# SSE: Multimodal Semantic Data Selection and Enrichment for Industrial-scale Data Assimilation

Maying Shen*, Nadine Chang*, Sifei Liu, Jose M. Alvarez
NVIDIA
{mshen, nadinec, sifeil, josea}@nvidia.com

## Abstract

*In recent years, the data collected for artificial intelligence has grown to an unmanageable amount. Particularly within industrial applications, such as autonomous vehicles, model training computation budgets are being exceeded while model performance is saturating – and yet more data continues to pour in. To navigate the flood of data, we propose a framework to select the most semantically diverse and important dataset portion. Then, we further semantically enrich it by discovering meaningful new data from a massive unlabeled data pool. Importantly, we can provide explainability by leveraging foundation models to generate semantics for every data point. We quantitatively show that our Semantic Selection and Enrichment framework (SSE) can a) successfully maintain model performance with a smaller training dataset and b) improve model performance by enriching the smaller dataset without exceeding the original dataset size. Consequently, we demonstrate that semantic diversity is imperative for optimal data selection and model performance.*

## 1. Introduction

The advent of successful artificial intelligence models have led to an exponential data growth in recent years. As we continue to collect, label, and train on ever expanding data, larger and more impressive models are being created, which recently culminated in foundation models. However, this insatiable data growth is now leading to two new challenges. 1) Our labeled datasets are so large that model performances have saturated. 2) The continuous stream of unlabeled data needs to be properly filtered in order to discover the most valuable data, like a needle in a haystack. Both of these challenges are especially pronounced within the industrial autonomous vehicle (AV) field, as data sizes are much larger due to big fleets that collect large-scale multi-sensor, temporal and 3D data. In this paper, we address

these two particular challenges in a task agnostical manner and showcase results in the context of AV.

The first challenge when working with an excessive data amount is that training on all labeled data is computationally unaffordable. Under a reasonable compute cost, studies have shown that more data does not equate to better performance. Instead, we need less data that is higher quality for optimal model performance [7]. This discovery has spurred an emerging new branch of literature: data selection/pruning. Hence, the second challenge consists of finding valuable data in a pool of unlabeled and unorganized data. A field of academic study that is related to this challenge is known as active learning, which primarily focuses on improving early stages of training. However, data pools in active learning are much smaller than the growing avalanche of data we see today. Therefore, we define this second challenge as data enrichment. Additionally, both data selection and enrichment would benefit from additional explainability, as knowing why data points were selected helps to verify their inherent value.

**Data Selection.** The emerging field of data selection aims to reduce an existing labeled dataset into a compact dataset. Recent studies achieve this by targeting high-quality and pruning low-quality data, mainly in classification and detection datasets. They identify high-quality datasets as either having a traditional object-balanced distribution, visual diversity, or semantic diversity. For traditional object distribution balancing, works refer to long-tail approaches that upsample rarer objects as overall object balancing strategy. Existing released datasets, such as NuScenes [3], incorporate a variety of these statistical dataset balancing approaches to produce a reasonably sized clean dataset. For visual diversity, other works maximize visual embedding coverage space [22].

Since balanced dataset statistics do not guarantee a meaningful content diversity, other works pivot to creating semantically diverse datasets. Within the context of AV, semantic diversity is defined as the overall scene diversity across one driving scenario, which may include the following: A description of the overall background scene,

*Equal contribution.

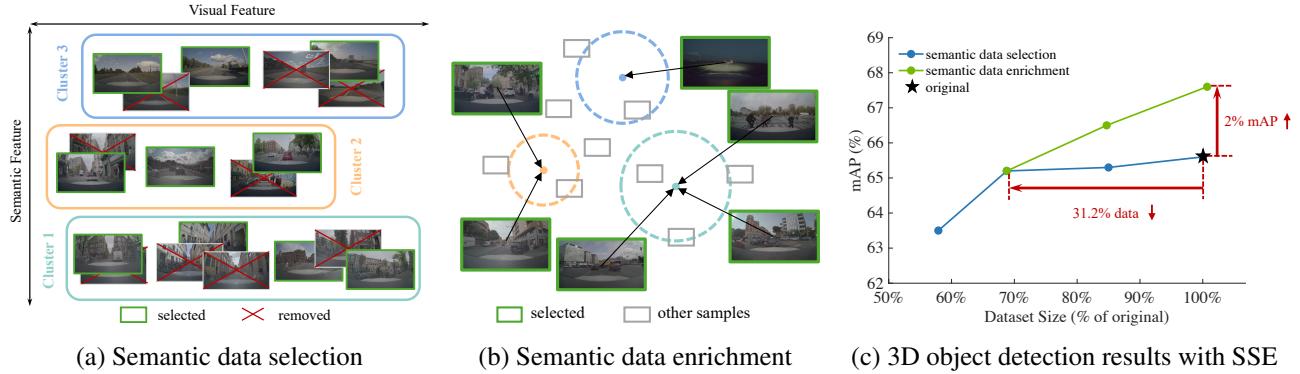|  (a) Semantic data selection | (b) Semantic data enrichment | (c) 3D object detection results with SSE |

Figure 1: We introduce our semantic data selection and enrichment framework (SSE) for autonomous vehicles. The framework generates semantic captions for each data point using a foundation model, capturing semantics including scene understanding (e.g., "crowded urban intersection") and crucial object interactions (e.g., "person about to cross in front of car"). (a) To create a compact dataset, we select the most semantically important portions of a curated and labeled dataset, removing visually similar scenes. (b) To enrich the dataset, we identify new important data points, which are semantically distant from our labeled dataset, from a growing unlabeled data pool. (c) With this approach, we maintain downstream 3D object detection performance using only 70% of the labeled dataset, and we can enhance model performance without increasing the original training dataset size by enriching the selected dataset.

dynamic objects, car behaviors, weather or presence of hazardous objects. Recent works leverage the world knowledge encapsulated by foundation models like CLIP [18] to encode images into stronger embedding spaces with more semantic information than previous vision embeddings, and then remove images with similar embeddings. Thus, they mostly focus on visual duplicate removal as opposed to semantically relevant data-point selection.

**Data Enrichment.** Within academia, parts of data enrichment fall under active learning. Similar to data pruning, active learning has been studied mostly for image classification. While some works consider 2D object detection, they often focus on early stages of the training with a small amount of data [5, 10]. Broadly, active learning uses specialized models, uncertainty estimation, or query by committee. Specialized models are trained to identify targeted objects (e.g. a "traffic light" detector) [23, 2], but this methodology does not scale and cannot respond to new targets. Uncertainty estimation utilizes the downstream model to measure data sample uncertainty to determine which new samples to select [11, 13, 30]. As AV includes several constantly evolving downstream models and multi-sensor data, relying on a single model and specialized, model-specific data for uncertainty estimation is impractical. Query by committee uses majority vote by model ensembles, which requires even more model development and computational costs [19]. Unlike in data pruning, these works also do not focus on data semantics.

**Common Issues.** Existing works have two main issues: First, they lack explainability, which is a shared issue across data selection and enrichment works. Describing the contents of selected data points is critical to safety and helps us understand why some data are more important than others. Second, existing works disjointly study data selection and enrichment. But realistically, both have to be considered simultaneously. We emphasize that within industry, data selection cannot be performed across the labeled and unlabeled data pools because the unlabeled pool increases and changes continuously. Thus, data selection is restricted to human-labeled datasets, which can be separately enriched with new data from the unlabeled data pool. To address the lack of explainability and disjoint processes, as shown in Fig. 1, we propose a joint semantic data selection and enrichment framework (SSE) that leverages multimodal language models (MLLMs) to generate explainable semantics in natural language.

**Our Framework (SSE).** Given an existing hand-curated and labeled dataset, we first focus on selecting the most relevant data without relying on labels. We define semantically diverse data points as the most relevant ones. Semantics are achieved by prompting MLLMs to describe each data sample in detail – from producing overall scene descriptions (intersection in city) to pointing out notable hazards (bicyclists potentially crossing). These captions are then embedded and clustered. As seen in Fig. 1, we aim to select the most semantically diverse samples. Therefore, we organize them into core semantic clusters and subsequently remove visually similar samples within each cluster. After data selection, we enrich our dataset with additional semantically diverse scenes from a new *unlabeled* data source( Fig. 1). Importantly, our methodology does not require any human annotated data labels.

**Contributions.** With our joint semantic data selection and data enrichment framework, we address the data challenges unique to industrial scale AV:

1. We leverage MLLMs to inject explainability into our process and remove reliance on downstream models.

2. We propose a joint framework to incorporate both data selection and enrichment(SSE) to address the two challenges arising from exponential data growth.

3. In addition to improving semantic selection and enrichment, we also provide a better semantic retrieval that uses our semantic embeddings.

4. We conduct thorough evaluations and analysis on large-scale industrial datasets and demonstrate downstream model performance improvements over baseline methodologies.

5. We experimentally show that a semantically tuned dataset with fewer but more semantically relevant class objects is more effective for model performance than a hand-curated and statistically tuned dataset with strictly more class objects.

## 2. Method

A high-level overview of our semantic selection and enrichment frame is shown in Fig. 1. An additional pseudocode is provided in Algorithm 1.

**Capturing Semantics within Scenes.** We define the semantics within a single scene as answers to the following description requests: a general scene description, a general description of what is happening, the important objects to consider while driving or the dynamic objects. All descriptions require world knowledge and high levels of scene understanding. Thus, we ask a number of questions regarding a single image to a pre-trained MLLMs and receive a paragraph of desired answers. By encoding this caption paragraph into a text embedding space via a sentence transformer, we capture high-order semantics for a single scene. This process is illustrated for the labeled dataset, $D$, and unlabeled data pool, $P$, in Lines 1 and 2, Lines 13 and 14.

**Semantic Data Selection.** Semantic embeddings enable identifying the various unique semantics within the dataset. As our goal is semantic diversity within a dataset, we perform $k$-means clustering on the semantic embeddings to group semantically similar scenes. To select a core group of semantic scenes, we remove visually similar scenes within each semantic clusters. To calculate visual similarity, we leverage another foundation model (CLIP) and its vision encoder to obtain visual embeddings for each sample. Within each semantic cluster, we calculate the pairwise cosine similarity between the visual embeddings for all samples and greedily remove those whose score exceeds a pruning threshold $\epsilon$. The final dataset, $D_s$, contains the remaining samples within each semantic cluster. More detailed pseu-

---

**Algorithm 1** Data Selection and Enrichment

---
**Require:** labeled dataset, $D$; unlabeled data pool, $P$
**Require:** prompt, $u$; pruning threshold, $\epsilon$
**Ensure:** $D_s$, reduced dataset after data selection and $D_e$, enriched dataset after data enrichment
1: Let $C^D = \{\text{MLLM}(u, i) : i \in D\}$
2: Let $T^D = \{\text{SenTrans}(c_i) : c_i \in C^D\}$
3: Let clusters $M = \text{kMeans}(T^D)$
4: Let $V^D = \{\text{CLIP}(i) : i \in D\}$
5: **for each** $m \in M$ **do**
6:     **for each** $i \in m$ **do**
7:         **if** $1 - \cos(v_i, v_j) < \epsilon$,
            for any other sample $j \in m, i \neq j$ **then**
8:             Remove $j$ from $m$
9:         **end if**
10:     **end for**
11: **end for**
12: $D_s = \{i : \exists i \in m\}$ for $m \in M$         ▷ Selection

13: Let $C^P = \{\text{MLLM}(u, i) : i \in P\}$
14: Let $T^P = \{\text{SenTrans}(c_i) : c_i \in C^P\}$
15: Define $O = \{\text{ImageClosestCentroid}(m) : m \in M\}$
16: $D_e = D_s$
17: **while** expanding $D_e$ **do**         ▷ Enrichment
18:     $D_e = D_e + \arg\max_{i \in P}(1 - \cos(t_i, t_o)) : o \in O$
19: **end while**

---

docode is shown from Line 3 to Line 12.

**Semantic Data Enrichment.** After semantic data selection, we can expand the dataset by adding new semantically meaningful data from an unlabeled data pool, $P$. Since we have a set of semantic clusters, we utilize the semantic embedding closest to the centroid of each cluster as a semantic anchor point. Using these semantic anchors, we identify the data points within $P$ that are most semantically removed from our current semantics. We identify these points by calculating the cosine similarity between semantic embeddings. Then, we incrementally add data until the desired amount is reached. Finally, we annotate these data points analogous to the labeled dataset. The semantic enrichment is detailed from Line 15 to Line 19.

## 3. Experiments

### 3.1. Implementation Details

**Data.** We utilize an internal large-scale experimental research dataset $D$ with $415,544$ unique scenes. Each scene is 1 driving timepoint recorded with 8 cameras (4 standard and 4 fisheye) and belongs to a full video sequence, which we define as a session. $D$ contains 2750 unique sessions. The images are reduced to $480 \times 960$ resolution for training and evaluation. For data selection, we select unique scenes
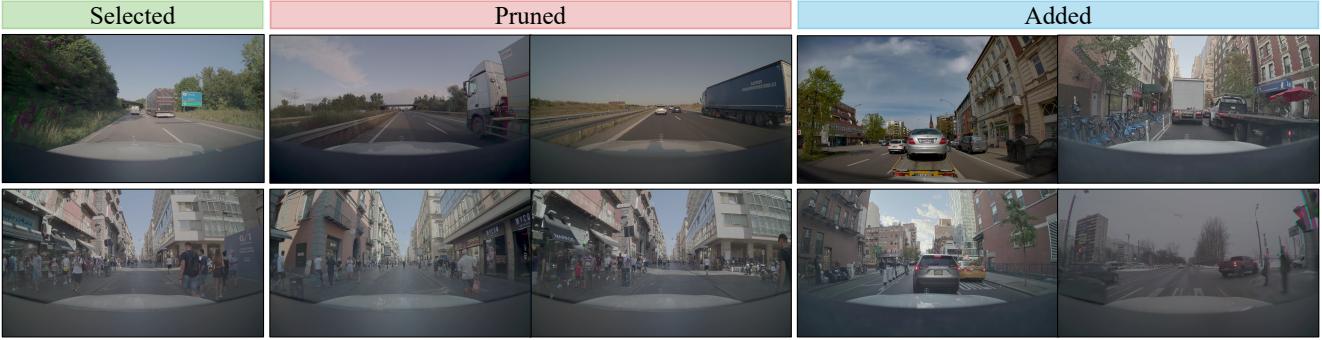
| Selected | Pruned | Added |
|:---:|:---:|:---:|

Figure 2: Examples of semantic selection and enrichment. The "Pruned" samples are visually and semantically similar to the "Selected" samples, not only a visual duplication. The "Added" samples add different semantics to existing data.

Table 1: Performance on downstream 3D detection with different <u>Data Selection</u> strategies. SSE achieves 30% data reduction while maintaining original mAP.

| Dataset Size (% of original) | Method | mAP | $\Delta$_original | $\Delta$_random |
|:---:|:---|:---:|:---:|:---:|
| 100% | Original Dataset | 65.6 | - | - |
| 80% | Random | 63.4 | -2.2 | - |
| | Long-tail [8] | 64.2 | -1.4 | +0.8 |
| | CLIP visual [18] | 64.9 | -0.7 | +1.5 |
| | **SSE (Ours)** | **65.1** | **-0.5** | **+1.7** |
| 70% | Random | 62.1 | -3.5 | - |
| | Long-tail [8] | 62.3 | -3.3 | +0.2 |
| | CLIP visual [18] | 64.4 | -1.2 | +2.3 |
| | **SSE (Ours)** | **65.2** | **-0.4** | **+3.1** |
| 60% | Random | 61.1 | -4.5 | - |
| | Long-tail [8] | 60.8 | -4.8 | -0.3 |
| | CLIP visual [18] | 62.9 | -2.7 | +1.8 |
| | **SSE (Ours)** | **63.5** | **-2.1** | **+2.4** |

from the dataset $D$. For data enrichment, we select unique scenes from an unlabeled data pool $P$ containing $662,325$ unique scenes from 5109 unique sessions.

**Evaluation Model.** We evaluate our semantic data selection and enrichment pipeline on the downstream Multi-Camera 3D Object detection model [17]. For the evaluations, we train the model for 20 epochs with data from all 8 cameras and an individual batchsize 8 on 32 GPUs in total. The learning rate is scheduled according to the one-cycle learning rate policy [25] with a maximum value at $5e^{-3}$. All final metrics are expressed as mean average precision (mAP), and per-class APs are reported in Appendix.

**Data Selection and Enrichment.** We use LLaVA1.6-34B-4bit [12] to generate semantics, because it can handle high-resolution images. Unless specified, our semantic selection and enrichment utilizes MPNet-base-v2 [26] for semantic caption encoding and $k = 300$ for clustering, the front camera image as input to MLLM, and the specialized AV prompt in Fig. 11 for generating image descriptions.

## 3.2. Main Results: Semantic Selection and Enrichment Performance

**Data Selection.** Our pipeline begins with data selection, where we aim to select a subset of semantically relevant data from an established, human-labeled internal dataset, $D$. We compare against three main baselines, as seen in Tab. 1. As an effective baseline, the random baseline randomly selects a portion of dataset to retain without any clustering. Next, following standard practice in dataset collection, we select a portion with the most balanced object distribution. We use Repeat Factor Sampling (RFS) [8], a strong standard baseline in long-tail detection for calculating a data sample's importance based on object count. Finally, our methodology leverages semantic scene understanding by converting generated image captions into text embeddings. To evaluate its effectiveness, we compare against the traditional approach of encoding images into visual embeddings. For this, we adopt CLIP [18], a powerful and widely adopted vision encoder, as a robust visual baseline. Once images are clustered based on both Visual and Semantic embeddings, data selection is performed by greedily pruning the most visually similar images within each cluster.

Industry applications have to address the tradeoff between dataset retention percentage and performance. Usually, they select the method with the most aggressive dataset retention and negligible performance change. Practically, this decision translates to the cheapest and most efficient model training for maximum performance effect. In Tab. 1, we observe that our semantic data selection selects the smallest required amount of data (70%) to achieve comparable performance to using the entire dataset. The difference of $0.4$ mAP is negligible. In comparison, the performance of other baselines drops noticeably when more data is pruned from the dataset. Note, that while object balancing seems effective for general long-tail detection settings, this strategy does not perform well when selecting relevant complex AV scenes.

Table 2: Performance on downstream 3D detection with different <u>Data Enrichment</u> methods. SSE improves 2 mAP when expanding the dataset to the original size. Note long-tail is impossible due to the lack of labels in data pool.

| Dataset Size (% of original) | Method | mAP | Δ_original | Δ_random |
|---|---|---|---|---|
| 100% | Original Dataset | 65.6 | - | - |
| 85% | Random | 64.5 | -1.1 | - |
| | CLIP visual [18] | 65.6 | 0.0 | +1.1 |
| | **SSE (Ours)** | **66.5** | **+0.9** | **+2.0** |
| 100% | Random | 65.2 | -0.4 | - |
| | CLIP visual [18] | 65.6 | +0.4 | +0.8 |
| | **SSE (Ours)** | **67.6** | **+2.0** | **+2.4** |



The image captures a bustling city street from the perspective of an ego vehicle's windshield. The road is shared with multiple pedestrians, some of whom are walking on the sidewalk while others are crossing the street nearby. There are several parked cars and motorcycles along the side of the road, adding to the urban feel of the scene. The weather appears to be clear, with ample sunlight illuminating the scene. The image conveys a typical day in a busy city with diverse road users and a mix of stationary and moving elements.
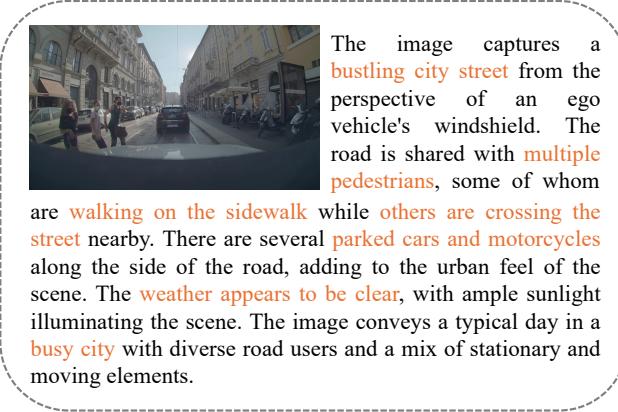
Figure 3: Semantic description with MLLMs. The highlighted phrases capture the relevant semantics.

**Do We Only Prune Identical Scenes?** In Fig. 2, we visually demonstrate that the pruned data is similar to the selected sample, but importantly *not* a visual duplication. For example, the first row shows that our pruning removed freeway scenes with trucks, which are similar to the selected data point.

**Data Enrichment.** Continuously collected unlabelled data is a potentially useful addition to an existing dataset. Here, we demonstrate the effect of applying our methodology on data enrichment for a carefully pruned and labeled dataset, $D_s$. We apply data enrichment on top of all previous methodologies and their respective $70\%$ retained datasets. As we do not have a labeled pool of data to work with, we cannot report a long-tail enrichment baseline. For the remaining enrichment results, the methodology used for data selection is also used for data enrichment. For example, in the CLIP visual baseline, we encode the unlabeled data with the CLIP vision encoder and apply our semantic selection method, as detailed in Sec. 2. Once data is selected from the unlabeled data pool $P$, we label the data following the same labeling guidelines used in our original dataset and expand $D$ for training.

We observe in Tab. 2 that by growing a dataset with se-

mantic embeddings only up to $85\%$ of the original dataset size, we are able to outperform the model trained on the original full dataset (0.9 mAP). By semantically expanding the dataset to the original dataset size, we are able to even push the improvement to 2 mAP.

**Explainability Enabled by Semantics.** Because our semantics are from generated captions, we can analyze why a data point is semantically important. The highlighted phrases in the generated captions shown in Fig. 3 succinctly capture the semantics within the given image. As a result of the generated caption for this sample, we are able to verify and understand that its semantic embedding refers to "bustling city street", "multiple pedstrains...walking on the sidewalk" and "several parked cars and motorcycles". Thus, we can inject explainability into every aspect of SSE.

### 3.3. Main Results: Semantic Importance

**How Semantically Diverse Are the Scenes in Enrichment?** In Fig. 2, we see that newly added data is visually and semantically different from the reference image in enrichment process. In the upper row, the new data points are non-highway scenes where the driving space is more restricted.

By leveraging the semantic embeddings, our methodology demonstrates two impactful results: First, through semantic selection, we can drastically reduce the size of a carefully human-curated and annotated dataset. Second, through semantic data enrichment, we obtain an effective measure to further grow this pruned dataset for further performance improvements.

**Do Our Semantic Clusters Capture True Semantics?** First, we qualitatively inspect whether our semantical clusters contain scenes that are semantically and not merely visually similar. The cluster in Fig. 5 contains scenes with pedestrians or cyclists near the ego car that are likely to cross the street. Significantly, these scenes are visually diverse, taken at different times of day and locations.

Because we cannot visually inspect every cluster for semantic confirmation, we quantitatively analyze the composition of the dataset in Fig. 4. The original dataset is hand-curated to contain scenes from a large number of unique driving sessions. These different driving sessions tend to be geographically diverse. Thus, scenes within a session are visually more similar than scenes across sessions. Note that in our semantic clusters (LLaVA generated as shown in Sec. 3.5) contain a high number of unique sessions, which indicates that semantic clusters indeed capture semantically similar but visually diverse scenes across sessions. In contrast, clusters formed with CLIP visual embeddings have fewer unique sessions, indicating that visual clusters gather more visually similar scenes from a single session. This trend is stable and holds when we cluster the whole dataset, post data selection, and post data enrichment.
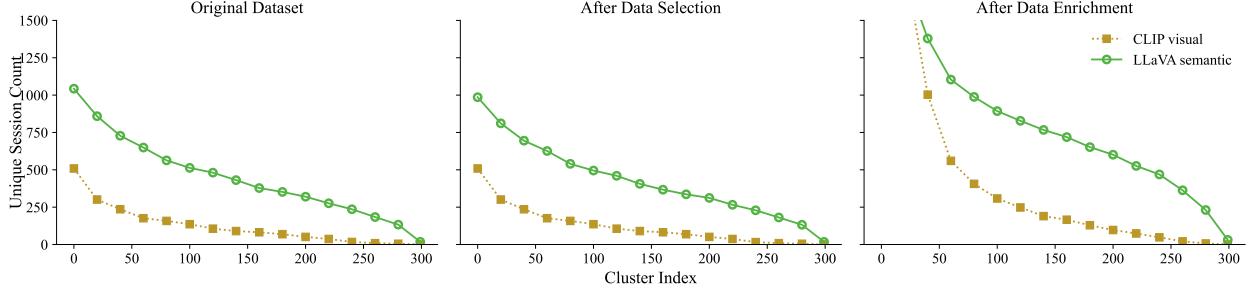
5

Figure 4: Number of unique driving video sessions in each cluster formed with different embeddings. Compared to clusters generated from visual embeddings, semantic clusters capture more semantically similar yet visually diverse scenes across sessions.



Figure 5: Visualization of samples in one of our semantic clusters. The scenes are visually different but semantically similar (Pedestrians/cyclists near the ego car and likely to cross the street in front).

**Semantically Diverse Data vs. Balanced Class Distribution Data.** We show how each methodology affects the number of unique ground truth objects. In our dataset, `car`, `truck`, `person`, and `bike with rider` cover the most to least represented objects. Interestingly, we observe that semantic selection and enrichment does not simply select and add images with rare objects, such as `person` and `bike with rider`, as seen in Fig. 6. Importantly, we see in Fig. 7 that this imbalance in object distribution after selection does not severely affect the per-class AP. This makes it even more surprising that after enrichment the rarer categories with fewer objects outperform their original counterpart ($+3.2$ AP for `person`, $+2.6$ AP for `bike with rider`). Furthermore, SSE performs consistently better than all baselines across all categories after enrichment. Thus, SSE semantically tunes a dataset and successfully demonstrates that **fewer but more high-quality objects** lead to a better performance in rare categories. Summarily, we show that semantically diversified data is more important than resampled and balanced data.
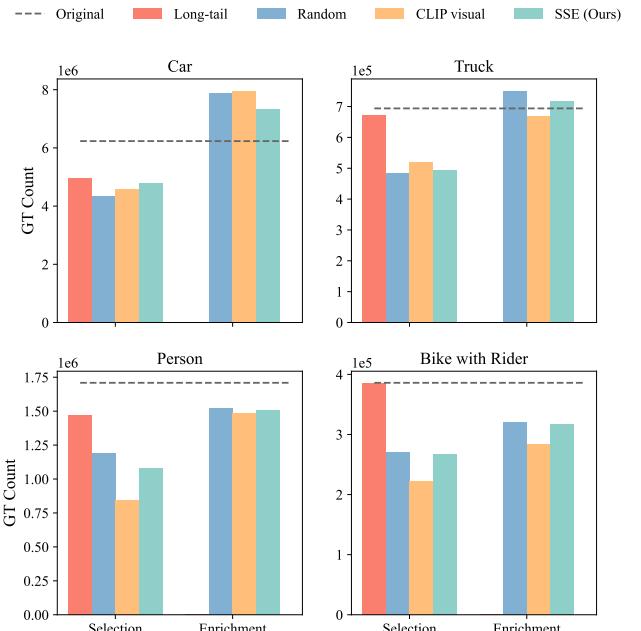


Figure 6: Number of unique objects after different selection and enrichment methods. SSE selects semantically diverse scenes that does not necessarily contain more rare objects.
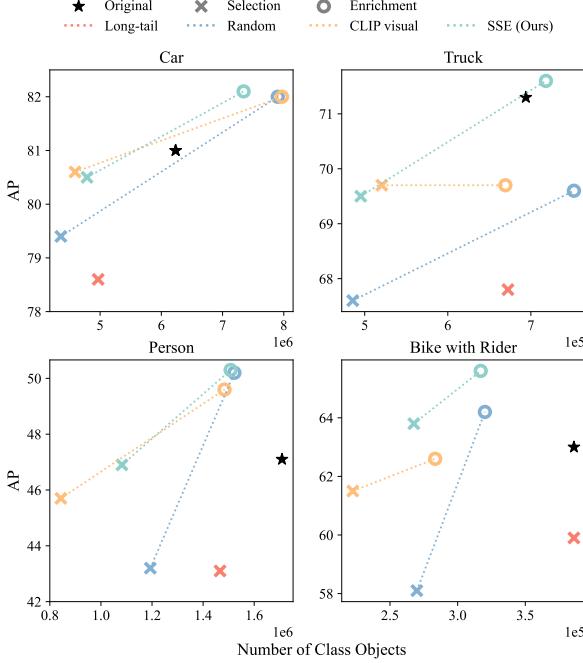
Figure 7: Per class detection accuracy as a function of object count in train data, across different methodologies. SSE semantically tunes the dataset and demonstrates that *less but more high-quality objects* lead to better performance in rare categories.



Figure 8: Fraction of data remaining vs pruning threshold $\epsilon$.



Figure 9: Study on the number of $k$ clusters.



Figure 10: Study on the effects of different MLLMs.

## 3.4. Analyzing Effects of Parameters on Data Selection and Enrichment

**Pruning Threshold $\epsilon$.** In data selection, we ensure visually diverse data within each semantic cluster through a pruning threshold $\epsilon$. Because the CLIP visual baseline uses the same visual embedding for clustering, we see that the pruning effects of $\epsilon$ are strongly pronounced, as shown in Fig. 8. The effects of using different MLLMs are addressed later in Sec. 3.5. For all studies, we closely match the percentage of data across methodologies for fair comparison.

**$k$ Selection.** Data selection and enrichment relies on clustering data embeddings with $k$-means. We randomly hold out 20% of the original dataset for validating $k$ clusters selection. We attempted to align the amount of retained data by setting similar pruning $\epsilon$. As observed in Fig. 9, we observe that semantic selection at $k = 300$ is able to retain the fewest samples at 70% with minimal performance change. Thus, we select $k = 300$ as our default number of clusters for all experiments. Importantly, we observe that all data selection models' performances across the val and train set are approximately similar. This indicates that the effects of $k$ are stable across retention thresholds and data splits.
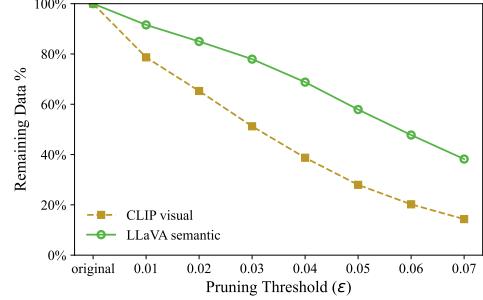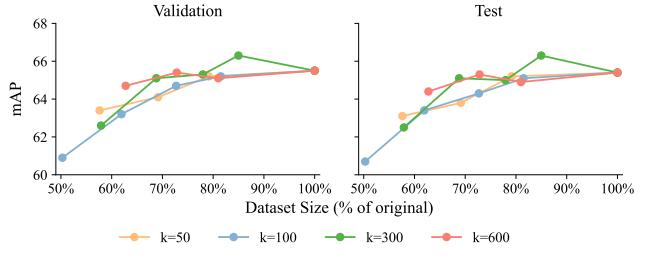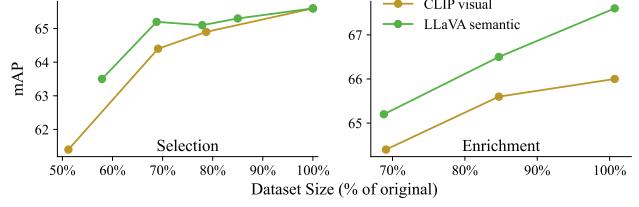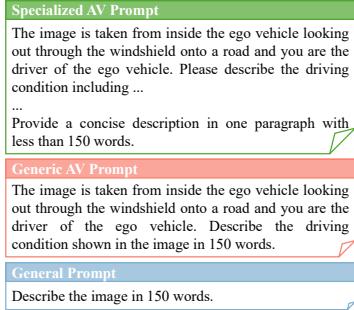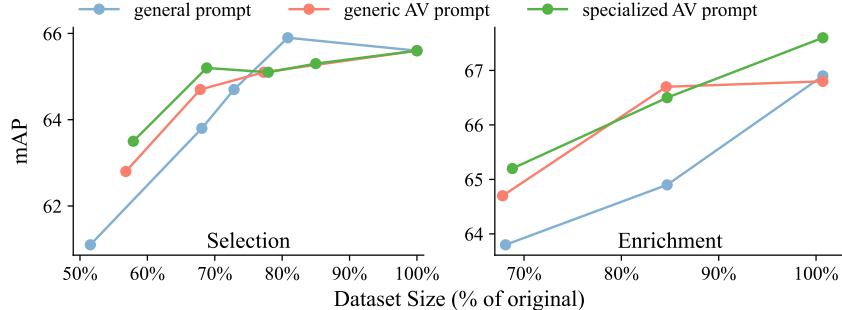
## 3.5. Studying the Effects of MLLMs

**Prompt Sensitivity.** The phenomenon of MLLM's sensitivity to prompts is well known and studied [15, 21, 31, 14]. We attempt to find the best prompts suitable for semantic data selection and enrichment. The results for using different prompts for the MLLM are in Fig. 11 for data selection and enrichment. In data selection, using an specialized AV prompt outperforms both the generic AV prompt and general prompt. The targeted questions within specialized AV prompt act as guildlines for selecting only the most critical scenes in the dataset. On the other hand, because we expand the dataset with data points that are more semantically diverse than our clusters, a more generic AV prompt is more befitting for a wider data search. Accordingly, we see that a generic AV prompt outperforms the rest. We report AV specialized prompt results as our main selection and enrichment results in Tabs. 1 and 2 because it best handles the data selection tradeoff in retention and performance.

**MLLMs.** On our semantic data selection and enrichment,

(a) Three prompts
(b) Dataset semantic selection and enrichment with different prompts.

Figure 11: Study of prompt sensitivity. Specialized AV prompt provides guidance for selecting driving critical scenes.



Query: Situations where need to slow down and drive carefully.

Figure 12: Retrieval examples using advanced visual embeddings vs semantic embeddings. Semantic retrieval finds semantically similar yet critically visually diverse scenes. Visual retrieval returns visually repetitive scenes and cannot show high-level scene understanding.

we compare the effects of different MLLMs for caption generation. Specifically we compare LLaVA with CLIP (vision language model). Significantly, we note that the performance trend in both semantic selection and enrichment is similar with both MLLMs, as seen in Fig. 10. We note that data retention for CLIP is notably lower because the default visual pruning process removes visually similar images based on CLIP embeddings.

### 3.6. Semantic Retrieval

Utilizing information rich captions as our semantic embeddings enables an additional capability: semantic retrieval. We set up a simple retrieval framework similar to standard retrieval tools powered by foundation models [18], where a similarity search is performed at a shared embedding space between a user prompt against a set of raw data points. We compare our retrievals against an advanced visual retrieval system, which uses CLIP vision encoder embeddings for patches of images [4]. As observed in Fig. 12, semantic retrieval is able to retrieve semantically similar scenes that are visually diverse. For example, semantic retrieval is able to break down the concept of "slow down and drive carefully" and successfully retrieve various scene in-

terpretations (starting from left to right images): vehicles cutting into lanes, bicyclist crossing, dense pedestrian traffic, and heavy rain. In comparison, CLIP retrievals are visually repetitive (mostly highway) and incapable of capturing high-level scene understanding beyond scenes with different traffic signs.

## 4. Related Works

**Data Pruning and Selection.** Selecting a good set of data points and pruning out bad data points are both considered even in the process of creating a dataset, prior to any model introduction. Public datasets, such as NuScenes, carefully balance the class and object count distributions and remove data points with common objects [3]. However, this balancing requires considerable manual effort, is unscalable, and is entirely subjective. Recent large-scale datasets, such as LAION [20], bypasses these issues by automatically removing visually similar images, as calculated through CLIP's visual similarity score. Beyond dataset creation, data selection is commonly applied to existing datasets in order to reduce compute costs. Many approaches for classification datasets attempt to select the most diverse set of data

8

points, either by maximizing visual space coverage [22] or maximizing gradient space coverage [1].

In contrast, recent works have attempted to inject semantics into their data selection process, particularly through the use of foundation models like CLIP. Within the context of fairness, FAIRdedup [24] uses CLIP to assign protected attributes for each image and attempts to balance the protected-attributes distribution. One analytical work shows that datasets pruned with optimal CLIP scores can improve model performance beyond normal datasets with the same distribution [27]. Another analytical work shows that under a computation budget, CLIP pruned dataset can lead to a higher quality dataset for better model performance [7]. However, both works use brute force grid search of CLIP score thresholds to show these analytical insights and do not provide a solution for data selection.

Finally, a series of studies (EL2N, GraNd, MoSo, Memorization, Forgetting score) assumes that hard samples are critical to datasets and aims to prune out easy samples [16, 28, 6, 29]. They often require model ensembles or use downstream model itself to find hard samples. Thus, they are not computationally adaptable to industrial applications nor do they provide any explainability. In our work, we are able to semantically select data points from a large-scale dataset and provide interpretable selection reasons.

**Active Learning.** Similar to data pruning, active learning has been studied mostly in image classification. While some works target 2D object detection, they primarily work with small amounts of data to improve early stages of training [5, 10, 33]. Within classification, active learning approaches range from uncertainty estimation over query by committee to training specialized models. Like data selection works, uncertainty estimation studies rely on the assumption that hard samples, for which the downstream model is uncertain in its predictions, are valuable for training [11, 13, 30, 32, 9]. Query by committee finds important data through majority vote of model ensembles, which requires even more models than uncertainty estimation [19]. Finally, works that train specialized models to find their targeted data points only apply to settings where desired data types are provided [23, 2]. All these works are reliant on a model in the loop or additional models to find new useful data. This reliance raises compute costs, making them unscalable for industrial applications.

## 5. Conclusion

We introduced semantic selection and enrichment (SSE) which identifies the most semantically diverse and significant parts of a dataset and enhances it further by uncovering new data from a vast pool of unlabeled information. Importantly, SSE offers explainability by utilizing foundation models to generate semantics for each data point. Through semantic data selection, we show that we can select a se-

mantical core group of data from an already carefully hand-curated, balanced, and human-labeled dataset. Our dataset post selection is only $70\%$ of the original dataset size but enables the downstream model to achieve similar performance to that from the original dataset. To expand this smaller semantic dataset back to $100\%$ of original dataset size, our enrichment process can successfully find and add more semantically meaningful data with increased downstream model performance. Crucially, we quantitatively show that the now semantically expanded dataset is able to improve rare class performance with even fewer rare but semantically important objects, demonstrating the strong advantage semantics imbue. Together, our semantic selection and enrichment contributes two main outcomes. First, we can reduce computational cost and achieve similar model performance. And second, a carefully handcrafted dataset can be semantically tuned without increasing the dataset size for better performance.

## 6. Acknowledgements

## References

[1] Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. *arXiv preprint arXiv:1906.03671*, 2019.

[2] Alessandro Bertoni and Tobias Larsson. Data mining in product service systems design: Literature review and research questions. *Procedia CIRP*, 64:306–311, 2017.

[3] 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. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020.

[4] Nadine Chang, Francesco Ferroni, Michael J Tarr, Martial Hebert, and Deva Ramanan. Thinking like an annotator: Generation of dataset labeling instructions. *arXiv preprint arXiv:2306.14035*, 2023.

[5] Jiwoong Choi, Ismail Elezi, Hyuk-Jae Lee, Clement Farabet, and Jose M Alvarez. Active learning for deep object detection via probabilistic modeling. In *ICCV*, pages 10264–10273, 2021.

[6] Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. *NeurIPS*.

[7] Sachin Goyal, Pratyush Maini, Zachary C Lipton, Aditi Raghunathan, and J Zico Kolter. Scaling laws for data filtering–data curation cannot be compute agnostic. In *CVPR*, pages 22702–22711, 2024.

[8] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *CVPR*, pages 5356–5364, 2019.

[9] Elmar Haussmann, Michele Fenzi, Kashyap Chitta, Jan Ivanecky, 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*, pages 1430–1435. IEEE, 2020.

[10] Chieh-Chi Kao, Teng-Yok Lee, Pradeep Sen, and Ming-Yu Liu. Localization-aware active learning for object detection. In *ACCV*, pages 506–522. Springer, 2019.

[11] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *NeurIPS*, 30, 2017.

[12] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023.

[13] Jeremiah Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax Weiss, and Balaji Lakshminarayanan. Simple and principled uncertainty estimation with deterministic deep learning via distance awareness. *NeurIPS*, 33:7498–7512, 2020.

[14] Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv preprint arXiv:2104.08786*, 2021.

[15] Moran Mizrahi, Guy Kaplan, Dan Malkin, Rotem Dror, Dafna Shahaf, and Gabriel Stanovsky. State of what art? a call for multi-prompt llm evaluation. *Transactions of the Association for Computational Linguistics*, 12:933–949, 2024.

[16] Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding important examples early in training. *NeurIPS*, 34:20596–20607, 2021.

[17] Trung Pham, Mehran Maghoumi, Wanli Jiang, Bala Siva Sashank Jujjavarapu, Mehdi Sajjadi, Xin Liu, Hsuan-Chu Lin, Bor-Jeng Chen, Giang Truong, Chao Fang, et al. Nvautonet: Fast and accurate 360 3d visual perception for self driving. *arXiv preprint arXiv:2303.12976*, 2023.

[18] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

[19] Soumya Roy, Asim Unmesh, and Vinay P Namboodiri. Deep active learning for object detection. In *BMVC*, volume 362, page 91, 2018.

[20] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022.

[21] Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. *arXiv preprint arXiv:2310.11324*, 2023.

[22] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. *arXiv preprint arXiv:1708.00489*, 2017.

[23] Xiaoling Shu and Yiwan Ye. Knowledge discovery: Methods from data mining and machine learning. *Social Science Research*, 110:102817, 2023.

[24] Eric Slyman, Stefan Lee, Scott Cohen, and Kushal Kafle. Fairdedup: Detecting and mitigating vision-language fairness disparities in semantic dataset deduplication. In *CVPR*, pages 13905–13916, 2024.

[25] Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial intelligence and machine learning for multi-domain operations applications*, volume 11006, pages 369–386. SPIE, 2019.

[26] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. Mpnet: Masked and permuted pre-training for language understanding. *Advances in neural information processing systems*, 33:16857–16867, 2020.

[27] Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. Beyond neural scaling laws: beating power law scaling via data pruning. *NeurIPS*, 35:19523–19536, 2022.

[28] Haoru Tan, Sitong Wu, Fei Du, Yukang Chen, Zhibin Wang, Fan Wang, and Xiaojuan Qi. Data pruning via moving-one-sample-out. *NeurIPS*, 36, 2024.

[29] Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. *arXiv preprint arXiv:1812.05159*, 2018.

[30] Matias Valdenegro-Toro. Deep sub-ensembles for fast uncertainty estimation in image classification. arxiv preprint arxiv: 191008168.(2019) doi: 10.48550. *arXiv*, 1910.

[31] Anton Voronov, Lena Wolf, and Max Ryabinin. Mind your format: Towards consistent evaluation of in-context learning improvements. *arXiv preprint arXiv:2401.06766*, 2024.

[32] Yeming Wen, Dustin Tran, and Jimmy Ba. Batchensemble: an alternative approach to efficient ensemble and lifelong learning. *arXiv preprint arXiv:2002.06715*, 2020.

[33] Donggeun Yoo and In So Kweon. Learning loss for active learning. In *CVPR*, pages 93–102, 2019.

## A. Additional Cluster Visualizations

We provide additional cluster visualization in Fig. 13. The cluster contains scenes where there are construction car or public service car near the ego car, likely with workers or pedestrians.

## B. Semantic Retrieval Visualizations

We provide additional examples of retrieval in Fig. 14. Semantic retrieval is able to understand the more complex semantics "car stops" and "breaking lights on", while CLIP retrieval focus on "lights on" and returns mostly night images with lights on.

## C. Detailed Per-Class AP

We provide detailed numbers for mAP and per-class AP for our main results in Tab. 3.

Figure 13: Visualization of samples in a cluster. The scenes are visually different but semantically similar. (Construction car/public service car, with potential worker/pedestrian)



Query: The car in front stops with breaking lights on.

Figure 14: Additional retrieval examples using advanced visual embeddings and semantic embeddings.

Table 3: Per-class AP of 3D object detection with different data strategies.

|  | Dataset Size (% of original) | Method | mAP | Per-Class AP | | | |
|---|---|---|---|---|---|---|---|
|  |  |  |  | Car | Truck | Person | Bike with Rider |
|  | 100% | Original | 65.6 | 81.0 | 71.3 | 47.1 | 63.0 |
| Selection | 70% | Random | 62.1 | 79.4 | 67.6 | 43.2 | 58.1 |
|  |  | Long-tail [8] | 62.3 | 78.6 | 67.8 | 43.1 | 59.9 |
|  |  | CLIP visual [18] | 64.4 | **80.6** | **69.7** | 45.7 | 61.5 |
|  |  | **SSD (Ours)** | **65.2** | 80.5 | 69.5 | **46.9** | **63.8** |
| Enrichment | 100% | Random | 66.5 | 82.0 | 69.6 | 50.2 | 64.2 |
|  |  | Long-tail [8] | - | - | - | - | - |
|  |  | CLIP visual [18] | 66.0 | 82.0 | 69.7 | 49.6 | 62.6 |
|  |  | **SSD (Ours)** | **67.6** | **82.1** | **71.6** | **50.3** | **65.6** |