

Differentiable SLAM-net: Learning Particle SLAM for Visual Navigation

Peter Karkus Shaojun Cai David Hsu
National University of Singapore

karkus@comp.nus.edu.sg

Abstract

Simultaneous localization and mapping (SLAM) remains challenging for a number of downstream applications, such as visual robot navigation, because of rapid turns, featureless walls, and poor camera quality. We introduce the Differentiable SLAM Network (SLAM-net) along with a navigation architecture to enable planar robot navigation in previously unseen indoor environments. SLAM-net encodes a particle filter based SLAM algorithm in a differentiable computation graph, and learns task-oriented neural network components by backpropagating through the SLAM algorithm. Because it can optimize all model components jointly for the end-objective, SLAM-net learns to be robust in challenging conditions. We run experiments in the Habitat platform with different real-world RGB and RGB-D datasets. SLAM-net significantly outperforms the widely adapted ORB-SLAM in noisy conditions. Our navigation architecture with SLAM-net improves the state-of-the-art for the Habitat Challenge 2020 PointNav task by a large margin (37% to 64% success). Project website: <http://sites.google.com/view/slamnet>

1. Introduction

Modern visual SLAM methods achieve remarkable performance when evaluated on suitable high-quality data [19]. However, in the context of downstream tasks, such as indoor robot navigation, a number of difficulties remain [8; 43; 12]. An imperfect navigation agent captures substantially different images from a human, as it may frequently face featureless walls and produce rapid turns (see Fig. 1). Further, despite advances in sensor technology, modern robots still often use cameras with noisy images, low frame rate, and a narrow field-of-view [32]. These factors make feature extraction and association difficult. Relocalization and loop closure can be challenging due to environmental changes and repetitive feature. Finally, integrating SLAM into a navigation pipeline is not trivial, because the map representation must be suitable for downstream planning, it may need to capture task-dependent information, and planning must be able to handle map imperfections.



Figure 1: Visual robot navigation is challenging for SLAM, *e.g.*, because the robot frequently faces featureless walls; it rotates quickly; the onboard camera produces noisy images; the frame rate is low; *etc.* The images were taken by our navigation agent in the Habitat environment.

This paper introduces the Differentiable SLAM Network (SLAM-net) together with a navigation architecture for downstream indoor navigation. The key idea of SLAM-net is to encode a SLAM algorithm in a differentiable computation graph, and learn neural network model components for the SLAM algorithm end-to-end, by backpropagating gradients through the algorithm. Concretely SLAM-net encodes the particle filter based FastSLAM algorithm [44] and learns mapping, transition and observation models. SLAM-net fills a gap in the literature on differentiable robot algorithms [55; 23; 34; 35].

The benefit of SLAM-net compared to unstructured learning approaches is that its encoded particle filter provides a strong prior for learning. The benefit over classic SLAM is that all components are learned, and they are directly optimized for the end-objective. Concretely, SLAM-net learns RGB and RGB-D observation models for the encoded FastSLAM algorithm, which previously relied on hand-crafted models and lidar sensors. Further, because of the task-oriented learning, feature extractors can learn to be more robust against domain specific challenges, *e.g.*, faced with downstream navigation; while on the flip side they may be less reusable across tasks.

We validate SLAM-net for localization with RGB and RGB-D input, as well as downstream robot navigation in previously unseen indoor environments. We use the Habitat

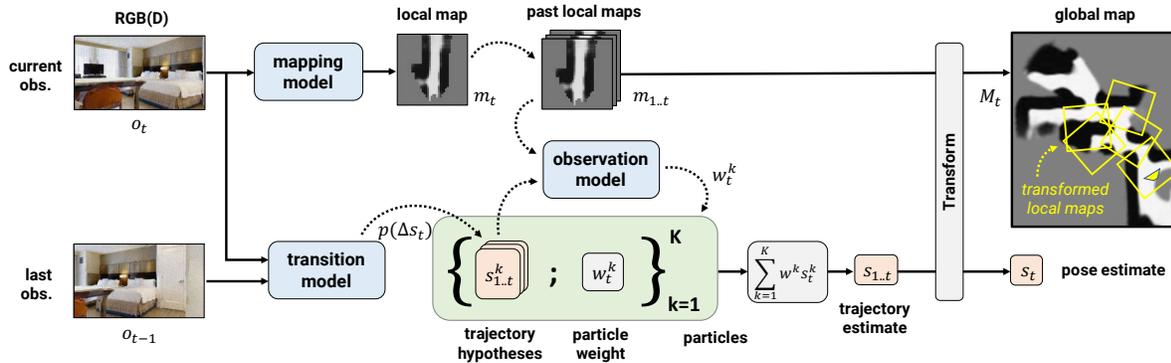


Figure 2: Differentiable SLAM-net. The global map is maintained by a collection of learned local grid maps. The trajectory is tracked by a particle filter. Particles represent trajectories and they are updated with learned neural network components: the mapping, transition, and observation models.

simulation platform [52] with three real-world indoor scene datasets. We additionally experiment with the KITTI visual odometry data [20]. SLAM-net achieves good performance under challenging conditions where the widely used ORB-SLAM [46] completely fails; and learned models transfer over datasets. For downstream navigation we propose an architecture similar to Neural SLAM [10], but with using our differentiable SLAM-net module. Our approach significantly improves the state-of-the-art for the CVPR Habitat 2020 PointNav challenge [25].

2. Related work

Learning based SLAM. Learning based approaches to SLAM have a large and growing literature. For example, CodeSLAM[7] and SceneCode[63] learn a compact representation of the scene; CNN-SLAM[57] learns a CNN-based depth predictor as the front-end of a monocular SLAM system; BA-net [56] learns the feature metric representation and a damping factor for bundle adjustment. While these works use learning, they typically only learn specific modules in the SLAM system. Other approaches do end-to-end learning but they are limited to visual odometry, *i.e.*, they estimate relative motion between consecutive frames without a global map representation [64; 41]. Our method maintains a full SLAM algorithm and learns all of its components end-to-end.

Classic SLAM. Classic SLAM algorithms can be divided into filtering and optimization based approaches [53]. Filtering-based approaches maintain a probability distribution over the robot trajectory and sequentially update the distribution with sensor observations [14; 4; 15; 44; 45]. Optimization-based approaches apply bundle adjustment on a set of keyframes and local maps; and they are popular for both visual [53; 47; 46; 38] and lidar-based SLAM [28]. Our approach builds on a filtering-based algorithm, Fast-

SLAM [44; 45]. The original algorithm (apart from a few adaptations [5; 40; 27]) works with a lidar sensor and hand-designed model components. Robot odometry information is typically used for its transition model, and either landmarks [44] or occupancy grid maps [21] are used for its observation model. In contrast, we learn neural network models for visual input by backpropagation through a differentiable variant of the algorithm. We choose this algorithm over an optimization based method because of the availability of differentiable particle filters [31; 34], and the suitability of the algorithm for downstream robot navigation.

Differentiable algorithms. Differentiable algorithms are emerging for a wide range of learning domains, including state estimation [24; 31; 34; 42], visual mapping [23; 36], planning [55; 33; 18; 49; 22; 62] and control tasks [2; 50; 17; 6]. Differentiable algorithm modules have been also composed together for visual robot navigation [23; 35]. This work introduces a differentiable SLAM approach that fills a gap in this literature. While Jatavallabhula et al. [30] have investigated differentiable SLAM pipelines, they focus solely on the effect of differentiable approximations and do not perform learning of any kind.

Visual navigation. A number of learning based approaches has been proposed for visual navigation recently [23; 3; 43; 35; 60; 10; 11; 51]. Modular approaches include CMP [23], DAN [35] and Neural SLAM [10]. However, CMP assumes a known robot location, circumventing the issue of localization. DAN assumes a known map that is given to the agent. Neural SLAM [10; 51] addresses the joint SLAM problem, but it relies solely on relative visual odometry without local bundle adjustment or loop closure, and thus it inherently accumulates errors over time. We propose a similar navigation architecture to Neural SLAM [10], but utilizing our Differentiable SLAM-net module in place of learned visual odometry.

3. Differentiable SLAM-net

3.1. Overview

The Differentiable SLAM-net architecture is shown in Fig. 2. Inputs are RGB(D) observations o_t , outputs are pose estimate s_t and global map M_t . SLAM-net assumes the robot motion is (mostly) planar. Poses are 2D coordinates with 1D orientation; the global map is a 2D occupancy grid.

Internally SLAM-net represents the global map as a collection of local maps, each associated with a local-to-global transformation. Local maps are $N \times N \times N_{\text{ch}}$ grids that, depending on the configuration, may encode occupancy and/or learned latent features. We add a local map for each observation, but without knowing the robot pose we do not know the correct local-to-global map transformation. Instead, the algorithm maintains a distribution over the unknown robot trajectory using particle filtering [16]. Our algorithm is based on FastSLAM [44; 45], and our differentiable implementation is built on PF-nets [34].

The algorithm works as follows. The particle filter maintains K weighted particles, where each particle represents a trajectory $s_{0:t}^k$. At $t = 0$ all particle trajectories are set to the origin; particle weights are constant, and the local map collection is empty. In each time step a *mapping model* predicts a local map m_t from the input observation o_t , and m_t is added to the collection. Particle trajectories are extended with samples from a probabilistic *transition model* that estimates the relative motion given o_t and o_{t-1} . Particle weights are then updated using an *observation model* which measures the compatibility of m_t and the past local maps $m_{1:t-1}$ assuming the particle trajectory $s_{0:t}^k$ was correct. The pose output is obtained by taking the weighted sum of particle trajectories. Optionally, a global occupancy grid map is obtained with simple 2D image transformations that combine local maps along the mean trajectory.

The key feature of SLAM-net is that it is end-to-end differentiable. That is, the mapping, transition and observation models are neural networks, and they can be trained together for the end-objective of localization accuracy (and/or global map quality). To make the algorithm differentiable we use the reparameterization trick [37] to differentiate through samples from the transition model; and we use spatial transformers [29] for differentiable map transformations. The rest of the operations of the algorithm, as presented, are already differentiable. While not used in our experiments, differentiable particle resampling could be incorporated from prior work [34; 65]. Further, due to limited GPU memory, to make use of the differentiable algorithm for learning our design choices on the local map representation and the formulation of the observation model are important.

Next we introduce each component of SLAM-net. Network architecture details are in the Appendix.

3.2. Transition model

The transition model is a CNN that takes in the concatenated current and last observations, o_t and o_{t-1} , and outputs parameters of Gaussian mixture models with separate learned mean and variance for the relative 2D pose and 1D orientation. The transition model is pre-trained to maximize the log-likelihood of true relative poses along the training trajectories. It is then finetuned together with the rest of the SLAM-net components optimizing for the end-objective.

3.3. Mapping model

The mapping model is a CNN with a pre-input perspective transformation layer. The input is observation o_t , the output is local map m_t . Local maps serve two purposes: to be fed to the observation model and aid pose estimation by closing loops; and to construct a global map for navigation.

Our local maps are $N \times N \times N_{\text{ch}}$ grids that capture information about the area in front of the robot. In one configuration local maps encode occupancy, *i.e.*, the probability of the area being occupied for the purpose of navigation. This model is trained with a (per cell) classification loss using ground-truth occupancy maps. In another configuration local maps encode learned latent features that have no associated meaning. This model is trained for the end-objective by backpropagating through the observation model. In both cases we found it useful to add an extra channel that encodes the visibility of the captured area. For depth input this is computed by a projection; for RGB input it is predicted by the network, using projected depth for supervision.

3.4. Observation model

The observation model is the most critical component of the architecture. It updates particle weights based on how “compatible” the current local map would be with past local maps if the particle trajectory was correct. Intuitively we need to measure whether local maps capture the same area in a consistent manner. Formally we aim to estimate a compatibility value proportional to the log-probability $\log p(m_t | m_{1:t-1}, s_{1:t}^k)$.

We propose a discriminative observation model that compares pairs of local maps with a learned CNN. The CNN takes in a current local map m_t and a past local map m_τ concatenated, and outputs a compatibility value. Importantly, the past local map is transformed to the viewpoint of the current local map according to the relative pose in the particle trajectory (s_t^k, s_τ^k) . We use spatial transformers [29] for this transformation. The overall compatibility is the sum of pairwise compatibility values along the particle trajectory. Compatibility values are estimated for all particles. Particle weights are then updated by multiplying with the exponentiated compatibility values, and they are normalized across particles. CNN weights are shared.

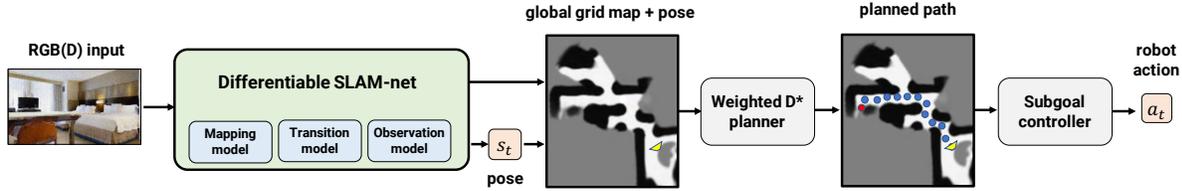


Figure 3: Visual navigation pipeline with the Differentiable SLAM-net, a path planner, and a motion controller.

For computational reasons, instead of comparing all local map pairs, we only compare the *most relevant* pairs. During training we pick the last 4–8 steps; during inference we dynamically choose 8 steps based on the largest overlapping view area (estimated using simple geometry).

3.5. Training and inference

Our training data consists of image-pose pair trajectories (depth or RGB images, and 2D poses with 1D orientation); and optionally ground-truth global occupancy maps for pre-training the mapping model. The end-to-end training objective is the sum of Huber losses for the 2D pose error and 1D orientation error.

We train in multiple stages. We first pre-train the transition model. We separately pre-train the observation model together with the mapping model for the end-objective, but in a low noise setting. That is, in place of the transition model we use ground truth relative motion with small additive Gaussian noise. Finally we combine all models and finetune them together for the end-objective. During finetuning we freeze the convolution layers and mixture head of the transition model. When the mapping model is configured to predict occupancy it is pre-trained separately and it is frozen during finetuning.

An important challenge with training SLAM-net is the computational and space complexity of backpropagation through a large computation graph. To overcome this issue, during training we use only short trajectories (4-8 steps), and $K = 32$ particles without resampling. During inference we use the full length trajectories, and by default $K = 128$ particles that are resampled in every step.

3.6. Implementation details

We implement SLAM-net in Tensorflow [1] based on the open-source code of PF-net [34]. We adopt the training strategy where the learning rate is decayed if the validation error does not improve for 4 epochs. We perform 4 such decay steps, after which training terminates, and the model with the lowest validation error is stored. The batch size is 16 for end-to-end training, and 64 for pre-training the mapping and transition models. We use Nvidia GeForce GTX 1080 GPUs for both training and inference.

For RGB input we configure local maps with 16 latent

feature channels that are not trained to predict occupancy. For RGB-D input local maps are configured with both latent channels and one channel that predicts occupancy. Further, with RGB-D data we only use depth as input to SLAM-net.

4. Visual Navigation with SLAM-net

We propose a visual navigation pipeline (Fig. 3) that combines SLAM-net with modules for path planning and motion control. In the pipeline SLAM-net periodically predicts the robot pose and a global occupancy grid map. The map and pose are fed to a 2D path planner that plans a path to the goal. The path is then tracked by a local controller the outputs robot actions.

Task specification. We follow the task specification of the Habitat 2020 PointNav Challenge [32]. A robot navigates to goals in previously unseen apartments using noisy RGB(D) input. The goal is defined by coordinates relative to the initial pose, but the robot location is unknown thereafter, and discrete robot actions generate noisy motion. Navigation is successful if the robot takes a dedicated *stop* action within 0.36 meters to the goal. Note that this success criteria places high importance on pose estimation accuracy.

Path planner. The challenge of path planning in the context of visual navigation is the imperfect partial knowledge of the map and the robot pose. To address this we adopt a weighted variant of the D* algorithm [39] with costs that penalize moving near obstacles. In each step when the map and pose are updated by the SLAM-net, the path is replanned. For planning we convert the occupancy grid map to an 8-connected grid where cells are assigned a cost. We threshold the map ($p > 0.5$ is an obstacle, $p \leq 0.5$ is free space) and define cell costs based on the distance to the nearest obstacle.

Additionally, we use a collision recovery mechanism. Upon detecting a collision an obstacle is registered to the map at the current estimated pose. The robot is then commanded to turn around (6 turn actions) and takes a step back (1 forward action). Collisions are not directly observed. We trigger collision recovery if the estimated pose does not change more than $3cm$ following a forward action. A similar mechanism was proposed in Chaplot et al. [10].

We also hard-coded an initial policy that makes the robot turn around in the beginning of an episode. The initial policy

terminates once the estimated rotation exceeds 370° .

Local controller. The planned path is tracked by a simple controller that chooses to turn or move forward aiming for the furthest straight-line traversable point along the path. The output of the controller are discrete actions.

5. Experiments

Our experiments focus on the following questions. 1) Can SLAM-net learn localization in previously unseen indoor scenes? How does it compare to existing methods? 2) Does SLAM-net enable downstream robot navigation? 3) What does SLAM-net learn and how do model components and hyper parameters affect performance? 4) Do learned models transfer to new datasets? 5) What are the limitations of SLAM-net if applied, *e.g.*, to autonomous driving data?

5.1. Datasets

Habitat. Our main experimental platform is the Habitat simulator [52] configured with different real-world datasets: Gibson [61], Replica [54], and Matterport [9]. The datasets contain a large number of 3D scans of real world indoor scenes, typically apartments. Habitat embeds the scenes in an interactive physics simulator for robot navigation. The simulator renders photorealistic but noisy first-person RGB(D) observations and simulates realistic robot motion dynamics. The camera has a horizontal FOV of 70° and a resolution of 640×360 . For SLAM-net we downscale images to 160×90 . Depth values are in the range of 0.1 to 10 meters. Unless stated otherwise, we use the Habitat Challenge 2020 configuration [32]: Gaussian noise for RGB images; the Redwood noise model for depth images [13]; the Locobot motion model for actuation noise [48]. To train and evaluate SLAM methods we generate a fixed set of trajectories. For navigation we let our method interact with the simulator.

Gibson data. Following Savva et al. [52] we use 72 scenes from the Gibson dataset for training and further split the original validation set to 7 scenes for validation and 7 scenes for testing. We use 36k of the provided navigation episodes for training (500 per scene). Given a start and goal pose we generate a trajectory with a navigation policy that switches between a shortest-path expert (30 steps) and random actions (40 steps). For evaluation we generate 105 trajectories (15 per test scene) using three different navigation policies: the shortest-path expert (**traj_expert**); the shortest-path expert mixed with random actions (**traj_exp_rand**); and our final navigation pipeline (**traj_nav**).

Replica and Matterport data. We use the Matterport and Replica datasets for transfer experiments without additional training. We generate trajectories for evaluation similarly to the Gibson data, using the shortest-path expert policy. We use the original validation split with a total of 170 and 210 trajectories for Replica and Matterport respectively.

Dataset	length[m]	#frames	#turns
Gibson (traj_expert)	7.4±3.8	51.1±24.7	22.6±11.6
Gibson (traj_exp_rand)	14.5±7.2	152.3±67.9	75.0±34.6
Gibson (traj_nav)	11.9±8.8	117.0±113.3	74.1±84.0
Matterport (traj_expert)	13.0±7.0	82.0±37.9	32.9±14.1
Replica (traj_expert)	8.0±2.8	53.9±17.7	23.5±9.17
KITTI-09	1680.3	1551	
KITTI-10	910.48	1161	

Table 1: Data statistics.

KITTI data. We conduct additional experiments with the KITTI real world driving data [20]. We use the full KITTI raw dataset for training, validation, and testing. Following the KITTI Odometry Split [20] the validation trajectories are 06 and 07, and the testing trajectories are 09 and 10. Since the original depth information for KITTI dataset is from sparse lidar, we use the completed depth data from the [58] as ground-truth depth.

Statistics. Tab. 1 provides statistics for each set of evaluation trajectories. We provide the mean and standard deviation of the trajectory length, number of frames, and number of turn actions (where applicable).

5.2. Baselines

Learned visual odometry. We use the transition model of SLAM-net as a learned visual odometry model. The model parameterizes a Gaussian mixture that predicts the relative motion between consecutive frames. When the model is used for visual odometry we simply accumulate the mean relative motion predictions.

ORB-SLAM. We adopt the popular ORB-SLAM [47; 46] as a classic baseline. The algorithm takes in RGB or RGB-D images, constructs a keypoint-based sparse map, and estimates a 6-DOF camera pose at every time-step. The algorithm relies on tracking features between consecutive frames. If there are not enough tracked key-points, the system is lost. When this happens we initialize a new map and concatenate it with the previous map based on relative motion estimated from robot actions. With RGB-D input re-initialization takes one time step, with RGB-only input it takes several steps.

We carefully tuned the hyperparameters of ORB-SLAM based on the implementation and configuration of Mishkin et al. [43], who tuned the algorithm for Habitat simulation data, although without sensor and actuation noise. For the main localization results in Tab. 2 we use the default velocity-based motion model as in [47]. In Tab. 4 we replace the motion model with relative motion estimated from actions, which gave better results.

Blind baseline. This baseline ignores observation inputs, and instead it accumulates the nominal robot motion based on the ground-truth (but noisy) motion model. This serves as a calibration of the performance of other methods.

Sensor		RGBD		RGBD		RGBD		RGB	
Trajectory generator		traj_expert		traj_exp_rand		traj_nav		traj_expert	
Metric	runtime↓	SR↑	RMSE↓	SR↑	RMSE↓	SR↑	RMSE↓	SR↑	RMSE↓
SLAM-net (ours)	0.06s	83.8%	0.16m	62.9%	0.28m	77.1%	0.19m	54.3%	0.26m
Learned visual odometry	0.02s	60.0%	0.26m	24.8%	0.63m	30.5%	0.47m	28.6%	0.40m
ORB-SLAM [46]	0.08s	3.8%	1.39m	0.0%	3.59m	0.0%	3.54m	X	X
- no obs. noise		7.5%	1.14m					18.0%	0.36m
- no obs. and act. noise		18.0%	1.01m					20.4%	0.31m
- high frame rate		30.4%	0.58m					X	X
- ideal condition		86.0%	0.18m					43.5%	0.42m
Blind baseline	0.01s	16.2%	0.80m	1.0%	4.13m	3.8%	1.50m	16.2%	0.80m

Table 2: Main SLAM results.

Sensor		RGBD	RGB
Metric		SR↑	SR↑
(1) All models (default)		77.1%	55.2%
(2) No learned transition model		43.8%	19.1%
(3) No observation model (VO)		30.5%	26.7%
(4) No joint training		66.7%	8.6%
(5) Occupancy map only		75.2%	23.8%
(6) Latent map only		70.5%	55.2%
(7) Occupancy + latent map		77.1%	44.8%
(8) Fixed comparisons (8)		44.8%	29.5%
(9) Dynamic comparisons (4)		73.3%	41.9%
(10) Dynamic comparisons (8)		77.1%	55.2%
(11) Dynamic comparisons (16)		77.1%	40.0%
(12) K=1 (VO)		30.5%	26.7%
(13) K=8		60.0%	35.2%
(14) K=32 (training)		72.4%	39.1%
(15) K=64		75.2%	46.7%
(16) K=128 (evaluation default)		77.1%	55.2%
(17) K=256		79.1%	44.8%
(18) K=512		82.9%	48.6%

Table 3: Ablation results.

6. Localization results

6.1. Main results for localization

Main localization results are summarized in Tab. 2. Visualizations and videos are in the supplementary material. Tab. 2 reports success rate (SR) that measures the percentage of episodes where the final pose error is below 0.36 meters (to enable successful downstream navigation); and root-mean-square-error (RMSE) which measures the absolute trajectory error as defined in [26]. Estimated trajectories are only aligned with the ground-truth at the beginning of each episode. We also report runtimes, measuring the average processing time per frame (including loading the data, RGBD sensor, traj_expert data).

SLAM-net learns successfully. We first observe that

SLAM-net successfully learned to localize in many episodes despite the challenging data. Comparing columns we see that an imperfect navigation policy can significantly increase the difficulty of localization. Comparing SLAM-net across sensor modalities we find that SLAM-net performs reasonably well with RGB-only input, and the depth sensor helps substantially (54.3% vs. 83.8% success for traj_expert data).

SLAM-net outperforms its alternatives. We find that SLAM-net outperforms both learned visual odometry and the classic ORB-SLAM by a large margin across all datasets and sensor modalities; and its runtime (on GPU) is slightly better than ORB-SLAM. ORB-SLAM performs particularly poorly. ORB-SLAM relies on temporal coherence between frames to track features, which is challenging here due to the combined effects of observation noise, sparse visual features, rapid turns (approx. $90^\circ/s$), low frame rate (approx. 3 fps), and narrow field of view (HFOV= 70°). We find that ORB-SLAM often fails to track features even with RGB-D input. With RGB-only input it fails in nearly all steps, hence we could not report a meaningful result. In contrast, SLAM-net does not rely explicitly on feature tracking, and it learns task-oriented models that can, *e.g.*, learn more robust feature extractors for the domain.

Why does ORB-SLAM fail? To better understand what makes our data challenging for ORB-SLAM, we evaluate it in a number of simplified settings for the traj_expert data. We remove different types of challenges in rows of Tab. 2: no observation noise, no action noise, high frame rate, and ideal condition. The high frame rate setting reduces the action step size in Habitat to achieve an equivalent 3 to ~ 30 fps increase in frame rate. The ideal condition setting removes all the above challenges together.

Our results show that ORB-SLAM only works well in ideal conditions, where its performance is comparable to SLAM-net in the hardest conditions. If we remove only one type of challenge the performance of ORB-SLAM remains significantly worse. The RGB-D results indicate that low frame rate has the largest impact. Results for RGB show a similar trend; however, the presence of observation noise

Dataset Sensor Metric	Replica RGBD		Matterport RGBD		Replica RGB		Matterport RGB	
	SR \uparrow	RMSE \downarrow	SR \uparrow	RMSE \downarrow	SR \uparrow	RMSE \downarrow	SR \uparrow	RMSE \downarrow
SLAM-net (ours)	78.8%	0.17m	49.5%	0.39m	45.3%	0.31m	23.3%	0.54m
Learned visual odometry	51.2%	0.31m	22.4%	0.75m	17.7%	0.67m	15.2%	0.93m
ORB-SLAM [46]	5.2%	1.46m	1.9%	2.90m	X	X	X	X
Blind baseline	7.7%	0.92m	5.7%	2.3m	7.7%	0.92m	5.7%	2.3m

Table 4: Transfer results.

Trajectory Metric	Kitti-09 RMSE \downarrow	Kitti-10 RMSE \downarrow
SLAM-net (ours)	83.5m	15.8m
SLAM-net (best of 5)	56.9m	12.8m
Learned visual odometry	71.1m	73.2m
ORB-SLAM-RGB [47]	7.62m	8.68m

Table 5: KITTI results.

makes feature tracking fail completely (hence we could not report meaningful results for this setting).

6.2. Ablation study

To better understand the workings of SLAM-net we perform a detailed ablation study. Results are summarized in Tab. 3. The table reports success rates for the Gibson traj_nav data, using SLAM-net in different conditions.

Which components are important? Line (2) of Tab. 3 replaces the learned transition model with a probabilistic model that matches the motion dynamics prior (without considering an observation). Line (3) removes the observation model and uses a single particle (equivalent to learned visual odometry). Line (4) evaluates the pre-trained transition and observation models without joint finetuning. We find that both model components are important. Finetuning is useful; and its benefits is much more significant for RGB input. Our RGBD model uses predicted occupancy in its local maps, which probably generalizes better to the noisier evaluation setting than the RGB model that only uses learned latent features.

Occupancy maps are useful for localization. Lines (5–7) use different channels in learned local maps, pre-trained occupancy predictions, learned latent features, or both. The RGBD model is comparable in all settings. Adding latent maps on top of occupancy only improves 1.9%, which indicates that 2D occupancy is sufficient for localization. The latent map configuration is 4.7% behind the occupancy maps, showing that the model was able to learn useful map features end-to-end, without any direct supervision.

We can learn better map features if occupancy prediction is difficult. Comparing the RGB models we find that the occupancy maps do not perform well here, but end-

to-end training allowed learning more effective features. The difference to RGBD can be explained by the substantially lower prediction accuracy of our occupancy map predictions.

Choosing what to compare matters. Lines (8-11) compare strategies for choosing which map-pose pairs to feed into our discriminative observation model. Line (8) uses the last 8 steps of the particle trajectory, while lines (9-11) chooses the most relevant past steps based on their estimated overlapping view area. As expected, dynamically choosing what to compare is useful. While one would expect more comparisons to be useful, over 8 comparisons do not improve performance. Since we trained with 8 comparisons, this result indicates that our model overfits to this training condition.

More particles at inference time are useful. Lines (12–18) vary the number of particles at inference time. Surprisingly, as little as 8 particles can already improve over the visual odometry setting (line 12). Increasing the number of particles helps, providing a trade-off between performance and computation. The effect for RGB is less pronounced, and improvement stops over 128 particles.

6.3. Transfer results

An important concern with learning based SLAM method is potential overfitting to the training environments. We take the SLAM-net models trained with the Gibson data, and evaluate them for the Replica and Matterport datasets, with no additional training or hyperparameter tuning. These datasets contain higher quality scenes and cover a wide range of smaller (Replica) and larger (Matterport) indoor environments. The robot and camera parameters remain the same.

Results are in Tab. 4. We observe strong performance for SLAM-net across all datasets and sensor modalities. Comparing alternative methods we observe a similar trend as for the Gibson data (traj_exp columns in Tab. 2). These results on photorealistic data are promising for sim-to-real transfer using real robot navigation data.

6.4. KITTI odometry results

To better understand the limitations of our approach we apply it to the KITTI odometry data, which contains long trajectories of autonomous driving. We do not expect a

SLAM component	SR \uparrow	SPL \uparrow
Ground-truth	90.7%	0.56
SLAM-net (ours)	65.7%	0.38
Learned visual odometry	32.4%	0.19

Table 6: Navigation results.

strong performance. SLAM-net is designed to enable indoor robot navigation which is reflected in a number of design decisions. First, our local maps ignore information far from the camera. Second, we do not have a dedicated loop closure mechanism for large loops. Third, a key benefit of our approach is the joint training of its components, however, this requires large and diverse training data. The KITTI data is relatively small for this purpose. Finally, images in the KITTI data are of high quality for which existing SLAM methods are expected to work well.

Our results are in Tab. 5. We report RMSE in meters after trajectory alignment. SLAM-net results are averaged over 5 seeds. The ORB-SLAM results are for RGB only input taken from [47]. As expected, SLAM-net does not perform as well as ORB-SLAM; nevertheless, it learns a meaningful model and outperforms learned visual odometry. Looking at predicted trajectories we find that SLAM-net occasionally fails to capture turns of the road (visualizations are in the Appendix). One reason is that there are no particles near the true trajectory, or the observation model gives a poor prediction. The training data contains only a limited number of turns, which makes learning from scratch difficult. Indeed our model starts to overfit after a few epochs, suggesting that more training data would improve the performance.

7. Navigation results

The motivation of our work is to enable visual robot navigation in challenging realistic conditions. Our navigation results are reported in Tab. 6 and Tab. 7, using RGB-D input. Videos are in the supplementary material. We report two key metrics following Anderson et al. [3]: success rates (SR) and success weighted path length (SPL).

In Tab. 6 we experiment with our navigation pipeline using different methods for localization and mapping, but keeping the planner and controller modules fixed. Navigation performance is strong with a ground-truth localization oracle, which validates our architecture and serves as an upper-bound for SLAM methods. The navigation architecture with SLAM-net significantly outperforms visual odometry, achieving 65.7% success. Our navigation architecture with visual odometry is conceptually similar to that of Active Neural SLAM [10] and Occupancy Anticipation [51]. Our results are consistent with that of Ramakrishnan et al. [51] in matching conditions. We did not compare with the classic ORB-SLAM method because of its poor performance in our

Rank	Method	SR \uparrow	SPL \uparrow
1	SLAM-net (ours)	64.5%	0.377
2	VO [59]	37.3%	0.266
3	OccupancyAnticipation [51]	29.0%	0.220

Table 7: Habitat 2020 PointNav Challenge leaderboard [25].

prior experiments.

Finally, we submitted our method to the Habitat Challenge 2020 evaluation server, which allows direct comparison with various alternative methods. Tab. 7 shows the top of the leaderboard for the PointNav task. SLAM-net achieves 64.5% success, significantly improving over the SOTA (VO [59], 37.3%). It also outperforms the challenge winner (OccupancyAnticipation [51], 29.0%) which was shown to be superior to Active Neural SLAM [10].

8. Conclusions

We introduced a learning-based differentiable SLAM approach with strong performance on challenging visual localization data and on downstream robot navigation, achieving SOTA in the Habitat 2020 PointNav task.

Together, our results provide new insights for understanding the strengths of classic and learning based SLAM approaches in the context of visual navigation. Our findings partially contradict the results of Mishkin et al. [43], who benchmarked classic and learned SLAM for navigation. While they found ORB-SLAM to be better than learning based SLAM in the same Habitat simulator, they used a noise-free setting and relative goals. As pointed out by Kadian et al. [32], this setting is not realistic in various aspects. Indeed, the public ORB-SLAM implementation of Mishkin et al. [43] failed entirely in our more realistic setting; while our learning-based approach achieved strong performance.

We believe that our work on differentiable SLAM may lay foundation to a new class of methods that learn robust, task-oriented features for SLAM. Future research may investigate alternative differentiable SLAM algorithms, *e.g.*, that build on an optimization-based method instead of particle filtering. While our initial results are promising, future work is needed to apply SLAM-net to real-world robot navigation. A particularly interesting application would be learning to relocalize with significant changes in the environment, a setting known to be challenging for current algorithms.

Acknowledgement

We would like to thank Rico Jonschkowski for suggesting to keep local maps and Gim Hee Lee for valuable feedback. This research/project is supported in part by the National Research Foundation, Singapore under its AI Singapore Program (AISG Award No: AISG2-RP-2020-016) and by the National University of Singapore (AcRF grant R-252-000-A87-114).

References

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, et al. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. 4
- [2] Brandon Amos, Ivan Jimenez, Jacob Sacks, Byron Boots, and J Zico Kolter. Differentiable MPC for end-to-end planning and control. In *Advances in Neural Information Processing Systems*, pages 8299–8310, 2018. 2
- [3] Peter Anderson, Angel Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, et al. On evaluation of embodied navigation agents. *arXiv preprint arXiv:1807.06757*, 2018. 2, 8
- [4] Ali Azarbayejani and Alex P Pentland. Recursive estimation of motion, structure, and focal length. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17(6):562–575, 1995. 2
- [5] Timothy D Barfoot. Online visual motion estimation using fastslam with sift features. In *2005 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 579–585. IEEE, 2005. 2
- [6] Mohak Bhardwaj, Byron Boots, and Mustafa Mukadam. Differentiable gaussian process motion planning. In *International Conference on Robotics and Automation (ICRA)*, pages 10598–10604, 2020. 2
- [7] Michael Bloesch, Jan Czarnowski, Ronald Clark, Stefan Leutenegger, and Andrew J Davison. Codeslam—learning a compact, optimisable representation for dense visual slam. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2560–2568, 2018. 2
- [8] Cesar Cadena, Luca Carlone, Henry Carrillo, Yasir Latif, Davide Scaramuzza, José Neira, Ian Reid, and John J Leonard. Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age. *IEEE Transactions on Robotics*, 32(6):1309–1332, 2016. 1
- [9] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *International Conference on 3D Vision (3DV)*, 2017. MatterPort3D dataset license agreement available at: http://kaldir.vc.in.tum.de/matterport/MP_TOS.pdf. 5
- [10] Devendra Singh Chaplot, Dhiraj Gandhi, Saurabh Gupta, Abhinav Gupta, and Ruslan Salakhutdinov. Learning to explore using active neural slam. In *International Conference on Learning Representations*, 2020. 2, 4, 8, 13
- [11] Devendra Singh Chaplot, Ruslan Salakhutdinov, Abhinav Gupta, and Saurabh Gupta. Neural topological slam for visual navigation. In *Computer Vision and Pattern Recognition (CVPR)*, pages 12875–12884, 2020. 2
- [12] Changhao Chen, Bing Wang, Chris Xiaoxuan Lu, Niki Trigoni, and Andrew Markham. A survey on deep learning for localization and mapping: Towards the age of spatial machine intelligence. *arXiv preprint arXiv:2006.12567*, 2020. 1
- [13] Sungjoon Choi, Qian-Yi Zhou, and Vladlen Koltun. Robust reconstruction of indoor scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015. 5
- [14] Javier Civera, Oscar G Grasa, Andrew J Davison, and JMM Montiel. 1-point ransac for ekf-based structure from motion. In *2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3498–3504. IEEE, 2009. 2
- [15] Andrew J Davison. Real-time simultaneous localisation and mapping with a single camera. In *null*, page 1403. IEEE, 2003. 2
- [16] Arnaud Doucet, Nando De Freitas, and Neil Gordon. An introduction to sequential monte carlo methods. In *Sequential Monte Carlo methods in practice*, pages 3–14. Springer, 2001. 3
- [17] Sebastian East, Marco Gallieri, Jonathan Masci, Jan Koutnik, and Mark Cannon. Infinite-horizon differentiable model predictive control. In *International Conference on Learning Representations*, 2019. 2
- [18] Gregory Farquhar, Tim Rocktäschel, Maximilian Igl, and Shimon Whiteson. TreeQN and ATreeC: Differentiable tree planning for deep reinforcement learning. In *International Conference on Learning Representations*, 2018. 2
- [19] Jorge Fuentes-Pacheco, José Ruiz-Ascencio, and Juan Manuel Rendón-Mancha. Visual simultaneous localization and mapping: a survey. *Artificial intelligence review*, 43(1):55–81, 2015. 1
- [20] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3354–3361. IEEE, 2012. 2, 5
- [21] Giorgio Grisetti, Cyrill Stachniss, and Wolfram Burgard. Improved techniques for grid mapping with rao-blackwellized particle filters. *IEEE transactions on Robotics*, 23(1):34–46, 2007. 2
- [22] Arthur Guez, Théophane Weber, Ioannis Antonoglou, Karen Simonyan, Oriol Vinyals, Daan Wierstra, et al. Learning to search with MCTSnets. In *International Conference on Machine Learning*, 2018. 2
- [23] Saurabh Gupta, James Davidson, Sergey Levine, Rahul Sukthankar, and Jitendra Malik. Cognitive mapping and planning for visual navigation. *arXiv preprint arXiv:1702.03920*, 2017. 1, 2

- [24] Tuomas Haarnoja, Anurag Ajay, Sergey Levine, and Pieter Abbeel. Backprop KF: Learning discriminative deterministic state estimators. In *Advances in Neural Information Processing Systems*, pages 4376–4384, 2016. 2
- [25] Habitat. Leaderboard for the CVPR Habitat 2020 PointNav challenge, 2020. URL <https://evalai.cloudcv.org/web/challenges/challenge-page/580/leaderboard/1631>. Accessed: 2020 November 16. 2, 8
- [26] Ankur Handa, Thomas Whelan, John McDonald, and Andrew J Davison. A benchmark for rgb-d visual odometry, 3d reconstruction and slam. In *International Conference on Robotics and Automation (ICRA)*, pages 1524–1531, 2014. 6
- [27] Jan Hartmann, Dariush Forouher, Marek Litza, Jan Helge Kluessendorff, and Erik Maehle. Real-time visual slam using fastslam and the microsoft kinect camera. In *ROBOTIK 2012; 7th German Conference on Robotics*, pages 1–6. VDE, 2012. 2
- [28] Wolfgang Hess, Damon Kohler, Holger Rapp, and Daniel Andor. Real-time loop closure in 2d lidar slam. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1271–1278. IEEE, 2016. 2
- [29] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *Advances in neural information processing systems*, pages 2017–2025, 2015. 3, 12, 13
- [30] Krishna Murthy Jatavallabhula, Ganesh Iyer, and Liam Paull. gradslam: Dense slam meets automatic differentiation. *arXiv preprint arXiv:1910.10672*, 2019. 2
- [31] Rico Jonschkowski, Divyam Rastogi, and Oliver Brock. Differentiable particle filters: End-to-end learning with algorithmic priors. *Proceedings of Robotics: Science and Systems*, 2018. 2
- [32] Abhishek Kadian, Joanne Truong, Aaron Gokaslan, Alexander Clegg, Erik Wijmans, Stefan Lee, Manolis Savva, Sonia Chernova, and Dhruv Batra. Are We Making Real Progress in Simulated Environments? Measuring the Sim2Real Gap in Embodied Visual Navigation. In *arXiv:1912.06321*, 2019. 1, 4, 5, 8
- [33] Peter Karkus, David Hsu, and Wee Sun Lee. QMDP-net: Deep learning for planning under partial observability. In *Advances in Neural Information Processing Systems*, pages 4697–4707, 2017. 2
- [34] Peter Karkus, David Hsu, and Wee Sun Lee. Particle filter networks with application to visual localization. In *Proceedings of the Conference on Robot Learning*, pages 169–178, 2018. 1, 2, 3, 4
- [35] Peter Karkus, Xiao Ma, David Hsu, Leslie Pack Kaelbling, Wee Sun Lee, and Tomás Lozano-Pérez. Differentiable algorithm networks for composable robot learning. In *Robotics: Science and Systems (RSS)*, 2019. 1, 2
- [36] Peter Karkus, Anelia Angelova, Vincent Vanhoucke, and Rico Jonschkowski. Differentiable mapping networks: Learning structured map representations for sparse visual localization. In *International Conference on Robotics and Automation (ICRA)*, 2020. 2
- [37] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 3
- [38] Georg Klein and David Murray. Parallel tracking and mapping for small ar workspaces. In *2007 6th IEEE and ACM international symposium on mixed and augmented reality*, pages 225–234. IEEE, 2007. 2
- [39] Sven Koenig and Maxim Likhachev. Fast replanning for navigation in unknown terrain. *IEEE Transactions on Robotics*, 21(3):354–363, 2005. 4
- [40] Gim Hee Lee, Friedrich Fraundorfer, and Marc Pollefeys. Rslam: Ransac sampling for visual fastslam. In *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 1655–1660. IEEE, 2011. 2
- [41] Ruihao Li, Sen Wang, Zhiqiang Long, and Dongbing Gu. Undeepvo: Monocular visual odometry through unsupervised deep learning. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 7286–7291. IEEE, 2018. 2
- [42] Xiao Ma, Peter Karkus, David Hsu, and Wee Sun Lee. Particle filter recurrent neural networks. In *AAAI Conference on Artificial Intelligence*, volume 34, pages 5101–5108, 2020. 2
- [43] Dmytro Mishkin, Alexey Dosovitskiy, and Vladlen Koltun. Benchmarking classic and learned navigation in complex 3d environments. *arXiv preprint arXiv:1901.10915*, 2019. 1, 2, 5, 8
- [44] Michael Montemerlo, Sebastian Thrun, Daphne Koller, Ben Wegbreit, et al. Fastslam: A factored solution to the simultaneous localization and mapping problem. *AAAI Conference on Artificial Intelligence*, 593598, 2002. 1, 2, 3
- [45] Michael Montemerlo, Sebastian Thrun, Daphne Koller, Ben Wegbreit, et al. Fastslam 2.0: An improved particle filtering algorithm for simultaneous localization and mapping that provably converges. In *IJCAI*, pages 1151–1156, 2003. 2, 3
- [46] Raul Mur-Artal and Juan D Tardós. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *IEEE Transactions on Robotics*, 33(5):1255–1262, 2017. 2, 5, 6, 7
- [47] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *IEEE Transactions on Robotics*, 31(5):1147–1163, 2015. 2, 5, 7, 8
- [48] Adithyavairavan Murali, Tao Chen, Kalyan Vasudev Alwala, Dhiraj Gandhi, Lerrel Pinto, Saurabh Gupta, and Abhinav Gupta. Pyrobot: An open-source robotics framework for research and benchmarking. *arXiv preprint arXiv:1906.08236*, 2019. 5

- [49] Junhyuk Oh, Satinder Singh, and Honglak Lee. Value prediction network. In *Advances in Neural Information Processing Systems*, 2017. 2
- [50] Masashi Okada, Luca Rigazio, and Takenobu Aoshima. Path integral networks: End-to-end differentiable optimal control. *arXiv preprint arXiv:1706.09597*, 2017. 2
- [51] Santhosh K Ramakrishnan, Ziad Al-Halah, and Kristen Grauman. Occupancy anticipation for efficient exploration and navigation. *arXiv preprint arXiv:2008.09285*, 2020. 2, 8
- [52] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, et al. Habitat: A platform for embodied AI research. In *International Conference on Computer Vision (ICCV)*, 2019. 2, 5
- [53] Hauke Strasdat, José MM Montiel, and Andrew J Davison. Visual slam: why filter? *Image and Vision Computing*, 30(2): 65–77, 2012. 2
- [54] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J. Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, Anton Clarkson, Mingfei Yan, Brian Budge, Yajie Yan, Xiaqing Pan, June Yon, Yuyang Zou, Kimberley Leon, Nigel Carter, Jesus Briales, Tyler Gillingham, Elias Mueggler, Luis Pesqueira, Manolis Savva, Dhruv Batra, Hauke M. Strasdat, Renzo De Nardi, Michael Goesele, Steven Lovegrove, and Richard Newcombe. The Replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019. Replica dataset license agreement available at: <https://github.com/facebookresearch/Replica-Dataset/blob/master/LICENSE>. 5
- [55] Aviv Tamar, Yi Wu, Garrett Thomas, Sergey Levine, and Pieter Abbeel. Value iteration networks. In *Advances in Neural Information Processing Systems*, 2016. 1, 2
- [56] Chengzhou Tang and Ping Tan. Ba-net: Dense bundle adjustment network. *arXiv preprint arXiv:1806.04807*, 2018. 2
- [57] Keisuke Tateno, Federico Tombari, Iro Laina, and Nassir Navab. Cnn-slam: Real-time dense monocular slam with learned depth prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6243–6252, 2017. 2
- [58] Jonas Uhrig, Nick Schneider, Lukas Schneider, Uwe Franke, Thomas Brox, and Andreas Geiger. Sparsity invariant cnns. In *2017 international conference on 3D Vision (3DV)*, pages 11–20. IEEE, 2017. 5
- [59] Unknown. VO. the method could not be identified at the time of submission. *Submission to the Habitat PointNav Challenge*, 2020. 8
- [60] Erik Wijmans, Abhishek Kadian, Ari Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis Savva, and Dhruv Batra. DD-PPO: Learning near-perfect pointgoal navigators from 2.5 billion frames. In *International Conference on Learning Representations*, 2020. 2
- [61] Fei Xia, Amir R. Zamir, Zhi-Yang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gibson env: real-world perception for embodied agents. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. Gibson dataset license agreement available at: http://svl.stanford.edu/gibson2/assets/GDS_agreement.pdf. 5
- [62] Ryo Yonetani, Tatsunori Tanai, Mohammadamin Berekatain, Mai Nishimura, and Asako Kanezaki. Path planning using neural a* search. *arXiv preprint arXiv:2009.07476*, 2020. 2
- [63] Shuaifeng Zhi, Michael Bloesch, Stefan Leutenegger, and Andrew J Davison. Scenecode: Monocular dense semantic reconstruction using learned encoded scene representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11776–11785, 2019. 2
- [64] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1851–1858, 2017. 2
- [65] Michael Zhu, Kevin Murphy, and Rico Jonschkowski. Towards differentiable resampling. *arXiv preprint arXiv:2004.11938*, 2020. 3

A. SLAM-net components

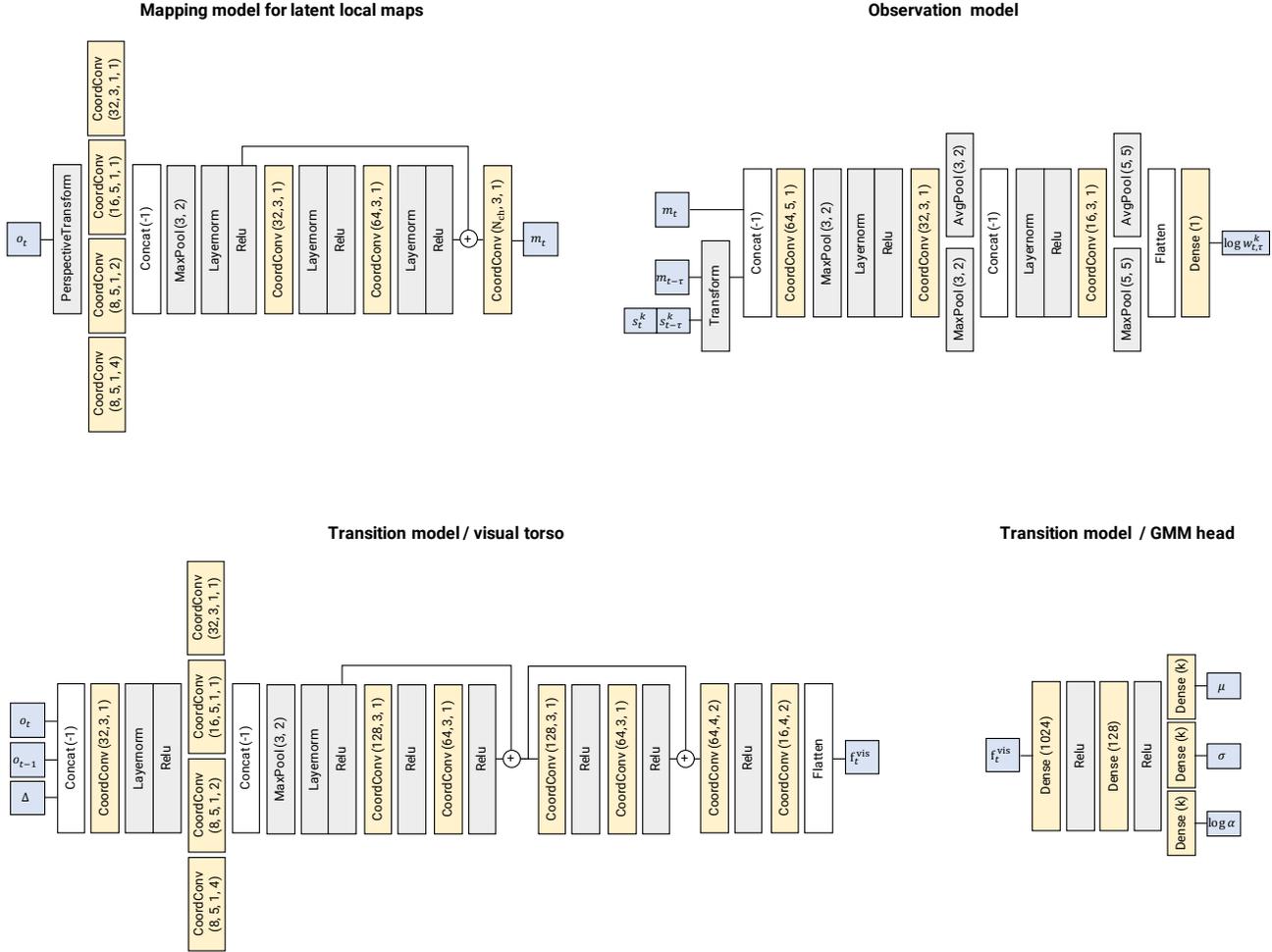


Figure 4: **Network architectures of SLAM-net components.** The notation is as follows. **Conv**(f, k, s, d): convolutional layer with f filters, $k \times k$ kernel, $s \times s$ strides, $d \times d$ dilation. **CoordConv**: convolutional layer that takes in extra channels that encode constant pixel coordinates. **MaxPool**(k, s): max-pooling layer with $k \times k$ kernel and $s \times s$ strides. **AvgPool**(k, s): average-pooling layer with $k \times k$ kernel and $s \times s$ strides. **Dense**(h) dense fully-connected layer with h output units. **Concat**(-1): concatenate input images along their channel axis. **Transform**: spatial image transformation using a spatial transformer network [29]. Given a relative pose and orientation it applies translational and rotational image transformations. **PerspectiveTransform**: geometric perspective transformation to a top-down view, given a known camera matrix and a target resolution. **Inputs and outputs.** o_t : input image at time t , either RGB or depth. m_t : local map at time t . s_{t-t}^k : state s along particle trajectory for particle k and time t . $\log w_{t-t}^k$: estimated compatibility of map-pose pairs at time t and $t - \tau$ for particle k . Δ : difference of input image pair $o_t - o_{t-1}$. f_{t-t}^{vis} : intermediate visual features in the transition model. μ : mean predictions for a Gaussian mixture model (GMM) with $k=3$ components. σ : predicted standard deviations. $\log \alpha$: mixture log-probabilities.

A.1. Observation model

In the particle filter the observation model estimates the log-likelihood $\log w_t^k$, the probability of the current observation o_t given a particle trajectory $s_{1:t}^k$ and past observations $o_{1:t-1}$. The particle weight is multiplied with the estimated log-likelihood.

In SLAM-net we decompose this function. A mapping model first predicts a local map m_t from o_t . The observation model

than estimates the compatibility of m_t with $m_{1:t-1}$ and $s_{1:t}^k$, by summing pair-wise compatibility estimates.

$$\log w_t^k \approx \sum_{\tau \in T} \log w_{t,\tau}^k \quad (1)$$

$$\log w_{t,\tau}^k = f_{\theta}^{\text{obs}}(m_t, s_t^k, m_{t-\tau}, s_{t-\tau}^k) \quad (2)$$

The network architecture of f_{θ}^{obs} is shown in Fig. 4. Weights are shared across particles and time steps. The important component of this model is the image transformation (Transform), that applies translational and rotations image transformations to $m_{t-\tau}$ given the relative poses defined by s_t^k and $s_{t-\tau}^k$. We use spatial transformer networks [29] that implement these transforms in a differentiable manner.

A.2. Mapping model

The mapping model component takes in the o_t observation (160×90 RGB or depth image) and outputs a local map m_t . Local maps are $40 \times 40 \times N_{\text{ch}}$ grids, which can be understood as images with N_{ch} channels. The local maps either encode latent features ($N_{\text{ch}} = 16$) or they are trained to predict the occupied and visible area. Each cell in the local map corresponds to a 12×12 cm area in front of the robot, *i.e.*, the local map covers a 4.8×4.8 m area.

The network architecture (configured for latent local maps) is shown in Fig. 4. The same network architecture is used for RGB and depth input. We apply a fixed perspective transform to the first-person input images, transforming them into 160×160 top-down views that cover a 4.8×4.8 m area. The perspective transformation assumes that the camera pitch and roll, as well as the camera matrix are known.

When SLAM-net is configured with local maps that predict occupancy we use a similar architecture but with separately learned weights. In case of RGB input our network architecture is the same as the mapping component in ANS [10], using the same ResNet-18 conv-dense-deconv architecture. We freeze the first three convolutional layer of the ResNet to reduce overfitting. In case of depth input our network architecture is similar to the one in Fig. 4, but it combines top-down image features with first-person image features. When SLAM-net is configured with both occupancy and latent local maps we use separate network components and concatenate the local map predictions along their last (channel) dimension.

For the KITTI experiments we use 40×40 local maps where cells are 70×70 cm, *i.e.*, a local map captures a 28×28 m area. We use an equivalent network architecture that is adapted to the wider input images.

A.3. Transition model

The transition model takes in the last two consecutive observations (o_t, o_{t-1}) and it outputs a distribution over the relative motion components of the robot. It can also take in the last robot action a_{t-1} when it is available.

Our transition model parameterizes a Gaussian mixture model (GMM) with $k=3$ mixture components, separately for each coordinate (x, y, yaw) and each discrete robot action a . The network architecture is shown in Fig. 4. The model first extracts visual features f_t^{vis} from a concatenation of o_t, o_{t-1} and their difference $\Delta = o_t - o_{t-1}$. The visual features are then fed to GMM heads for each combination of robot action a and x, y, yaw motion coordinates. Each GMM head uses the same architecture with independent weights.

B. Additional figures

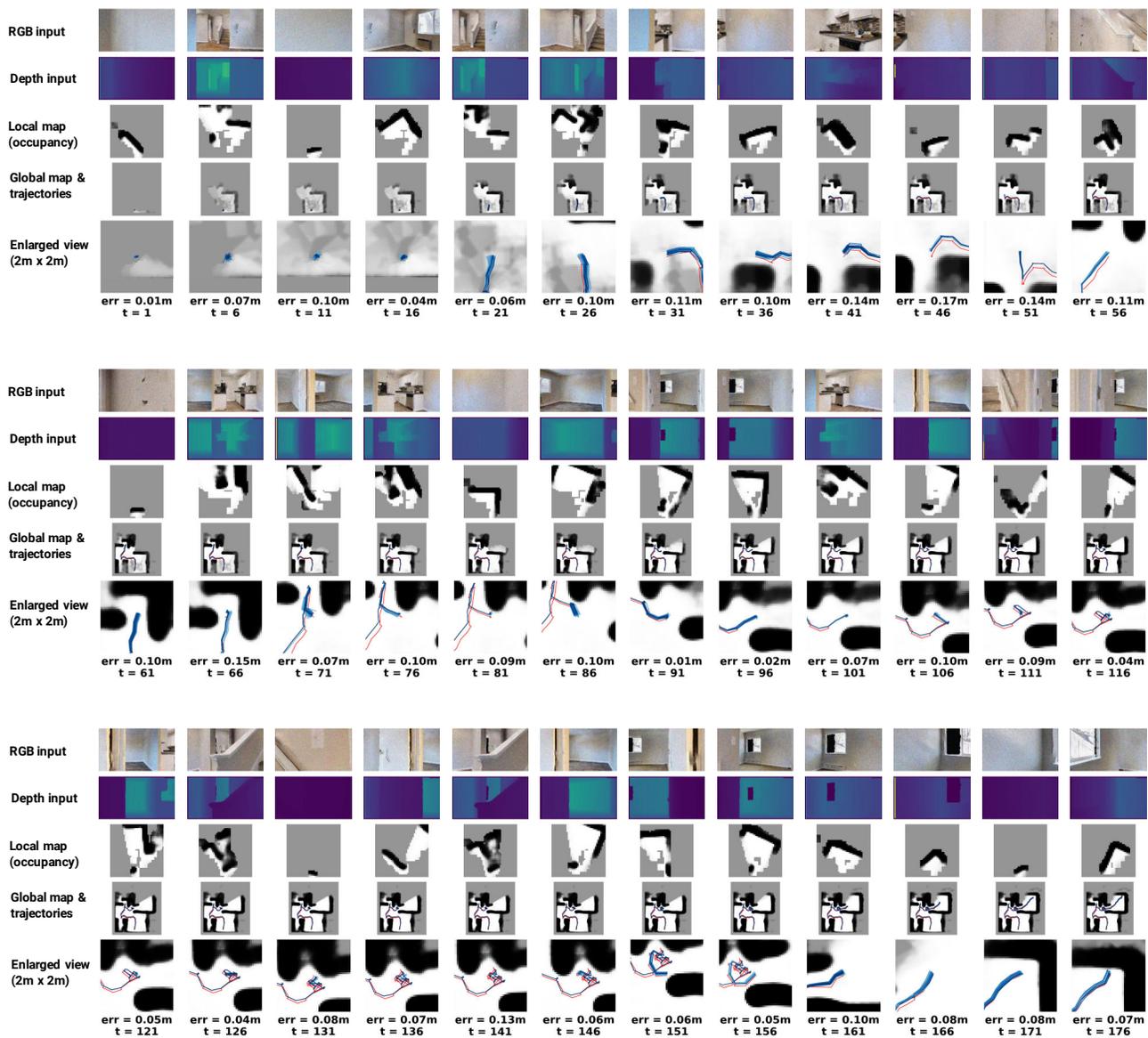


Figure 5: **SLAM-net with RGBD sensor.** Trajectory sample from the Gibson data `traj_exp_rand` test set. In this configuration SLAM-net local maps predict occupancy (third row). The enlarged view (bottom row) shows a 2×2 m window of the predicted global map. The true trajectory is in red; particle trajectories are in blue. The shade of blue indicate particle weights (darker color means larger particle weight). Notice that particle trajectories maintain multi-modal trajectories. Low weight particles are dropped after resampling. The figure is best viewed using zoom.

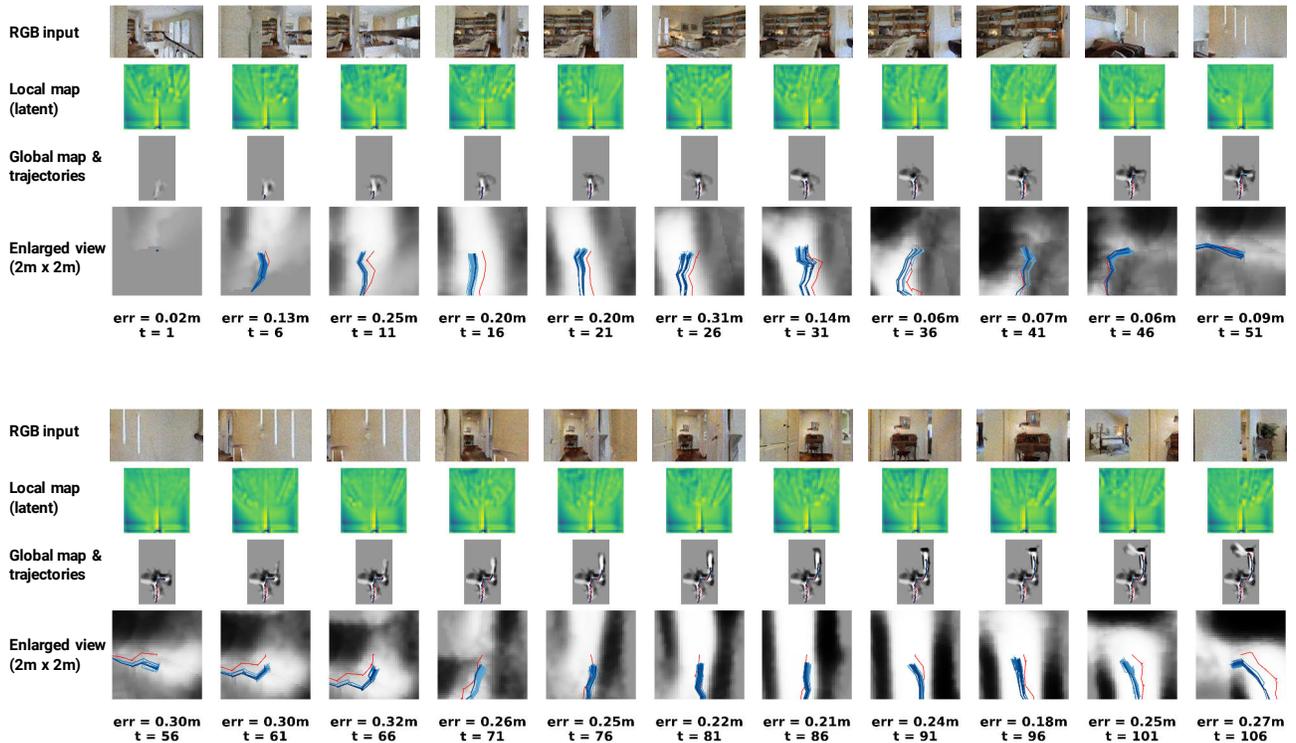


Figure 6: **SLAM-net with RGB only sensor.** Trajectory sample from the Gibson data traj_exp test set. In this configuration SLAM-net local maps have 16 latent feature channels. We visualize normalized features of a single channel (second row). The enlarged view (bottom row) shows a 2×2 m window of the predicted global map. The true trajectory is in red; particle trajectories are in blue. The shade of blue indicate particle weights (darker color means larger particle weight). Notice that particle trajectories maintain multi-modal trajectories. Low weight particles are dropped after resampling. The figure is best viewed using zoom.

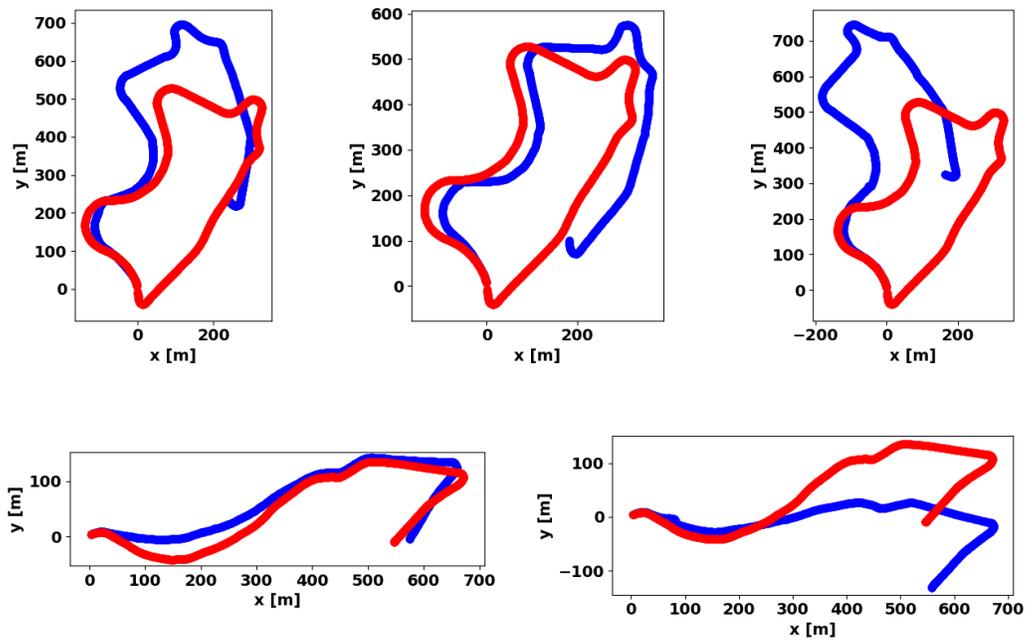


Figure 7: **KITTI trajectories.** The figure shows SLAM-net predictions (before alignment) with different random seeds for the Kitti-09 trajectory (top row) and the Kitti-10 trajectory (bottom row). The predicted trajectory is in blue, the true trajectory is in red.