

330 A Details on experimental setup

331 A.1 Setup for LiDAR and camera input

332 **Setup for LiDAR input:** To learn robust features for the calibration of different transforms, we
333 train our network to recover the true calibration from a variety of initial guesses $\hat{\mathbf{T}} = \mathbf{T}_p \mathbf{T}$. Where
334 we perturb the ground truth \mathbf{T} using uniformly sampled \mathbf{T}_p , modeled by a translation and angle-
335 axis vector. We sample these vectors uniformly on the 3-sphere using [19] and scale them with a
336 uniformly distributed scalar. For training, the translation and rotation magnitudes are in the range
337 $[0, 1.5]\text{m}$ and $[0, 15]^\circ$, respectively.

338 To account for different intensity profiles registered by different LiDAR models, and cater to the
339 fact that some LiDAR drivers don't provide intensity readings, we augment the intensity channel of
340 our LiDAR data to minimize our dependence on intensity information. We apply uniform random
341 scalar perturbations in the range $[0, 1.0]$ to the intensity channel of the LiDAR points, meaning that
342 in some samples the intensity information is close to dropped out.

343 Lastly, to process the LiDAR data using our sparse 3-D CNN, we voxelize the points using isotropic
344 voxels of 2 cm per side. We found this resolution to be reasonable since it is in the range of the
345 measurement error reported by most LiDAR manufacturers.

346 **Setup for camera input:** We've found experimentally that the camera feature extractor fails to
347 learn generalizable geometric features unless spatial augmentations are applied. In our training, we
348 performed random crop augmentations of $[512, 256]$ pixels in width and height to the input image,
349 and updated the camera intrinsic parameters accordingly.

350 A.2 Model setup

351 Both image and point cloud feature extractors in our model follow the U-Net [14] architecture with
352 5 layers of coarse-to-fine features, each layer being a factor of 2 finer than the previous layer. We
353 use features from 3 layers for our alignment, the 1/16-scale, the 1/4-scale, and the 1-scale. The
354 feature dimensionalities at these layers are 128, 128, and 32, respectively. To help adaptation, at
355 each pyramid level, the features from both domains are passed through a shared 2-layer MLP with
356 the same input and output dimensions at each layer and Leaky ReLU activation.

357 With recent advancements in 3-D convolutions [20], we've chosen the spconv [21] implementation
358 of sparse 3-D CNNs for our LiDAR feature extractor, as it efficiently handles the large point clouds,
359 and effectively treats the sparsity pattern of the data.

360 Learning features in the image domain is relatively straightforward, we've found that using the same
361 U-Net [14] architecture (similar to Pixloc [11]) with the 2-D convolutional VGG [22] backbone was
362 sufficient for learning image features.

363 We initialize the visual extractor using weights from a pre-trained PixLoc [11] model trained on
364 the CMU Seasons dataset [23], and the LiDAR extractor using only the sparse 3D CNN weights
365 of a SphereFormer [24] model pre-trained on Semantic KITTI [25]. For the pose optimization, we
366 trained with $M = 5$ iterations at every pyramid level.