

Multi-Robot Distributed Semantic Mapping in Unfamiliar Environments through Online Matching of Learned Representations

Stewart Jamieson^{1,2}, Kaveh Fathian², Kasra Khosoussi², Jonathan P. How², Yogesh Girdhar³

Abstract—We present a solution to multi-robot distributed semantic mapping of novel and unfamiliar environments. Most state-of-the-art semantic mapping systems are based on supervised learning algorithms that cannot classify novel observations online. While unsupervised learning algorithms can invent labels for novel observations, approaches to detect when multiple robots have independently developed their own labels for the same new class are prone to erroneous or inconsistent matches. These issues worsen as the number of robots in the system increases and prevent fusing the local maps produced by each robot into a consistent global map, which is crucial for cooperative planning and joint mission summarization. Our proposed solution overcomes these obstacles by having each robot learn an unsupervised semantic scene model online and use a multiway matching algorithm to identify consistent sets of matches between learned semantic labels belonging to different robots. Compared to the state of the art, the proposed solution produces 20-60% higher quality global maps that do not degrade even as many more local maps are fused.

OPEN SOURCE SOFTWARE

An implementation of this solution is contributed in the “Sunshine” 3D semantic mapping ROS package, a general purpose single- and multi-robot semantic mapping system: <https://gitlab.com/warplab/ros/sunshine>.

I. INTRODUCTION

Semantic mapping is a relatively young field that was initially motivated by giving robots a spatial awareness of nearby terrains, objects, and activities [1]. Semantic maps describe the world using a set of classes, and have been used with great effectiveness in solving many field robotics problems such as mission summarization [2], object-based SLAM [3]–[5], and context-aware planning [6], [7]. Great progress has been made in improving the accuracy of semantic mapping systems by leveraging deep learning models trained on large datasets [8]–[12]. However, most of these systems do not support *a priori* unknown classes, which are an essential part of scientific exploration. For example, in an underwater exploration scenario an explicit goal is novelty detection, e.g., discovering new species or

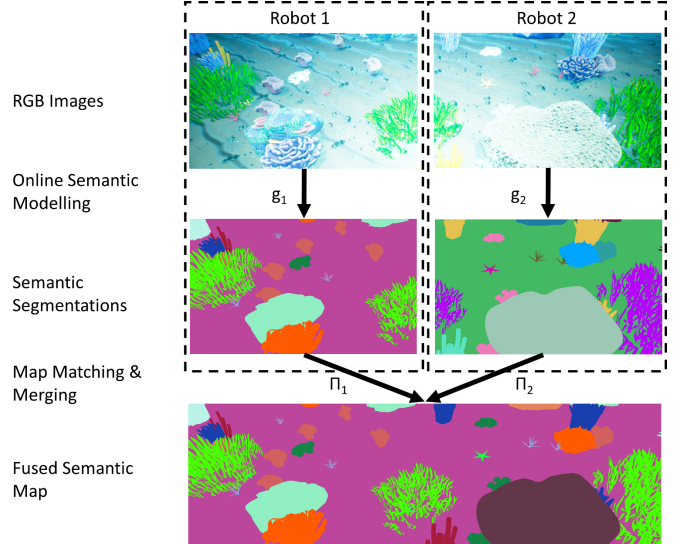


Fig. 1. When robots develop different individual semantic models g_i , they solve the correspondences Π_i into a global (shared) semantic language in order to fuse their results. Images are from environment #2 (see Fig. 3).

geological phenomena. This complication makes it crucial to find algorithms that use unsupervised learning to develop new semantic representations *online*, so that the semantic mapping system can detect and classify novel observations.

Another important aspect of scientific exploration is that the task often involves large and unknown environments, which can be very time consuming to cover with a single robot. Furthermore, many large-scale phenomena of scientific interest, such as mass migrations, feeding events, and geological activities like volcanism, are transient and dynamic and can thus be easily missed or insufficiently covered by a single robot. These issues necessitate using a team of robots working in parallel for *distributed* mapping and exploration.

When learning semantic representations online using unsupervised algorithms, the learned models are egocentric. Fig. 1 demonstrates this can be an issue: one robot has learned a different permutation of the same semantic labels (represented with colors) learned by the other robot. In addition to the unknown permutation between corresponding labels, some labels learned by one robot may not correspond to any label observed by another robot, or a single label may represent the union of multiple labels learned by another robot. Therefore, multi-robot distributed semantic mapping with learned representations requires correctly estimating the total number of distinct labels, and associating and fusing labels that correspond to the same semantic category.

¹S. Jamieson is with the MIT-WHOI Joint Program in Oceanography/Applied Ocean Science and Engineering sjamieson@whoi.edu

²S. Jamieson, K. Fathian, K. Khosoussi, and J. P. How are with the Department of Aeronautics and Astronautics at the Massachusetts Institute of Technology (MIT) [{sjamieson, kavehf, kasra, jhow}@mit.edu](mailto)

³Y. Girdhar is with the Applied Ocean Physics and Engineering Department at the Woods Hole Oceanographic Institution (WHOI) yogi@whoi.edu

*This work was supported in part by NSF-NRI Award Number 1734400, by ARL DCIST under Cooperative Agreement Number W911NF-17-2-0181, and by ONR under BRC award N000141712072.

This work presents a novel system for multi-robot distributed semantic mapping that addresses the previously described issues. Each robot uses an online semantic 3D mapping system to model its own observations and create a high quality semantic map. The robots explore the target environment in parallel, sharing their learned semantic maps and models with each other and with the human operator whenever communication constraints permit. Finally, a multiway matching algorithm that can run on any robot estimates the total number of unique phenomena observed across the robot team, finds matches between the same phenomena labeled differently by various robots, and fuses the local maps to obtain a consistent global map across all robots.

II. RELATED WORKS

Spatiotemporal topic models (STMs) [13]–[15] are a class of unsupervised learning algorithms that has been specifically augmented for realtime semantic mapping in novel environments [16], [17]. BNP-ROST is an STM that adaptively develops semantic labels online as new phenomena are observed [18], and has been used for 3D semantic mapping in bandwidth-limited environments without any pre-training [19]. STMs can create high quality semantic maps in realtime by leveraging the spatial and temporal correlations between observations and using efficient sampling algorithms to discover good semantic representations online [19], [20]. While there has been progress in other unsupervised semantic image segmentation approaches [21], including deep-learning based methods [22], by not leveraging these correlations they produce lower quality maps of large environments.

Doherty et al. [23] showed that repeatedly matching the topic models of two robots with the Hungarian algorithm [24], merging them together, and distributing the merged model would result in both robots converging to a single set of good semantic labels, even if this was done at a low frequency. To our knowledge, [23] is the only prior work that has explored multi-robot distributed semantic mapping with representations learned online. The present work improves upon [23] with a novel solution for matching many ($N \gg 2$) semantic maps that is more robust to major variations across what each robot sees, and does not rely on distributing the merged topic model throughout the entire robot team. Eliminating the need to distribute the merged model halves communication bandwidth usage, and makes the approach more robust to transient connection failures.

Multiway matching algorithms are a class of data association techniques that leverage the transitivity property (cycle consistency) to rectify wrong correspondences and construct a unified representation from the partial and noisy observations of multiple agents. Multiway matching leads to superior accuracy compared to classical pairwise approaches (e.g., [24]), however it has combinatorial complexity. State-of-the-art methods consider approximations of this problem via convex relaxations [25]–[28], spectral relaxations [29], [30], or graph clustering [31] to obtain a solution in polynomial time. In this work we use CLEAR [32], a spectral clustering based approach, to match topic models because it

is one of the most efficient multiway matching algorithms with leading performance in both precision and recall.

III. PROBLEM SETUP

Let us denote the environment to be mapped as $E \subset \mathbb{R}^3$. We partition E into a grid of disjoint cells (boxes) $\mathcal{B} = \{b_i\}$ such that $E = \cup_i b_i$. An oracle (e.g., the human operator) could assign to each box a distribution over human-defined semantic labels \mathcal{Z}^H that represent its semantic contents. We assume that if these boxes are sufficiently small, each box can be effectively represented by a single dominant label. We thus define the ground truth semantic segmentation $f : \mathcal{B} \rightarrow \mathcal{Z}^H$, where $f(b) \in \mathcal{Z}^H$ represents the human label for $b \in \mathcal{B}$. The model f and set \mathcal{Z}^H are unknown *a priori* because the operator did not know what would be found during the mission, and communication bandwidth limitations prevent the operator from seeing most of the observations until the mission is over and the robots are recovered. The goal of the robot team is to construct a fused semantic map $g : \mathcal{B} \rightarrow \mathcal{Z}^G$, where \mathcal{Z}^G is a shared set of learned semantic labels, such that $g(b)$ is “similar” to $f(b)$. A metric to evaluate this similarity will be presented in Section V.

We assume that the team consists of N autonomous robots. By timestep t , the n^{th} robot has collected its own set of localized image observations and used them to build, in an unsupervised manner, a local semantic map $g_{n,t} : \mathcal{B} \rightarrow \mathcal{Z}_t^n$, where \mathcal{Z}_t^n is the set of semantic labels the robot has developed to describe its own observations.¹ Due to the egocentricity of unsupervised semantic mapping, robots which observe different phenomena, or the same but in a different order, will likely develop disparate semantic models, as in Fig. 1. In order to construct a fused semantic map, these N unique semantic models must first be fused to use a common set of labels. Therefore, the team must construct a set of global semantic labels \mathcal{Z}_t^G and a set of correspondences $\Pi_t = \{\Pi_{n,t} : \mathcal{Z}_t^n \rightarrow \mathcal{Z}_t^G\}_{n=1}^N$ that translate individual robots’ labels \mathcal{Z}_t^n into \mathcal{Z}_t^G . Given Π_t , the individual semantic maps can be fused into a single global map $g_t : \mathcal{B} \rightarrow \mathcal{Z}_t^G$ for time t (see Fig. 1). Each label can be computed as $g_t(b) = \Pi_{n_b^*,t}(g_{n_b^*,t}(b))$ where n_b^* is the index of the robot that most recently visited and observed the cell $b \in \mathcal{B}$.

IV. CONSISTENT ONLINE TOPIC MATCHING

We present an algorithm to produce a fused semantic map from semantic maps built by robots with disparate semantic models. *Consistency* ensures that the fused model created by any agent is the same across all agents, and *online* indicates that our algorithm will be run during the mission whenever a new matching is required. Typically, we expect a human operator would run our algorithm at a central node whenever they require an updated global map, however it is computationally lightweight enough that it may be run by any robot that has collected other robots’ semantic maps.

The stages of the proposed system are presented graphically in Fig. 2 and in detail by the following subsections.

¹In practice, we assume $g_{n,t}$ is learned from RGB-D image observations $\{o_{n,\tau}\}_{\tau=1}^t$ paired with estimated camera poses $\{x_{n,\tau} \in \text{SE}(3)\}_{\tau=1}^t$.

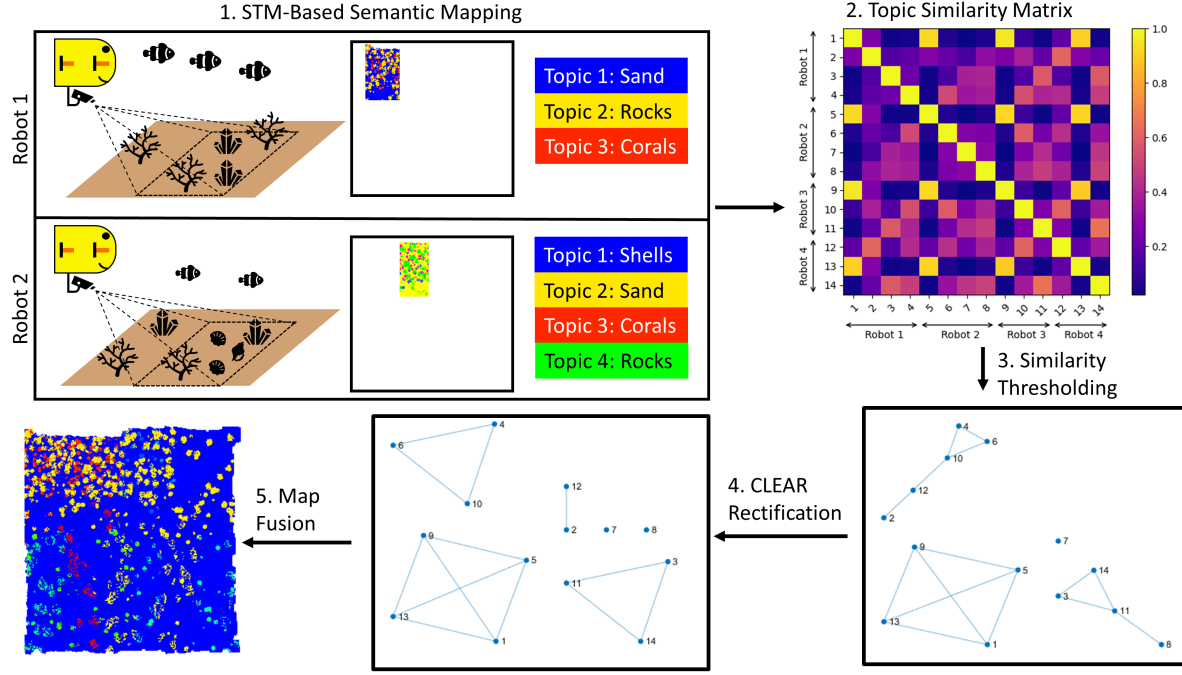


Fig. 2. Our proposed system is composed of five stages. Each robot first learns an individual semantic model online which it uses to describe its own observations. When a fused map is required, a topic similarity matrix is constructed by using a similarity metric for topic descriptors to compute all pairwise similarity scores. A noisy association graph is produced by treating each topic as a vertex and using similarities above a specified threshold as edges. The noisy graph is rectified using the CLEAR multiway matching algorithm to produce a consistent topic matching, which has the form of a cluster graph. These topic matches are used to fuse the individual maps into one consistently labelled global map.

A. Online STM-Based Semantic Mapping

The proposed approach has each robot construct a spatiotemporal topic model online as the basis for its individual semantic map, as in [19]. Nonetheless, it would be straightforward to adapt this system to use any similar unsupervised online semantic mapping module.

While an exhaustive description of the BNP-ROST model used is left to [18], a few details are relevant here. First, the model requires no pre-training, although it can be bootstrapped with topics for common phenomena. Second, the model is tuned by varying the feature vocabulary, the spatiotemporal grid cell size, and three scalar hyperparameters. Typically, the vocabulary is domain-specific while the cell size and the scalar hyperparameters are tuned for a specific mission.² All robots on the same mission use the same hyperparameters. Finally, each robot’s topic model uses a stochastic process to develop an ever-evolving set of topics online based on its own visual observations. At time t the n^{th} robot has $K_{n,t}$ topics, so $\mathcal{Z}_t^n = \{1, \dots, K_{n,t}\}$. Each topic $z_k \in \mathcal{Z}_t^n$ is characterized by a semantic “descriptor” $\phi_k \in \Delta^V$, a distribution over “words” in the predefined vocabulary of size V [16], where $\Delta^V = \{p \in \mathbb{R}_+^V : \|p\|_1 = 1\}$. Each grid cell b_i of the environment E is labelled with a single maximum likelihood topic, which may change as the model evolves; thus, we will start using the terms “topic” and “label” interchangeably.

²While they may be tuned with a dataset representative of the target environment, these hyperparameters encode only abstract information about the domain and thus tend to generalize well to novel environments.

B. Computing Topic Similarity

We require a similarity metric that measures how similar two topics are to each other in order to identify when multiple robots have developed any semantically equivalent topics. Since each descriptor ϕ_k represents a probability mass function, it is natural to consider similarity metrics that operate on discrete probability distributions. Total Variation Distance (TVD) [33] measures the largest possible difference in probability that two topics assign to the same set of words. Thus, the Topic Overlap (TO),

$$\text{TO}(\phi_1, \phi_2) = 1 - \text{TVD}(\phi_1, \phi_2), \quad (1)$$

is a similarity metric that represents the total probability mass which both ϕ_1 and ϕ_2 assign similarly. It can be computed using the identity $\text{TVD}(\phi_1, \phi_2) = \frac{1}{2} \|\phi_1 - \phi_2\|_1$.

Another metric commonly used for comparing topic descriptors is Cosine Similarity (CS) [34], which computes the cosine of the angle between two descriptors as

$$\text{CS}(\phi_1, \phi_2) = \frac{\phi_1 \cdot \phi_2}{\|\phi_1\| \|\phi_2\|}. \quad (2)$$

The CS metric is about as efficient to compute as TO, but assigns a higher score when ϕ_1 and ϕ_2 are very similar and a lower score when they are very dissimilar. This may be preferable because, as the stochastic nature of topic models means that the descriptors $\{\phi_k\}$ fluctuate constantly, two topics with the same semantic meaning are likely to have slightly different descriptors at any given time.

Each similarity metric presented above is symmetric and bounded by $[0, 1]$, where a score of 0.0 indicates two topics

have no words in common, and 1.0 indicates that two topics are exactly the same. We use the chosen similarity metric s to construct the *pairwise similarity graph*, a weighted and undirected graph in which vertices are topics and edge weights are the similarity of the adjoining vertices; it is represented in Fig. 2 by the topic similarity matrix.

C. Constructing the Noisy Association Graph

In practice, similarity metrics are “noisy” in that topics which a human would judge to have the same semantic meaning may not have a similarity score of 1.0, and topics with very different semantic meanings may not have a similarity of 0.0. The pairwise similarity graph is simplified by removing edges with weights below some $\sigma \in (0, 1)$ that represents low similarity, and setting weights above σ to 1, resulting in the unweighted *noisy association graph*.

In general, a good threshold σ for considering two topics to be “sufficiently similar” will depend on factors including the similarity metric s , the topic model hyperparameters, and the subjective opinion of the human operator. It is difficult to choose σ analytically because the expected topic growth rate and average inter-topic similarity are complicated functions of the topic model hyperparameters. A simple solution for choosing σ is to collect a validation set of topics developed by robots in past missions, for which the human operator can infer their semantic meanings, and then tune σ low enough that the algorithm merges as many topics with the equivalent meanings as possible but high enough that it does not match distinct topics. In the training dataset used to choose the topic model hyperparameters, to be described in Section V, a threshold of $\sigma = 0.75$ was found to work well with both topic similarity metrics. However, even if the threshold σ is chosen well, it may not be obvious from this graph how many unique topics should be used in the final map, or which sets of topics would form a consistent matching.

D. Rectifying the Noisy Association Graph

When visualized as a graph where vertices represent topics and edges represent matches, a consistent matching has the structure of a *cluster graph*. This is a graph composed of disjoint fully-connected components, so that any two topics in the same component are matched and no topics are matched between components. In this cluster graph, the number of distinct topic labels developed by the entire robot team is equal to the number of disjoint components.

CLEAR [32] is a spectral clustering algorithm that estimates the closest cluster graph to a noisy association graph (see Fig. 2). A key reason for choosing CLEAR is that it is one of the fastest algorithms to perform multiway matching with high accuracy. The Laplacian L of the noisy association graph is a matrix defined in terms of the graph adjacency matrix A and degree matrix D as $L = D - A$, where

$$[A]_{ij} = \begin{cases} 1, & s(\phi_i, \phi_j) \geq \sigma \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

$$[D]_{ij} = \begin{cases} \sum_{k=1}^N [A]_{ik}, & i = j \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

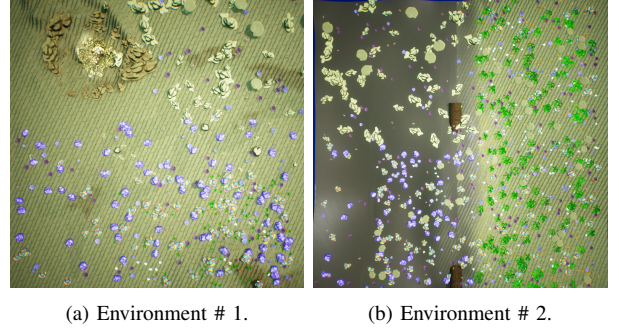


Fig. 3. Top-down views of the two simulated test environments used in the experiments. Each map is approximately 250m×250m, and contains a rich variety of coral species, seaweed, and rocks.

CLEAR uses a special normalization of the Laplacian based on the degree matrix plus identity, denoted by L_{nm} , to identify clusters of semantic labels in the noisy association graph with high pairwise similarity. The number of eigenvalues of L_{nm} less than 0.5 is a robust estimate of the number of global labels, $|\mathcal{Z}_t^G|$. CLEAR then uses the eigenvectors of L_{nm} to find a consistent set of label correspondences Π_t .

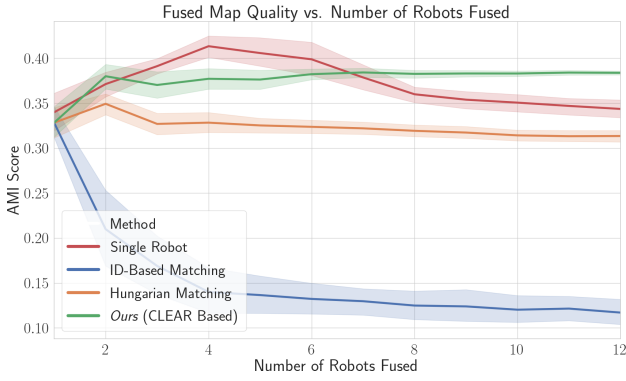
V. EXPERIMENTAL METHODOLOGY

The proposed consistent online topic matching system was evaluated using semantic mapping experiments in two unique high-resolution 3D simulated coral reef environments produced in the Unreal Engine [35] using the Automatic Coral Generator package [36]. A top-down view of each simulated environment is presented in Fig. 3. In each experiment, a team of 12 simulated robots traversed one of the two environments and each collected 250 RGB-D observations using the AirSim plugin [37]. AirSim also provided a ground truth semantic segmentation for each image. Each robot was given noiseless localization information and ran a new release of the BNP-ROST [18] spatiotemporal topic model called “Sunshine”, which is conceptually identical to the system presented in [19] but redesigned for ease of use and with optimized code to produce higher quality maps with less processing power. Throughout the experiment, sets of 1 to 12 robots’ local maps were randomly chosen and fused together using the approach described in Section IV. Each experiment was repeated 24 times to control for between-run variation in the topic models produced by each robot.

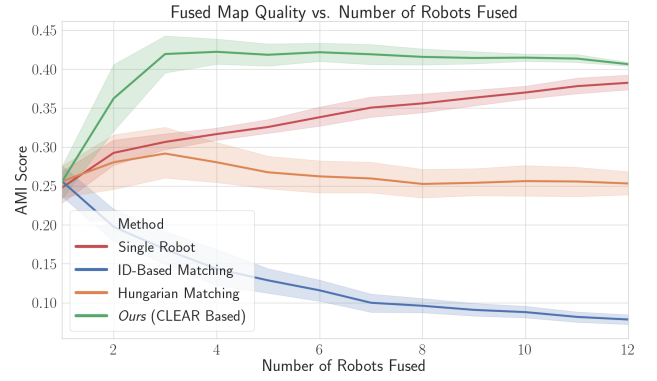
The topic model hyperparameters were set using a Bayesian Optimization algorithm [38] to find the values that resulted in the highest map quality, evaluated using a third simulated reef environment which was similar to Environment #1. The topic model vocabulary was the same one used in [19]. All code, instructions, hyperparameters, and datasets required to reproduce these experiments, as well as instructions to generate similar test environments, are available in the Sunshine repository.

A. Evaluating Semantic Map Quality

A useful fused semantic map $g(x)$ is one that is, at every location x , a *good predictor* of the ground truth label $f(x)$ defined by the human operator. We measure this predictive



(a) In environment #1, sparser and less varied phenomena (coral species) were spread throughout a uniformly sandy reef. The prevalence of sand in both shaded and well lit conditions tended to cause the single robot to develop two sand topics, reducing its performance.



(b) In environment #2, more variation between robots in the phenomena (coral species) and terrains observed meant that the Hungarian algorithm's assumptions were violated, leading to reduced performance compared to the proposed CLEAR-based approach.

Fig. 4. The AMI scores between the fused maps and corresponding ground truth maps extracted from the simulator demonstrate that the multi-robot distributed mapping system with CLEAR-based label matching outperforms all other multi-robot approaches. Error bars represent the 95% confidence interval of the mean score. The score of a single-robot that explored the same total area as the corresponding fused maps is shown in red, for comparison.

strength using Adjusted Mutual Information (AMI) [39], a normalized variant of Mutual Information (MI). MI represents the number of bits of information contained in one random variable that describe another, and is defined for discrete random variables $U \in \mathcal{U}, V \in \mathcal{V}$ as

$$\text{MI}(U, V) = \sum_{u \in \mathcal{U}} \sum_{v \in \mathcal{V}} p(u, v) \log \left(\frac{p(u)p(v)}{p(u, v)} \right). \quad (5)$$

Given the semantic maps f and g and a finite set of cells \mathcal{B} , we define the random variable B as a cell chosen uniformly at random from \mathcal{B} , and thus define the random variables $Y_f = f(B)$ and $Y_g = g(B)$ as the robot team's semantic label and the ground truth semantic label for B , respectively. The joint probability of these random variables is:

$$p(y_f, y_g) = \frac{1}{|\mathcal{B}|} \sum_{b \in \mathcal{B}} \mathbb{1}_{f(b)=y_f \wedge g(b)=y_g}, \quad (6)$$

which gives the probabilities $p(y_f), p(y_g)$ through marginalization. Denoting the entropy of random variables Y_f and Y_g as $H(Y_f)$ and $H(Y_g)$, the AMI is computed as

$$\text{AMI}(Y_f, Y_g) = \frac{\text{MI}(Y_f, Y_g) - \mathbb{E}[\text{MI}(Y_f, Y_g)]}{\max\{H(Y_f), H(Y_g)\} - \mathbb{E}[\text{MI}(Y_f, Y_g)]}. \quad (7)$$

A perfect AMI score of 1 indicates that the robot team's semantic map contains all information required to reproduce the ground truth semantic map of the same area. Conversely, a score of 0 indicates that the team's semantic map contains no more information about the corresponding part of the ground truth map than a randomly generated map is expected to. This is because the normalization subtracts out the expected mutual information $\mathbb{E}[\text{MI}(Y_f, Y_g)]$ between the labelings Y_f and Y_g , computed according to [39]. If a fused semantic map has a high AMI score with respect to the ground truth, then there is a consistent correspondence between each semantic label used by the robot team and each label used to produce the ground truth map. Thus, for high AMI scores, the human operator only needs to look at a few example images for any label in the fused map in order to determine that label's human-interpretable meaning.

B. Baseline Comparisons

The local maps were also fused using the ID-based matching and Hungarian matching approaches described in [23] to get baseline performance metrics. ID-based matching assumes that every robot observed the same phenomena in the same order, and so the first topic learned by one robot corresponds to the first topic learned by every other robot, and likewise for the remaining topics. The Hungarian approach only assumes that every robot observed the same phenomena, and finds the maximum similarity permutation between the first robot's topics and each additional robot's. This is a *sequential*, not multiway, matching approach because it compares topics belonging to a pair of robots at a time instead of considering all of the topics together. We are not aware of any previous baselines that used a multiway matching algorithm or did not assume that every robot observed the same phenomena.

Separate from the multi-robot experiments, a single robot was used to explore the same environments and independently build its own semantic map. It used the same topic model hyperparameters as the robots in the multi-agent experiment, but did not require any topic matching or map fusion. The single robot required Nt seconds to explore the same area that N robots explored in t seconds; the quality of the map it produced after exploring the same area as the N fused robots is reported in the results as "Single Robot".

VI. RESULTS & DISCUSSION

Figure 4 shows the performance of the proposed system using CLEAR, compared to using baseline matching solutions and to using a single robot. The performance was measured in terms of AMI of the fused map with the ground truth semantic map produced by AirSim. While the performance of the other matching algorithms declines as more robots are fused, the performance of the proposed matching solution *increases* or stays steady. This happens because CLEAR leverages redundant edges in the noisy association graph, added by additional robots, to help compensate for incorrect

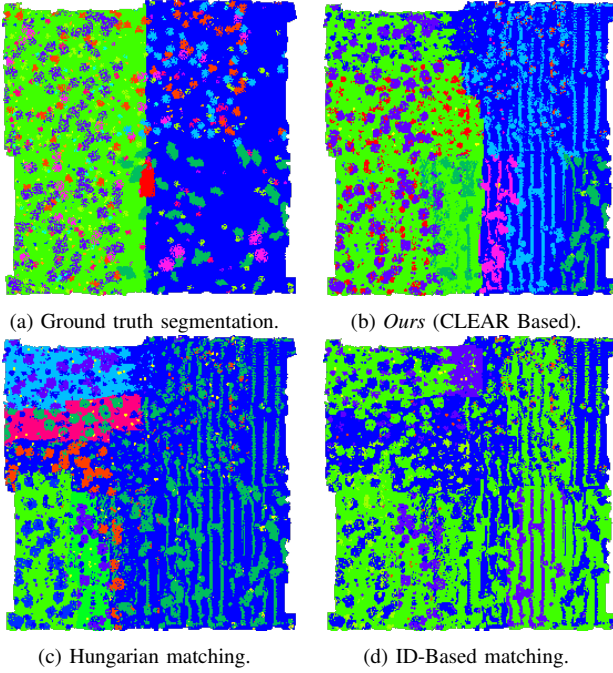


Fig. 5. Sample fused maps from each multi-robot matching approach with all 12 robots, alongside the ground truth segmentation, for Environment #2. Note that each map has been manually colored (1 color per label) with the same palette to ease comparison. Our approach most accurately captures the variation in terrain and coral species present in each quadrant.

edges. The number of incorrect edges at each vertex grows slower than the number of correct edges as topics are added, so our system is able to find a better solution when more robots' maps are fused.

In the first test environment, Fig. 4a, the fused map quality of the proposed approach goes from about 10% lower than the semantic map produced by a single robot to about 10% higher as more robots are fused. This is excellent performance considering that the robot team was able to map the entire environment in 1/12th the amount of time. Compared to the Hungarian matching approach, the proposed system achieves 23% higher AMI scores on average; as shown in Fig. 5, this is primarily because CLEAR is better suited to recognize when different robots have observed distinct phenomena. In the second test environment, Fig. 4b, this difference was magnified as there was very little in common between what any pair of robots observed. Table I summarizes numerical results for the map quality after fusing all 12 local maps together with each matching algorithm and for various similarity and distance metrics. The other figures shown used Cosine similarity for CLEAR matching and Euclidean (L2) distance as the Hungarian cost metric.

As seen in Fig. 6, as the team explores environment #2 the fused map quality is mostly constant after each robot has collected 125 images, i.e., covered half of its assigned area. This suggests that environment #2 would be most efficiently explored using 24 robots; in general, the optimal number will depend on the size and complexity of the environment.

The results presented are expected to generalize well to real world environments. Previous versions of BNP-ROST

TABLE I. Semantic Mapping Performance with 12 Robots.

Matching Alg.	Metric	MEAN AMI SCORE (STD. DEV.)	
		Env. #1	Env. #2
ID-Based	N/A	0.117 (0.035)	0.078 (0.016)
Hungarian	L1 Distance	0.297 (0.006)	0.216 (0.004)
	L2 Distance	0.313 (0.016)	0.253 (0.039)
	Cosine Distance	0.304 (0.011)	0.203 (0.015)
	TO Similarity	0.250 (0.002)	0.341 (0.003)
CLEAR	Cosine Similarity	0.384 (0.006)	0.406 (0.007)
Single Robot (No Matching)		0.344 (0.026)	0.382 (0.024)

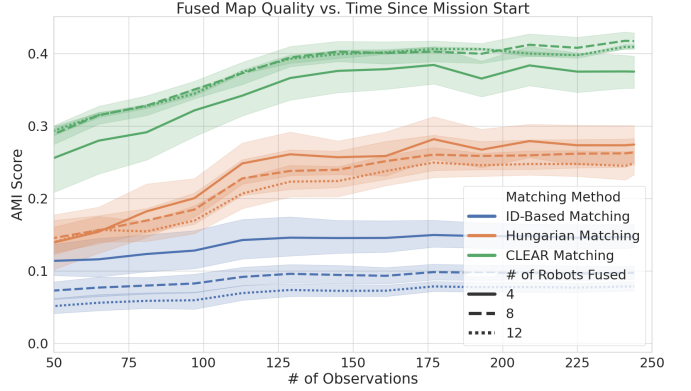


Fig. 6. The fused map quality varies throughout each experiment; shown here is how the fused map quality (AMI) changes as the robots explore environment #2. The increase in performance across all methods is caused by each robot's topic model improving over time with more data.

have demonstrated the ability to produce useful semantic maps of real-world environments while running on real robot hardware [19]. Furthermore, the compressed semantic maps and topic descriptors shared between robots are very small (typically around 10 to 100 kB) so they can be transmitted between robots in even severely bandwidth-constrained environments, like the deep sea [19], [40].

VII. CONCLUSIONS

We have presented a novel multi-robot distributed semantic mapping system that produces accurate semantic maps even when fusing maps from *many* robots and when each robot is building its unsupervised semantic model online with *no pre-training*. The proposed topic matching approach results in 20-60% higher map quality than pairwise Hungarian matching, with the largest gains in mapping complex and diverse environments, while also using less communication bandwidth than the previous state-of-the-art [23]. The fused maps are suitable for the human operator to use for mission summarization and informative path planning. We find that the fused maps approximate the quality of the best single-robot maps, hence further performance increases will likely come from improving the STM-based online semantic mapping component. The presented system for accurate topic matching over low-bandwidth enables novel multi-robot distributed autonomous exploration capabilities, such as cooperative-adaptive path planning and distributed reward learning, which will be explored in future work.

REFERENCES

- [1] D. F. Wolf and G. S. Sukhatme, "Semantic Mapping Using Mobile Robots," *IEEE Transactions on Robotics*, vol. 24, no. 2, pp. 245–258, Apr. 2008, conference Name: IEEE Transactions on Robotics.
- [2] G. Flaspohler, N. Roy, and Y. Girdhar, "Feature discovery and visualization of robot mission data using convolutional autoencoders and Bayesian nonparametric topic models," in *IEEE International Conference on Intelligent Robots and Systems*, 2017, pp. 1–8, arXiv: 1712.00028 ISSN: 21530866.
- [3] S. L. Bowman, N. Atanasov, K. Daniilidis, and G. J. Pappas, "Probabilistic data association for semantic SLAM," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. Singapore, Singapore: IEEE, May 2017, pp. 1722–1729, tex.ids: Bowman2017a.
- [4] K. Himri, P. Ridao, N. Gracias, A. Palomer, N. Palomeras, and R. Pi, "Semantic SLAM for an AUV using object recognition from point clouds," *IFAC-PapersOnLine*, vol. 51, no. 29, pp. 360–365, 2018, tex.ids: 2018 publisher: Elsevier.
- [5] J. Zhang, M. Gui, Q. Wang, R. Liu, J. Xu, and S. Chen, "Hierarchical Topic Model Based Object Association for Semantic SLAM," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 11, pp. 3052–3062, Nov. 2019.
- [6] M. Everett, J. Miller, and J. P. How, "Planning Beyond the Sensing Horizon Using a Learned Context," *arXiv:1908.09171 [cs]*, Oct. 2019, arXiv: 1908.09171.
- [7] S. Jamieson, J. P. How, and Y. Girdhar, "Active Reward Learning for Co-Robotic Vision Based Exploration in Bandwidth Limited Environments," in *2020 International Conference on Robotics and Automation (ICRA)*. Paris, France: IEEE, May 2020.
- [8] N. Sunderhauf, T. T. Pham, Y. Latif, M. Milford, and I. Reid, "Meaningful Maps With Object-Oriented Semantic Mapping," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. Vancouver, BC: IEEE, Sept. 2017, pp. 5079–5085, citation Key Alias: Sunderhauf.
- [9] Q.-H. Pham, B.-S. Hua, D. T. Nguyen, and S.-K. Yeung, "Real-time Progressive 3D Semantic Segmentation for Indoor Scene," *arXiv:1804.00257 [cs]*, Apr. 2019.
- [10] A. Rosinol, M. Abate, Y. Chang, and L. Carlone, "Kimera: an Open-Source Library for Real-Time Metric-Semantic Localization and Mapping," *arXiv:1910.02490 [cs]*, Oct. 2019, arXiv: 1910.02490.
- [11] Y. Nakajima and H. Saito, "Efficient Object-Oriented Semantic Mapping With Object Detector," *IEEE Access*, vol. 7, pp. 3206–3213, 2019.
- [12] M. Grinvald, F. Furrer, T. Novkovic, J. J. Chung, C. Cadena, R. Siegwart, and J. Nieto, "Volumetric Instance-Aware Semantic Mapping and 3D Object Discovery," *IEEE Robotics and Automation Letters*, vol. 4, no. 3, pp. 3037–3044, July 2019.
- [13] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet Allocation," *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003, arXiv: 1111.6189v1 ISBN: 9781577352815.
- [14] L. Cao and L. Fei-Fei, "Spatially coherent latent topic model for concurrent segmentation and classification of objects and scenes," *Proceedings of the IEEE International Conference on Computer Vision*, 2007, ISBN: 9781424416318.
- [15] X. Wang and E. Grimson, "Spatial Latent Dirichlet Allocation," in *Neural Information Processing Systems*, 2007, pp. 1–8.
- [16] Y. Girdhar, P. Giguère, and G. Dudek, "Autonomous adaptive exploration using realtime online spatiotemporal topic modeling," *The International Journal of Robotics Research*, vol. 33, no. 4, pp. 645–657, Apr. 2014.
- [17] Y. Girdhar, Walter Cho, M. Campbell, J. Pineda, E. Clarke, and H. Singh, "Anomaly detection in unstructured environments using Bayesian nonparametric scene modeling," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*, May 2016, pp. 2651–2656.
- [18] Y. Girdhar and G. Dudek, "Modeling curiosity in a mobile robot for long-term autonomous exploration and monitoring," *Autonomous Robots*, vol. 40, no. 7, pp. 1267–1278, 2016, arXiv: 1509.07975 Publisher: Springer US.
- [19] Y. Girdhar, L. Cai, S. Jamieson, N. McGuire, G. Flaspohler, S. Suman, and B. Claus, "Streaming Scene Maps for Co-Robotic Exploration in Bandwidth Limited Environments," in *2019 International Conference on Robotics and Automation (ICRA)*. Montréal, Canada: IEEE, May 2019, pp. 7940–7946.
- [20] Y. Girdhar and G. Dudek, "Gibbs Sampling Strategies for Semantic Perception of Streaming Video Data," *arXiv:1509.03242 [cs]*, Sept. 2015, arXiv: 1509.03242.
- [21] M. Thoma, "A Survey of Semantic Segmentation," pp. 1–16, 2016, arXiv: 1602.06541.
- [22] X. Xia and B. Kulis, "W-Net: A Deep Model for Fully Unsupervised Image Segmentation," 2017, arXiv: 1711.08506.
- [23] K. Doherty, G. Flaspohler, N. Roy, and Y. Girdhar, "Approximate Distributed Spatiotemporal Topic Models for Multi-Robot Terrain Characterization," in *Intelligent Robots and Systems (IROS)*, 2018.
- [24] H. W. Kuhn, "The Hungarian method for the assignment problem," *Naval Research Logistics Quarterly*, vol. 2, no. 1-2, pp. 83–97, Mar. 1955.
- [25] D. Pachauri, R. Kondor, and V. Singh, "Solving the multi-way matching problem by permutation synchronization," in *Advances in Neural Information Processing Systems*, 2013, pp. 1860–1868.
- [26] Y. Chen, L. Guibas, and Q. Huang, "Near-Optimal Joint Object Matching via Convex Relaxation," in *Proceedings of the 31st International Conference on Machine Learning*, vol. 32. Beijing, China: JMLR: W&CP, 2014, p. 9.
- [27] N. Hu, Q. Huang, B. Thibert, and L. Guibas, "Distributable Consistent Multi-object Matching," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. Salt Lake City, UT, USA: IEEE, June 2018, pp. 2463–2471.
- [28] J.-G. Yu, G.-S. Xia, A. Samal, and J. Tian, "Globally consistent correspondence of multiple feature sets using proximal Gauss–Seidel relaxation," *Pattern Recognition*, vol. 51, pp. 255–267, Mar. 2016.
- [29] X. Zhou, M. Zhu, and K. Daniilidis, "Multi-image Matching via Fast Alternating Minimization," in *2015 IEEE International Conference on Computer Vision (ICCV)*. Santiago, Chile: IEEE, Dec. 2015, pp. 4032–4040.
- [30] E. Maset, F. Arrigoni, and A. Fusiello, "Practical and Efficient Multi-view Matching," in *2017 IEEE International Conference on Computer Vision (ICCV)*. Venice: IEEE, Oct. 2017, pp. 4578–4586.
- [31] J. Yan, Z. Ren, H. Zha, and S. Chu, "A constrained clustering based approach for matching a collection of feature sets," June 2016, pp. 3832–3837, arXiv: 1606.03731.
- [32] K. Fathian, K. Khosoussi, Y. Tian, P. Lusk, and J. P. How, "CLEAR: A Consistent Lifting, Embedding, and Alignment Rectification Algorithm for Multi-View Data Association," *arXiv:1902.02256 [cs]*, July 2019, arXiv: 1902.02256.
- [33] D. A. Levin, Y. Peres, and E. L. Wilmer, *Markov Chains and Mixing Times, second edition*, 2nd ed. American Mathematical Society, 2017.
- [34] N. Aletras and M. Stevenson, "Measuring the Similarity between Automatically Generated Topics," in *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics, volume 2: Short Papers*. Gothenburg, Sweden: Association for Computational Linguistics, 2014, pp. 22–27.
- [35] Epic Games, "Unreal Engine," Oct. 2019.
- [36] Kelint, "Automatic Coral Generator," Apr. 2016, publisher: Unreal Engine Marketplace.
- [37] S. Shah, D. Dey, C. Lovett, and A. Kapoor, "AirSim: High-Fidelity Visual and Physical Simulation for Autonomous Vehicles," *arXiv:1705.05065 [cs]*, July 2017.
- [38] A. Cully, K. Chatzilygeroudis, F. Allocati, and J.-B. Mouret, "Limbo: A Flexible High-performance Library for Gaussian Processes modeling and Data-Efficient Optimization," *Journal of Open Source Software*, vol. 3, no. 26, p. 545, June 2018.
- [39] N. X. Vinh, J. Epps, and J. Bailey, "Information theoretic measures for clusterings comparison: is a correction for chance necessary?" in *Proceedings of the 26th International Conference on Machine Learning*, Montreal, Canada, 2009, p. 8.
- [40] J. W. Kaeli, J. J. Leonard, and H. Singh, "Visual summaries for low-bandwidth semantic mapping with autonomous underwater vehicles," in *2014 IEEE/OES Autonomous Underwater Vehicles (AUV)*. IEEE, Oct. 2014, pp. 1–7.