# LIP-Loc: LiDAR Image Pretraining for Cross-Modal Localization

Sai Shubodh Puligilla [*†]     Mohammad Omama[‡]     Husain Zaidi[§]     Udit Singh Parihar[†]

Madhava Krishna[†]

## Abstract

*Global visual localization in LiDAR-maps, crucial for autonomous driving applications, remains largely unexplored due to the challenging issue of bridging the cross-modal heterogeneity gap. Popular multi-modal learning approach Contrastive Language-Image Pre-Training (CLIP) [32] has popularized contrastive symmetric loss using batch construction technique by applying it to multi-modal domains of text and image. We apply this approach to the domains of 2D image and 3D LiDAR points on the task of cross-modal localization. Our method is explained as follows: A batch of N (image, LiDAR) pairs is constructed so as to predict what is the right match between N X N possible pairings across the batch by jointly training an image encoder and LiDAR encoder to learn a multi-modal embedding space. In this way, the cosine similarity between N positive pairings is maximized, whereas that between the remaining negative pairings is minimized. Finally, over the obtained similarity scores, a symmetric cross-entropy loss is optimized. To the best of our knowledge, this is the first work to apply batched loss approach to a cross-modal setting of image & LiDAR data and also to show Zero-shot transfer in a visual localization setting. We conduct extensive analyses on standard autonomous driving datasets such as KITTI and KITTI-360 datasets. Our method outperforms state-of-the-art recall@1 accuracy on the KITTI-360 dataset by 22.4%, using only perspective images, in contrast to the state-of-the-art approach, which utilizes the more informative fisheye images. Additionally, this superior performance is achieved without resorting to complex architectures. Moreover, we demonstrate the zero-shot capabilities of our model and we beat SOTA by 8% without even training on it. Furthermore, we establish the first benchmark for cross-modal localization on the KITTI dataset.*

## 1. Introduction

Visual localization serves as a crucial aspect in the domain of mobile robotics, playing a pivotal role in applications such as autonomous vehicles and Simultaneous Localization and Mapping (SLAM). Our focus in this work is to address the complex issue of determining the pose of an image within an expansive 3D map. This task of localization gets quite challenging due to large, occluded, dynamic scenes and repetitive places.

Outdoor and indoor visual localization approaches may differ in their pipelines because of the different nature of their challenges; our work addresses outdoor visual localization. However, they can both broadly be categorized into 2 steps: global localization, which gives a rough estimate of the pose, and then, local localization, which gives a more accurate version that further involves the use of PnP [16, 18, 21] in a RANSAC [22] setting. In our paper, we deal with the problem of global localization. Global localization can be achieved using Global Navigation Satellite Systems (GNSSs), however, it does not always give reliable estimates. This is particularly prevalent in urban environments where high-rise buildings can interfere with the signal quality, leading to inaccuracies. Other contributing factors can include multipath effects, where signals bounce off multiple surfaces before reaching the receiver, and atmospheric conditions, which can alter the signal speed. Therefore, Light Detection And Ranging (LiDAR)-based [24,44] and vision-based approaches [1, 5, 11] are seen as established sensor modalities in the vision community to estimate the pose of a robot accurately. While LiDAR modality is robust to variation in illumination and can detect objects at long distances with high accuracy, they are generally expensive and are especially prone to failure modules such as degenerate places like tunnels and suffers from issues such as surface reflections and interference. Vision modality-based methods use 2D images and extract features from them using methods such as NetVLAD [1] to eventually match them with a query image for localization. While vision-based approaches have seen large success, there are still important limitations, such as dynamic environments, illumination, or weather changes.

Given the complementary nature of these sensor modal-

---
[*]Corresponding author: p.saishubodh@gmail.com, Project page: https://shubodhs.ai/liploc

[†]Robotics Research Center, KCIS, IIIT Hyderabad

[‡]University of Texas at Austin

[§]Microsoft

ities, their advantages can be combined in a multi-modal fashion through the fusion of 2D and 3D data [19, 28, 35], which enhances localization accuracy significantly. However, this fusion is not straightforward because of the heterogeneity gap between these two modalities, and therefore, it remains an unexplored and largely unsolved problem. Further, the multi-modal approach does not solve the problem of localizing a sensor of one modality in a map of another, something that we refer to as 'cross-modal localization'. Our task revolves around the novel application of a contrastive loss based on batch construction approach to the distinct domains of 2D images and 3D LiDAR point clouds, specifically in the context of cross-modal localization. This involves creating a shared embedding space for both 2D and 3D data, enabling the localization process to occur even when only one modality is available at a given time.

The practical motivation of our work lies in its utility in autonomous navigation scenarios. Imagine a setting where a detailed and expensive LiDAR map has been constructed using a resource-intensive setup initially. In subsequent navigation instances, however, our approach enables localization solely via 2D images, thereby eliminating the need for resource-heavy operations. This feature is particularly advantageous as it mitigates resource constraints and provides an economical and efficient solution for repeat localizations. Furthermore, our method showcases its versatility by being applicable even when the initial map is constructed in a different modality, such as 3D. Thus, our work facilitates cost-effective and efficient navigation by capitalizing on the power of cross-modal localization.

The main contributions of our paper are as follows:

- **Batched Loss Approach:** This work is the first of its kind to apply the batched contrastive approach in a cross-modal setting involving image and LiDAR data, establishing a novel direction in metric learning for autonomous driving applications.

- **Superior Performance with Simpler Methods:** We demonstrate that our method outperforms state-of-the-art (*AECMLoc*) [51] recall@1 accuracy on the KITTI-360 [23] dataset by 22.4% using only perspective images and standard Vision Transformer [8] architecture for the encoders, contrasting with the state-of-the-art approaches that rely on more informative fisheye images and complex architectures.

- **Zero-shot Analyses and Benchmark Establishment:** We conduct exhaustive analyses on standard autonomous driving datasets such as KITTI [10] and KITTI-360 [23] and establish the first benchmark for cross-modal localization on the KITTI dataset.

The remainder of this paper is organized in the following manner: Section II provides an overview of previous research in the domain of visual localization. Our proposed methodology, encompassing batch construction, contrastive loss, and architecture design, is detailed in Section III. Experiments conducted on public datasets and their results are exhibited in Section IV. In Section V, we demonstrate the zero-shot capability of our model. The paper concludes with Section VI.

## 2. Related Work

Localization of a robot involves understanding where it is in the world using a pre-existing map. Generally, this has been done using the same type of sensor that was used to create the map, such as images with images or 3D scans with 3D scans, i.e., between the same corresponding modalities. Our work expands upon this by showing that you can localize using different types of sensors than those used to create the map, like using simple cameras to localize in a map built from expensive 3D scans, which is more flexible and cost-effective. We review the literature on both of these approaches in this section.

### 2.1. Same modal localization

The standard pipeline for localization approaches begins with the acquisition of reference data, which is typically a large 3D map. The first step in the pipeline is to retrieve prior information for which traditional methods include the bag of words approach [27], but more recent work has leveraged deep learning techniques. For instance, Arandjelovic *et al.* [1] extended the Vector of Locally Aggregated Descriptors (VLAD) [15] approach, introducing a differentiable generalized VLAD layer that can be integrated into any CNN architecture, i.e., NetVLAD [1], a CNN-based image retrieval algorithm, retrieves the most similar images, or reference images, from an image database. This stage is called global localization, where we can directly take the pose based on the most similar reference images for a given query image. But this pose can be refined further as follows. Once the reference images are retrieved, local feature extraction and matching are performed. Feature extraction has traditionally relied on methods such as Scale-Invariant Feature Transform (SIFT) [50], Speeded Up Robust Features (SURF) [2], and Oriented FAST and Rotated BRIEF (ORB) [37], which are designed to extract local features from the image. However, with the advent of deep learning, recent years have seen a shift towards the use of convolutional neural networks (CNNs) such as ResNet [12], ConvNets [17] and more recently, Vision Transformers ViT [8] as a means to extract features. These neural networks have shown remarkable performance in feature extraction, replacing handcrafted features with learned representations. Following local feature matching, the Pose from N Points

(PnP) [16, 18] algorithm along with RANSAC is leveraged to estimate the 6-DoF pose. This whole process to obtain finer estimate of pose is termed local localization. In our current work, we apply our batched loss approach to the global localization problem. No single localization method currently exists that can universally adapt to a wide range of environments, such as urban landscapes, rural settings, nighttime conditions, warehouses, varying weather, foggy conditions, and busy marketplaces [6, 26, 30, 38, 42, 46]. Most state-of-the-art methods specialize in one or a few of these scenarios and lack the generalizability to perform well across all. Our paper aims to take a step in this direction towards creating a more general localization method.

So far, we have discussed about image-to-image based retrieval methods. Now let us discuss the same modal retrieval for point clouds. Point clouds provide a robust way to represent scenes under varying lighting conditions and seasonal changes and have the ability to maintain the structural integrity of scenes. The initial developments in this field were PointNet [31], and EdgeConv [48], which process unordered points to extract permutation-invariant features. Building upon this, PointNetVLAD [44] introduced an end-to-end trainable model, which merges the strengths of PointNet and NetVLAD for point cloud-based place recognition. When a query point cloud is provided, the task involves retrieving the most similar sub-maps from this database, i.e., this is 3D-3D based retrieval. However, the cost and weight of LiDAR technology can be prohibitive for large-scale applications. Here is where our work has significant motivation. Consider a scenario where an exhaustive and expensive map has been previously constructed using a combination of sensor setups. With our approach, you can localize within this map multiple times using only simple and cost-effective sensors, such as RGB cameras. This eliminates resource constraints and allows for flexibility even if the map has not been pre-built in that specific sensor modality.

## 2.2. Cross modal localization

Cross modal localization has gained traction in recent years [3, 25, 29, 41]. Early works which localize between the domains of LiDAR and 2D images include [3, 25, 41]. These approaches use CNN based approaches to do pose regression with standard translation or rotation-based losses. These methods typically need some initial estimate to work and only works as a local localization task. 2D3D-MatchNet [9] proposes a deep network to jointly learn the 2D image and 3D point cloud keypoint descriptors.

More recently, Cattaneo *et al*. [4] proposed a teacher-student training approach on Oxford Robotcar using triplet loss and jointly trains a 2D network for the images and 3D network for the point clouds. Similarly, AdaFusion [19] presents a visual-LiDAR descriptor fusion in a weighted
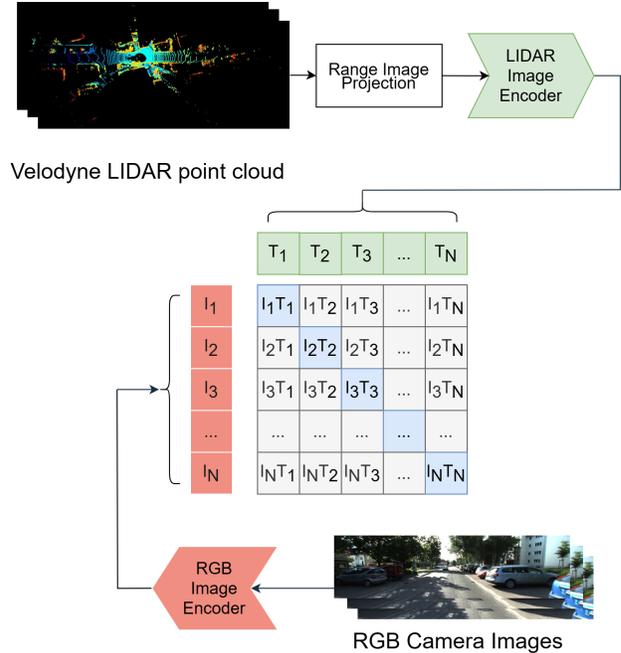


Figure 1. Batched Contrastive Learning Architecture

way using a pairwise margin-based loss. A similar approach i3dLoc [49] uses Generative Adversarial Networks (GANs) on the domains of 360-images and 2D range projections to deal with inconsistent environmental conditions on a custom dataset. Recently, LiDARCLIP [14] uses principles from CLIP to produce joint text, LiDAR and camera embeddings, although not for the task of global localization. The work closest to ours is AECMLoc [51] which is the first work to address the task of cross-modal localization on KITTI-360 dataset. It uses a spherical convolution based 2D network and another PointNet based 3D network with attention enhancement with a triplet loss and achieves reasonable cross-modal localization accuracy. In contrast to traditional methods which employ a triplet loss with select negative samples through hard negative mining, we introduce the use of a contrastive loss incorporating numerous negative samples within a batch construction method.

## 3. Methodology

In this section, we discuss in detail our methodology where we first discuss how the data is constructed in a batched manner and how contrastive loss is imposed upon it. Then we discuss the details about encoder architecture and training.

### 3.1. Batch Construction and Contrastive Loss

Here, we discuss the batch construction procedure and the contrastive symmetric loss in detail. This approach

was first introduced as multi-class N-pair loss [40] and then popularized by InfoNCE [45] and CLIP [32] under various names. In the context of our paper, we call it batched contrastive loss.

Firstly, our batch construction method is explained as follows: A batch of $N$ (image, LiDAR) pairs is constructed so as to predict what is the right match between $N \times N$ possible pairings across the batch by jointly training an image encoder and LiDAR encoder to learn a multi-modal embedding space. In this way, we would have $N$ positive pairings and $(N^2 - N)$ negative pairings as shown in Fig 1.

To put it formally, let us say $x_i$ represents a 2D image, $x_i^+$ represents a LiDAR sample, and $f(x)$ represents an embedding vector for $x$, simply written as $f$. In a batch size of $N$, we have $N$ such pairs $\left\{ \left( x_1, x_1^+ \right), \cdots, \left( x_N, x_N^+ \right) \right\}$. First, let us consider for one 2D image (the same explanation works vice-versa for LiDAR given the symmetric nature of loss). Consider an $(N+1)$-tuplet of one 2D image and $N$ LiDAR samples i.e. $S_i = \left\{ x_i, x_1^+, x_2^+, \cdots, x_N^+ \right\}$: The anchor here is $x_i$ while $x_i^+$ is a positive example to the anchor and $x_j^+, j \neq i$ are the negative examples. In other words, every 2D image would have 1 positive and $N-1$ negative LiDAR examples (and vice-versa). Therefore, our loss can be finally expressed as:

$$
\mathcal{L}_{\text{batched}} \left( \left\{ (x_i, x_i^+) \right\}_{i=1}^N ; f \right) = \frac{1}{N} \sum_{i=1}^N \log \left( 1 \right.
$$
$$
\left. + \sum_{j \neq i} \exp \left( f_i^\top f_j^+ - f_i^\top f_i^+ \right) \right)
\tag{1}
$$

Each $i$ in the outer summation would correspond to every row in the Fig. 1. The above loss can equivalently be expressed as standard softmax loss as follows (for full theoretical details, refer to [40]):

$$
\mathcal{L}_{\text{batched}} \left( \left\{ (x_i, x_i^+) \right\}_{i=1}^N ; f \right) =
$$
$$
-\frac{1}{N} \sum_{i=1}^N \log \left( \frac{\exp \left( f_i^\top f_i^+ \right)}{\exp \left( f_i^\top f_i^+ \right) + \sum_{j \neq i} \exp \left( f_i^\top f_j^+ \right)} \right)
\tag{2}
$$

This loss is used to train our dual encoders as explained in the next section.

## 3.2. Contrastive encoder training

Our methodology uses contrastive learning to jointly train two encoders: one for 2D images from the camera and the other for 3D points from the LiDAR. We generate range images from the LiDAR points. In our experimentations, we found that using range images when passed
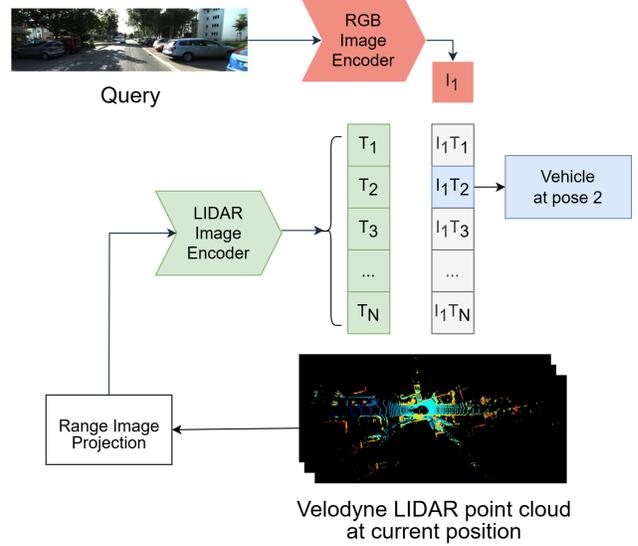


Figure 2. Inference pipeline showing that query camera image is calculating similarity with all the lidar range images in the database



Figure 3. Qualitative visualization: 2D query, 3D prediction, 3D ground truth.

through standard ResNet or ViT architectures gave us better performance in comparison with specialized or sophisticated architectures, which use 3D point clouds as input. The training process ensures corresponding camera and LiDAR images are closely aligned by learning a shared multi-modal embedding space and, thus, enables efficient cross-modal localization.

Both encoders follow the same two-step process: feature extraction and projection. For feature extraction, we employ the established ViT [8] model. We have also experimented with other models such as ResNet [12] which too gives superior performance to baseline (ResNet backbone beats baseline by more than 10%), however amongst

these two, ViT is much better. We deliberately opted for a standard architecture, aiming to highlight the effectiveness of our batched loss method without resorting to more advanced architectures. The feature extraction step uses the pre-trained *vit_small_patch16_224* model to transform the input images into intermediate feature vectors, and then, the feature projection step transforms these high-dimensional vectors into a shared and lower-dimensional embedding space. The projection step is aimed at bolstering the task-specific features and preserving the compact nature of embeddings to ensure computational efficiency. We found *vit_small_patch16_224* to be our best model amongst other variants, henceforth, whenever not explicitly mentioned, it can be assumed that we are referring to this model by 'LIP-Loc'.

The joint training of the encoders ensures both the camera and LiDAR encoders update their parameters during each training step simultaneously to maximize the cosine similarity between the embeddings of the corresponding camera and LiDAR pairs and to minimize the similarity between those of non-corresponding pairs. This is achieved by optimizing a symmetric cross-entropy loss based on a batch construction procedure, as explained previously.

## 4. Experiments and Results

We have shown the effectiveness of our contrastive learning via Batch Construction in two different datasets, KITTI and KITTI 360. Our ablation results show the improvement of batched contrastive loss over triplet loss in terms of network recall and GPU memory footprint. With the incorporation of cropping in the image field of view and distance thresholding in LiDAR space, we are able to achieve better generalization, as shown quantitatively in the table.

### 4.1. Dataset and Preprocessing

For both KITTI and KITTI 360, our training pairs consist of synchronized LiDAR scans and camera images taken in the same snapshot of the world.

**KITTI dataset**

Evaluation is performed on the KITTI odometry dataset which consists of Velodyne HDL-64E LiDAR scans and a color stereo camera rig. Our experiment id and corresponding training sequences are shown in Table 1. The evaluation sequences are 08 and 09.

**KITTI 360 dataset**

It comes with Velodyne HDL-64E and raw images from the perspective camera. This has about 80k frames of lidar and image pairs over a distance of 73.7 Km, along with precise vehicle pose information require for evaluation. We have used sequences 3,4,5,6,7,9,10 for training and sequence 0 for evaluation. Of course, in the Zero-shot setting, we did not use any of this. To not confuse the standard train-test evaluation with Zero-shot, we separately dedicate

| Experiment id | Training sequences |
|---|---|
| exp_large | KITTI(5): 03, 04, 05, 06, 07 |
| exp_larger | KITTI(8): 00, 01, 02, 03, 04, 05, 06, 07 |
| exp_largest | KITTI(18): 00, 01, 02, 03, 04, 05, 06, 07, 11, 12, 13, 15, 16, 17, 18, 19, 20, 21 |
| exp_360 | KITTI_360(7): 03, 04, 05, 06, 07, 09, 10 |

Table 1. Training sequences for different experiments.

Section 5 for Zero-shot results. We have applied distance threshold in lidar scans and field of view threshold on lidar range image.

### 4.2. Objective and Evaluation metrics

We further clearly elaborate the objective here. During training time, we take a batch of N (image, LiDAR) pairs, where N is typically 32 as illustrated in Fig 1 which we train for 50 number of epochs. It is hereby reiterated that hyperparameter tuning is not necessary and our method works out of the box as mentioned in next section in even a zero shot setting. The inference procedure is explained in Fig 2, where given an RGB query, we pass it through its corresponding encoder to extract the embedding and similarly for reference LiDAR range images to extract corresponding LiDAR embeddings. Then, the reference with highest similarity score with query is picked as final prediction whose pose is finally extracted for the global localization problem. The same process is done vice-versa for 3D to 2D localization task.

We used the Recall@1 evaluation metric with a distance threshold of 20 meters. For a given query camera image embedding we find the closest LiDAR range image embedding in joint space, if the distance between the retrieved camera and LiDAR is less than 20 meters, then we consider it as a True Positive. Thus,

$$\text{Recall@1} = \frac{\text{Number of True LiDAR and camera matches}}{\text{Total number of query camera images}}$$

### 4.3. Comparative Analysis: Triplet Versus Batched Contrastive

Our experimental findings align with the premise that an increase in negative samples in metric learning contributes to better generalization, a claim also made by foundational works such as CLIP [32]. We observe the same in the results we obtained, as illustrated in Table 2.

The traditional triplet approach tends to require more GPU memory as the number of negative samples increases [40]. This can create limitations when working with larger datasets or when more robust learning is required.

In contrast, the batched contrastive approach presents a more efficient and scalable solution. It has demonstrated improved performance with the addition of more data [40], making it particularly suitable for metric learning. The scalability of this approach allows for the management of larger datasets, making it a key consideration for improving the efficiency of intermediate modules in the learning process. Table 2 shows comparison of the batched contrastive loss against standard triplet loss for lidar-camera alignment on the KITTI dataset. Please note that we use ResNet50 for this analysis and not ViT. We see that with batched construction, we can get more negative samples with lesser GPU footprint leading to better evaluation results.

| Exp_ID/Metric | triplet loss | batched loss |
|---|---|---|
| Seq 8 r@1 | 0.215 | 0.295 |
| Seq 9 r@1 | 0.232 | 0.309 |
| GPU Utilization | ∼8214MB×2 | ∼8214MB |
| No of -ve Samples | 1 | 31 |
| Batch Size | 32 | 32 |

Table 2. Recall@1 comparison between Triplet and Batched Contrastive Training for *exp_larger* setting i.e. no of sequences is 8.

## 4.4. Data Preprocessing for Better Generalization

A lower quantity of data can potentially reduce the model's ability to generalize effectively. To mitigate this, we advocate the incorporation of intelligent pre-processing techniques designed to boost generalization. In our study, we experiment with distance threshold cropping for LiDAR data and field of view (FoV) cropping for LiDAR range images. The horizontal field of view of a LiDAR range image is greater than the camera field of view, so we crop the LiDAR range image such that both sensors have common overlapping information present.

The rationale for utilizing distance cropping stems from the observation that objects at a greater distance while being accurately captured by LiDAR, may not be equally discernible through camera imaging. In contrast, if we constrain the LiDAR data too much to a more confined area, we risk losing sight of more distant meaningful, and relevant information for our model.

Our empirical analysis in Table 3 strongly indicates that a distance threshold of 50 meters for LiDAR cropping provides the most beneficial outcome for our model's performance. Adjusting this threshold to either higher or lower distances tends to degrade the overall results, underscoring the critical role of this specific parameter in optimizing the data preprocessing stage for better generalization.

## 4.5. Scaling Data for Better Generalization

The robustness of metric learning is directly proportional to the volume of data available for processing [32]. This

| Sequence/ExpID | Seq 8 | Seq 9 |
|---|---|---|
| exp_larger ‖ No Threshold | 0.295 | 0.309 |
| exp_larger ‖ (50m Threshold) | 0.325 | 0.370 |
| exp_larger ‖ (30m Threshold) | 0.317 | 0.308 |
| exp_largest ‖ No Threshold | 0.484 | 0.457 |
| exp_largest ‖ (50m Threshold) | 0.540 | 0.495 |

Table 3. Ablation of thresholding on lidar scans with increase in training sequences

concept is clearly exemplified in our empirical results, as presented in Table 4. As we progressively increase the number of sequences, there is a noticeable upswing in accuracy across both sequences.

To conclude, our results emphasize the importance of leveraging scalable techniques and large, varied datasets in metric learning, as this approach can notably enhance model accuracy and its generalization capabilities.

| Sequence/ExpID | Seq 8 | Seq 9 | Sequences |
|---|---|---|---|
| exp_large | 0.278 | 0.260 | 5 |
| exp_larger | 0.547 | 0.525 | 8 |
| exp_largest | 0.805 | 0.780 | 18 |

Table 4. Increasing in number of sequences and mixture of datasets leads to better generalization

## 4.6. Setting New Benchmark on KITTI-360 Cross-Modal Place Recognition Task

In this section, we demonstrate our approach's superiority to the prior state-of-the-art *AECMLoc* [51] on the KITTI-360 dataset (Tables 5 and 6). Unlike the former, which relied on fisheye images requiring complex preprocessing due to distortion and uses a specialized architecture, we employ perspective images. Figure 4 demonstrates LIPLoc's superiority over AECMLoc in both Zero-shot and standard same-dataset training by plotting recall values from $k$'s value of 1 to 20. LIP-Loc overwhelmingly beats the baseline by about 20% in both 2D-3D as well as 3D-2D settings; and at Recall@20, we almost reach 97% accuracy, meaning that our method would be a robust retrieval method for further fine grained localization. In their paper, AECM-Loc show that these kind of 95+% recall@20 values are observed only in same modal localization, i.e. 2D to 2D or 3D to 3D. It is interesting to note that our method reaches similar values while being a cross modal method. Furthermote, notice that baseline's 3D to 2D reduces by 15% compared to 2D to 3D, whereas ours is almost the same, demonstrating the versatility of our method. Note that we beat the baseline method without even training on KITTI-360 by 8%. The next section addresses this in detail.

| 2D to 3D | recall@1 | recall@5 | recall@20 |
|---|---|---|---|
| AECMLoc | 0.462 | 0.660 | 0.782 |
| LIP-Loc | 0.686 | 0.868 | 0.966 |
| Zero-shot LIP-Loc | 0.540 | 0.770 | 0.919 |

Table 5. Baseline comparison of Recall values for 2D query to 3D database localization

| 3D to 2D | recall@1 | recall@5 | recall@20 |
|---|---|---|---|
| AECMLoc | 0.311 | 0.472 | 0.710 |
| LIP-Loc | 0.6982 | 0.8745 | 0.9665 |
| Zero-shot LIP-Loc | 0.574 | 0.809 | 0.946 |

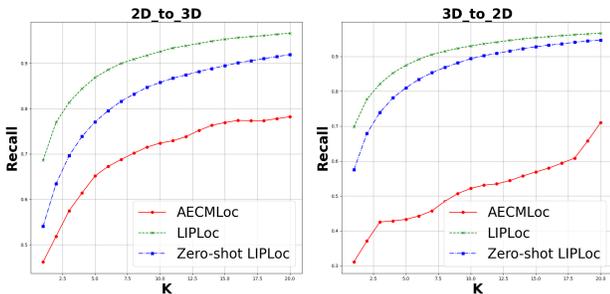Table 6. Baseline comparison of Recall values for 3D query to 2D database localization



Figure 4. Recall@k plot for various values of k for baseline comparison with our LIP-Loc and Zero-shot LIP-Loc.

## 5. Zero-shot Transfer

### 5.1. Standard Definition of Zero-shot

Zero-shot learning in the context of computer vision refers to the ability of a model to generalize to classes not seen during training [20]. CLIP redefines this term and extends it to refer to generalisation to unseen datasets.

CLIP attempts to emphasize the task-learning capabilities of models through zero-shot transfer; however, since popular computer vision datasets are inclined towards generic image classification rather than task-specific evaluations, their analysis on these datasets primarily serve as assessment to domain generalization and robustness to distribution shift. We also focus on the latter in our paper i.e. domain generalization.

### 5.2. Zero-shot Transfer for Localisation

In the context of visual localization, we define zero-shot transfer as the model's capability to estimate the pose of an object in unseen datasets. Note that as an early work in this area, we are referring to coarse estimate of pose i.e. global localization problem. To the best of our knowledge, zero-shot transfer has never been applied to visual localization
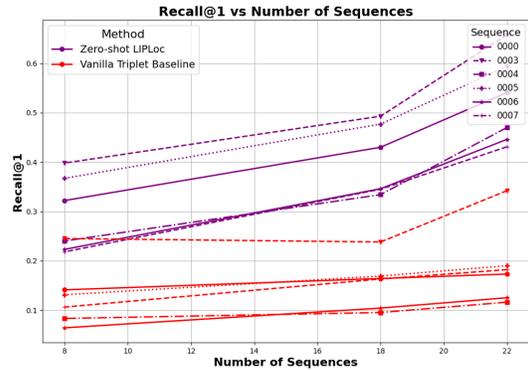


Figure 5. Zero-shot Scaling: Here, x-axis corresponds to the number of sequences of KITTI. y-axis refers to recall@1 on different sequences of KITTI-360. Note that our batched loss scales with increase in dataset size, while triplet loss which most standard methods use shows very marginal increase.

previously.

It is worth mentioning that while our work is inspired by CLIP's application on computer vision tasks, CLIP itself was inspired by GPT-1 [33] and GPT-2 [34] which have studied zero-shot transfer over the course of pre-training and "unexpected" task-learning capabilties of language models.

In our paper, we train the model on KITTI dataset and evaluate it on KITTI-360, the baseline for which is AECM-Loc which is exclusively trained on KITTI-360. Both these datasets are captured from different camera modalities: KITTI-360 employs 180-degree FOV fisheye cameras, while KITTI uses 90-degree FOV cameras. In fact, when KITTI-360 input data is converted to perspective images, there is heavy perspective distortion because of which training directly on perspective results in under-performance as AECMLoc has reported. More importantly, KITTI 360 has no trajectory overlap with KITTI. All these factors make KITTI-360 an interesting candidate for out-of-distribution dataset. We do acknowledge that KITTI-360 might not be analogous to how ImageNet Rendition [13] or ImageNet Sketch [47] was for ImageNet, but all these factors could make it equivalent to ImageNetV2 [36].

### 5.3. Results and Analysis: Scaling & Robustness

To be truly Zero-shot, we report our accuracy values without doing any customization of hyperparameters for KITTI-360 nor do we form any validation sets on KITTI-360. The trained model truly works out of the box and outperforms AECMLoc which is trained on KITTI-360. We further compare the Zero-shot capabilities of our approach with the 'vanilla-triplet' baseline which uses the standard triplet loss formulation which most standard methods use.

Our best model which we refer to as "Zero-shot LIPLoc" is trained on full KITTI dataset and beats AECMLoc by 8%. It reaches comparable accuracy to AECMLoc when trained on about 80% of the KITTI sequences.

It is common in computer vision that models scale with dataset size. However, such phenonemon is non-trivial in localization and hence we have shown in the previous section that increase in number of sequence increases accuracy. In fact, it is specific sometimes to such an extent that there could be a separate model fitted every sequence of dataset [7] and these models don't even work on other sequences of the same dataset, let alone on another dataset.

Therefore, we go one more step ahead and show in Figure 5 that as we scale up training on original KITTI data, the accuracy on KITTI-360 progressively improves although Zero-shot LIP-Loc has never seen KITTI-360; demonstrating domain generalization. The average increase in Zero-shot LIP-Loc is 23% whereas the that in baseline is 6%, clearly showing the scalability of constrastive formulation. The more common triplet formulation which is used in most standard methods only increases by few percentage points. We have tried bigger ViT models as encoders, but the accuracy saturates. This is because we are dealing with much smaller datasets, the model may become overparameterized.

In CLIP, as they train on internet scale data, they admit that it is ambiguous what exactly results in accuracy increase: data, model or loss function? But in our case, we exploit the benefit of working with smaller dataset and clearly explained how each factor contributed to our training.

How well does a Zero-shot model work on out of distribution datasets? Typically in deep learning when models are trained and evaluated on same dataset like ImageNet, they exploit spurious correlations because of which robustness gap arises. CLIP does this robustness analysis by testing on 7 natural distribution shift datasets [43]. Such a benchmark does not exist for visual localization, therefore our work motivates the building of such a benchmark.

Within the scope of this paper, we do the robustness analysis between sequences of KITTI and KITTI-360 in Fig 6. We firstly consider models which are trained on subset of KITTI data as per Table 1. Then we evaluate on the rest of KITTI data as test dataset whose recall we plot on x-axis and then evaluate on full KITTI-360 sequences whose recall we plot on y-axis. An ideal robust model would perform equally well on both these test sets, i.e. y=x line. When our curve is closest to the robust model plot, we notice in Fig 6 that when there is only about 10-20% accuracy gap between seen validation dataset and unseen data, proving the robustness of LIP-Loc. This holds true even when the model is trained on a smaller dataset.

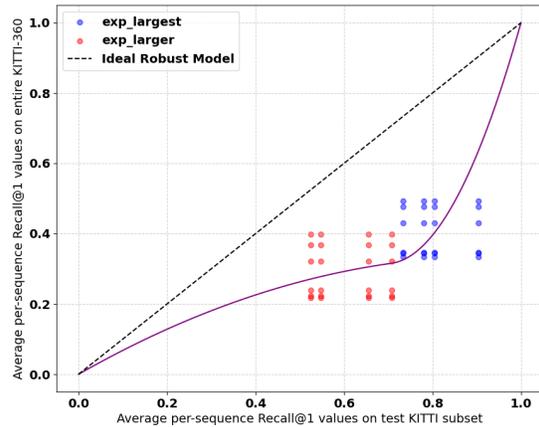Note that this is only the first step towards establishing



Figure 6. Zero-shot Robustness: When recall on KITI subset is around 50%, our zero-shot model reaches upto 40% on KITTI-360; when it is around 80%, our model reaches upto 50% recall, strongly demonstrating robustness.

robustness of Zero-shot localization models. There needs to be a dedicated benchmark along with standardised metrics to truly evaluate these models. Since there is no equivalent of internet scale (text, image) data for localization, how we train our localization models and evaluate them is an open question, as we discuss in the next section.

## 6. Future Work

Our approach showcases the potential of batched contrastive learning in bridging the cross-modal heterogeneity gap. By achieving superior performance without relying on complex architectures or fisheye images, our method offers a simpler yet highly effective solution for cross-modal localization. Additionally, we establish the first benchmark for cross-modal localization on the KITTI dataset, providing a foundation for future research.

Our reflections on Zero-shot transfer open a set of thought provoking questions: What is the internet scale equivalent for localization? Can a larger version of dataset like KITTI act as like one? Or would it involve synthetic dataset? What is the equivalent of natural distribution scale datasets for localization? Would recall@K be the right metric for evaluating such systems or do we need better metrics? One of the weaknesses of CLIP is its task learning capabilities, for instance CLIP struggles to find out the closest objects in an image. Could combining depth encoder with text encoder solve this problem and other task generalization problems?

All of these are predominantly open questions which we hope our work will motivate the readers to address in their own work.

# References

[1] Relja Arandjelovic, Petr Gronat, Akihiko Torii, Tomas Pajdla, and Josef Sivic. Netvlad: Cnn architecture for weakly supervised place recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(6):1437–1451, Jun 2018. 1, 2

[2] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. Surf: Speeded up robust features. In Aleš Leonardis, Horst Bischof, and Axel Pinz, editors, *Computer Vision – ECCV 2006*, pages 404–417, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg. 2

[3] D. Cattaneo, M. Vaghi, A. L. Ballardini, S. Fontana, D. G. Sorrenti, and W. Burgard. Cmrnet: Camera to lidar-map registration. In *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pages 1283–1289, 2019. 3

[4] D. Cattaneo, M. Vaghi, S. Fontana, A. L. Ballardini, and D. G. Sorrenti. Global visual localization in lidar-maps through shared 2d-3d embedding space. *2020 IEEE International Conference on Robotics and Automation (ICRA)*, May 2020. 3

[5] Zetao Chen, Adam Jacobson, Niko Sunderhauf, Ben Upcroft, Lingqiao Liu, Chunhua Shen, Ian Reid, and Michael Milford. Deep learning features at scale for visual place recognition. *2017 IEEE International Conference on Robotics and Automation (ICRA)*, May 2017. 1

[6] Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu, Zichong Yang, Kuei-Da Liao, Tianren Gao, Erlong Li, Kun Tang, Zhipeng Cao, Tong Zhou, Ao Liu, Xinrui Yan, Shuqi Mei, Jianguo Cao, Ziran Wang, and Chao Zheng. A survey on multimodal large language models for autonomous driving. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) Workshops*, pages 958–979, January 2024. 3

[7] Siyan Dong, Qingnan Fan, He Wang, Ji Shi, Li Yi, Thomas Funkhouser, Baoquan Chen, and Leonidas Guibas. Robust neural routing through space partitions for camera relocalization in dynamic indoor environments, 2020. 8

[8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. 2, 4

[9] Mengdan Feng, Sixing Hu, Marcelo H Ang, and Gim Hee Lee. 2d3d-matchnet: Learning to match keypoints across 2d image and 3d point cloud. *2019 International Conference on Robotics and Automation (ICRA)*, May 2019. 3

[10] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. *2012 IEEE Conference on Computer Vision and Pattern Recognition*, Jun 2012. 2

[11] Albert Gordo, Jon Almazán, Jerome Revaud, and Diane Larlus. End-to-end learning of deep visual representations for image retrieval. *International Journal of Computer Vision*, 124(2):237–254, Jun 2017. 1

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. 2, 4

[13] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization, 2021. 7

[14] Georg Hess, Adam Tonderski, Christoffer Petersson, Kalle Åström, and Lennart Svensson. Lidarclip or: How i learned to talk to point clouds, 2023. 3, 14

[15] Hervé Jégou, Matthijs Douze, Cordelia Schmid, and Patrick Pérez. Aggregating local descriptors into a compact image representation. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 3304–3311, 2010. 2

[16] Laurent Kneip, Davide Scaramuzza, and Roland Siegwart. A novel parametrization of the perspective-three-point problem for a direct computation of absolute camera position and orientation. In *CVPR 2011*, pages 2969–2976, 2011. 1, 3

[17] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*, NIPS'12, page 1097–1105, Red Hook, NY, USA, 2012. Curran Associates Inc. 2

[18] Zuzana Kukelova, Martin Bujnak, and Tomas Pajdla. Real-time solution to the absolute pose problem with unknown radial distortion and focal length. In *2013 IEEE International Conference on Computer Vision*, pages 2816–2823, 2013. 1, 3

[19] Haowen Lai, Peng Yin, and Sebastian Scherer. Adafusion: Visual-lidar fusion with adaptive weights for place recognition. *IEEE Robotics and Automation Letters*, PP:1–8, 10 2022. 2, 3

[20] Christoph H. Lampert, Hannes Nickisch, and Stefan Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 951–958, 2009. 7

[21] Viktor Larsson, Zuzana Kukelova, and Yinqiang Zheng. Making minimal solvers for absolute pose estimation compact and robust. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2335–2343, 2017. 1

[22] Karel Lebeda, Jiri Matas, and Ondrej Chum. Fixing the locally optimized ransac. 09 2012. 1

[23] Yiyi Liao, Jun Xie, and Andreas Geiger. Kitti-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, page 1–1, 2022. 2

[24] Zhe Liu, Shunbo Zhou, Chuanzhe Suo, Peng Yin, Wen Chen, Hesheng Wang, Haoang Li, and Yunhui Liu. Lpd-net: 3d point cloud learning for large-scale place recognition and environment analysis. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2831–2840, 2019. 1

[25] Xudong Lv, Boya Wang, Dong Ye, and Shuo Wang. Lc-cnet: Lidar and camera self-calibration using cost volume network, 2020. 3

[26] Will Maddern, Geoffrey Pascoe, Chris Linegar, and Paul Newman. 1 year, 1000 km: The oxford robotcar dataset. *The International Journal of Robotics Research*, 36(1):3–15, 2017. 3

[27] D. Nister and H. Stewenius. Scalable recognition with a vocabulary tree. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 2161–2168, 2006. 2

[28] Amadeus Oertel, Titus Cieslewski, and Davide Scaramuzza. Augmenting visual place recognition with structural cues. *IEEE Robotics and Automation Letters*, 5(4):5534–5541, Oct 2020. 2

[29] Mohammad Omama, Pranav Inani, Pranjal Paul, Sarat Chandra Yellapragada, Krishna Murthy Jatavallabhula, Sandeep Chinchali, and Madhava Krishna. Alt-pilot: Autonomous navigation with language augmented topometric maps. *arXiv preprint arXiv:2310.02324*, 2023. 3

[30] Sai Shubodh Puligilla, Satyajit Tourani, Tushar Vaidya, Udit Singh Parihar, Ravi Kiran Sarvadevabhatla, and K. Madhava Krishna. Topological mapping for manhattan-like repetitive environments. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6268–6274, 2020. 3

[31] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation, 2017. 3

[32] 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. 1, 4, 5, 6

[33] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018. 7

[34] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019. 7

[35] Sebastian Ratz, Marcin Dymczyk, Roland Siegwart, and Renaud Dube. Oneshot global localization: Instant lidar-visual pose estimation. *2020 IEEE International Conference on Robotics and Automation (ICRA)*, May 2020. 2

[36] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet?, 2019. 7

[37] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. Orb: An efficient alternative to sift or surf. In *2011 International Conference on Computer Vision*, pages 2564–2571, 2011. 2

[38] Torsten Sattler, Will Maddern, Carl Toft, Akihiko Torii, Lars Hammarstrand, Erik Stenborg, Daniel Safari, Masatoshi Okutomi, Marc Pollefeys, Josef Sivic, Fredrik Kahl, and Tomas Pajdla. Benchmarking 6dof outdoor visual localization in changing conditions, 2018. 3

[39] Mustafa Shukor, Guillaume Couairon, and Matthieu Cord. Efficient vision-language pretraining with visual concepts and hierarchical alignment, 2022. 14

[40] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. 4, 5, 6

[41] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume, 2017. 3

[42] Hajime Taira, Masatoshi Okutomi, Torsten Sattler, Mircea Cimpoi, Marc Pollefeys, Josef Sivic, Tomas Pajdla, and Akihiko Torii. Inloc: Indoor visual localization with dense matching and view synthesis, 2018. 3

[43] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. Measuring robustness to natural distribution shifts in image classification, 2020. 8

[44] Mikaela Angelina Uy and Gim Hee Lee. Pointnetvlad: Deep point cloud based retrieval for large-scale place recognition. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Jun 2018. 1, 3

[45] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding, 2019. 4

[46] Johanna Wald, Torsten Sattler, Stuart Golodetz, Tommaso Cavallari, and Federico Tombari. Beyond controlled environments: 3d camera re-localization in changing indoor scenes, 2020. 3

[47] Haohan Wang, Songwei Ge, Eric P. Xing, and Zachary C. Lipton. Learning robust global representations by penalizing local predictive power, 2019. 7

[48] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein, and Justin M. Solomon. Dynamic graph cnn for learning on point clouds, 2019. 3

[49] Peng Yin, Xu Lingyun, Ji Zhang, Howie Choset, and Sebastian Scherer. i3dloc: Image-to-range cross-domain localization robust to inconsistent environmental conditions. 07 2021. 3

[50] Zhengyou Zhang, Rachid Deriche, Olivier Faugeras, and Quang-Tuan Luong. A robust technique for matching two uncalibrated images through the recovery of the unknown epipolar geometry. *Artificial intelligence*, 78(1-2):87–119, 1995. 2

[51] Zhipeng Zhao, Huai Yu, Chenwei Lyv, Wen Yang, and Sebastian Scherer. Attention-enhanced cross-modal localization between 360 images and point clouds, 2022. 2, 3, 6

# Supplementary Material

## A. Qualitative Results & Semantic Breakdown of KITTI360

Fig 7 and Fig 8 demonstrate the qualitative results of 2D to 3D localization and 3D to 2D localization respectively on KITTI-360 dataset using the best model LIP-Loc. The 3D scans are shown in top view. As can be seen in Fig 7, the 3D scans that we are able to predict are very close to ground truth scans, which can help in navigation in an environment where 3D information is not available at test time. Similarly, Fig 8 shows that the RGB images retrieved via 3D-2D localization are similar to ground truth, thus this could further result in downstream application like finding a finer pose through perspective-n-point algorithm.

Let us discuss about semantic breakdown of KITTI-360 dataset now. KITTI-360 is a diverse suburban dataset with 37 label classes, including 24 "instance" classes and 13 "stuff" classes. They define a category and within a category come many classes. For example, the category "flat" contains "classes" like road, sidewalk, parking etc; construction contains building, garage, wall, fence etc. They additionally do a statistical analysis over the distribution of the semantic labels, through which they plot 2D semantic labels over frames and 3D semantic labels over points and bounding boxes. When it is done over frames or points, they find that the highest distribution is of classes vegetation, sky, terrain, car and road while when done over bounding boxes reflects that the highest distribution is for classes car, pedestrian, rider, building and bicycle. There are some predominantly "downtown" scenes, i.e. those with buildings/houses, and many objects like trees, bicycles are common, such as sequences 0000, 0002, 0009. There are also some predominantly "highway" scenes, i.e. i.e. those with open areas, continuous vegetation, roads and cars such as 0003, 0004, 0005. Fig-5 of main paper reported recall@1 values on these sequences. Overall, without any training, our "Zero-shot LIP-Loc" performs well in all sets of diverse conditions having a recall of around 0.5 for most sequences and reaching a maximum of 0.658 for 1 sequence. There is no clear correlation between the accuracy on a sequence and its semantic distribution, i.e. whether it is highway or downtown. For example, if we look at highway scenes such as sequences 0003 and 0005 have recall of 0.658 and 0.594, whereas other highway sequences like 0004 have low recall value like 0.470 whereas downtown scene like 0000 has recall of 0.541. This could mean that our "Zero-shot LIP-Loc" model is not learning spurious correlations, in other words, it is not fitting to certain distribution, rather it is learning in a generalized way. As opposed to our baseline AECMLoc which has tested only on 0000 which has one kind of distribution predominantly, we have tested on 6 sequences each of which differs and we get good recall values for each and do not get abnormally poor values anywhere, which suggests that our approach is robust to distribution shift. With that being said, we have to point out that KITTI-360 does not give clear per-sequence breakdown of semantics, and there is a necessity for a benchmark to do thorough analysis and demonstrate the true zero-shot effectiveness of approaches like ours.

## B. Architecture: Different Encoders & Bigger Models

| Different Encoder | Seq 8 | Seq 9 |
|---|---|---|
| exp_large (resnet) | 0.179 | 0.147 |
| exp_larger (resnet) | 0.295 | 0.309 |
| exp_largest (resnet) | 0.484 | 0.457 |
| trip_larger_vanila (resnet) | 0.215 | 0.232 |
| exp_large (ViT) | 0.278 | 0.260 |
| exp_larger (ViT) | 0.547 | 0.525 |
| exp_largest (ViT) | 0.805 | 0.780 |
| trip_vanila_larger (ViT) | 0.279 | 0.282 |

Table 7. Recall@1 on Different encoders

| Bigger Model | Seq 8 | Seq 9 |
|---|---|---|
| exp_largest_resnet101 | 0.477 | 0.462 |
| exp_large_resnet101 | 0.178 | 0.152 |
| exp_largest_vit_base_patch16_224 | 0.777 | 0.720 |
| exp_large_vit_base_patch16_224 | 0.238 | 0.230 |

Table 8. Recall@1 on Bigger Models

| Combined Models | Seq 8 | Seq 9 |
|---|---|---|
| exp_largest (ViT) | 0.805 | 0.780 |
| exp_combined_vit (thresh 50) | 0.817 | 0.741 |
| exp_combined_vit (thresh 100) | 0.785 | 0.758 |
| exp_combined_vit_fewshot | 0.811 | 0.773 |
| exp_combined_vit_base_patch16_224 | 0.827 | 0.805 |

Table 9. Recall@1 for model trained on both KITTI and KITTI-360 and inference on KITTI-360 (no overlap): test on 0000 of KITTI-360, while train on the rest.

In this section, we report additional experiments on experimenting with different encoders and bigger models. We report in the main paper that *vit_small_patch16_224* is the final model we have chosen. Here we discuss more on why we have chosen that model and what results we have obtained for other models. Please do note that in this supplementary too, wherever not explicitly mentioned, we are referring to ViT's model of *vit_small_patch16_224* model and ResNet's model of *ResNet50*.

Figure 7. Visualization of 2D to 3D localization

Table 7 reports recall values on different encoders. We observed in our experiments that ViT models have a significant accuracy improvement over ResNet. To ensure the comparison is fair, we pick models which have roughly same number of parameters, i.e. *ResNet50* which has 25M parameters (25,557,032) and *vit_small_patch16_224* which has 22M parameters (22,050,664). Even with 3M less parameters, we notice a rise of over 30% accuracy in `exp_largest` case. This pattern can be observed in training over smaller sequences too, such as `exp_larger` and `exp_large`. Even in the triplet vanilla case, we can see a marginal 5% improvement, clearly demonstrating that the edge of ViT over the standard ResNet models.

In Table 8, we report the recall values of bigger models such as *resnet101* and *vit_base_patch16_224*. Although these models have significantly higher parameters, such as 87M for the latter, we do not observe any much change in accuracy. In fact, it dropped marginally. This could be due to the fact that localization datasets are much smaller compared to internet-scale datasets like CLIP and bigger models result in overparametrization, thus dropping accuracy.

Does adding more data improve accuracy for these big-

ger models? In the main paper, we have discussed the standard train-test setting, where we tested on 0000 sequence of KITTI-360 while its training was on rest of sequences of KITTI-360 (0000) or KITTI (Seq 8 and Seq 9), this was "LIP-Loc". Other setting was when we trained on KITTI data and evaluated on KITTI-360, called as "Zero-shot LIP-Loc". `exp_combined` refers to the a third setting, where we train on all sequences of KITTI and KITTI-360 excluding test sequences of KITTI-360 (i.e. 0000) and KITTI (i.e. 8 and 9) on which we test. To get back to our question of whether adding more data will improve accuracy for bigger models, see last row of Table 9 whose accuracy improved over `exp_largest_vit_base_patch16_224` of 8 by 5%. This further reaffirms that if we scale the model, we need to scale the data in order to improve the accuracy. Do note that this fact is not as established in visual localization as much as in computer vision or language models, it is still an open question as to how much role big data will play for localization, hence these analyses play crucial role.

We also try a few-shot experiment here wherein we give just 1% of data of KITTI-360 when compared to the `exp_combined_vit` experiment. To

Figure 8. Visualization of 3D to 2D localization



Figure 9. Semantic Breakdown of KITTI-360 evaluation sequences: The top row represents downtown with cars, buildings whereas bottom row represents highway with more greenery, wide roads. Our "Zero-shot LIP-Loc" model performs well in these diverse conditions without even being trained on this data.

be clear, `exp_combined_vit` uses all sequences of KITTI and KITTI-360 (excluding test sequences), whereas `exp_combined_vit_fewshot` uses all sequences of KITTI but just 1% of KITTI-360. This is a very captivating result: We receive almost same or marginally improve upon the accuracy as the other experiment despite using significantly very less dataset.

It is also worth noting that the combined experiments don't improve significantly from `exp_largest`, unless we use a bigger model like `exp_largest_vit_base_patch16_224`, which is also 2% improvement. Future experiments have to be done to
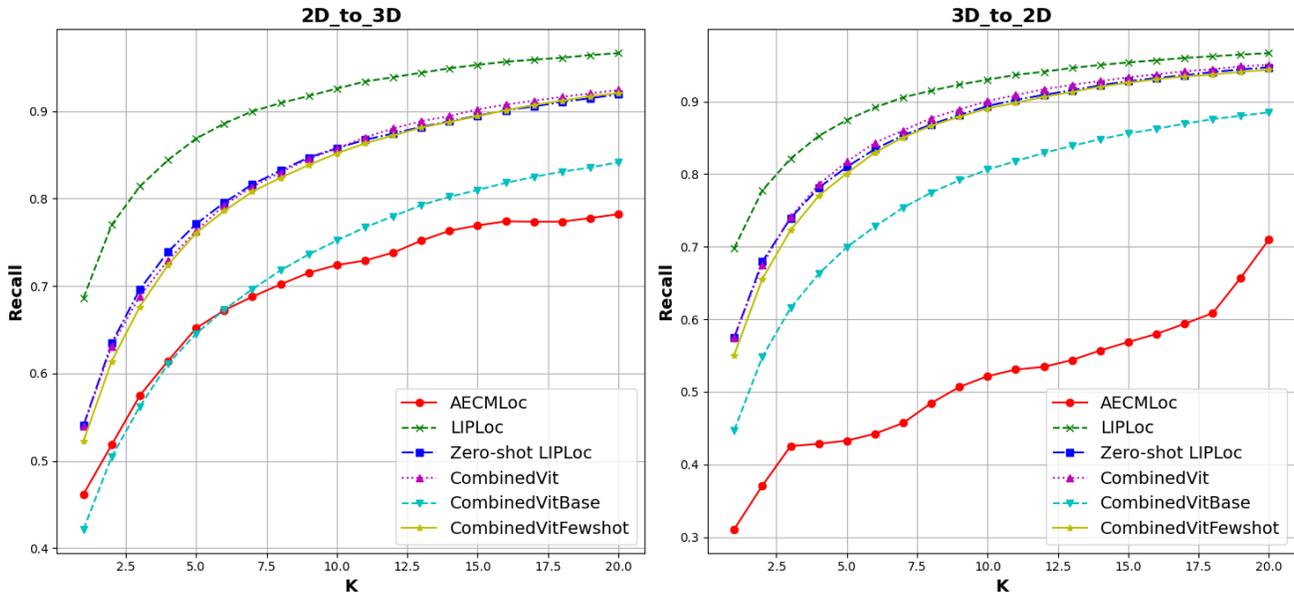
Figure 10. Recall@K curves on KITTI-360 sequence 0000: Combined Models comparison with our best models and baseline AECMLoc

establish even clearer understanding.

So far, we have discussed evaluation on KITTI dataset. Now let us discuss about evaluation on KITTI-360 dataset by looking at Fig 10. Previously in the main paper, we reported LIPLoc, Zero-shot LIPLoc and AECMLoc (baseline). Here we additionally add the plots of combined experiments, which as described above, merges the training sequences of KITTI and KITTI-360 and trains a single model using full data. Do note that all of our models beat the SOTA AECMLoc. But amongst our models themselves, we rather see ambigious or counterintuitive results.

Firstly to clarify, when we use the term "LIP-Loc", we are referring to standard train/test paradigm, for example when reporting LIP-Loc on KITTI-360, we mean we trained on certain train split of KITTI-360 and evaluating on its test splits; similarly when reporting on KITTI, we mean we trained on train split of KITTI and evaluating on its corresponding test splits. "Zero-shot LIP-Loc" on other hand are trained on full KITTI data but has not seen any KITTI-360 data on which we evaluate. Whereas combined models are trained on train split of KITTI and train split of KITTI-360. Therefore, please keep these nuances in mind when interpreting the result. With that being said, since we are evaluating on test split of KITTI-360, we would expect combined models to significantly outperform Zero-shot LIP-Loc. However, that's not the case here: Amongst the combined models, all the standard ViT model `CombinedVit` and bigger model `CombinedVitBase` and the few shot model `CombinedVitFewshot` give similar recall compared to Zero-shot LIP-Loc and subpar performance com-

pared to LIP-Loc. This further proves that Zero-shot LIP-Loc has generalized very well.

As future work, it will be interesting to see an analysis between zero-shot and few-shot LIP-Loc. This raises many open questions: In computer vision problems which CLIP deals with, few-shot is clearly defined because it is talking about classification categories. However it is not well defined in visual localization context, which further asserts the necessity of establishment of a well thought benchmark. We encourage the reader to address these open questions and ask the question, "Can big data solve the localization problem?"

CLIP admits that it is not good at task generalization for tasks such as finding close objects in an image or counting the number of objects in an image. Extending our work along the lines of the recent work LiDARCLIP [14] which connects CLIP's embedding space to LiDAR point cloud domain could result in an approach which uses text features to query the right set of points in the LiDAR scan, explicitly identify distance and location of the objects and applying clustering in 3D space to count number of objects (for example) and correlate them with image features to identify the class and appearance of an object. This is especially helpful in extreme low visibility conditions where RGB camera will not work well and LiDAR can help identify objects close to the ego vehicle.

## C. Architecture: Hierarchical Design

In models as a follow up to CLIP, many models such as ViCHA [39] propose architectural improvement such as

hierarchical alignment. What this essentially proposes is that aligning the two encoders at various levels by adding multiple losses at various layers of text and image encoder. They claim that this helps in convergence faster and results in superior performance. In our experiments, we have hierarchically aligned image and lidar encoders at various layers and report it in first half of the table 10. We have tried two experiments: One that aligns only at final layers, the other that aligns throughout the encoder, as ViCHA argues that aligning at the beginning could result in confusing the model. However, in our experiments we did not observe any noticeable improvement, although ViCHA's observation of alignment at final layers could be verified in the case of visual localization as well.

The second half of the table pertains to the following. In standard CLIP setting, there is no relation between any consecutive images in a batch, as they are just (image, text) pairs. However, in our localization setting, the images are sequential. Therefore, we attempted the question: Can we achieve higher accuracy by grouping together adjacent images and having additional encoder for groups of images which results in secondary loss? The last 3 rows of Table 10 correspond to these experiments. Do note that these experiments are with *ResNet* architecture. Our results actually deteriorated during our experiments. There is a simple rationale for this: The training of deep models works so well because of randomization of samples in a batch, especially in the case of our batch construction technique. When we contruct groups within the batch and ensure the images within the group are consecutive but groups themselves are random, we are asking for a tradeoff: will the additional hierarchical loss improve accuracy more than reducing randomization will decrease it? We have found in our experiments that the answer is no, even for smalller group sizes such as 4.

This section concludes that sticking to non-complicated architectures works the best since the power of CLIP model subsumes any minor architectural improvement.

| Advanced Architectures | Seq 8 | Seq 9 |
|---|---|---|
| exp_large (ViT) | 0.278 | 0.260 |
| hier_align_large_vit (final layers) | 0.275 | 0.258 |
| hier_align_large_vit (all layers) | 0.218 | 0.197 |
| exp_large (resnet) | 0.179 | 0.147 |
| exp_larger (resnet) | 0.295 | 0.309 |
| exp_largest (resnet) | 0.484 | 0.457 |
| hier_group_shuffle_large_resnet | 0.170 | 0.1495 |
| hier_group_shuffle_larger_resnet | 0.239 | 0.212 |
| hier_group_shuffle_largest_resnet | 0.378 | 0.346 |

Table 10. Architecture: Hierarchical Design