

BioSLAM: A Bio-inspired Lifelong Memory System for General Place Recognition

Peng Yin^{1,*}, Member, IEEE, Abulikemu Abuduweili^{2,*}, Shiqi Zhao¹
Lingyun Xu², Changliu Liu², Member, IEEE and Sebastian Scherer², Senior Member, IEEE

Abstract—We present BioSLAM, a lifelong SLAM framework for learning various new appearances incrementally and maintaining accurate place recognition for previously visited areas. Unlike humans, artificial neural networks suffer from catastrophic forgetting and may forget the previously visited areas when trained with new arrivals. For humans, researchers discover that there exists a memory replay mechanism in the brain to keep the neuron active for previous events. Inspired by this discovery, BioSLAM designs a gated generative replay to control the robot’s learning behavior based on the feedback rewards. Specifically, BioSLAM provides a novel dual-memory mechanism for maintenance: 1) a dynamic memory to efficiently learn new observations and 2) a static memory to balance new-old knowledge. When the agent is encountered with different appearances under new domains, the complete processing pipeline can help to incrementally update the place recognition ability, robust to the increasing complexity of long-term place recognition.

We demonstrate BioSLAM in three incremental SLAM scenarios: 1) a 120km city-scale trajectories with LiDAR-based inputs, 2) a multi-visited 4.5km campus-scale trajectories with LiDAR-vision inputs, and 3) an official Oxford dataset with 10km visual inputs under different environmental conditions. We show that BioSLAM can incrementally update the agent’s place recognition ability and outperform the state-of-the-art incremental approach, Generative Replay, by 24% in terms of place recognition accuracy. To our knowledge, BioSLAM is the first memory-enhanced lifelong SLAM system to help incremental place recognition in long-term navigation tasks.

Index Terms—Lifelong SLAM, Incremental Place Recognition, Continuous Localization

I. INTRODUCTION

AN essential capability for long-term robotics autonomy in the open world without human assistance is lifelong Simultaneous Localization and Mapping (SLAM) [1]. In the context of lifelong SLAM, the system needs to consider work in long-term navigation in large-scale environments and diverse environmental conditions, as depicted in Fig. 1. Current SLAM methods are mainly conducted under single-type environments, where the environmental conditions (such as illuminations, weather, seasons, etc.) are consistent. Recent works attempt to relax the *single-type* assumption to accommodate diverse environments by leveraging domain adaptation

Peng Yin and Shiqi Zhao are with the Department of Mechanical Engineering, City University of Hong Kong, Hong Kong 518057, China. {pengyin@cityu.edu.hk, ryanzhao9459@gmail.com}. Abulikemu Abuduweili, Lingyun Xu, Changliu Liu, and Sebastian Scherer are with the Robotics Institute, Carnegie Mellon University, Pittsburgh, PA 15213, USA. (abulikea, lingyun2, cliu6, basti@andrew.cmu.edu).

*Authors Peng Yin and Abulikemu Abuduweili contributed equally.

Corresponding author: Peng Yin (pengyin@cityu.edu.hk)

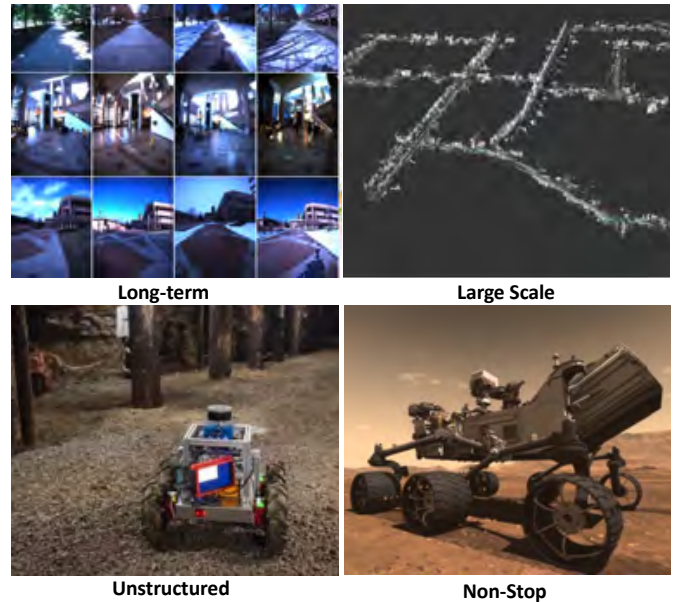


Fig. 1: **Challenges in Real-world Robotic Localization.** For real-world field applications, robotic localization usually encounters the following challenges: 1) changing appearance under long-term traveling, 2) diverse geometric differences under large-scale areas, 3) mixture structure/unstructured environments, and 4) non-stop restriction for long-term autonomy.

techniques [2], [3] into model learning with deep neural networks. However, the domain knowledge under new scenarios can affect the localization accuracy in previous learned areas, an effect known as “catastrophic forgetting”.

In real-world long-term navigation [4], the robot may encounter complicated 3D environments, such as campus areas, open streets, residential blocks, commercial buildings, etc., and each place has its unique patterns in place recognition. The robot platform can’t collect datasets under all scenarios at once and train the localization module in a supervised manner. A naive solution for incremental observations is to source additional data for model adaptation with a new scenario; however, this adaptation is not feasible when the goal is to ensure the uninterrupted and long-term operation of the robot since it causes catastrophic forgetting of previous knowledge. Moreover, changes in environments can be sudden, e.g., rapid illumination and weather changes, while it may take too long for traditional learning-based approaches to react to the changes. Given the above consideration, the main challenges for lifelong place recognition include:

- *Various environmental conditions*: the appearances of the

same area under different environmental conditions will be represented with different patterns.

- *Diverse scenarios*: the robot platform will encounter different 3D environments in large-scale navigation tasks, and most areas are a combination of different types.
- *Non-stop training*: the robot will accumulate new datasets, and model fine-tuning is usually required to improve localization performance for new scenarios.

Though the topic of long-term SLAM [5]–[9] has been well studied in the past decade, in this work, we narrow down the scope to the lifelong place recognition in long-term navigation and proposed BioSLAM, that can continuously re-localize to new environments without sacrificing recognizing ability in previously seen environments. In our previous work [10], we notice that cross-domain appearance differences will significantly affect the recognizing performance; the recognizing module encounters the catastrophic forgetting problem, where it is only robust to the most recently trained scenarios. In contrast, humans and animals do not suffer from catastrophic forgetting, and short-term and long-term memory mechanisms exist within the hippocampus [11] and the front lobe of the brain [12], which plays the main role in lifelong knowledge updating. Recently, new evidence from fMRI studies in humans [13] finds that the hippocampus may ‘act as a librarian to retrieve the cortical books of memory’, i.e., the hippocampus can index the memories for fast retrievals. Inspired by the biological mechanism, we design two memory zones for BioSLAM, namely static memory zones (SMZ) for historical memory encoding with low frequency and dynamic memory zone (DMZ) for quickly memory reply, and propose a dual-memory selection mechanism to balance the short-term adaptation for new observations and long-term memory retention for historical knowledge. Specifically, BioSLAM also develops a sleeping cycle for memory consolidation within SMZ, which is also inspired by a similar mechanism in the hippocampus [14]. Based on the above mechanism, BioSLAM has the ability to achieve long-term place recognition.

The evaluation methods [15] for the traditional place recognition using supervised learning approaches do not apply to lifelong systems. The adaptation capability reflects the performance of lifelong systems concerning new observations and the long-term memory retention of previously visited areas. In this work, we formulate two metrics, namely adaptation efficiency (AE) analysis, and retention ability (RA) analysis, and perform an extensive evaluation using three long-term datasets: 1) *ALITA Urban Dataset*, which is focused on changing geometric patterns, and 2) *ALITA Campus Dataset*, which is focused on changing illumination patterns, 3) *Oxford RobotCar Dataset*, which is an official long-term datasets with different environmental conditions. The contributions of this paper are as follows:

- BioSLAM provides a systematic framework to learn about ever-changing environments without interruption. This framework enables incremental place feature learning for long-term autonomy.
- Within BioSLAM, we develop a rewarding mechanism with a dynamic and static memory zone, which contains the task-oriented external reward and curiosity-oriented internal

reward, which can quickly adapt new patterns and maintain memorization for long-term memory retention.

The rest of the paper will introduce the related works for place recognition and lifelong incremental learning in section II. Section III gives the structural overview of BioSLAM. Section IV and section V explain the details of the general place feature learning and bio-inspired lifelong memory, respectively. The experiment setup and qualitative/quantitative analysis are given in section VI and section VII.

II. RELATED WORKS

There are two essential modules in lifelong navigation: 1) navigation and 2) lifelong learning. The navigation task usually contains the place recognition (PR) or Loop closure detection (LCD) module as stated in [16]–[18], which mainly serves as the data association for large-scale re-localization and map optimization in SLAM tasks. Lifelong learning, also known as continual, incremental, or sequential learning, aims at incrementally building up knowledge from a sequential data stream [19], [20], which is essential for long-term localization where robots will encounter many infinite environments. In the following subsections, we will mainly introduce the related works in visual/LiDAR navigation and recent lifelong learning works from a robotics perspective.

A. Long-term Navigation

In long-term navigation, the place recognition targets identifying the exact areas under different perspectives and environmental conditions [16].

The traditional geometry descriptors (e.g., scale-invariant feature transform (SIFT) [21] and oriented FAST and rotated BRIEF (ORB) [22]) are widely used in visual place recognition because of their invariant properties to scale, orientation and illumination changes. Based on these handcrafted features, FAB-MAP [23] build a Bag-of-visual-words (BoW) architecture to achieve large-scale visual re-localization [8], [24]. iBoW-LCD [25] uses an incremental BoW scheme based on binary descriptors to retrieve matched images more efficiently. Shan An *et al.* introduces FILD++ [26], an incremental loop closure detection approach via constructing a hierarchical small-world graph. With the booming of deep learning, new convolutional neural network (CNNs) features, provide significant improvements in feature/semantic extraction. NetVLAD [27] combined the CNN features and a differentiable VLAD [28] layer to enable deep learning for visual place recognition; and based on NetVLAD, recent deep learning approaches [29], [30] further improve the recognition accuracy with different networks.

Despite the success of existing place recognition methods, the non-learning-based approaches are sensitive to parameter tuning under different scenarios; and learning-based techniques are trained in a supervised learning manner, restricting their generalization ability within the offline training datasets. However, in real-world and long-term tasks [31], [32], the data stream is infinite with the combination of different areas under varying environmental conditions; meanwhile, robotic systems

cannot stop and wait for the network model to update for newly encountered scenarios.

Except for the above place descriptor-based SLAM system, there are other remarkable works for long-term navigation. [5] provides the Experience-based approach for ever-changing environments, and the robot can switch between different experience traces while maintaining the robustness of the localization system. Recently, [33] extended the experience-based long-term navigation for the UGV routine-following task, and [34], [35] further extended it to the teach-and-repeat visual navigation task for UAV systems. [36] developed the linear regression-based supervised change prediction mechanism to handle the predictable changes in the long-term navigation task. In [6], the author provides a map summarizing framework for lifelong visual navigation, where the multiple visited maps are incorporated into a joint map which shows better generalization ability for changing environments. On the other hand, dynamic changes within the 3D environments may also cause localization failures in the long-term navigation task. [37] provides a frequency-enhanced map monitor mechanism, which can detect the regular changeable appearance in the long-term navigation task. Besides the experience traces, [38] also provides a lifelong navigation approach based on a particle filter with a hidden Markov model; this method can also long-term UGV localization over the parking areas under ever-changing parking spaces. [39] construct graph pruning for Lifelong SLAM, which can balance the graph size and mapping performance for long-term navigation requirements.

In this work, we target lifelong localization, where the place observations will be viewed only once in the sequential order [19]. Instead of focusing on short-term localization or fixed pattern localization [15] in most existing place recognition methods or using the experience-based/frequency-based mechanisms to keep the long-term robust appearance features, we focus on how to provide the lifelong training procedure for the learnable place descriptor.

B. Lifelong Learning for Robotics

Lifelong learning, also known as continual learning, aims at providing incrementally updated knowledge in ever-changing environments. Though this area has been studied for a long time, most approaches are still restricted to simulation or toy datasets [19] and can not be applied in real robotic applications [40]–[42]. As mentioned in [20], the fundamental challenge for lifelong learning is not necessarily finding solutions that work in the real world but rather finding stable algorithms that can learn in the real world and overcome the catastrophic forgetting problem. Recent works can be roughly divided into four families: dynamic architectures, regularization-based, rehearsal, and generative replay approach.

Dynamic architecture-based methods either 1) add additional parameters to the models, such as LwF [43], which use shared early feature extraction layers and fixed task layers; or 2) use model adaptation to avoid catastrophic forgetting, such as PackNet [44], which defines the mask layer to protect weights when learning new tasks. Regularization-based methods in the context of lifelong learning can add constraints to

avoid overfitting to new tasks and keep inference ability for the previous mission, such as Elastic Weight Consolidation (EWC) [45] and Synaptic Intelligence (SI) [46]. Airloop [47] proposed a lifelong learning method for visual loop closure detection utilizing a euclidean-distance knowledge distillation loss on images. InCloud [48] proposed an angular distillation loss to encourage the network to preserve the structure of its embedding space. However, the above methods must deal with specific network structures and can quickly converge to undesired local optima for complex tasks. Rehearsal-based methods, on the other hand, use memory replays to enhance the knowledge from the previous tasks or processes such as iCaRL [49], GEM [50], which use a small subset of the previous dataset to balance the knowledge distribution for different tasks. Instead of maintaining the knowledge based on past data samples, generative replay [51] combines the actual raw data and generated artificial data for model updating. In [52], the authors use a dual teacher-student generative replay method for incremental learning, where the teacher network is frozen to guide new networks, and the networks will switch roles when the student network surpasses the teacher. Meanwhile, for the memory mechanism in lifelong navigation, [53] introduced a long-short-mechanism to transfer the new observation to the long-term memory while maintaining robust localization performance under changing environments.

In the long-term localization task, the robots will encounter many areas under different viewpoints and environmental conditions. Hence, the ideal approach would be tackling the real-world localization problem in an embodied platform: an autonomous agent that can efficiently and incrementally update its localization ability with limited computation resources. Similar to [53], BioSLAM also utilizes the dual-memory mechanism to adapt to new observations while maintaining memorization ability for long-term knowledge; the significant difference is that BioSLAM introduced the rewarding mechanism (internal & external rewards) and the memory consolidation for the dual memory system.

III. PROBLEM FORMULATION & SYSTEM OVERVIEW

This paper presents BioSLAM, an incremental place recognition method that enables continual learning of place feature extraction modules across various domains. Various domains may encompass different weather conditions, road areas, or input signal modalities. The approach consists of two fundamental modules: 1) a general place feature extraction module for encoding place features under different domains, and 2) a bio-inspired lifelong memory system for continual learning place recognition from different domains. In this section, we will first formulate the lifelong localization problem, then introduce the two key modules in the BioSLAM system.

A. Problem Formulation

We define a sequence of place observations under domain D (i.e., visual or LiDAR) as $O^D = \{O_1, \dots, O_M\}$, and a query of observations under the same domain as $Q^D = \{Q_1, \dots, Q_N\}$. The task of traditional place recognition is to learn a feature extraction function \mathcal{F} with parameter θ to help each frame in

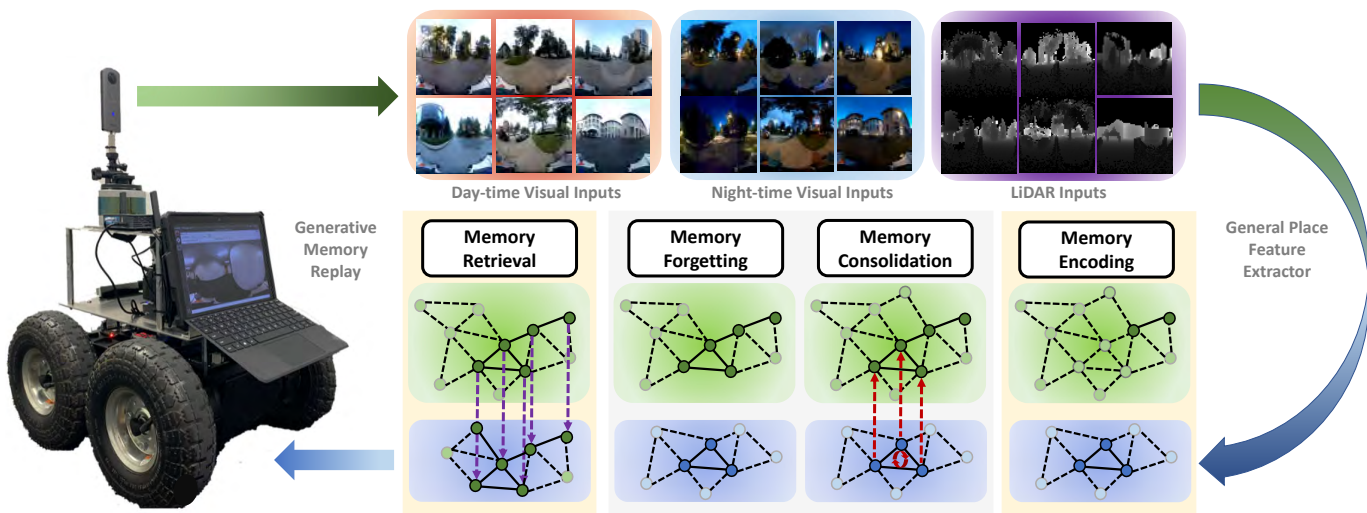


Fig. 2: **Lifelong Localization System Framework.** The lifelong localization system contains the following modules: 1) the general place feature extraction network, which extracts the place feature from different domains; 2) the bio-inspired lifelong memory system can provide short-term and long-term assistance to capture new knowledge and maintain old knowledge.

Q^D find positive (or neighbor) places from the reference sets O^D . Let $d(\cdot, \cdot)$ denote the difference matrix (i.e., Euclidean distance). The objective is to make the feature differences of positive (or neighbor) places smaller than negative (or far-away) places by feature extraction function \mathcal{F}_θ .

$$\begin{aligned} \mathcal{L}(Q_k^D) &= d(\mathcal{F}_\theta(Q_k^D), \mathcal{F}_\theta(O_{\approx}^D)) - d(\mathcal{F}_\theta(Q_k^D), \mathcal{F}_\theta(O_{\neq}^D)) \\ \theta^* &= \arg \min_{\theta} \sum_{k=1}^N \mathcal{L}(Q_k^D) \end{aligned} \quad (1)$$

where Q_k^D is the current k -th query, O_{\approx}^D is the positive reference in a neighbor range near Q_k^D , and O_{\neq}^D is the negative reference away from Q_k^D .

In lifelong localization, the environmental domains ($D_1, \dots, D_t, \dots, D_T$) can vary under different environmental conditions or sensor modalities (e.g. incrementally learning from different illuminations or weathers). Thus, both references set O^{D_t} and the query set Q^{D_t} are obtained incrementally. As depicted in Fig. 2, the lifelong localization problem is to incrementally learn and update the feature extraction function \mathcal{F}_θ , that can quickly adapt its feature extraction ability in the newest domain $\{O^{D_T}, Q^{D_T}\}$, and also in parallel maintain the feature distinguish ability for previous domains $\{O^{D_t}, Q^{D_t}\}_{t=1, \dots, T-1}$. Since we are considering lifelong learning [19], raw data of different domains are fed sequentially for one-time usage and cannot be stored for offline training. Thus, when optimizing feature extraction function \mathcal{F}_{θ_T} at the current domain $\{O^{D_T}, Q^{D_T}\}$, we cannot access the previous raw data from $\{O^{D_t}, Q^{D_t}\}_{t=1, \dots, T-1}$. For time step T , lifelong localization can be formulated as,

$$\begin{aligned} \mathcal{L}^t(Q_k^{D_t}) &= d(\mathcal{F}_\theta(Q_k^{D_t}), \mathcal{F}_\theta(O_{\approx}^{D_t})) - d(\mathcal{F}_\theta(Q_k^{D_t}), \mathcal{F}_\theta(O_{\neq}^{D_t})) \\ \theta_T &= \arg \min_{\theta} \sum_{t=1}^T \sum_{k=1}^N \mathcal{L}^t(Q_k^{D_t}) \end{aligned} \quad (2)$$

B. General Place Feature Extraction

For the lifelong purpose of long-term localization, we developed a General Place Feature Extraction or General Place Learning (GPL) network based on our previous works in visual [10] and LiDAR-based [54] localization, which can be referred to as the feature extraction function \mathcal{F}_θ as in Eq. (2). We use the shared spherical convolution network to simultaneously achieve LiDAR and visual place localization. The spherical harmonic-based convolution can help the learned descriptor have the viewpoint-invariant propriety for the same place recognition. The major difference between the current GPL and our previous works is that GPL does not contain any domain-transfer module, which has been used to reduce the feature differences for the same areas under different domains [10]. This modification is because we want to evaluate the adaptation ability for the same network. On the other hand, no task-specific network layers are used in the dynamic architecture-based lifelong modules, as stated in section II-B. We want to avoid uncertain parameters and only focus on how memory mechanisms can help incremental learning for real-world applications.

C. Bio-inspired Lifelong Memory

Inspired by the memory system in human-being and other mammal animals [13], we propose a dual-memory, dynamic memory (DM) zone and static memory (SM) zone, enhanced lifelong learning mechanism to deal with catastrophic forgetting in lifelong localization. As studied in [55], to create long-term memories in our brain, we have so-called *sleep circle* during our sleep: 1) the brain can *encoding* our daily observation into the *hippocampus* zone with decay along the time; 2) and a consolidation mechanism is triggered between the hippocampus and the neocortex to store essential memory traces and forget the rest traces; 3) then humans can retrieve the relative memory traces based on the consolidated ones in the neocortex. In BioSLAM, we re-build the ‘*sleep circle*’ for the

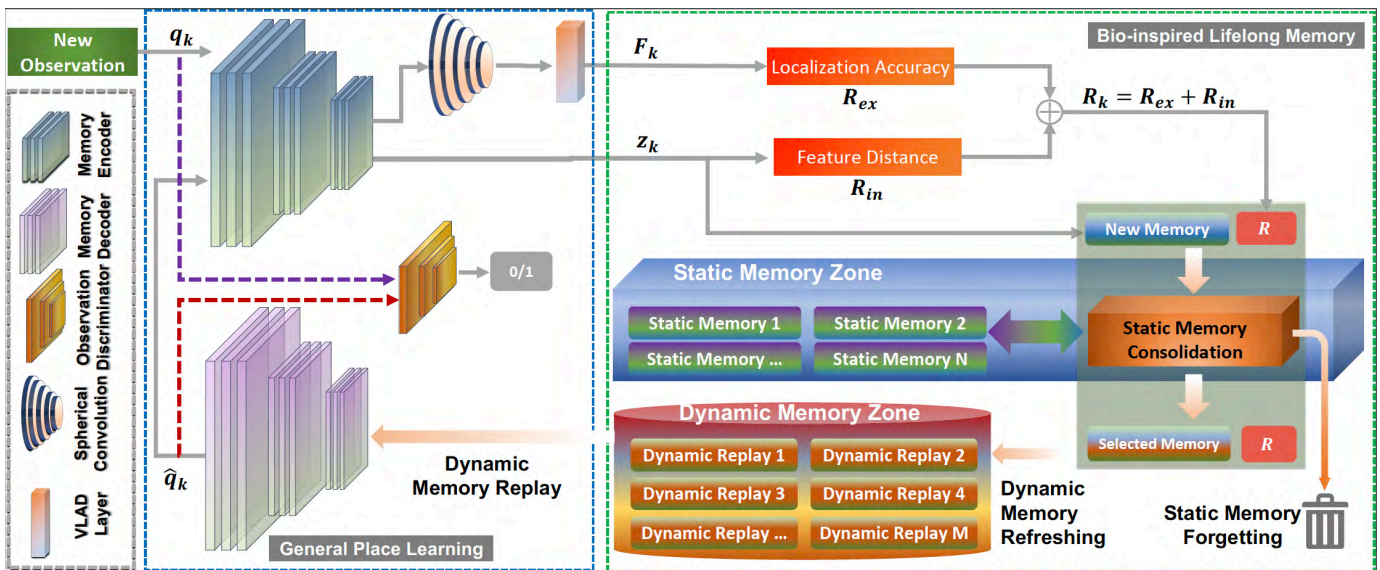


Fig. 3: **BioSLAM Network Structure.** The structure of BioSLAM includes the General Place Learner (GPL) network, the rewarding mechanism to guide the memory storing and consolidation, and the dual-memory module with static-/dynamic-memory zones. The whole procedure includes: 1) new observations q_k are fed into the networks only one time and sequentially; 2) In GPL, the memory encoder converts the inputs q_k to encoded memory z_k , followed by spherical convolution and VLAD layer to generate place feature F_k . 3) the rewarding mechanism will estimate the external reward \mathcal{R}_{ex} and internal rewards \mathcal{R}_{in} to guide the memory operations; 4) the dual-memory will conduct memory storing/consolidation and replay the important (high-rewarded) memories. 5) the memory decoder generates the synthetic samples from the replayed memories; 6) updating the GPL network module with both raw (real) samples and synthetic samples.

lifelong localization task. As we can see in Fig. 2, the memory system of BioSLAM also includes the place feature ‘encoding’ procedure for new observations, the memory ‘consolidation’ controlled by a rewarding mechanism to filter out necessary traces for more extended storage, and the ‘retrieved’ memory to re-enhance the long-term place recognition ability. Based on the above architecture, the BioSLAM system constructs two major modules, the *General Place Learning* (GPL) module and the *Bio-inspired Lifelong Memory* (BiLM) module, which will be investigated in section IV and section V.

IV. GENERAL PLACE LEARNING

As shown in Fig. 3, the general place learning (GPL) (blue dashed box) system mainly contains two sub-modules: a place memory encoding module (upper part of the blue dashed box) and a generative memory reply module (lower part). All samples under different domains are fed into the system sequentially once during the lifelong learning procedure. The GPL system uses symmetric encoder-decoder networks to encode *new observations* to ‘memory codes’ and decode ‘memory codes’ into the *synthetic observations*. In this section, we introduce the design of the encoder, the decoder, and the place feature learning within the GPL system.

A. Place Memory Encoding and Place Feature Extraction

GPL applies the encoder module \mathcal{E} to convert raw sensor observations into the ‘memory codes’ with VGG [56]-based networks, and the memory codes are basic materials in the

BiLM system. Viewpoint differences and environmental appearance changes in lifelong localization will affect the final localization performance in real-world applications. Based on the orientation-equivalent property of spherical harmonics, we utilize the spherical convolution [10], [54] on top of the encoder module to provide viewpoint-invariant feature (also called place descriptor) to reduce the viewpoint differences in long-term re-localization.

The GPL system encodes panorama camera and 3D local point cloud with the same encoding network structure \mathcal{E} . For the visual inputs, we convert the raw image to $[H \times W]$ spherical perspectives. For the LiDAR inputs, instead of a single scan, we generate dense local 3D maps using the similar voxel mapping mechanism in our previous work [57] and map the points onto the spherical projections, which have the same omnidirectional view as a panorama camera. Then the preprocessed (visual or Lidar) inputs q_k are fed into the encoder module for generating memory codes and place features. The memory code z_k are encoded from q_k :

$$z_k = \mathcal{E}(q_k) \quad (3)$$

To extract the viewpoint-invariant place feature from z_k , we utilize the spherical convolution based on the spherical harmonics [58]. In theory, Spherical convolution can avoid space-varying distortions in Euclidean space by convolving spherical signals in the harmonic domain. Let f is the signal on spherical harmonic, which satisfy the viewpoint-equivariant [59] property with the signal \mathcal{E} ,

$$[f \star_{SO(3)} [H_R \mathcal{E}]](q_k) = [H_R [f \star_{SO(3)} \mathcal{E}]](q_k) \quad (4)$$

where $H_R (R \in SO(3))$ is the rotation operator for spherical signals. $f \star_{SO(3)} \mathcal{E}$ denotes the spherical convolution between f and \mathcal{E} . Practically, the spherical convolution is computed in three steps. We first expand f and $\mathcal{E}(q_k)$ to their spherical harmonic basis, then compute the point-wise product of harmonic coefficients, and finally invert the spherical harmonics.

Let V be the unsupervised VLAD layer [27]. Then the feature extraction function $\mathcal{F}_\theta = V \circ [f \star_{SO(3)} \mathcal{E}]$ is the viewpoint-invariant function, where θ are learnable parameters of the feature extraction function. For details about the viewpoint-invariant analysis, please refer to works [10], [54]. Given the data sample q_k , the place feature F_k can be denoted as

$$F_k = \mathcal{F}_\theta(q_k) = V \circ [f \star_{SO(3)} z_k] \quad (5)$$

The above procedure is relevant to the biological ‘encoding’ procedure within the ‘sleep cycle’ as we mentioned section III. The extracted ‘memory codes’ z_k will be used for later ‘memory consolidation’ in next section V and ‘retrieval’ for generative replay in next section IV-B

B. Memory Decoding and Memory Replay

As depicted in Fig. 3, we generate synthetic samples from stored place memory codes by memory decoder, which is called memory replay. In particular, the retrieved memory codes will contain the latent place information under previous domains as shown in Fig. 4, which enforces the memory decoder to extract a portion of history samples under all previous domains to maintain the localization performance. Thus, we can combine the raw real samples with (decoded) synthetic samples to train the place feature extraction function.

To ensure the generalization ability of synthetic samples, we use a deep generative adversarial network (GANs) to mimic the distribution of raw samples when generating synthetic samples. The GANs-based generative model defines a zero-sum minimax game with the memory decoder \mathcal{G} and the discriminator \mathcal{D} as stated in [60], the objective function is thereby defined by,

$$\begin{aligned} \min_{\mathcal{G}} \max_{\mathcal{D}} \mathcal{L}_{gan}(\mathcal{G}, \mathcal{D}) = \\ \min_{\mathcal{G}} \max_{\mathcal{D}} \mathbf{E}_{q \sim P_{data}} [\log \mathcal{D}(q)] + \mathbf{E}_{z' \sim P_z} [\log(1 - \mathcal{D}(\mathcal{G}(z')))] \end{aligned} \quad (6)$$

where P_z is the retrieved memory buffer from the BiLM system, and P_{data} is the new observed data samples. \mathcal{L}_{gan} denotes the generator and discriminator losses [60]. The detailed generative play strategy for lifelong learning can be found in [61], [62].

Given the combination of new raw data $q_k \sim P_{data}$ and retrieved synthetic samples $q_k \sim \mathcal{G}(P_z)$, we can obtain input tuple sets $(q_k, \{o_k^{pos}\}, \{o_k^{neg}\})$, where for each query sample q_k we have a set of potential positives (close-by samples) $\{o_k^{pos}\}$ and the set of negatives (far away samples) $\{o_k^{neg}\}$. The localization loss metric is defined by:

$$\begin{aligned} L_{loc}(q_k) = \\ \max_{i,j} (\|\mathcal{F}(q_k) - \mathcal{F}(o_k^{pos}_i)\|^2 + \alpha - \|\mathcal{F}(q_k) - \mathcal{F}(o_k^{neg}_j)\|^2, 0) \\ \mathcal{L}_{loc} = \mathbf{E}_{q_k \sim P_{data}} [L_{loc}(q_k)] + \mathbf{E}_{z' \sim P_z} [L_{loc}(\mathcal{G}(z'))] \end{aligned} \quad (7)$$

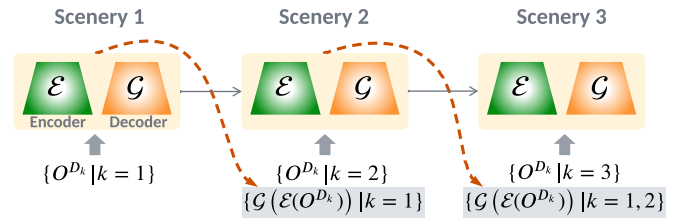


Fig. 4: **Generative Memory Reply within the GPL system.** In the lifelong localization, new observations O^{D_t} under domains D_t will be streamed into the BioSLAM system sequentially. GPL’s generative memory replay module can generate synthetic samples $\mathcal{G}(\mathcal{E}(O^{D_t}))$ from stored memories.

the above equation is the triplet loss version of Eq. (1), where $L_{loc}(q_k)$ is the localization loss metric for single query q_k , ($\alpha > 0$) is a margin to control the feature difference threshold, and \mathcal{L}_{loc} is the localization loss for the joint $\{P_{data}, P_z\}$ sets.

To keep the consistency of memory encoding-decoding, we further use a reconstruction loss between encoder and decoder modules. For the encoded memory code z' and the decoded memory $\mathcal{E}(\mathcal{G}(z'))$, we have

$$\mathcal{L}_{rec} = \mathbf{E}_{z' \sim P_z} [\|\mathcal{E}(\mathcal{G}(z')) - z'\|] \quad (9)$$

The joint loss metric for the generative memory replay enhanced place recognition can be written as,

$$\mathcal{L}_{joint} = \mathcal{L}_{loc} + \mathcal{L}_{rec} + \mathcal{L}_{gan}(\mathcal{G}, \mathcal{D}) \quad (10)$$

The major difference between our work and the traditional generative replay [61] is that BioSLAM can manage the retrieved memory based on their long-term behavior instead of treating all data on the same manifold distribution. The next section will deeply investigate the lifelong memory system.

V. BIO-INSPIRED LIFELONG MEMORY

In our BioSLAM, as shown in the greed dashed box of Fig. 3, the bio-inspired lifelong memory system mainly contains two modules: 1) *Rewarding mechanism* to measure the memory codes’ importance (or reward calculation), the importance score will be used to manage memory consolidation and selection in the dual-memory module. 2) *Dual-memory* module to cooperate with the rewarding mechanism for long-/short- term memory storage and importance-retrieval with limited space usage; which includes static memory zone and dynamic memory zone.

A. Rewarding mechanism

When the memory system encounters a new place ‘memory codes’ z , we define the hybrid reward to control the learning behavior: an external reward \mathcal{R}_{ex} which indicates localization ability, and the internal reward \mathcal{R}_{in} which presents the intrinsic familiarity on observations.

1) *External Reward*: The external reward is related to the learning difficulty of data samples, which indicates their distinguishing ability in the place recognition task. In the standard learning paradigm, all samples with different difficulties are equally considered during the model optimization. However,

humans and animals usually spend more energy and time to learn harder concepts. Inspired by animal training [63] and curriculum learning [64], it is practical to differentiate the data samples into different difficulty levels, such as “easy”, “medium”, and “hard”. For lifelong localization, the feature extraction function \mathcal{F}_θ may require more ‘energy’ or iterations to learn “hard” samples. Hence, we encourage “hard” samples to have a higher chance of re-training by defining the triplet loss to measure the sample’s difficulty and distinguishability. Based on the place recognition loss metric L_{loc} , we define the external reward for each input as follows:

$$\mathcal{R}_{ex}(q_k) = L_{loc}(q_k) \quad (11)$$

which means that if the sample q_k has a higher loss than other samples, it will require more iterations in model training. “Hard” samples tend to have a high value of \mathcal{R}_{ex} , and therefore the dual-memory module selects them for memory replay with a higher probability, as described in section V-B4.

2) *Internal Reward*: The internal reward is related to the robustness of the feature extraction model when given a sample. Let $A(q_k)$ denote the data augmentation (i.e. random rotation and random translation) for sample q_k . The internal reward \mathcal{R}_{in} for sample q_k is defined as the cosine distance between its feature and the feature over data augmentation:

$$\mathcal{R}_{in}(q_k) = 1 - \frac{\mathcal{E}(q_k) \cdot \mathcal{E}(A(q_k))}{\|\mathcal{E}(q_k)\|_2 \cdot \|\mathcal{E}(A(q_k))\|_2} \quad (12)$$

The internal reward \mathcal{R}_{in} also indicates the network’s familiarity with the observations. In large-scale place recognition, similar place patterns (street view, buildings, trees) can be frequently visited with different views; the encoder \mathcal{E} has a robust representation and low internal reward of frequently visited places. Therefore, the internal reward \mathcal{R}_{in} can guide the dual memory system to focus more on unfamiliar areas.

The final reward for q_k can be obtained by combining the external reward with the internal reward,

$$\mathcal{R}_k = \mathcal{R}_{ex}(q_k) + \mathcal{R}_{in}(q_k), \quad (13)$$

Based on this rewarding mechanism, we can evaluate all the queries and obtain a set of memory trace $m_k = (z_k, p_k, \mathcal{R}_k)$, where p_k is the estimated location of q_k obtained from our previous re-localization system [10], [65]. The memory traces m_k are then used as the main factor in memory operations described in section V-B.

B. Dual-Memory & Memory Operations

The memory of human beings is highly connected with long-term memory (the neocortex) and short-term memory (the hippocampus) mechanisms within our brains. BioSLAM also constructs such paired dual-memory mechanisms,

- **Static Memory** M_S is similar to the long-term memory of human beings, which stores the selected memory traces by memory consolidation.
- **Dynamic Memory** M_D is similar to the short-term memory of human beings, which is a quick access memory with a portion of pre-stored historical memory traces.

Dynamic memory is automatically refreshed from the static memory and connected with the memory decoder.

In general, dynamic memory has a smaller buffer size (e.g. 1000), while static memory has a larger buffer size (e.g. 4000). Based on the dual-memory structure, we construct two important operations for static memory: memory consolidation and forgetting, and two important operations for dynamic memory: memory refreshing and memory replay.

As shown in Fig. 3, the running mechanisms of the dual-memory module can be summarized as follows. Consider an agent employing the BioSLAM algorithm, continuously navigating an environment. Upon encountering a new observation, the memory encoder encodes the observation into memory codes. Memory consolidation enables static memory to store the memory code, while forgetting allows static memory to discard unimportant samples and retain significant and diverse ones. Next, dynamic memory obtains selected memories from static memory via memory refreshing and replays them to a memory decoder using memory replay. The memory decoder generates synthetic samples from the replayed memory codes, which BioSLAM employs to update the agent’s place recognition function in combination with the new real samples.

1) *Static Memory Consolidation*: As stated in [66], memory consolidation is defined as a time-dependent process by which recently learned experiences are transformed into long-lasting forms to extend the long-term memory circle. In the long-term and large-scale place recognition, the observations may vary not only in the spatial domain (Euclidean distance) but also in the feature domain (feature distance). Furthermore, real-world navigation tasks involve an unlimited data stream. Memory consolidation is essential to abstract concise representations and guarantees memory efficiency. To provide memory consolidation for static memory, we construct a feature-spatial code $c_k = [z_k, p_k]$ for memory trace m_k , which can capture both spatial and feature properties.

At time step T and observations $q_k^T \subset Q^{D_T}$, the obtained new memory traces $m_k^T = (z_k, p_k, \mathcal{R}_k)^T$ typically consist of a large number of samples. To obtain a diverse and smaller subset of memory traces (abstraction), we use K-means-based unsupervised clustering. K-means clustering partition the $\{m_k\}^T$ into K sets $\mathbf{S}^T = \{S_1^T, S_2^T, \dots, S_K^T\}$ by feature-spatial code $c_k = [z_k, p_k]$ to minimize the following,

$$\mathbf{S}^T = \arg \min_{\mathbf{S}^T} \sum_{i=1}^k \frac{1}{|S_i^T|} \sum_{c_x, c_y \in S_i^T} \|c_x - c_y\|^2 \quad (14)$$

$$\mu_i^T = \frac{1}{|S_i^T|} \sum_{c_i \in S_i^T} c_i$$

where μ_i^T is the cluster centroid for cluster S_i^T . The memory traces within a single cluster have similar characteristics and storing all of them is redundant and memory-intensive. To optimize memory usage and retrieval efficiency, downsampling is applied within each cluster S_i^T :

$$(\tilde{S}_i^T, \tilde{\mu}_i^T) = \text{downsampling}(S_i^T, \mu_i^T), \quad |\tilde{S}_i^T| < N_{max} \quad (15)$$

Where N_{max} is the (predefined) maximum number of samples in each cluster. After the sampling process, a set of new

clusters $\tilde{\mathbf{S}}^T = \{\tilde{S}_1^T, \tilde{S}_2^T, \dots, \tilde{S}_K^T\}$ and their corresponding centroids $\tilde{\mu}^T = \{\tilde{\mu}_1^T, \tilde{\mu}_2^T, \dots, \tilde{\mu}_K^T\}$ are generated from a subset of the traces $\{m_k\}^T$.

BioSLAM can integrate the newly generated clusters $\tilde{\mathbf{S}}^T$ with the existing clusters from previous steps, resulting in a total of $\mathbf{S}^{(M_s)} = \{\tilde{\mathbf{S}}^1, \tilde{\mathbf{S}}^2, \dots, \tilde{\mathbf{S}}^T\}$ clusters in static memory with corresponding centroids $\mu^{(M_s)} = \{\tilde{\mu}^1, \tilde{\mu}^2, \dots, \tilde{\mu}^T\}$. If the total number of clusters $C^{(M_s)} = |\mu^{(M_s)}|$ exceeds the maximum threshold K_{max} , similar clusters are merged to prevent memory overflow, as described in the Memory Forgetting mechanism section V-B2. The consolidation mechanism is shown in Algorithm 1.

Algorithm 1: Memory Consolidation

Input: Static memory M_S , new memory traces $\{m_k\}$, maximum number of clusters K_{max}

Output: Updated static memory

- 1 Construct feature-spatial codes $\{c_k\} = \{z_k, p_k | m_k\}$;
 - 2 Calculate clusters $\{S_i^T\}_{i=1}^K$ and centroids $\{\mu_i^T\}_{i=1}^K$ for $\{c_k\}$ based on Eq. (14);
 - 3 Downsample within clusters, based on Eq. (15), to generate smaller clusters $\{\tilde{S}_i^T\}_{i=1}^K$ and centroids $\{\tilde{\mu}_i^T\}_{i=1}^K$;
 - 4 Append new clusters $\{\tilde{S}_i^T\}_{i=1}^K$ and centroids $\{\tilde{\mu}_i^T\}_{i=1}^K$ to M_S ;
 - 5 Calculate the total cluster number $C^{(M_s)}$ in M_S ;
 - 6 **if** $C^{(M_s)} > K_{max}$ **then**
 - 7 Memory Forgetting with Algorithm 2;
 - 8 Update static memory M_S ;
-

2) *Static Memory Forgetting*: As mentioned in the previous section, the space allocated to static memory is limited in long-term lifelong learning. Memory forgetting, a crucial operation in static memory, is implemented to eliminate redundant memory clusters that are too similar to the existing ones. If the number of current clusters exceeds the maximum limit of clusters $C^{(M_s)} > K_{max}$, the memory forgetting mechanism removes $K^* = C^{(M_s)} - K_{max}$ clusters. To accomplish this, we first calculate the cluster similarity based on the distance matrix $d_{(i,j)}$ between every pair of cluster centroids.

$$d_{(i,j)} = \|\mu_i - \mu_j\|, \quad \forall \mu_i, \mu_j \in \mu^{(M_s)} \quad (16)$$

We search for cluster pairs (i^*, j^*) with the smallest distance between their centroids and remove one of the clusters i^* from each pair. This removal procedure is repeated K^* times.

$$(i^*, j^*) = \arg \min_{i,j} d_{(i,j)} \quad (17)$$

This method allows us to maintain a diverse set of memory clusters while removing any redundant clusters. The memory forgetting mechanism is presented in Algorithm 2.

3) *Dynamic Memory Refreshing*: Dynamic memory is similar to the short-term memory of humans, brief and storage-limited. To effectively replay important memory traces from dynamic memory, we need to refresh it periodically and clone relevant memory traces from static memory to dynamic memory. This is done through memory refreshing mechanisms,

Algorithm 2: Memory Forgetting

Input: Static memory M_S , maximum number of clusters K_{max}

- 1 Load clusters $\mathbf{S}^{(M_s)}$ and centroids $\mu^{(M_s)}$ from static memory M_S ;
 - 2 Calculate the number of forgettable clusters $K^* = |\mu^{(M_s)}| - K_{max}$;
 - 3 Calculate the distance matrix $d_{(i,j)}$ between every pair of cluster centroids on Eq. (16);
 - 4 **while** repeat K^* times **do**
 - 5 Find most similar cluster pairs (i^*, j^*) based on Eq. (17);
 - 6 Remove cluster i^* from static memory M_S and distance matrix $d_{(i,j)}$;
-

where dynamic memory M_d obtains memory traces $\{m_k\}$ from static memory M_s through importance sampling.

$$M_d = \text{importance_sampling}(\{m_k\}, \{w_k\}) \quad (18)$$

$$m_k = (z_k, p_k, \mathcal{R}_k) \sim M_s, \quad w_k = \gamma^{n(m_k)} \cdot \mathcal{R}_k$$

where importance weights w_k are determined by the reward \mathcal{R}_k and the time-decaying factor $\gamma^{n(m_k)}$. Here, γ ($0 \leq \gamma \leq 1$) is a predefined decay parameter, and $n(m_k)$ denotes the replayed time (or revisited time) for the trace m_k . Traces with higher rewards are assigned higher sampling weights as they have lower localization ability and robustness, requiring more attention from BioSLAM. Moreover, new traces are assigned higher sampling weights since the network's ability to learn samples with many occurrences has reached an upper limit, and storing samples replayed many times is unnecessary. The decaying mechanisms also encourage dynamic memory to be more curious about new traces. These reward decay mechanisms are inspired by the decaying factor in human memory [67], which suggests that repeated learning of the same things decreases the boost in memorization.

4) *Dynamic Memory Replay*: During lifelong learning, dynamic memory replays selected memory traces to the memory decoder for generating synthetic samples, as stated in section IV-B. Importance sampling is used to obtain replayed memories $\{z'_k\}$ from dynamic memory M_d in the same manner as the refreshing memory mechanisms.

$$\{z'_k\} = \text{importance_sampling}(\{z_k\}, \{w_k\}) \quad (19)$$

$$m_k = (z_k, p_k, \mathcal{R}_k) \sim M_d, \quad w_k = \gamma^{n(m_k)} \cdot \mathcal{R}_k$$

We use memory decoder \mathcal{G} to generate replayed samples $\{\hat{q}_k\}$ from memories $\{z'_k\}$,

$$\hat{q}_k = \mathcal{G}(z'_k) \quad (20)$$

Both new observations $\{q_k\}$ and generated samples $\{\hat{q}_k\}$ are used to train the General Place Learner (GPL) network by minimizing the total loss Eq. (10). The overall lifelong learning algorithm of BioSLAM is shown in Algorithm 3.

VI. EXPERIMENT SETUP AND CRITERIA

In this section, we will introduce the experimental setup for lifelong localization. Unlike traditional localization tasks,

Algorithm 3: Lifelong Learning with BioSLAM

Input: Initial place feature extraction model \mathcal{F}_θ with parameters θ , Initial static memory $M_s = \emptyset$ and dynamic memory $M_d = \emptyset$

```

1 for  $T = 1, 2, \dots$  do
2   Obtain observation set  $\{q_k\}$ ;
3   while repeat until converge do
4     Generate replayed samples  $\{\hat{q}_k\}$  from dynamic
      memory based on Eq. (19) and (20);
5     Calculate loss  $\mathcal{L}_{joint}$  using real samples  $\{q_k\}$ 
      and replayed samples  $\{\hat{q}_k\}$  based on Eq. (10);
6     Calculate gradient  $\frac{d\mathcal{L}_{joint}}{d\theta}$  then optimize  $\mathcal{F}_\theta$ 
      with parameters  $\theta$  by gradient descend;
7   Calculate rewards  $\{\mathcal{R}_k\}$  of observations  $\{q_k\}$ 
      based on Eq. (13);
8   Static memory  $M_s$  consolidation based on
      Algorithm 1;
9   Dynamic memory  $M_d$  refreshing based on Eq. (18)

```

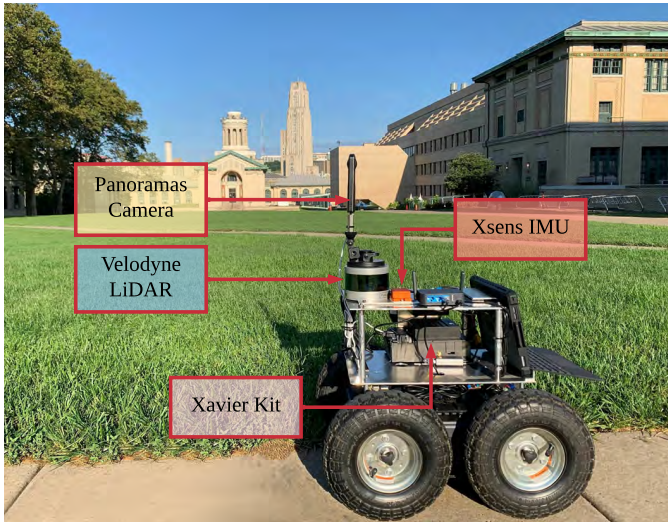


Fig. 5: **Data Collection Platform.** The platform can record the omnidirectional visual inputs, Velodyne VLP-16 LiDAR inputs, and Xsens MTI IMU data on an Nvidia Jetson AGX Xavier. We utilize the LiDAR odometry [68] to generate the relative odometry for each trajectory and GNSS or Generalized-ICP [69] to estimate the relative transformation between different trajectories.

lifelong localization requires recorded data that includes either long-term differences or large-scale geometric differences. For these reasons, we built our own data collection platform and created our own lifelong localization datasets, including the ALITA Urban dataset and the ALITA Campus dataset [70]. Then we evaluated the performance of BioSLAM on the official Oxford RobotCar dataset.

A. Data Collection Platform

Fig. 5 shows our data collection platform, which includes an omnidirectional camera, a Velodyne VLP-16 LiDAR device, an inertial measurement unit (Xsense MTI 30, 0.5° error

TABLE I: Comparison between different datasets.

Dataset	Domains	Scales (km)
ALITA Urban	Areas: Street, Residential, Terrain	120×1
ALITA Campus	Input Modality: Lidar, Visual	4.5×8
Oxford RobotCar	Weather and Road Conditions	≈ 10.0

in roll/pitch, 1° error in yaw, $550mW$), and an embedded GPU device (Nvidia Xavier, 8G memory). To collect time-synced LiDAR projection and omnidirectional images, we first generated dense 3D maps through well-known LiDAR odometry [68]. Then project the point cloud within a certain distance (default is $30m$) to the spherical projections, which have the same perspective as the omnidirectional images. We will revisit the same area under large-scale and long-term assumptions in lifelong localization. To provide the relative ground truth position between different visits: to outdoor environments, we rely on the GNSS system and Generalize-ICP [69] to estimate the relative transformation; For indoor environments, we mainly rely on Generalize-ICP. Please note that we cannot guarantee the meter-level global absolute localization, but we can provide accurate relative localization, which is enough for the lifelong localization task. Based on the collected datasets, we have hosted a General Place Recognition Competition for long-term place recognition. For more details on the data collection platform and the datasets, please refer to our dataset paper [70] and competition site (<http://gprcompetition.com/>).

B. Lifelong Localization Datasets and Learning Settings

We intend to analyze the lifelong learning performance under *large-scale*, *multi-modal*, and *long-term* three perspectives. To this end, our localization datasets include three tracks:

- *ALITA Urban* dataset: as shown in Fig. 6, is targeted towards large-scale lifelong learning performance. We collected 50 trajectories within the city of Pittsburgh, focusing only on LiDAR inputs within a short-term drive, as our main concern is large-scale localization. The total trajectory distance for this dataset is 110 km. 729 query and database frames are selected from 11 trajectories to construct the test set.
- *ALITA Campus* dataset: as shown in Fig. 7, is targeted towards multi-modal lifelong learning performance. Both LiDAR inputs and visual inputs are collected, and we picked up 10 trajectories within Carnegie Mellon University. Each trajectory is revisited 8 times under different day and night-time conditions to meet the long-term requirements. The test set contains 739 query and database frames.
- *Oxford RobotCar* dataset [71]: is a widely used public dataset aimed at long-term lifelong learning performance. This dataset covers over 1000 km of driving from May 2014 to December 2015, providing long-term observations. As our main concern is long-term localization, we use a subset of the Oxford RobotCar dataset that includes long-term behavior and different weather conditions.

We use the lifelong learning procedure depicted in Fig. 2 to feed the sequential data stream for all datasets. In our experiments, a feature extraction function \mathcal{F}_θ continually learns from different domains (D_1, \dots, D_t, \dots) , where the domain of each dataset is different. For the ALITA Urban dataset, domains are different areas, such as commercial buildings,

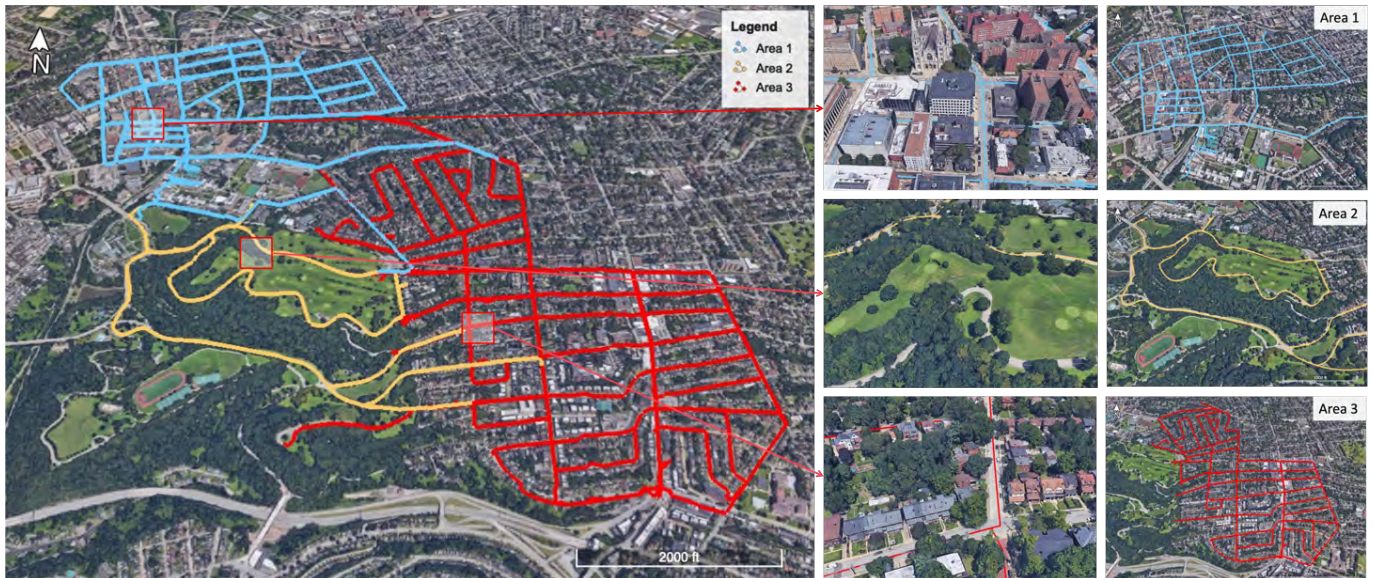


Fig. 6: **ALITA Urban Dataset.** The dataset includes 50 trajectories (110 km) within the city of Pittsburgh. The dataset includes three areas (colored in blue, yellow, and red) covering commercial buildings, parks, and residential areas.

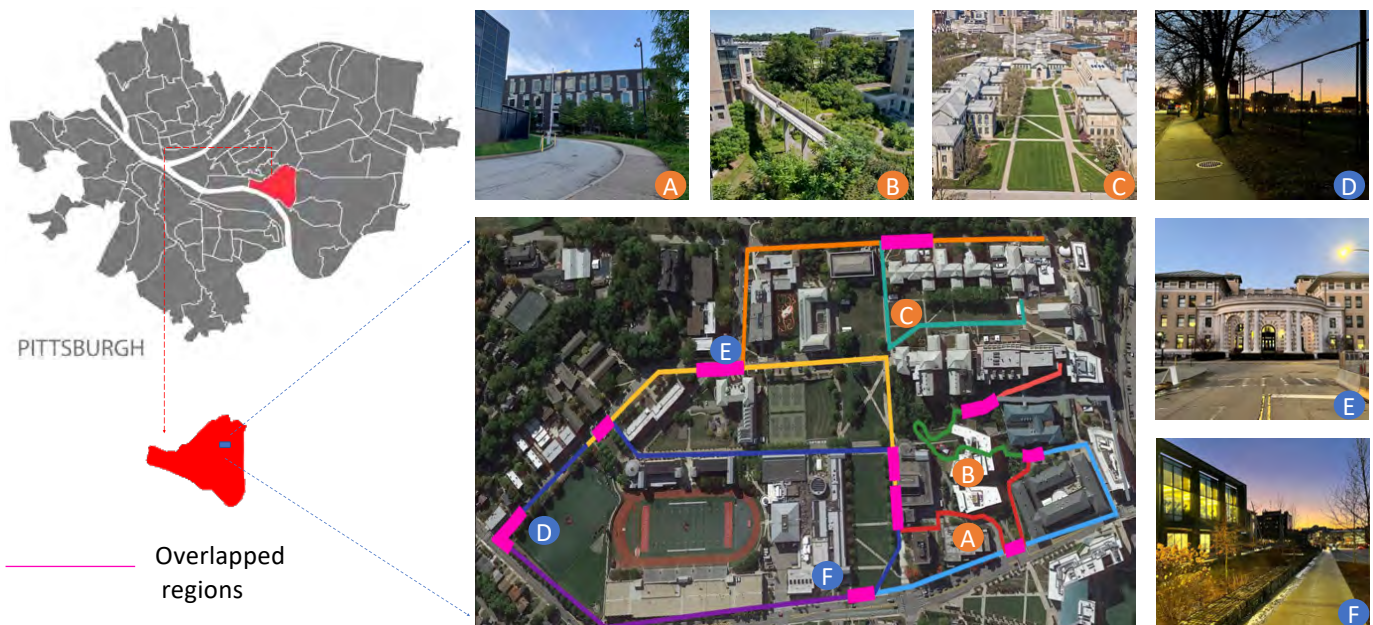


Fig. 7: **ALITA Campus Dataset.** For the dataset, omnidirectional camera and LiDAR data are recorded for 2D-to-2D and 2D-to-3D place recognition within CMU. The datasets are generated during 08/2021 ~ 10/2021, which are mainly taken from normal day-light (2pm ~ 5pm) and dawn-light (5am ~ 6am or 7pm ~ 8pm).

parcs, and residential areas. For the ALITA Campus dataset, domains are different input modalities, such as LiDAR input, daylight visual input, and night light visual input. For the Oxford RobotCar dataset, domains are different weather and road conditions, such as sunny days, daylight with roadworks, night, snowy days, cloudy days, and rainy days. It's worth noting that each data sample is only fed into the system once, and BioSLAM does not save a copy of that data. The comparison between different datasets is presented in Table I.

C. Performance Evaluation

During lifelong place learning, we mainly evaluate the online localization performance through the *Weighted Recall (WR)* of top-6 retrievals over lifelong learning, which is defined by $WR = \sum_{k=1}^6 \omega_k r_k$, $\omega_1 = 0.5, \omega_k = 0.1$ for $k \neq 1$. r_k ($1 \leq k \leq 6$) denotes Recall@k, which measures the percentage of correctly localized queries using top-k elements returned from the database. To be considered a correct match, the query image must have at least one of the top k retrieved reference images within the predefined neighbor range (e.g.,

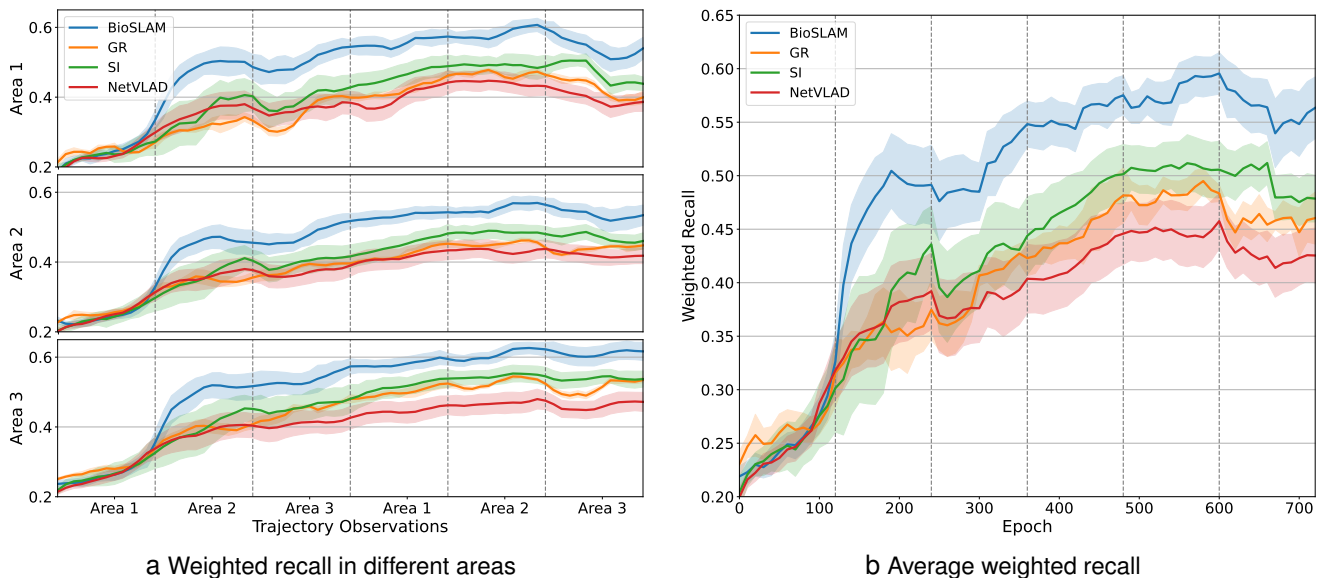


Fig. 8: **Comparison of weighted recall w.r.t. training epochs on ALITA Urban dataset.** Place recognition methods incrementally trained on trajectory observations from 3 different areas. The shaded region shows the standard deviation.

3m) from the ground truth position.

In lifelong place recognition, our focus is on the WR curve during incremental learning by continually feeding a data stream from different domains, as opposed to evaluating the WR on fixed place recognition models for descriptors, as is the case with classical place recognition. We also introduce unique matrices for evaluating the performance of lifelong place recognition, namely Adaptation Efficiency (AE) for measuring fast learning ability in the current new domain, and Retention Ability (RA) for measuring the ability to memorize without forgetting the previous domains when training samples from previous domains are no longer available. When optimizing the feature extraction function \mathcal{F}_θ in the current domain D_T , the adaptation efficiency is defined as the WR on observations and queries from the current domain $\{O^{D_T}, Q^{D_T}\}$. Meanwhile, the retention ability is defined as the WR on observations and queries from the previous domains $\{O^{D_t}, Q^{D_t}\}_{t=1, \dots, T-1}$.

D. Baselines

As our place recognition task involves different sensor modalities, non-learning/learning methods, and lifelong/non-lifelong methods, it is not feasible to cover all the relevant state-of-the-art methods. Therefore, we focus on performance comparisons from a 2D perspective and exclude Point-like 3D methods [72]. To compare our proposed method, we select well-known non-learning methods (Bag-of-wards (BOW) [73] and CoHOG [74]), learning-based methods (NetVLAD [27] and RegionVLAD [29]) and lifelong-based methods (Generative Replay (GR) [75] and Synaptic Intelligence (SI) [46]). All of the learning-based baselines and our proposed method were incrementally trained on the same sequential data during the experiments. Among the above methods, GR and SI are the most related and essential baselines to BioSLAM. Although both BioSLAM and GR use memory replay, BioSLAM has

more efficient and reasonable memory replay mechanisms. This is due to two reasons: first, BioSLAM replays samples according to their reward (importance), while GR replays them randomly and evenly. Second, BioSLAM has static memory, which refreshes the dynamic memory buffer to keep diverse and important memory traces, making it easier to adapt to new trajectory observations.

VII. EXPERIMENT ANALYSIS

In this section, we analyze the lifelong place recognition results on large-scale *City* areas, multi-modal *Campus* scenarios, and long-term *Oxford RobotCar* datasets.

A. Large-scale City Place Recognition

In the localization task under city-scale environments, robots may encounter various 3D geometric structures, such as open streets, bridges, parks, and residential areas. We evaluated the performance of BioSLAM in a large-scale lifelong learning scenario using the ALITA Urban dataset. We divided the 50 trajectories covering a distance of 120km within the city into three different areas based on their geometric properties: area 1 for commercial buildings, area 2 for parks, and area 3 for residential districts. BioSLAM and the baseline methods learn the place features incrementally by continually feeding the trajectory observations from different areas.

Fig. 8a illustrates the weighted recall curve of trajectory observations within area 1, area 2, and area 3, respectively. The vertical dotted line indicates the epoch when the trajectory observations switched from one area to another, e.g., the training observations switched from area 1 to area 2 at epoch 120. Fig. 8b presents the average weighted recall curve of all trajectories during training. It is evident that BioSLAM outperforms other methods during training and achieves at least a 14% improvement in terms of final average recall. When the

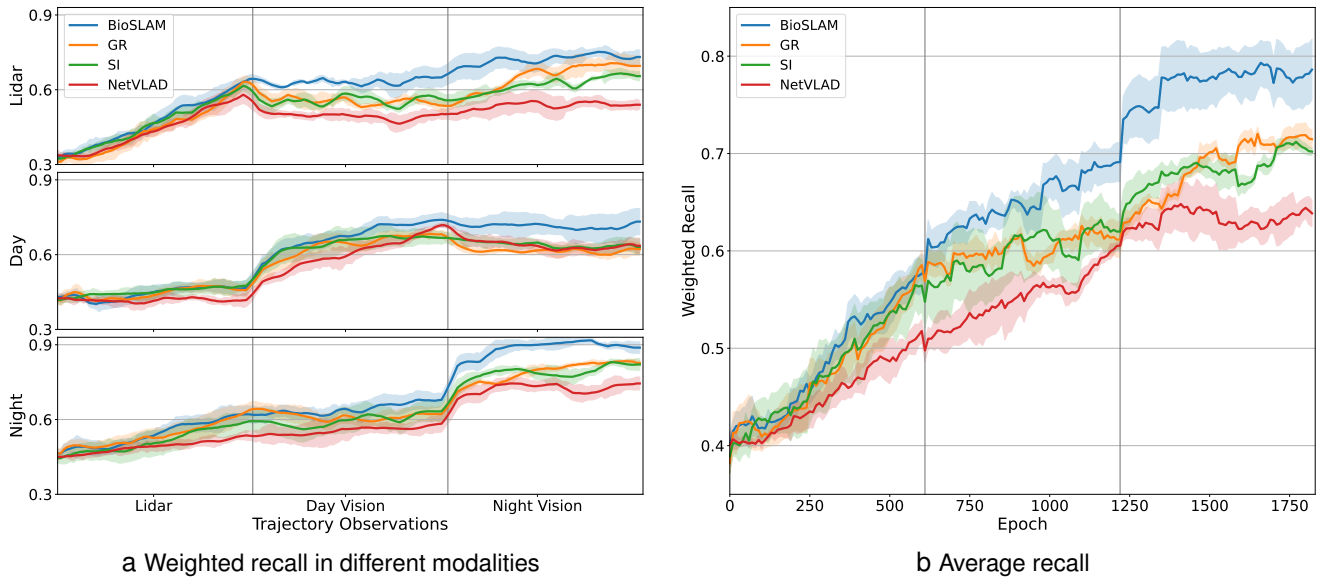


Fig. 9: **Comparison of weighted recall w.r.t. training epochs on ALITA Campus dataset.** Place recognition methods incrementally trained on trajectories from Lidar, day-time visual, and night-time visual inputs.

training observations switched from area 2 to area 3 at epoch 240, BioSLAM’s performance drop on previous trajectories is much smaller than that of other methods, as shown in Fig. 8b. This is because BioSLAM retrains necessary previous knowledge by replaying related memory traces. While GR only replays randomly, BioSLAM replays essential and highly rewarded memory traces, leading to a higher convergence rate and better final performance than other baselines.

B. Multi-modal Campus Place Recognition

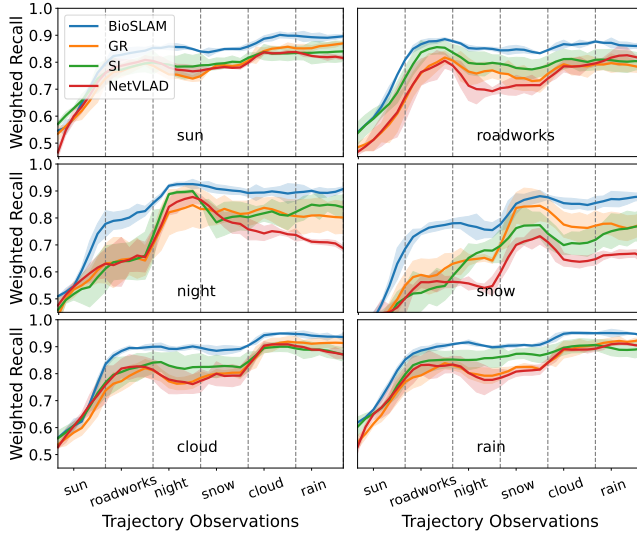
Training place recognition models separately on independent modalities, such as Lidar and visual, can be inefficient as no information is shared between them. In this study, we demonstrate the effectiveness of BioSLAM in more practical settings where the model benefits from solving place recognition using multiple modalities, such as Lidar, day-light vision, and night-light vision. Firstly, knowledge gained from one modality can help to better and more quickly understand other modalities since the modalities are not completely independent in place recognition. Secondly, generalization across multiple modalities may lead to the acquisition of more universal knowledge that applies to unseen modalities, a phenomenon that is also observed in infants’ learning [76], [77]. We evaluate the performance of BioSLAM in multi-modal lifelong learning scenarios using the ALITA Campus dataset. Both BioSLAM and the baseline methods incrementally learn the place features from different modalities (Lidar, day-light vision, and night-light vision) by continually feeding trajectory observations from each modality.

Fig. 9a presents a performance comparison between BioSLAM and baselines on different modalities of the ALITA Campus dataset. The vertical dotted line indicates the epoch when the trajectory observations switched from one modality to another, e.g., the training observations switched from Lidar

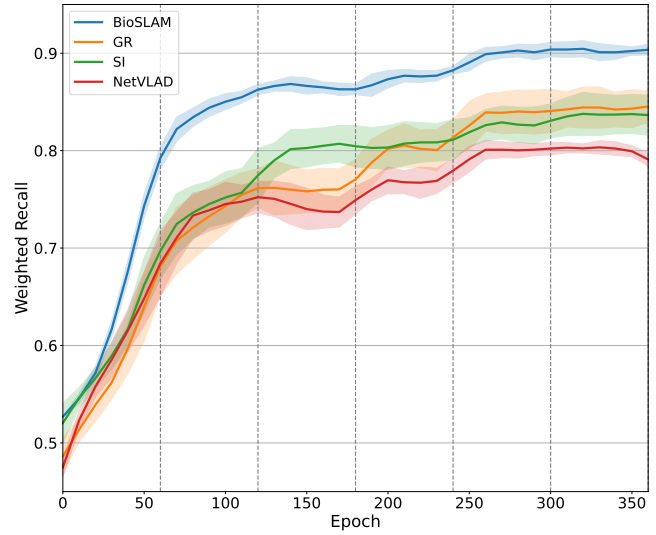
to day-light visual image at epoch 600. During the initial 0 ~ 600 epochs, we train the place recognition model on the Lidar inputs, and all methods show performance improvement across all modalities. This result supports the idea that the knowledge gained from one modality can aid in better and faster understanding of other modalities. Between 600 ~ 1200 epochs, we train the model on daylight visual images. The performance of all methods improves in daylight visual images, but the Lidar performance drops around epoch 600 due to the significant difference in the Lidar and visual images. However, BioSLAM exhibits a smaller performance drop than other methods, indicating its ability to learn new observations in a new modality without forgetting previous modalities. Between 1200 and 1800 epochs, we train the model on night-light visual images. Here, we observe an increase in performance in night-light visual images for all methods, and BioSLAM also shows increased performance in Lidar, thanks to its efficient replay mechanisms. Fig. 9b displays the comparison of average weighted recall between BioSLAM and the baselines on different modalities of the ALITA Campus dataset. Initially, the performance of all methods is comparable. However, as new observations from new modalities are added, BioSLAM exhibits faster and better convergence compared to the baselines. Moreover, BioSLAM outperforms the other methods by 10% in terms of final average recall.

C. Long-term Place Recognition

In long-term localization, robots need to navigate through various weather and road conditions such as heavy rain, night, direct sunlight, snow, and also account for changes in road and building structures over time. To evaluate BioSLAM’s performance in long-term lifelong learning, we use the Oxford RobotCar dataset. We use a subset of the original Oxford RobotCar dataset that includes diverse weather and road



a Weighted recall in different weather and road conditions



b Average recall

Fig. 10: **Comparison of weighted recall w.r.t. training epochs on Oxford dataset.** Place recognition methods incrementally trained on trajectories from different weather and road conditions, such as sun, roadworks, night, snow, cloud, and rain.

conditions such as sun, roadworks, night, snow, cloud, and rain. BioSLAM and the baseline methods incrementally learn the place features by continually feeding the trajectory observations captured in various weather and road conditions.

Fig. 10a depicts the weighted recall curve of trajectory observations captured under various weather and road conditions. The vertical dotted line indicates the epoch when the trajectory observations switched from one condition to another, e.g., the training observations switched from sun-light vision to roadworks at epoch 60. As shown, BioSLAM addresses catastrophic forgetting by retraining necessary previous knowledge and replaying related memory traces. When the observations switch from night light to snowy day at epoch 120, BioSLAM’s performance drop on previous trajectories (night light) is significantly smaller than other methods. A similar phenomenon is observed when the observations switch from snow to cloud at epoch 160.

Fig. 10b shows the average weighted recall curve of all trajectories during training. As demonstrated, BioSLAM outperforms other methods during training and achieves at least an 8% improvement in final average recall compared to different baselines. Although both BioSLAM and GR utilize memory mechanisms to recall previous samples, BioSLAM exhibits more efficient learning due to its rewarding and memory selection mechanisms. The comparison between BioSLAM and GR in different buffer sizes (for dynamic memory) on the Oxford RobotCar dataset is shown in Fig. 11. It is evident that BioSLAM consistently outperforms GR across various buffer sizes, showcasing its superior efficiency, especially in small memory settings.

D. Adaptation, Retention and Generalization Ability

To clearly show the lifelong property of BioSLAM, in this section, we will show the Adaptation Efficiency (AE)

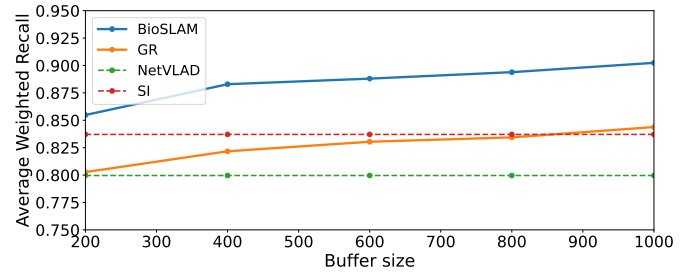


Fig. 11: **Comparison of BioSLAM and GR in different buffer sizes on Oxford RobotCar dataset.**

and Retention Ability (RA) of ALITA Urban (large-scale localization) and Oxford RobotCar (long-term localization) datasets. AR is measured as the WR in the current domain. RA is measured as the WR on the previous domains.

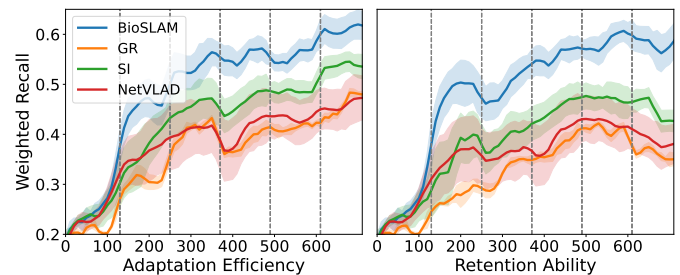


Fig. 12: **Comparison of BioSLAM and baselines in AE (left) and RA (right) on ALITA Urban dataset.**

Fig. 12 illustrates the performance of different methods on the ALITA Urban dataset in terms of AE and RA, while Fig. 13 shows the same metrics for the Oxford RobotCar dataset. The AE metric represents the fast adaptation ability to a new domain, and BioSLAM achieves the highest AE

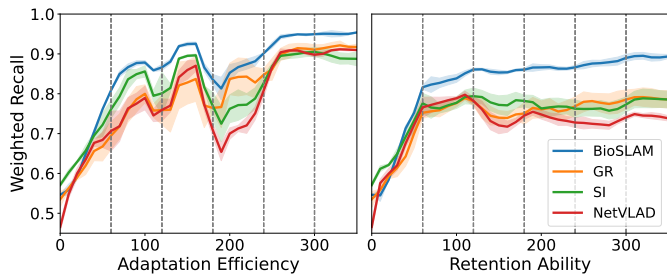


Fig. 13: Comparison of BioSLAM and baselines in AE (left) and RA (right) on Oxford RobotCar dataset.

compared to the baselines on both datasets, demonstrating its fast-learning capability. On the other hand, RA represents the memorizing ability without forgetting previous domains. BioSLAM surpasses the other methods in terms of RA, thanks to its efficient rewarding and dual-memory mechanisms. In the Oxford dataset, the RA of BioSLAM almost monotonically increases, whereas other methods can exhibit a decrease in some domains (e.g., epoch 120-180 for training with night-light observations). BioSLAM’s RA remains stable as it leverages its efficient replay mechanisms to replay important samples from previous domains such as sun and roadworks to overcome catastrophic forgetting.

In addition to the lifelong learning metrics AR and RA, we also evaluate the generalization ability of the final trained model on a fixed test set, similar to classic place recognition tasks. Fig. 14 displays the comparison between BioSLAM and other baselines in terms of top- k recall on the test set of the ALITA Urban dataset. Despite some non-lifelong learning methods being designed or trained for offline evaluation, BioSLAM still outperforms baselines on classic (non-lifelong learning) offline test set evaluation. We then report the generalization ability of the final trained model on the fixed test set in terms of weighted recall, as illustrated in Table II.

On the ALITA Urban dataset, BioSLAM outperforms a state-of-the-art lifelong learning method GR by 7.6%. On the ALITA Campus dataset, BioSLAM outperforms a lifelong learning method GR by 15.6%. Note that this paper primarily focuses on incremental and lifelong learning scenarios, and thus the most crucial evaluation metric is the recall curve over incremental learning (as shown in Fig. 8 and 9). While some non-lifelong learning methods (i.e., CoHOG) may perform well on recall for a fixed test set, these methods cannot learn incrementally, and therefore their performance is limited. Consequently, lifelong learning methods have a higher potential in wider and dynamic real-world environments.

E. Ablation Study

As outlined in section V, BioSLAM stands out from previous lifelong learning methods due to several novel mechanisms, including (1) external reward \mathcal{R}_{ex} to indicate localization performance; (2) internal reward \mathcal{R}_{in} to indicate the robustness of feature representation; (3) Static memory consolidation to abstract concise memory traces, where clustering (Eq. (14)) is a key component; (4) Dynamic memory refreshing to effectively replay diverse and important memo-

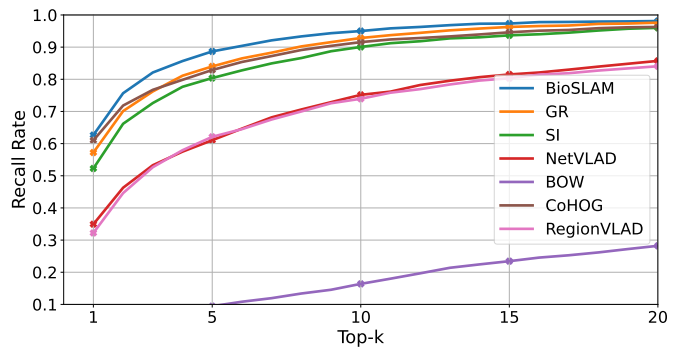


Fig. 14: Comparison of BioSLAM and baselines in terms of recall@ k on ALITA Urban dataset.

TABLE II: Comparison of weighted recall of trained model and fixed test set on ALITA Urban, ALITA Campus, and Oxford RobotCar datasets. We use **boldface** for highlighting the best results.

Weighted Recall (%)		Urban	Campus	Oxford
Non-learning	BOW	5.7	60.1	63.2
	CoHOG	70.1	85.1	90.3
Learning based (not lifelong)	RegionVLAD	45.3	75.1	81.4
	NetVLAD	47.8	72.2	82.9
Lifelong learning	SI	65.1	73.7	85.7
	GR	68.4	76.1	87.9
	BioSLAM	73.6	90.2	94.2

ries, with the time-decay mechanism (Eq. (18)) for importance weight being critical. To evaluate the effectiveness of these mechanisms, we compared BioSLAM with the following variants: (1) w/o \mathcal{R}_{ex} , which does not apply external reward and leads to $\mathcal{R}_k = \mathcal{R}_{in}(q_k)$; (2) w/o \mathcal{R}_{in} , which does not apply internal reward and results in $\mathcal{R}_k = \mathcal{R}_{ex}(q_k)$; (3) w/o consolidation-clustering, which does not use clustering in static memory consolidation, resulting in Algorithm 1 directly storing all memory traces in static memory; (4) w/o time-decay, which does not use the time decay factor in importance sampling, which is equivalent to setting $\gamma = 1$ in dynamic memory refresh Eq. (18). Note that these variants follow the control variates method, covering all important mechanisms of BioSLAM without overlapping functionalities.

TABLE III: Ablation study of BioSLAM. We use **boldface** for highlighting the best results.

Weighted Recall \ Dataset	Urban	Campus	Oxford
BioSLAM	59.8	79.0	90.2
w/o \mathcal{R}_{ex}	53.6	69.7	81.7
w/o \mathcal{R}_{in}	47.9	68.2	80.4
w/o consolidation-clustering	56.1	74.2	86.4
w/o time-decay	50.9	71.2	83.1

Section VII-E presents the results of the ablation study conducted on the Urban dataset. It is evident that BioSLAM outperforms its variants, and the removal of any of its components leads to a significant drop in performance. To facilitate a clear comparison of the results obtained from different variants, we present the ablation study conducted on different datasets in terms of the final average recall during lifelong learning, as depicted in Table III. Notably, the larger drop in performance is observed when removing internal rewards, highlighting the

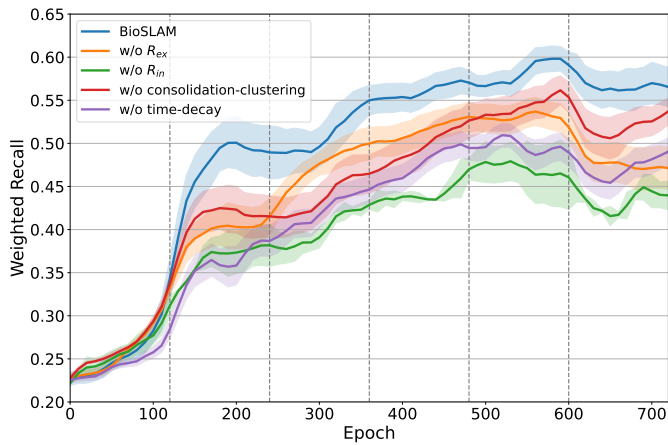


Fig. 15: **Ablation study of BioSLAM.** Comparison between BioSLAM and its variants on the ALITA Urban dataset.

significance of the internal reward as an indicator of feature representation robustness for retrieving memories. While the performance of the “w/o cluster-consolidation” variant is close to that of BioSLAM, the latter is more memory-efficient by clustering and downsampling.

F. BioSLAM Feature Property

In this section, we analyze the learned features of BioSLAM, denoted by $\mathcal{F}_\theta(q_k)$, using similarity matrix and Principle Component Analysis (PCA) on ALITA Urban and Campus datasets. The similarity matrix M_{sim} is computed by taking the cosine similarity between the reference O^D and query Q^D features, where $M_{\text{sim}}(i, j) = \cos(\mathcal{F}(O_i^D), \mathcal{F}(Q_j^D))$. A high-contrast similarity matrix indicates that the learned feature representation, \mathcal{F} , has a strong ability to express and discriminate between places, as evidenced by high similarity values for similar places and low values for dissimilar places.

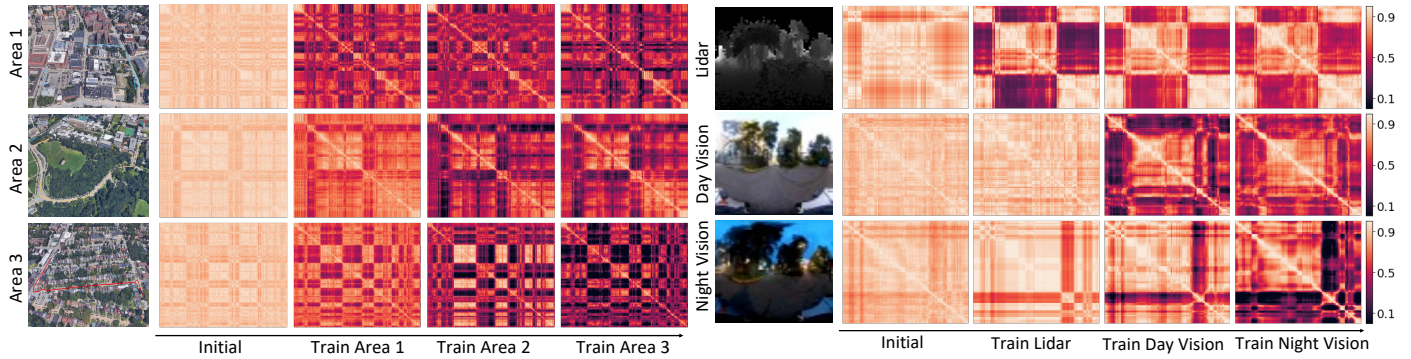
The similarity matrix of BioSLAM during lifelong learning on the ALITA Urban dataset is presented in Fig. 16a, where the left column illustrates the sampled trajectories from area 1, area 2, and area 3, respectively. The three right columns display the similarity matrices of the corresponding trajectories (from left to right) after incrementally learning from area 1, area 2, and area 3, respectively. Notably, after learning from a new area (e.g., area 2), the similarity matrix of the previous area (e.g., area 1) shows almost no decay. This observation indicates that BioSLAM still has a strong expression ability on past trajectories even when learning from different areas. Fig. 16b presents the similarity matrix of BioSLAM on the ALITA Campus dataset. The left column shows the sampled observations from Lidar, day-light visual image, and night-light visual image. The three right columns display the similarity matrices after incremental learning from each modality. Notably, the similarity matrix of the previous modality (e.g., Lidar) shows almost no decay after learning from a new modality (e.g., day-light visual image). This indicates that BioSLAM can remember past modalities when learning from different new modalities.

We utilized PCA to project the BioSLAM learned features into a 2D space for visualization. The PCA visualization of observations from different areas on the Urban dataset is displayed in Fig. 17a, with sub-figures from left to right representing the initial step and incremental learning from area 1, area 2, and area 3, respectively. The results show that BioSLAM training enables the clustering of observations within the same area and facilitates discrimination between clusters of different areas during incremental learning. Similarly, the PCA visualization of learned features from different modalities and trajectories on the ALITA Campus dataset is presented in Fig. 17b, with sub-figures from left to right representing the initial step and incremental training on Lidar, day-light visual, and night-light visual images. Notably, BioSLAM can differentiate not only different modalities (as seen from the three clusters from left to right) but also different trajectories within each modality. Moreover, BioSLAM can find cross-domain relationships between place observations from different modalities. For instance, the PCA results of Lidar and night-time visual domains for trajectory 1 are relatively close and located in the lower part of the PCA visualization results.

G. BioSLAM Memory Activity

In this section, we provide a visualization of the memory buffer during lifelong learning on the Oxford RobotCar dataset. For this experiment, we used a buffer size of 1000 samples for the dynamic memory (DM) and 5000 samples for the static memory (SM). For SM, we store the samples on the hard disk since it requires a large capacity and is not accessed frequently. For DM, we store the data in RAM because it requires frequent access and does not need a large capacity.

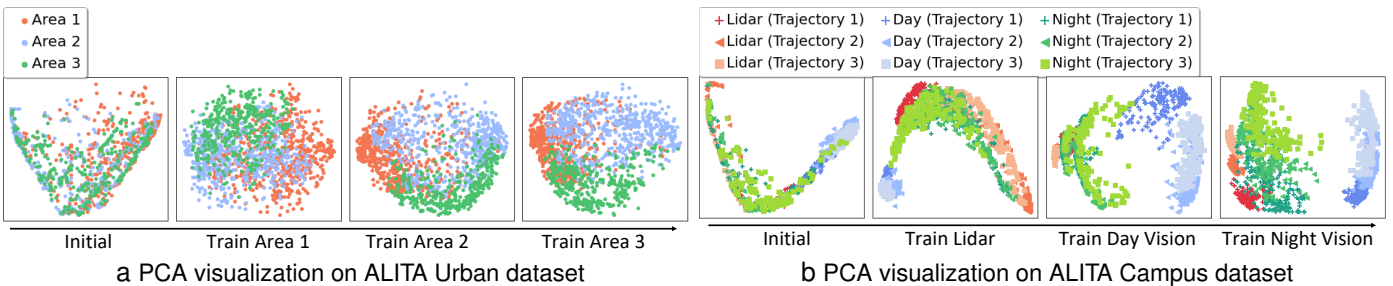
As described in section V-B, static memory stores concise and diverse observations based on their feature and spatial properties, which are achieved through memory consolidation. For observations from a current domain, memory consolidation mechanisms cluster the observations into K clusters (where $K = 20$ in our experiment) and then downsample to keep a maximum of N_{max} samples per cluster (where $N_{\text{max}} = 50$ in our experiment). Therefore, at most $K \times N_{\text{max}}$ samples were stored in the static memory for a domain or trajectory ($K \times N_{\text{max}} = 10000$ in our experiment). If the total number of clusters for all domains exceeds the maximum threshold K_{max} (where $K_{\text{max}} = 100$ in our experiment) or the total number of samples exceeds the buffer size $K_{\text{max}} \times N_{\text{max}}$, the memory forgetting mechanism is triggered to reduce the number of samples in static memory to keep at most $K_{\text{max}} \times N_{\text{max}}$ samples. The value of $K_{\text{max}} \times N_{\text{max}}$ is equivalent to the buffer size of the static memory, which is 5000 in our experiments. Fig. 18a shows the number of samples in static memory from different domains, such as weather and road conditions, during lifelong learning on the Oxford RobotCar dataset. In the experiment, we receive the observations from a new domain every 60 epochs. As depicted, the static memory incrementally stores samples from different domains throughout the lifelong learning process. If the total number of samples in the static memory exceeds the buffer



a Similarity matrix from different areas on Urban dataset

b Similarity matrix from different modalities on Campus dataset

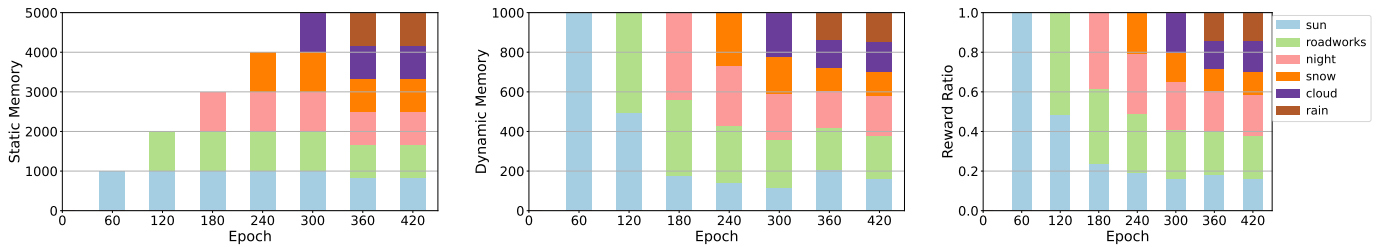
Fig. 16: **Similarity matrix during lifelong learning.** (a) ALITA Urban dataset. The left column represents the trajectories from area 1, area 2, and area 3. The right three columns display the corresponding similarity matrices with incremental learning. (b) ALITA Campus dataset. The left column represents Lidar, day-time visual, and night-time visual observations. The right three columns represent the corresponding similarity matrices with incremental learning.



a PCA visualization on ALITA Urban dataset

b PCA visualization on ALITA Campus dataset

Fig. 17: **PCA visualization over training.** (a) ALITA Urban dataset. Visualization of observations from different areas. (b) ALITA Campus dataset. Visualization of observations from different trajectories and different modalities.



a Number of samples in static memory

b Number of samples in dynamic memory

c Reward ratio of different domains

Fig. 18: **Static Memory, Dynamic memory zone, and Reward ratio (normalized) of different trajectories during lifelong learning training on Oxford RobotCar dataset.**

size, the memory forgetting mechanism deletes some similar samples in the memory.

Dynamic memory M_d selects memory traces from static memory with Memory Refreshing. Memory refreshing is based on importance sampling, where the sampling weight is proportional to its reward value. These selected memory traces are then sent to the memory decoder for training through memory replay. Fig. 18b illustrates the number of samples in dynamic memory from different domains during lifelong learning on the Oxford RobotCar dataset. To better understand the sample ratio between different domains, we also plot the normalized average reward for each domain in Fig. 18c. The average reward of a domain can be computed as the average of

all samples' reward in the domain $R^{D_t} = \frac{1}{|D_t|} \sum_{q_k \in D_t} \mathcal{R}(q_k)$. The normalized average reward of a domain is then $\bar{R}^{D_t} = R^{D_t} / \sum_{s \in [1, T]} R^{D_s}$. The plot shows that the dynamic memory zone holds more memory traces (samples) from domains with higher normalized average reward. Given a domain, a higher reward indicates worse performance, and BioSLAM uses higher sampling weights to retrieve more memory traces from the high-rewarded domain to achieve better performance.

H. Incremental Learning Property

To demonstrate the incremental learning property of BioSLAM, we evaluate its performance using a confidence

score during lifelong learning. The confidence score is calculated as a function of the place recognition loss L_{loc} , where $confidence(q_k) = \frac{L_{max} - L_{loc}(q_k)}{L_{max}}$, and L_{max} is the normalization constant set to the maximum loss value. A higher confidence score indicates a higher place recognition ability. We calculate the confidence score for all observations on real trajectories. Supplemental movie 1¹ presents the confidence score of BioSLAM during lifelong learning on the ALITA Campus dataset. Similarly, supplemental movie 2² presents the confidence score of BioSLAM during lifelong learning on the ALITA Urban dataset. As shown in the videos, the proposed BioSLAM framework enables incremental improvements in place recognition ability across all areas during lifelong learning.

TABLE IV: Comparison of GPU memory (Megabyte) of different methods.

Method	NetVLAD	SI	GR	BioSLAM
GPU Memory (MB)	1261	1265	1695	1695

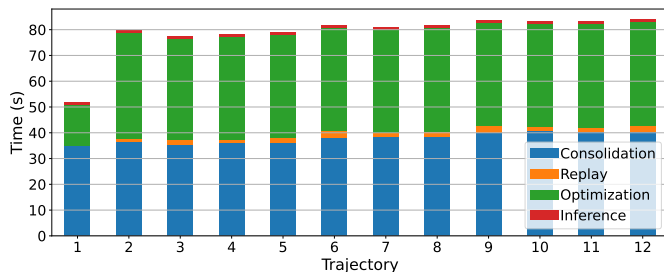


Fig. 19: Time usage (inference, optimization, replay, consolidation time) of the BioSLAM lifelong learning on ALITA Urban dataset during incremental learning.

I. Run-time Analysis

In this section, we present the memory and time usage of BioSLAM and compare it to another method for lifelong learning. Our experiments were conducted on an Ubuntu 18.04 system, with an Nvidia RTX 2080 Ti (12 GB) GPU, Intel Core i9-7900x processors, and 64 GB memory. We report the total memory usage on ALITA Urban datasets in IV. As can be seen, the memory usages of BioSLAM are acceptable for current embedded system structures.

Fig. 19 displays the time usage of the BioSLAM lifelong learning procedure on the ALITA Urban dataset, where new trajectory segments are incrementally fed. Each segment is approximately 2km in distance and composed of around 200 observation frames. The data inference procedure takes less than 1 second per trajectory segment, demonstrating efficiency in real-world inference. Memory consolidation and forgetting, on average, take approximately 40 seconds for each trajectory, which is fast enough to analyze the newly captured memory traces and update the memory system. Memory replay takes

2 seconds to generate replayed samples for place recognition training. Additionally, place recognition optimization takes 40 seconds to optimize a 2km new trajectory in one epoch. For multi-epoch training, BioSLAM runs inference, replay, and optimization multiple times but only runs memory consolidation once. Typically, training a trajectory segment about 50 times achieves convergence results. Therefore, the total learning time for a 2km trajectory segment is $40s + (1 + 2 + 40) * 50s = 2190s \approx 36min$. Given that the distance between neighbor keyframes is 10m, BioSLAM can learn 100m of new areas in around 1.8 minutes.

One important property of BioSLAM is that the time required for memory operations is not affected by differences in spatial or temporal scale. This is made possible by our memory forgetting mechanism, which allows us to maintain the searching space of M_S and M_d and keep memory traces up-to-date for localization. As a result, BioSLAM can be used for long-term place recognition tasks on low-cost robotic systems that use NVIDIA embedded systems.

VIII. DISCUSSION & LIMITATIONS

BioSLAM provides robust and efficient lifelong learning for large-scale and long-term place recognition tasks. The framework adopts a dual-memory mechanism, where long-lasting memory traces are stored in the static memory M_S , while generative memories are retrieved from dynamic memory M_D , enabling efficient learning of new types of observations and maintaining the lifelong memorization ability for old knowledge. This mechanism ensures robust place recognition under diverse conditions. In addition, BioSLAM demonstrates strong adaptability to changes in domains, as shown in the evaluation of the ALITA Campus dataset. The performance will not drop significantly when shifting from one domain to another, such as LiDAR, day-time, and night-time visual domains. This adaptability enables robots to achieve long-term autonomy in real-world scenarios. Supplemental movie 3³ demonstrates the mechanism and performance of BioSLAM in lifelong learning. The upper part of the video illustrates the BioSLAM framework, which comprises the General Place Learner, the rewarding mechanism, and a dual-memory module. The lower part of the video displays the confidence scores during lifelong learning.

However, BioSLAM also has some limitations for lifelong navigation. Firstly, it cannot offer sub-meter level localization ability like traditional visual SLAM systems. This limitation arises from the triplet loss applied in place descriptor learning, which cannot support feature-level alignment for 6D pose estimation. A potential solution could be to combine meter-level place descriptors and sub-meter-level features into a joint SLAM system. However, transferring both descriptors and features into the same lifelong learning framework would be another open challenge. Secondly, while BioSLAM enables the place recognition network module to learn observations incrementally under diverse domains, it cannot achieve cross-domain place recognition via direct transfer without learning when the target domains significantly differ (e.g., summer vs.

¹<https://youtu.be/eNrwUw7BWuE>

²<https://youtu.be/K8pSDJ5rLYs>

³<https://youtu.be/oa61retjo-U>

winter), as the appearance differences between such domains exceed the network's distinguishing ability. One potential solution is to combine an experience-based approach with BioSLAM. However, determining the required number of experiences for a given area and how to fuse/delete the experience and relative place descriptors presents another significant challenge for lifelong navigation.

In this work, BioSLAM provides a memory system for lifelong place recognition. On the other hand, it also provides a new option for other lifelong learning tasks. Recall the network structures as shown in Fig. 3. The functional modules related to the place recognition task are mainly the place descriptor extraction \mathcal{F}_θ and the relative external reward \mathcal{R}_{ex} in the rewarding mechanism. For other tasks, such as 3D segmentation, researchers can replace the place descriptor extraction network (i.e., the spherical convolution and VLAD layer) with a 3D U-Net [78], utilize segmentation loss metric instead of triplet loss, and not need to replace the entire blocks in the lifelong memory system. However, unlike place recognition which can leverage self-supervised mechanisms without human labeling, providing accurate segmentation annotations will be a significant challenge. Additionally, another potential option is to develop a parallel hybrid lifelong learning system for multiple tasks since the encoder module \mathcal{E} can be shared.

IX. CONCLUSION

The real-world robots will encounter diverse environmental changes under long-term autonomy. In the place recognition task, the robots continuously observe new scenarios, which are unbounded under variant conditions. In this work, we proposed BioSLAM, a lifelong place recognition method, to alleviate the above problem. BioSLAM combines a general place learning (GPL) system and a bio-inspired lifelong memory (BiLM) system. The GPL system utilizes a viewpoint-invariant place descriptor and a generative replay module to achieve the 'memory encoding' and 'memory replay' for continual place feature learning. The BiLM system provides a dual-memory mechanism controlled by a rewarding mechanism to guide the 'memory consolidation,' 'memory forgetting,' and 'memory replay' to enhance the memorization of long-term traces. We investigate the large-scale and long-term place recognition ability in experiments with city-scale 3D point-cloud maps, campus-scale visual-LiDAR hybrid inputs, and long-term city-scale visual inputs. Both results show that BioSLAM can significantly balance the place learning ability for new observations and maintain the memorization ability for historical observations.

Our method can be applied to low-cost mobile robots with current embedded devices, as it has a lightweight memory system that does not require saving massive streaming datasets. Another interesting direction for future work is to enable memory sharing between client agents and the cloud server. In this case, the server can be synced with data from all kinds of scenarios by various robots to update a more general place recognition. Finally, the BioSLAM system can be applied to other perception tasks by modifying objective functions in the rewarding mechanism according to specific requirements.

IEEE Transactions on Robotics (T-RO) paper, presented at ICRA 2024, Yokohama, Japan. Cite as T-RO paper.

X. ACKNOWLEDGMENT

This research is supported by grants from NVIDIA and utilized NVIDIA SDKs (CUDA Toolkit, TensorRT, and Omniverse). This research is supported by the ARL grant NO.W911QX20D0008 and partially by the National Science Foundation (NSF) under Grant No. 2144489. Any opinions, findings, conclusions, or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of ARL and NSF.

REFERENCES

- [1] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, "Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age," *IEEE Transactions on robotics*, vol. 32, no. 6, pp. 1309–1332, 2016.
- [2] D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superpoint: Self-supervised interest point detection and description," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 224–236.
- [3] Y. Latif, R. Garg, M. Milford, and I. Reid, "Addressing challenging place recognition tasks using generative adversarial networks," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 2349–2355.
- [4] T. Naseer, W. Burgard, and C. Stachniss, "Robust visual localization across seasons," *IEEE Transactions on Robotics*, vol. 34, no. 2, pp. 289–302, 2018.
- [5] W. Churchill and P. Newman, "Experience-based navigation for long-term localisation," *The International Journal of Robotics Research*, vol. 32, no. 14, pp. 1645–1661, 2013.
- [6] P. Mühlfellner, M. Bürki, M. Bosse, W. Derendarz, R. Philippsen, and P. Furgale, "Summary maps for lifelong visual localization," *Journal of Field Robotics*, vol. 33, no. 5, pp. 561–590, 2016.
- [7] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardós, "Orb-slam: A versatile and accurate monocular slam system," *IEEE Transactions on Robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [8] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. M. Montiel, and J. D. Tardós, "Orb-slam3: An accurate open-source library for visual, visual-inertial, and multimap slam," *IEEE Transactions on Robotics*, vol. 37, no. 6, pp. 1874–1890, 2021.
- [9] T. Qin, P. Li, and S. Shen, "Vins-mono: A robust and versatile monocular visual-inertial state estimator," *IEEE Transactions on Robotics*, vol. 34, no. 4, pp. 1004–1020, 2018.
- [10] P. Yin, L. Xu, J. Zhang, H. Choset, and S. Scherer, "i3dloc: Image-to-range cross-domain localization robust to inconsistent environmental conditions," in *Proceedings of Robotics: Science and Systems (RSS '21)*. Robotics: Science and Systems 2021, 2021.
- [11] E. I. Moser, E. Kropff, and M.-B. Moser, "Place cells, grid cells, and the brain's spatial representation system," *Annu. Rev. Neurosci.*, vol. 31, pp. 69–89, 2008.
- [12] J. Stretton and P. Thompson, "Frontal lobe function in temporal lobe epilepsy," *Epilepsy research*, vol. 98, no. 1, pp. 1–13, 2012.
- [13] G. Berdugo-Vega and J. Graeff, "Inquiring the librarian about the location of memory," *Cognitive Neuroscience*, vol. 0, no. 0, pp. 1–3, 2022.
- [14] J. G. Klinzing, N. Niethard, and J. Born, "Mechanisms of systems memory consolidation during sleep," *Nature neuroscience*, vol. 22, no. 10, pp. 1598–1610, 2019.
- [15] M. Zaffar, S. Garg, M. Milford, J. Kooij, D. Flynn, K. McDonald-Maier, and S. Ehsan, "VPR-bench: An open-source visual place recognition evaluation framework with quantifiable viewpoint and appearance change," *International Journal of Computer Vision.*, vol. 129, no. 7, pp. 2136–2174, 2021.
- [16] S. Lowry, N. Sünderhauf, P. Newman, J. J. Leonard, D. Cox, P. Corke, and M. J. Milford, "Visual place recognition: A survey," *IEEE Transactions on Robotics*, vol. 32, no. 1, pp. 1–19, 2016.
- [17] T. Barros, R. Pereira, L. Garrote, C. Premevida, and U. J. Nunes, "Place recognition survey: An update on deep learning approaches," *arXiv preprint arXiv:2106.10458*, 2021.
- [18] X. Zhang, L. Wang, and Y. Su, "Visual place recognition: A survey from deep learning perspective," *Pattern Recognition*, vol. 113, p. 107760, 2021.

- [19] M. Delange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and T. Tuytelaars, "A continual learning survey: Defying forgetting in classification tasks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [20] T. Lesort, V. Lomonaco, A. Stoian, D. Maltoni, D. Filliat, and N. D. Rodríguez, "Continual learning for robotics: Definition, framework, learning strategies, opportunities and challenges," *Information Fusion*, vol. 58, pp. 52–68, 2020.
- [21] K. Mikolajczyk and C. Schmid, "A performance evaluation of local descriptors," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 10, pp. 1615–1630, 2005.
- [22] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "ORB: An efficient alternative to SIFT or SURF," in *2011 International Conference on Computer Vision*, 2011, pp. 2564–2571.
- [23] M. Cummins and P. Newman, "Fab-map: Probabilistic localization and mapping in the space of appearance," *The International Journal of Robotics Research*, vol. 27, no. 6, pp. 647–665, 2008.
- [24] Y. Tian, Y. Chang, F. Herrera Arias, C. Nieto-Granda, J. P. How, and L. Carlone, "Kimera-multi: Robust, distributed, dense metric-semantic slam for multi-robot systems," *IEEE Transactions on Robotics*, vol. 38, no. 4, pp. 2022–2038, 2022.
- [25] E. Garcia-Fidalgo and A. Ortiz, "ibow-lcd: An appearance-based loop-closure detection approach using incremental bags of binary words," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 3051–3057, 2018.
- [26] S. An, H. Zhu, D. Wei, K. A. Tsintotas, and A. Gasteratos, "Fast and incremental loop closure detection with deep features and proximity graphs," *Journal Of Field Robotics*, vol. 39, no. 4, pp. 473–493, 2022.
- [27] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "Netvlad: Cnn architecture for weakly supervised place recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5297–5307.
- [28] R. Arandjelovic and A. Zisserman, "All about vlad," in *2013 IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 1578–1585.
- [29] A. Khaliq, S. Ehsan, Z. Chen, M. Milford, and K. McDonald-Maier, "A holistic visual place recognition approach using lightweight cnns for significant viewpoint and appearance changes," *IEEE Transactions on Robotics*, vol. 36, no. 2, pp. 561–569, 2020.
- [30] L. Hui, H. Yang, M. Cheng, J. Xie, and J. Yang, "Pyramid point cloud transformer for large-scale place recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6098–6107.
- [31] L. Kunze, N. Hawes, T. Duckett, M. Hanheide, and T. Krajník, "Artificial intelligence for long-term robot autonomy: A survey," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 4023–4030, 2018.
- [32] A. Francis, A. Faust, H.-T. L. Chiang, J. Hsu, J. C. Kew, M. Fiser, and T.-W. E. Lee, "Long-range indoor navigation with prm-rl," *IEEE Transactions on Robotics*, vol. 36, no. 4, pp. 1115–1134, 2020.
- [33] K. MacTavish, M. Paton, and T. D. Barfoot, "Visual triage: A bag-of-words experience selector for long-term visual route following," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 2065–2072.
- [34] M. Warren, M. Greeff, B. Patel, J. Collier, A. P. Schoellig, and T. D. Barfoot, "There's no place like home: Visual teach and repeat for emergency return of multi-robot uavs during gps failure," *IEEE Robotics and automation letters*, vol. 4, no. 1, pp. 161–168, 2018.
- [35] M. Paton, K. MacTavish, L.-P. Berczi, S. K. van Es, and T. D. Barfoot, "I can see for miles and miles: An extended field test of visual teach and repeat 2.0," in *Field and Service Robotics*. Springer, 2018, pp. 415–431.
- [36] S. Lowry and M. J. Milford, "Supervised and unsupervised linear learning techniques for visual place recognition in changing environments," *IEEE Transactions on Robotics*, vol. 32, no. 3, pp. 600–613, 2016.
- [37] T. Krajník, J. P. Fentanes, J. M. Santos, and T. Duckett, "Fremen: Frequency map enhancement for long-term mobile robot autonomy in changing environments," *IEEE Transactions on Robotics*, vol. 33, no. 4, pp. 964–977, 2017.
- [38] G. D. Tipaldi, D. Meyer-Delius, and W. Burgard, "Lifelong localization in changing environments," *The International Journal of Robotics Research*, vol. 32, no. 14, pp. 1662–1678, 2013.
- [39] G. Kurz, M. Holoch, and P. Biber, "Geometry-based graph pruning for lifelong slam," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 3313–3320.
- [40] S. Chen, J. Wu, Q. Lu, Y. Wang, and Z. Lin, "Cross-scene loop-closure detection with continual learning for visual simultaneous localization and mapping," *International Journal of Advanced Robotic Systems*, vol. 18, no. 5, p. 17298814211050560, 2021.
- [41] N. Vödisch, D. Cattaneo, W. Burgard, and A. Valada, "Continual slam: Beyond lifelong simultaneous localization and mapping through continual learning," in *The International Symposium of Robotics Research*. Springer, 2022, pp. 19–35.
- [42] H. Blum, F. Milano, R. Zurbrügg, R. Siegwart, C. Cadena, and A. Gawel, "Self-improving semantic perception for indoor localisation," in *Conference on Robot Learning*. PMLR, 2022, pp. 1211–1222.
- [43] Z. Li and D. Hoiem, "Learning without forgetting," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 12, pp. 2935–2947, 2017.
- [44] A. Mallya and S. Lazebnik, "Packnet: Adding multiple tasks to a single network by iterative pruning," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 7765–7773.
- [45] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska et al., "Overcoming catastrophic forgetting in neural networks," *Proceedings of the national academy of sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.
- [46] F. Zenke, B. Poole, and S. Ganguli, "Continual learning through synaptic intelligence," in *International Conference on Machine Learning*. PMLR, 2017, pp. 3987–3995.
- [47] D. Gao, C. Wang, and S. Scherer, "Airloop: Lifelong loop closure detection," in *2022 International Conference on Robotics and Automation (ICRA)*, 2022, pp. 10 664–10 671.
- [48] J. Knights, P. Moghadam, M. Ramezani, S. Sridharan, and C. Fookes, "Includ: Incremental learning for point cloud place recognition," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 8559–8566.
- [49] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "icarl: Incremental classifier and representation learning," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2017, pp. 2001–2010.
- [50] D. Lopez-Paz and M. Ranzato, "Gradient episodic memory for continual learning," *Advances in neural information processing systems*, vol. 30, 2017.
- [51] G. M. Van de Ven and A. S. Tolias, "Generative replay with feedback connections as a general strategy for continual learning," *arXiv preprint arXiv:1809.10635*, 2018.
- [52] Y. Choi, M. El-Khamy, and J. Lee, "Dual-teacher class-incremental learning with data-free generative replay," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3543–3552.
- [53] F. Dayoub and T. Duckett, "An adaptive appearance-based map for long-term topological localization of mobile robots," in *2008 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2008, pp. 3364–3369.
- [54] P. Yin, F. Wang, A. Egorov, J. Hou, J. Zhang, and H. Choset, "Seqspherevlad: Sequence matching enhanced orientation-invariant place recognition," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020, pp. 5024–5029.
- [55] R. Stickgold, "Sleep-dependent memory consolidation," *Nature*, vol. 437, no. 7063, pp. 1272–1278, 2005.
- [56] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *3rd International Conference on Learning Representations (ICLR 2015)*. Computational and Biological Learning Society, 2015.
- [57] P. Yin, L. Xu, J. Zhang, and H. Choset, "Fusionvlad: A multi-view deep fusion networks for viewpoint-free 3d place recognition," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 2304–2310, 2021.
- [58] C. Esteves, C. Allen-Blanchette, A. Makadia, and K. Daniilidis, "Learning so (3) equivariant representations with spherical cnns," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 52–68.
- [59] T. S. Cohen, M. Geiger, J. Köhler, and M. Welling, "Spherical cnns," in *6th International Conference on Learning Representations, ICLR, 2018*.
- [60] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," *Advances in neural information processing systems*, vol. 27, 2014.
- [61] H. Shin, J. K. Lee, J. Kim, and J. Kim, "Continual learning with deep generative replay," in *ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS 30 (NIPS 2017)*, ser. Advances in Neural Information Processing Systems, I. Guyon, U. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30, 2017.

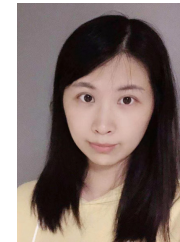
- [62] G. M. van de Ven, H. T. Siegelmann, and A. S. Tolia, "Brain-inspired replay for continual learning with artificial neural networks," *Nature Communications*, vol. 11, p. 4069, 2020.
- [63] K. A. Krueger and P. Dayan, "Flexible shaping: How learning in small steps helps," *Cognition*, vol. 110, no. 3, pp. 380–394, 2009.
- [64] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proceedings of the 26th annual international conference on machine learning*, 2009, pp. 41–48.
- [65] P. Yin, S. Zhao, H. Lai, R. Ge, J. Zhang, H. Choset, and S. Scherer, "Automerger: A framework for map assembling and smoothing in city-scale environments," *IEEE Transactions on Robotics*, pp. 1–19, 2023.
- [66] J. Byrne, *Learning and memory: a comprehensive reference*. Academic Press, 2017.
- [67] M. G. Berman, J. Jonides, and R. L. Lewis, "In search of decay in verbal short-term memory," *Journal of Experimental Psychology: Learning, Memory, and Cognition*, vol. 35, no. 2, p. 317, 2009.
- [68] J. Zhang and S. Singh, "Loam: Lidar odometry and mapping in real-time," in *Robotics Science and Systems*, vol. 2, no. 9. Berkeley, CA, 2014, pp. 1–9.
- [69] A. Segal, D. Haehnel, and S. Thrun, "Generalized-icp," in *Robotics: science and systems*, vol. 2, no. 4. Seattle, WA, 2009, p. 435.
- [70] P. Yin, S. Zhao, R. Ge, I. Cisneros, R. Fu, J. Zhang, H. Choset, and S. Scherer, "Alita: A large-scale incremental dataset for long-term autonomy," *arXiv preprint arXiv:2205.10737*, 2022.
- [71] W. Maddern, G. Pascoe, C. Linegar, and P. Newman, "1 year, 1000 km: The oxford robotcar dataset," *The International Journal of Robotics Research*, vol. 36, no. 1, pp. 3–15, 2017.
- [72] M. A. Uy and G. H. Lee, "Pointnetvlad: Deep point cloud based retrieval for large-scale place recognition," in *2018 IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4470–4479.
- [73] D. Gálvez-López and J. D. Tardós, "Bags of binary words for fast place recognition in image sequences," *IEEE Transactions on Robotics*, vol. 28, no. 5, pp. 1188–1197, October 2012.
- [74] M. Zaffar, S. Ehsan, M. Milford, and K. McDonald-Maier, "Cohog: A light-weight, compute-efficient, and training-free visual place recognition technique for changing environments," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 1835–1842, 2020.
- [75] H. Shin, J. K. Lee, J. Kim, and J. Kim, "Continual learning with deep generative replay," *Advances in neural information processing systems*, vol. 30, 2017.
- [76] D. A. Baldwin, E. M. Markman, and R. L. Melartin, "Infants' ability to draw inferences about nonobvious object properties: Evidence from exploratory play," *Child Development*, vol. 64, no. 3, pp. 711–728, 1993.
- [77] M. H. Bornstein and M. E. Arterberry, "The development of object categorization in young children: hierarchical inclusiveness, age, perceptual attribute, and group versus individual analyses," *Developmental psychology*, vol. 46, no. 2, p. 350, 2010.
- [78] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, "3d u-net: learning dense volumetric segmentation from sparse annotation," in *International conference on medical image computing and computer-assisted intervention*. Springer, 2016, pp. 424–432.



Abulikemu Abuduweili received his B.S. in Electrical Engineering from Peking University, China, in 2017, and his M.S. in Electrical Engineering from Peking University, China, in 2020. He is currently a PhD student at Electrical Engineering and Robotics at Carnegie Mellon University, USA. His research interests include Deep Learning, Robotics, and Computer Vision.



Shiqi Zhao received his Bachelor's degree from Dalian University of Technology, Dalian, China, in 2018, and his Master's degree from the University of California San Diego, U.S., in 2020. He is currently working as a research assistant at City University of Hong Kong. His research interests include Place Recognition, 3D Perception, and Deep Learning.



Lingyun Xu received her Bachelor's degree from Huazhong University of Science and Technology, Wuhan, China, in 2011, and her Ph.D. degree from the University of Chinese Academy of Sciences, Beijing, in 2018. She is research Post-doctoral with the Department of the Robotics Institute, Carnegie Mellon University, Pittsburgh, USA. Her research interests include Visual Tracking, SLAM, Place Recognition, 3D Perception, and Autonomous Driving. Dr. Xu has served as a Reviewer for several IEEE Conferences ICRA, IROS, CVPR, and ICCV.



Changliu Liu received the B.S. degree in mechanical engineering and the B.S. degree in economics from Tsinghua University, China, in 2012, the M.S. degree in mechanical engineering, and the M.A. degree in mathematics from the University of California, Berkeley, U.S.A., in 2014 and 2016 respectively, and the Ph.D. degree in mechanical engineering from the University of California, Berkeley, U.S.A., in 2017. She is an assistant professor at the Robotics Institute at Carnegie Mellon University. Her research interests lie in designing and verifying intelligent systems with applications to manufacturing and transportation. She received NSF Career Award, Amazon Research Award, and Ford URP Award.



Peng Yin received his Bachelor's degree from Harbin Institute of Technology, Harbin, China, in 2013, and his Ph.D. degree from the University of Chinese Academy of Sciences, Beijing, in 2018. He is currently an Assistant Professor at City University of Hong Kong, China. His research interests include LiDAR SLAM, Place Recognition, 3D Perception, and Reinforcement Learning. Dr. Yin has served as a Reviewer for several IEEE Conferences ICRA, IROS, ACC, RSS.



Sebastian Scherer received his B.S. in Computer Science, M.S. and Ph.D. in Robotics from CMU in 2004, 2007, and 2010. Sebastian Scherer is an Associate Research Professor at the Robotics Institute at Carnegie Mellon University. His research focuses on enabling autonomy for unmanned rotorcraft to operate at low altitude in cluttered environments. He is a Siebel scholar and a recipient of multiple paper awards and nominations, including AIAA@Infotech 2010 and FSR 2013. His research has been covered by the national and internal press including IEEE

Spectrum, the New Scientist, Wired, der Spiegel, and the WSJ.