpoint trajectory snippets, within the set of annotated 155
 156 scenes, that satisfy the given natural-language description. 156

157 Early efforts in scenario mining leveraged explicit hu- 158
 159 man knowledge. These methods encode traffic laws, domain 159

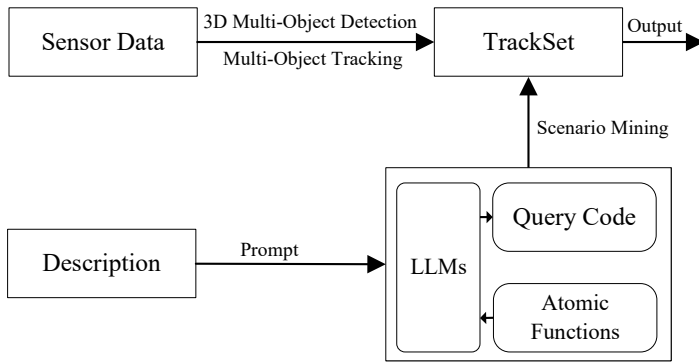


Figure 2. Overview of the RefAV framework: RefAV harnesses large language models to transmute natural-language descriptions into executable, structured code and then runs these scripts automatically to mine the dataset. The trajectory labels consumed by the queries originate from upstream 3D multi-object detection and tracking modules, while handcrafted prompts and a library of twenty-eight atomic functions—each expressing a particular object state—steer the LLM to generate the final retrieval program.

148 expertise, and parameters from accident databases into formal
 149 languages and ontologies. By defining a logical scenario
 150 space with parameters (e.g., road curvature, number
 151 of vehicles, weather) and their valid ranges, scenarios can
 152 be generated through techniques like combinatorial testing
 153 to cover a wide array of predefined conditions. A significant
 154 initiative in this area is the PEGASUS project [20],
 155 which established a systematic, knowledge-driven workflow
 156 for defining scenarios. Formal languages such as ASAM
 157 OpenSCENARIO [1] have become industry standards for
 158 describing the dynamic content of driving scenarios. The
 159 primary advantage of knowledge-based methods is the high
 160 degree of control and interpretability, making them ideal for
 161 testing system compliance with known rules. However, their
 162 main limitation is that they are bound by existing knowledge
 163 and manual effort, often failing to uncover novel modes and
 164 lacking the behavioral complexity of real-world traffic.

165 With the advent of large-scale, real-world driving datasets,
 166 data-driven methods have become prominent. These ap-
 167 proaches mine vast logs of sensor data to extract realistic
 168 scenarios or learn generative models of traffic behavior. The
 169 typical pipeline involves data acquisition from datasets like
 170 the Waymo Open Motion Dataset [24], nuScenes [4] or Ar-
 171 goverse2 [30], using scenario identification by mining the
 172 data for events that exceed a certain criticality threshold. By
 173 articulating the task in a native database query language,
 174 a bespoke domain-specific language (DSL), or a general-
 175 purpose programming language, the problem is recast as
 176 one of label retrieval. Erwin de Gelder et al. [10] present a

177 label-based scenario-mining system for autonomous driving
 178 that operates on datasets pre-annotated either automatically
 179 or by hand. Although the labels are applied in a semi-manual
 180 fashion, they remain coarse-grained; as a result, the frame-
 181 work is rigid and scales poorly—supporting richer, more
 182 nuanced scenes would require an unwieldy proliferation
 183 of tags. Motional’s scenario mining pipeline [19] adopts a
 184 continual-learning paradigm: the system cyclically discovers
 185 scenes, annotates them (both manually and automatically),
 186 retrain its models with the expanded data, and then performs
 187 automatic evaluation. Acting as the data-sourcing engine,
 188 it maintains a tag vocabulary whose compositions encode
 189 basic spatio-temporal relations between the ego vehicle and
 190 surrounding traffic participants. Given 3D trajectories, the
 191 ego path, and the HD map, it automatically labels both the
 192 ego and other actors, stores the tags in a relational database,
 193 and exploits SQL for efficient retrieval. The strength of data-
 194 driven methods lies in their ability to produce highly realistic
 195 scenarios grounded in real-world behavior. The coverage of
 196 the source data is inherently limiting their primary drawback;
 197 discovering truly novel edge cases remains a "needle-in-a-
 198 haystack" problem, and the generated scenarios are often
 199 descriptive rather than actively challenging.

200 To overcome the limitations of passive methods, the most
 201 recent trend in scenario mining involves the application of
 202 large foundation models [28]. As outlined in the compre-
 203 hensive survey by Gao et al. [8], this new paradigm leverages
 204 the power of Large Language Models (LLMs), Vision-Language
 205 Models (VLMs), and Diffusion Models to generate scenarios
 206 from high-level, often semantic, inputs. For instance, a user
 207 can provide a natural language prompt like, "Create a chal-
 208 lenging scenario where a truck illegally overtakes a bicycle
 209 on a rainy night," and the model generates the corresponding
 210 scene parameters for the simulator.

211 Works such as ChatScene [34] have demonstrated the abil-
 212 ity of LLMs to understand complex spatial and behavioral
 213 relationships to produce diverse and contextually rich scenar-
 214 ios. This approach holds immense promise for bridging the
 215 gap between abstract human knowledge and concrete simu-
 216 lation data. While still an emerging area, the key challenges
 217 include ensuring the physical plausibility and controllability
 218 of generated scenarios and managing the significant computa-
 219 tional resources required by these large models. RefAV [6] is
 220 a large-scale scenario-mining framework that houses 10,000
 221 distinct natural-language queries describing the complex
 222 multi-agent interactions present in the 1,000 driving logs
 223 of the Argoverse 2 sensor suite. It exposes 28 handcrafted
 224 atomic functions capable of recognizing trajectory states, ex-
 225 pressing relational predicates between a target agent and its
 226 surrounding entities, and supporting basic Boolean logic. At
 227 its core, the framework feeds the natural-language query, the
 228 atomic function inventory, and carefully engineered prompts
 229 into an LLM, which synthesizes an executable script com-

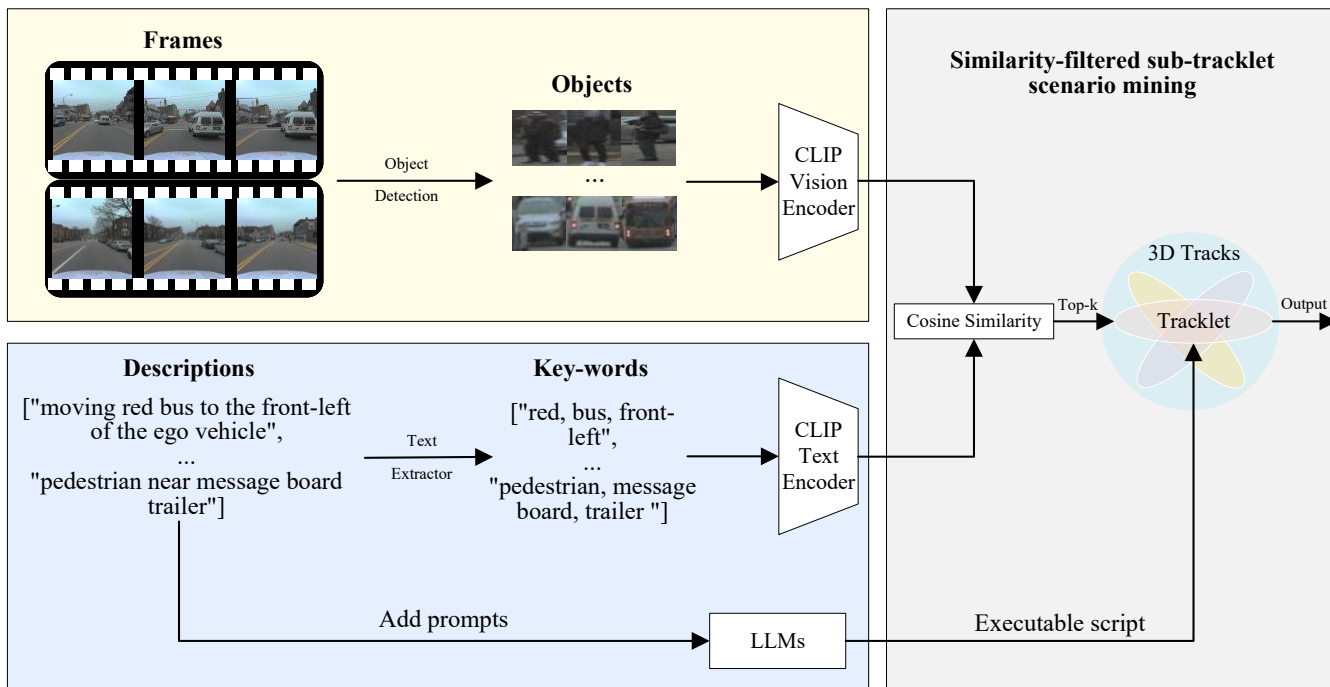


Figure 3. Framework of the multimodal semantics assisted scenario mining.

230 posed of these atoms; running the script over the dataset
 231 then retrieves the trajectories that satisfy the query, as shown
 232 in Figure 2. Our work builds on RefAV to deliver robust
 233 scenario mining assisted by multimodal semantics.

234 In conclusion, the field of scenario mining is evolving
 235 from static, knowledge-driven methods towards more dynamic
 236 and intelligent approaches. Current research increasingly
 237 focuses on hybrid methods, such as using data-driven
 238 models to create a realistic basis for subsequent scenario
 239 mining. Key future directions include richer multimodal
 240 annotations, more faithful modeling of complex multi-agent
 241 interactions, and greater explainability of the critical scenarios
 242 that are retrieved.

243 2.2. Video retrieval

244 Text–video retrieval is a task closely aligned with, and analogous
 245 to, scenario mining: both seek to locate semantically coherent
 246 segments within large-scale, long-sequence data streams based
 247 on natural-language descriptions. Multimodal fusion–based
 248 approaches constitute the most prevalent and emblematic
 249 paradigm within the text–video retrieval literature. We may
 250 glean valuable insights for scenario mining by studying
 251 advances in video retrieval. Since 2020, research on text–
 252 video retrieval has advanced through successive innovations
 253 in cross-modal alignment and temporal modeling:

254 in 2020, Gabeur et al. proposed a multi-modal Transformer
 255 that jointly encodes visual modalities and explicitly models
 256 temporal dependencies via cross-modal attention, optimizing
 257 language embeddings together with video features [7];
 258 in 2021, Wang et al. released T2VLAD, which introduces
 259 shared semantic centers to perform computationally efficient
 260 global-local alignment for fine-grained comparison [27];
 261 Gorti et al. presented X-Pool, enabling text to attend
 262 selectively to semantically relevant frames and thereby
 263 filtering visual noise for improved accuracy [11]; in 2023,
 264 Wu et al. developed Cap4Video, leveraging zero-shot
 265 video-generated captions for data augmentation, cross-modal
 266 interaction, and an auxiliary inference branch, pushing
 267 performance on multiple benchmarks [31]; in 2024, Wang
 268 et al. introduced T-MASS, a stochastic text-embedding
 269 strategy that treats queries as deformable semantic masses,
 270 employing a similarity-aware radius and support-text
 271 regularization to boost expressiveness and set new records
 272 on five datasets [26]; entering 2025, Zhang et al. proposed
 273 TokenBinder, a two-stage framework adopting a one-to-many
 274 coarse-to-fine alignment paradigm inspired by comparative
 275 judgment and equipped with a Focused-view Fusion Network
 276 for cross-attention, achieving state-of-the-art results
 277 across six benchmarks [33], while Bian et al. introduced
 278 the SMA framework, which performs selective multi-grained
 279 alignment at both

279 video-sentence and object-phrase levels with token aggre-
280 gation and similarity-aware keyframe selection, attaining
281 strong performance on MSR-VTT, ActivityNet, and beyond
282 [3].

283 The motivation driving these video-retrieval tasks is akin
284 to that of scenario mining: both seek to locate a contiguous
285 scene that matches a natural-language description. Scenario
286 mining, however, often requires finer-grained retrieval tar-
287 geting the trajectory of a specific agent or a set of interact-
288 ing agents. In almost all video-retrieval work, retrieval is
289 achieved by directly aligning or contrasting visual and tex-
290 tual semantics, which constitutes a latent opportunity for
291 scenario mining: cross-modal semantic matching can be
292 leveraged for coarse retrieval, after which detailed scenario
293 mining can be confined to the trajectory subset thus obtained.

294 3. Methods

295 This section first details the multimodal semantics assisted
296 scenario mining pipeline, and then presents the work’s two
297 further robustness-oriented contributions: a Fault-Tolerant It-
298 erative Code Generation (FT-ICG) mechanism and enhanced
299 prompting for spatial-relation functions.

300 3.1. CLIP-based Multimodal Semantic Filter

301 Within the original RefAV pipeline, scenario mining is per-
302 formed exclusively via scripts assembled from atomic func-
303 tions, a design that overlooks the direct correspondence be-
304 tween natural-language queries and raw image frames. In
305 RefAV, the pipeline begins with 3D object detection to ex-
306 tract each target’s class, heading, velocity, and related at-
307 tributes, assigns corresponding labels, and then waits for
308 an LLM-generated script to query them. The correctness
309 of those attributes rests wholly on the upstream detector’s
310 performance, and the elongated processing chain renders the
311 connection between natural-language queries and raw sensor
312 data highly indirect—raising the likelihood that crucial infor-
313 mation will be missed or inaccurately captured. Inspired by
314 advances in text–video retrieval, we enhance RefAV with a
315 coarse-to-fine filtering stage: natural-language descriptions
316 first delimit tracklets that are likely to contain the target
317 situation, and LLM-generated code then probes only those
318 candidates. This hierarchical procedure sharply reduces false
319 positives and scene ambiguities, yielding results that hew
320 more closely to ground truth.

321 Specifically, on the visual side, we enumerate every object
322 visible in the nine synchronized camera views and encode
323 each one offline with a pretrained CLIP image encoder, so
324 that every embedding captures a localized slice of the frame
325 at its timestamp. On the language side we employ spaCy—an
326 industrial-grade NLP library—to extract colour, entity, and
327 spatial-relation words from each query; these discrete key-
328 words succinctly convey the sentence semantics and, when
329 embedded by the CLIP text encoder, have been shown by

Xie et al [32]. to be more discriminative than full-sentence
330 encodings. We rank frames by cosine similarity to the key-
331 word embeddings, retain the top-k matches, and record their
332 timestamps to assemble candidate tracklets—subsets of the
333 full 3D trajectories. This design enables raw RGB imagery to
334 be compared directly against the textual description, thereby
335 discarding tracklets whose semantics deviate markedly from
336 the query; confining the LLM-generated scripts to mine only
337 within these semantically aligned candidates greatly dimin-
338 ishes the risk of false positives. Executing the original RefAV
339 scripts on this pruned search space markedly improves the
340 HOTA-T metric while reducing inference consumption. 341

342 3.2. Fault-Tolerant Iterative Code Generation

A significant challenge in the practical application of LLMs
343 for code generation is the propensity for the generated code
344 to contain errors. These errors can range from simple syntax
345 mistakes to more complex logical flaws or incorrect usage of
346 the provided atomic functions, all of which lead to runtime
347 exceptions. Such failures can terminate the scenario min-
348 ing process prematurely, resulting in missed scenarios and
349 reduced overall system reliability. The pseudocode for the

Algorithm 1: Fault-Tolerant Iterative Code Genera- tion

Input: Natural-language query $NLQuery$; set of
atomic functions \mathcal{A} ; maximum iterations K

Output: Executable Python code $ValidCode$

```

Prompt ←
COMPOSE( $NLQuery$ , DESCRIBE( $\mathcal{A}$ ));
for  $i \leftarrow 1$  to  $K$  do
  try
    Code ← LLMGENERATE( $Prompt$ );
    PYTHONEXEC( $Code$ );
    Break;
  catch (RuntimeError  $\varepsilon$ )
    ErrorMessage ← MESSAGE( $\varepsilon$ );
    IterationPrompt ← "This is the
code generated last time:
{Code}, with the error
message: {ErrorMessage}. Please
avoid code runtime errors.";
    Prompt ←
COMPOSE( $NLQuery$ , IterationPrompt);

```

350 fault-tolerant iterative code generation mechanism is shown
351 in Algorithm 1. Algorithm 1 proceeds as follows: first, the
352 natural-language scenario query is concatenated with a de-
353 scription of the available atomic-function library to form
354 an initial prompt, giving the LLM full context about which
355 functions are permissible and how they behave so that its
356

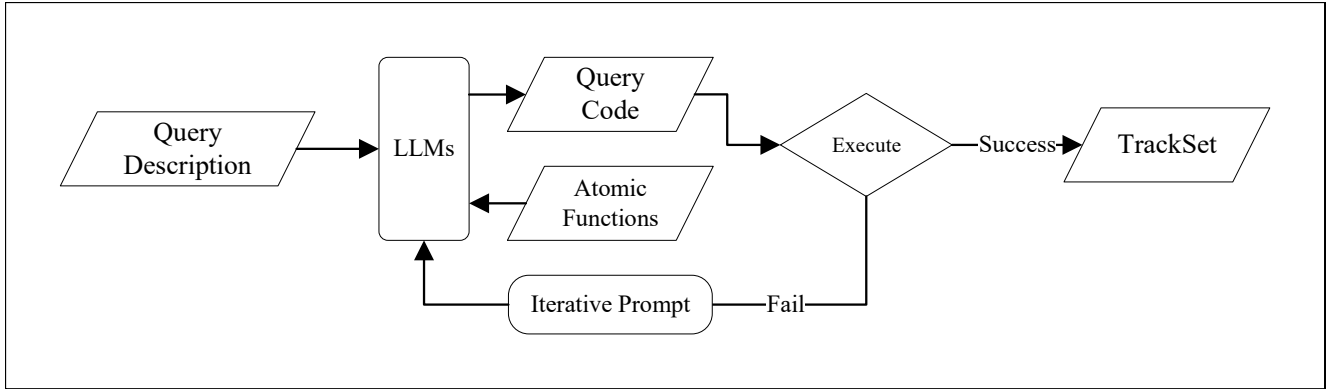


Figure 4. Framework of the Fault-Tolerant Iterative Code-Generation mechanism: whenever the executable script raises a runtime exception, the error trace is fed back to the LLM, which regenerates a revised script; this loop repeats for a preset number of iterations, thereby substantially reducing mining failures caused by execution errors.

357 first attempt is well informed. The model then generates a
 358 Python snippet that chains the atomic functions to imple-
 359 ment the query logic and retrieve the desired scenarios. This
 360 code is executed immediately in a sandboxed environment;
 361 if it runs without error, the resulting track set is accepted
 362 and the pipeline terminates successfully. If, however, exe-
 363 cution raises an exception—such as a ‘NameError’ for an
 364 undefined variable, a ‘TypeError’ due to incorrect argument
 365 counts, or any other syntactic or logical fault—the error is
 366 caught and its message recorded. The system then feeds both
 367 the faulty code and the accompanying error message back to
 368 the LLM in a new prompt that explicitly instructs the model
 369 to correct the identified problem. Armed with this feedback,
 370 the LLM produces a revised snippet, which is executed and
 371 validated again. This refinement loop continues, with the
 372 model iteratively “learning” from each failure, until the code
 373 executes cleanly or a maximum of K iterations is reached.
 374 In our implementation $K = 5$; if a runnable solution still has
 375 not been produced after five attempts, the query is flagged
 376 for manual review. This upper bound prevents pathologi-
 377 cal infinite-loop behavior while still allowing most errors
 378 to be resolved through a handful of feedback cycles. This
 379 iterative approach treats the LLM not as a single-shot code
 380 generator but as an entity capable of learning from explicit
 381 feedback on its errors. By providing the context of the pre-
 382 vious failure, the LLM is guided towards a correct solution.
 383 This significantly increases the success rate of code genera-
 384 tion, thereby enhancing the robustness and coverage of the
 385 scenario mining pipeline, allowing it to handle a broader
 386 spectrum of queries and code complexities without manual
 387 intervention. This process mirrors a human programmer’s
 388 debugging cycle, iteratively refining code based on observed
 389 errors.

3.3. Enhanced Prompting for Spatial Relational Functions 390 391

Beyond syntactic correctness, the semantic accuracy of the 392
 generated code is paramount. LLMs often fail to correctly 393
 interpret and parameterize functions that describe the relative 394
 spatial relationships between multiple objects in a specific 395
 domain. For example, a query like "a cyclist to the left of a 396
 bus" requires the LLM to correctly assign the ‘cyclist’ and 397
 ‘bus’ tracks to the appropriate parameters of a function like 398
 has objects in relative direction(). An incorrect assignment 399
 could lead the system to search for "a bus to the left of a 400
 cyclist," fundamentally misinterpreting the query. A more 401
 comprehensive prompt-engineering strategy can effectively 402
 suppress ambiguities and hallucinations in large models. [5] 403
 To mitigate such semantic errors, Enhanced Prompting for 404
 Spatial Relational Functions is introduced. This involves 405
 augmenting the initial prompt provided to the LLM with 406
 specific instructions that clarify the argument semantics for 407
 these critical functions. Before the LLM attempts to generate 408
 code involving functions that define relative positions or 409
 orientations, it receives the following guiding information: 410

If you use has objects in relative direction(), being crossed 411
 by(), heading in relative direction to() functions, direction 412
 parameter specifies the orientation of related candidates rel- 413
 ative to track candidates. The facing toward() and heading 414
 toward() functions indicate that the track candidates param- 415
 eter is oriented toward the related candidates parameter. 416

This explicit instruction serves as a form of contextual 417
 disambiguation. It clearly defines the roles of track candi- 418
 dates (often the primary subject of the relation) and related 419
 candidates (the reference object) within the context of each 420
 specified function. For directional functions like has objects 421
 in relative direction, it clarifies which entity’s perspective 422
 defines the direction. For orientational functions like facing 423
 toward, it specifies which entity is performing the action of 424

425 facing. By providing this upfront clarification, the LLM is
 426 better equipped to map the natural language description of
 427 spatial relationships to the correct functional representation
 428 and parameter assignment. This leads to a higher fidelity
 429 in translating complex spatial queries, ultimately improving
 430 the semantic accuracy and relevance of the mined scenarios.
 431 This addresses the challenge that code might run correctly
 432 but perform the wrong semantic operation if the LLM mis-
 433 understands these subtle but critical distinctions.

434 4. Experiments

435 4.1. Implementation Details

436 The experiments were conducted using the Argoverse 2
 437 dataset. The dataset provides rich multi-modal information,
 438 including RGB camera frames, LiDAR point clouds, HD
 439 Maps, and 3D track annotations for 26 object categories.

440 The primary metric is HOTA-Temporal. It is a spatial
 441 tracking metric that considers only the scenario-relevant
 442 objects during the precise timeframe when the scenario is
 443 occurring. HOTA[18] was introduced to provide a unified
 444 evaluation of multi-object tracking by jointly accounting for
 445 detection, association, and localization—three facets that to-
 446 gether reflect human intuition of tracking quality. Secondary
 447 metrics include HOTA, Timestamp F1, and Log F1. Times-
 448 tamp F1 treats the video as a sequence of frames, labeling
 449 each timestamp as “scenario” or “non-scenario.” Precision
 450 and recall are computed from the comparison of predicted
 451 and ground-truth frame labels. Log F1 simplifies the task
 452 to a single binary decision per log. After aggregating true
 453 positives, false positives, and false negatives across all logs,
 454 a conventional F1-score is produced.

455 In our setup, we adopt the pretrained ViT-B/32 size of
 456 CLIP—both its image and text encoders—as the backbone
 457 of the multimodal semantic filter. We employ YOLOv8-l
 458 [15] as the object detector. The Qwen2.5-VL-7B model
 459 [2] was deployed locally on a workstation outfitted with an
 460 NVIDIA RTX 4090 GPU, whereas the Gemini model [25]
 461 was accessed remotely via API calls. For 3D object detection
 462 and tracking, we utilized the track obtained directly from
 463 the LT3D method [21]. We set K in Algorithm 1 to 5. For
 464 the generated code, if the number of iterations of the fault
 465 tolerance mechanism exceeds the K value, we manually edit
 466 the generated code, manually modify the reported errors, and
 467 fill in the correct track candidates, related candidates, and
 468 direction parameters. Method evaluation is conducted on the
 469 validation set.

470 4.2. Comparative Experiments

471 As Table 1-3 demonstrate, our method outperforms the base-
 472 line under both upstream 3D tracking pipelines—Le3DE2E
 473 and TransFusion. We report results with three distinct LLMs
 474 acting as code generators, and the largest performance gain

3D Track	Method	HOTA-T	HOTA	TS-F1	Log-F1
Le3DE2E	RefAV*	33.27	36.72	61.94	58.12
	Our method	44.54	44.71	70.37	71.47
TransFusion	RefAV*	30.06	31.27	59.96	59.31
	Our method	44.11	44.67	69.44	68.66

Table 1. With Qwen2.5-VL-7B as the LLM, comparison of our method and baseline across two distinct 3D tracking pipelines—Le3DE2E and TransFusion. * represents the baseline reproduced in our implementation.

3D Track	Method	HOTA-T	HOTA	TS-F1	Log-F1
Le3DE2E	RefAV*	40.17	40.33	66.70	62.71
	Our method	48.30	49.52	72.30	73.41
TransFusion	RefAV*	35.50	35.93	59.89	59.13
	Our method	47.07	47.19	69.79	70.93

Table 2. With Gemini 2.5 Flash as the LLM, comparison of our method and baseline across two distinct 3D tracking pipelines—Le3DE2E and TransFusion.

3D Track	Method	HOTA-T	HOTA	TS-F1	Log-F1
Le3DE2E	RefAV*	42.73	44.27	69.84	66.13
	Our method	52.10	51.07	74.21	70.45
TransFusion	RefAV*	38.76	39.22	60.36	60.31
	Our method	47.37	47.79	69.73	71.66

Table 3. With Gemini 2.5 Pro as the LLM, comparison of our method and baseline across two distinct 3D tracking pipelines—Le3DE2E and TransFusion.

475 arises when Qwen2.5-VL-7B is used. This is likely because
 476 Qwen2.5-VL-7B, relative to Gemini Flash and Gemini Pro,
 477 exhibits a weaker innate understanding of spatial relations
 478 and atomic functions; our pipeline compensates for this
 479 shortcoming. The semantics-assisted filtering stage confines
 480 Qwen’s search to a much smaller candidate subset, the FT-
 481 ICG loop produces more robust and executable code, and
 482 the spatially informed prompts help Qwen correctly interpret
 483 and invoke the atomic functions. Collectively, these elements
 484 drive the substantial improvement over the original RefAV
 485 baseline when Qwen serves as the LLM. When Gemini Flash
 486 and Gemini Pro are used as the LLM, our method likewise
 487 outperforms the RefAV baseline—an advantage attributable
 488 to its coarse-to-fine mining cascade, the greater robustness
 489 of the generated code, and a deeper, more accurate handling
 490 of the atomic-function semantics.

491 As reported in Table 4, we benchmark the end-to-end infer-
 492 ence time for generating and executing a single query under
 493 RefAV and under our framework. Our pipeline delivers
 494 a notable speed-up, attributable to its coarse-to-fine mining
 495 strategy: CLIP features for all images and query tokens are

496 pre-extracted offline, making their embeddings immediately
497 available at inference time. The system therefore performs
498 only a lightweight cosine-similarity check to complete the
499 coarse retrieval stage, swiftly pruning trajectories that are se-
500 mantically irrelevant; by sparing the subsequent fine-grained
501 miner from comparing against obviously incorrect tracks,
502 the overall runtime is substantially reduced.

LLMs	Method	Time(s)
Qwen2.5-VL-7B	RefAV*	47.3
	Our method	19.4
Gemini Flash	RefAV*	48.9
	Our method	17.7
Gemini Pro	RefAV*	42.7
	Our method	18.7

Table 4. A comparison of the inference time required by RefAV and our approach to generate and execute a single query.

MSA	FT-ICG	EP-SRF	HOTA-T	HOTA	TS-F1	Log-F1
✓			41.06	39.98	67.31	67.04
	✓		34.71	39.32	62.77	58.09
		✓	35.37	39.21	62.93	60.11
✓	✓		41.12	40.38	66.68	68.45
✓		✓	43.27	44.13	70.12	69.93
✓	✓	✓	44.54	44.71	70.37	71.47

Table 5. With Qwen2.5-VL-7B as the LLM, performance comparison under different configurations.

503 4.3. Ablation study

504 To demonstrate the broad performance gains delivered by our
505 three contributions to scenario mining, we perform ablation
506 studies under three distinct LLM configurations. In the table,
507 the multimodal semantics-assisted filter is denoted **MSA**,
508 Fault-Tolerant Iterative Code Generation appears as **FT-ICG**,
509 and Enhanced Prompting for Spatial-Relational Functions is
510 labeled **EP-SRF**; all results are reported using the Le3DE2E
511 3D tracker.

512 Across all three LLMs, the multimodal semantic filter
513 (MSA) raises performance consistently—most conspicu-
514 ously on the TS-F1 metric—demonstrating that the CLIP-
515 based filter effectively selects image frames whose content
516 aligns with the query keywords, thereby producing more ac-
517 curate timestamps and tracklets. Restricting scenario mining
518 to these subsets of the full 3D tracks improves multi-agent
519 retrieval precision and minimizes false positives. The FT-
520 ICG mechanism likewise yields uniform gains, particularly
521 in HOTA-T, underscoring the practical benefit of resolv-
522 ing runtime code errors: each iterative refinement produces
523 scripts with higher correctness and executability, which in

MSA	FT-ICG	EP-SRF	HOTA-T	HOTA	TS-F1	Log-F1
✓			44.34	45.98	70.97	66.45
	✓		44.13	45.07	70.44	60.66
		✓	41.77	40.93	69.73	60.95
✓	✓		44.20	46.37	72.01	67.17
✓		✓	45.97	46.63	71.60	69.31
✓	✓	✓	48.30	49.52	72.30	73.41

Table 6. With Gemini Flash as the LLM, performance comparison under different configurations.

MSA	FT-ICG	EP-SRF	HOTA-T	HOTA	TS-F1	Log-F1
✓			47.28	47.97	67.70	66.95
	✓		44.10	46.37	65.90	60.32
		✓	43.74	45.62	69.90	59.13
✓	✓		51.98	51.10	74.05	69.98
✓		✓	50.56	49.47	71.30	70.01
✓	✓	✓	52.10	51.07	74.21	70.45

Table 7. With Gemini Pro as the LLM, performance comparison under different configurations.

turn lifts the HOTA-T score. Subsequent incorporation of
EP-SRF provides additional enhancements—most notably in
HOTA-Temporal, Timestamp-F1, and Log-F1—highlighting
the critical role of semantically precise parameterization
of spatial-relation functions and revealing untapped LLM
potential that can be unlocked through better prompt en-
gineering. The consistency of these improvements across
different LLM backbones indicates that our approach tackles
fundamental challenges in LLM-driven code generation and
interpretation rather than exploiting model-specific quirks.

5. Conclusion

In this paper we presented a robust, multimodal scenario-
mining framework that augments the RefAV pipeline with
CLIP-based semantic filtering, a fault-tolerant iterative code-
generation loop, and relation-explicit prompt engineering;
through a coarse-to-fine retrieval strategy that first constrains
the search space via image-text similarity and then refines
results with LLM-composed atomic-function scripts, the
proposed method simultaneously mitigates error propaga-
tion from upstream tracking, suppresses LLM runtime fail-
ures, and corrects common mis-parameterizations of spatial-
relation functions. Comprehensive experiments on the Ar-
goverse 2 benchmark—covering two distinct 3D tracking
backbones and three heterogeneous LLMs—demonstrate
consistent gains across all evaluation metrics. These re-
sults confirm that tightly coupling vision-language align-
ment with error-aware code synthesis delivers substantial
practical benefits for large-scale autonomous-driving data
mining, and they suggest a clear path toward even richer
multimodal integration and adaptive, self-refining prompting
in future work.

555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610

References

[1] ASAM e.V. ASAM OpenSCENARIO XML, 2022. Online; accessed 13 Jul 2025. 3

[2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. 7

[3] Ziyi Bian, Cong Jiang, Fangzhi Zhu, and Zheng Zhang. Selective multi-grained alignment for text-video retrieval. In *Companion Proceedings of the ACM on Web Conference 2025*, pages 873–877, 2025. 5

[4] 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, 2020. 3

[5] Lucas Choi and Ross Greer. Beyond general prompts: Automated prompt refinement using contrastive class alignment scores for disambiguating objects in vision-language models, 2025. 6

[6] Cainan Davidson, Deva Ramanan, and Neehar Peri. Refav: Towards planning-centric scenario mining. *arXiv preprint arXiv:2505.20981*, 2025. 2, 3

[7] Valentin Gabeur, Chen Sun, Karteek Alahari, and Cordelia Schmid. Multi-modal transformer for video retrieval. In *European Conference on Computer Vision*, pages 214–229. Springer, 2020. 4

[8] Yuan Gao, Mattia Piccinini, Yuchen Zhang, Dingrui Wang, Korbinian Moller, Roberto Brusnicki, Baha Zarrouki, Alessio Gambi, Jan Frederik Totz, Kai Storms, Steven Peters, Andrea Stocco, Bassam Alrifaae, Marco Pavone, and Johannes Betz. Foundation models in autonomous driving: A survey on scenario generation and scenario analysis, 2025. 3

[9] Erwin de Gelder, Jeroen Manders, Corrado Grappiolo, Jan-Pieter Paardekooper, Olaf Op den Camp, and Bart De Schutter. Real-world scenario mining for the assessment of automated vehicles. In *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, pages 1–8, 2020. 2

[10] Erwin de Gelder, Jeroen Manders, Corrado Grappiolo, Jan-Pieter Paardekooper, Olaf Op den Camp, and Bart De Schutter. Real-world scenario mining for the assessment of automated vehicles. In *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, page 1–8. IEEE, 2020. 1, 3

[11] Satya Krishna Gorti, Noël Vouitsis, Junwei Ma, Keyvan Golestan, Maksims Volkovs, Animesh Garg, and Guangwei Yu. X-pool: Cross-modal language-video attention for text-video retrieval. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5006–5015, 2022. 4

[12] Ross Greer, Bjørk Antoniussen, Andreas Møgelmoose, and Mohan Trivedi. Language-driven active learning for diverse open-set 3d object detection. In *Proceedings of the Winter Conference on Applications of Computer Vision*, pages 980–988, 2025. 1

[13] Elmar Haussmann, Michele Fenzi, Kashyap Chitta, Jan Ivanek, Hanson Xu, Donna Roy, Akshita Mittel, Nicolas Koumchatzky, Clement Farabet, and Jose M Alvarez. Scalable active learning for object detection. In *2020 IEEE intelligent vehicles symposium (iv)*, pages 1430–1435. IEEE, 2020. 1

[14] Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. spaCy: Industrial-strength Natural Language Processing in Python. 2020. 2

[15] Muhammad Hussain. Yolo-v1 to yolo-v8, the rise of yolo and its complementary nature toward digital manufacturing and industrial defect detection. *Machines*, 11(7):677, 2023. 7

[16] Aryan Keskar, Srinivasa Perisetla, and Ross Greer. Evaluating multimodal vision-language model prompting strategies for visual question answering in road scene understanding. In *Proceedings of the Winter Conference on Applications of Computer Vision*, pages 1027–1036, 2025. 1

[17] Junyi Li, Jie Chen, Ruiyang Ren, Xiaoxue Cheng, Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. The dawn after the dark: An empirical study on factuality hallucination in large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10879–10899, Bangkok, Thailand, 2024. Association for Computational Linguistics. 2

[18] Jonathon Luiten, Aljoša Ošep, Patrick Dendorfer, Philip Torr, Andreas Geiger, Laura Leal-Taixé, and Bastian Leibe. Hota: A higher order metric for evaluating multi-object tracking. *International Journal of Computer Vision*, 129(2):548–578, 2020. 7

[19] Motional. Technically speaking: Mining for scenarios to help better train our avs, 2022. Accessed 14 Jul 2025. 3

[20] PEGASUS Project. Scenario Description: Requirements & Conditions. In *Proc. PEGASUS Symposium*, Aachen, Germany, 2017. Online; accessed 13 Jul 2025. 3

[21] Neehar Peri, Achal Dave, Deva Ramanan, and Shu Kong. Towards long-tailed 3d detection, 2023. 7

[22] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 2

[23] Vipula Rawte, Amit Sheth, and Amitava Das. A survey of hallucination in large foundation models, 2023. 2

[24] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Sheng Zhao, Shuyang Cheng, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo open dataset, 2020. 3

[25] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, 611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667

- 668 Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap,
669 and Angeliki Lazaridou. Gemini: A family of highly capable
670 multimodal models, 2025. 7
- 671 [26] Jiamian Wang, Guohao Sun, Pichao Wang, Dongfang Liu,
672 Sohail Dianat, Majid Rabbani, Raghuvver Rao, and Zhiqiang
673 Tao. Text is mass: Modeling as stochastic embedding for text-
674 video retrieval. In *Proceedings of the IEEE/CVF conference*
675 *on computer vision and pattern recognition*, pages 16551–
676 16560, 2024. 4
- 677 [27] Xiaohan Wang, Linchao Zhu, and Yi Yang. T2vlad: global-
678 local sequence alignment for text-video retrieval. In *Proceed-*
679 *ings of the IEEE/CVF conference on computer vision and*
680 *pattern recognition*, pages 5079–5088, 2021. 4
- 681 [28] Yuping Wang, Shuo Xing, Cui Can, Renjie Li, Hongyuan Hua,
682 Kexin Tian, Zhaobin Mo, Xiangbo Gao, Keshu Wu, Sulong
683 Zhou, Hengxu You, Juntong Peng, Junge Zhang, Zehao Wang,
684 Rui Song, Mingxuan Yan, Walter Zimmer, Xingcheng Zhou,
685 Peiran Li, Zhaohan Lu, Chia-Ju Chen, Yue Huang, Ryan A.
686 Rossi, Lichao Sun, Hongkai Yu, Zhiwen Fan, Frank Hao
687 Yang, Yuhao Kang, Ross Greer, Chenxi Liu, Eun Hak Lee,
688 Xuan Di, Xinyue Ye, Liu Ren, Alois Knoll, Xiaopeng Li,
689 Shuiwang Ji, Masayoshi Tomizuka, Marco Pavone, Tianbao
690 Yang, Jing Du, Ming-Hsuan Yang, Hua Wei, Ziran Wang,
691 Yang Zhou, Jiachen Li, and Zhengzhong Tu. Generative ai
692 for autonomous driving: Frontiers and opportunities, 2025. 3
- 693 [29] Hiroki Watanabe, Lukas Tobisch, Julia Rost, Johannes Wall-
694 ner, and Günther Prokop. Scenario mining for development of
695 predictive safety functions. In *2019 IEEE International Con-*
696 *ference on Vehicular Electronics and Safety (ICVES)*, pages
697 1–7, 2019. 1, 2
- 698 [30] Benjamin Wilson, William Qi, Tanmay Agarwal, John Lam-
699 bert, Jagjeet Singh, Siddhesh Khandelwal, Bowen Pan, Rat-
700 nesh Kumar, Andrew Hartnett, Jhony Kaesemodel Pontes,
701 Deva Ramanan, Peter Carr, and James Hays. Argoverse 2:
702 Next generation datasets for self-driving perception and fore-
703 casting, 2023. 2, 3
- 704 [31] Wenhao Wu, Haipeng Luo, Bo Fang, Jingdong Wang, and
705 Wanli Ouyang. Cap4video: What can auxiliary captions do
706 for text-video retrieval? In *Proceedings of the IEEE/CVF*
707 *conference on computer vision and pattern recognition*, pages
708 10704–10713, 2023. 4
- 709 [32] Dong Xie, Linhu Liu, Shengjun Zhang, and Jiang Tian. A
710 unified multi-modal structure for retrieving tracked vehicles
711 through natural language descriptions. In *2023 IEEE/CVF*
712 *Conference on Computer Vision and Pattern Recognition*
713 *Workshops (CVPRW)*, pages 5419–5427, 2023. 5
- 714 [33] Bingqing Zhang, Zhuo Cao, Heming Du, Xin Yu, Xue
715 Li, Jiajun Liu, and Sen Wang. Tokenbinder: Text-video
716 retrieval with one-to-many alignment paradigm. In *2025*
717 *IEEE/CVF Winter Conference on Applications of Computer*
718 *Vision (WACV)*, pages 4957–4967. IEEE, 2025. 4
- 719 [34] Jiawei Zhang, Chejian Xu, and Bo Li. Chatscene: Knowledge-
720 enabled safety-critical scenario generation for autonomous
721 vehicles, 2024. 3