

# Using Variational Encoding and Generative Adversarial Learning for Inferring Missing Knowledge for Mobile Robots

Francesco Riccio  
Dept. of Engineering in  
Computer Science,  
Management and Automation,  
Sapienza University of Rome,  
Email: riccio@diag.uniroma1.it

Roberto Capobianco  
Dept. of Engineering in  
Computer Science,  
Management and Automation,  
Sapienza University of Rome,  
Email: capobianco@diag.uniroma1.it

Daniele Nardi  
Dept. of Engineering in  
Computer Science,  
Management and Automation,  
Sapienza University of Rome,  
Email: nardi@diag.uniroma1.it

**Abstract**—In artificial intelligence, look-ahead capabilities are key for successful planning and decision making. Unfortunately, robots have a limited perception horizon, which is additionally affected by significant noise. Moreover, the availability of world models is reduced and mostly limited to partial simulations. As a result, robots have to deal with unknown and missing knowledge during task execution, leading to costly re-planning routines. In literature, approaches to spatial knowledge prediction – beyond the sensory horizon – assume static environments and repetitive patterns (e.g. rectangular rooms). However, these assumptions are unrealistic. Hence, we propose a novel methodology that allows a robot to represent unknown spatial knowledge in dynamic and unstructured environments. In particular, we exploit generative adversarial networks and latent representations to (1) learn a distribution of spatial landmarks observed during task execution and to (2) generate missing information in real-time. In this paper, we describe the proposed approach and the obtained experimental results.

## I. INTRODUCTION

Spatial knowledge is fundamental to enable successful task execution in autonomous robots [23, 17], due to their planning and decision-making requirements. Several approaches assume that a map of the environment is given, or it can be generated before the deployment of a robot [9, 4] as a compact representation of its appearance and landmarks. In these cases, planning routines can be easily executed. However, often times this is not possible and, consequently, robots have to explore the surrounding environment while completing their tasks. This is the case, for example, in search and rescue scenarios [13], door-to-door delivery [14], planetary exploration [18], visual inspection [1] and mining [16].

In cases where a map, or a model of the world, is unknown or partially known, the robot has to build its model during its mission. In this paper, hence, we attack the problem of enabling an autonomous robot to explore an unknown environment, and represent portion of it not yet perceived through sensors. Several approaches consider the problem of estimating unknown parts of environment. These methods usually rely on frontier-based [17, 11] or gain-based [15, 24] techniques. The former exploit shape and location of frontiers, while the

latter reason about the expected information gain of visiting a given area. Only few contributions attempt to explicitly reconstruct the portion of the environment not (yet) observable by the robot, either with pre-trained models [19, 2], spoken instructions [6], or structure prediction, based on geometric features [3]. These methods require a pre-trained model of the world that tells the robot how to classify the environment and the expected structure. However, such approaches are often inaccurate and hardly generalize to dynamic environment.

In this paper, we extend and improve a previous formalization of GUESS [20] to enable a robot to *guess* the structure of the environment visited by the robot. GUESS is designed to refine and improve map prediction during the robot operation, by relying upon a Variational AutoEncoder (VAE) [12] paired with a Generative Adversarial Network (GAN) [8]. The VAE is used to learn a latent representation of the structure of the environment, while the GAN is used to generate expected future observations in the form of latent features. Then the VAE decoder collects such samples and reconstruct the generated observation to the original feature representation.

Our ultimate goal is to provide a robot with the ability to represent portions of the world beyond its sensory horizon, in order to support reasoning in partially known environments. In this paper, we describe our approach, and we demonstrate how it can be used to predict laser scans of a mobile robot. Our contribution is to extend our novel approach and to improve its prediction accuracy in inferring missing knowledge.

## II. GUESS

GUESS is a deep iterative algorithm based on variational autoencoding [12] and adversarial generative learning [8]. In this contribution we improve the GUESS [20] architecture and its accuracy in generating future observations. However, we preserve the most desirable feature of being agnostic with respect to the type of data and the application of deployment. In particular, we evaluate the GUESS architecture in inferring spatial knowledge modeled through 2D laser scans.

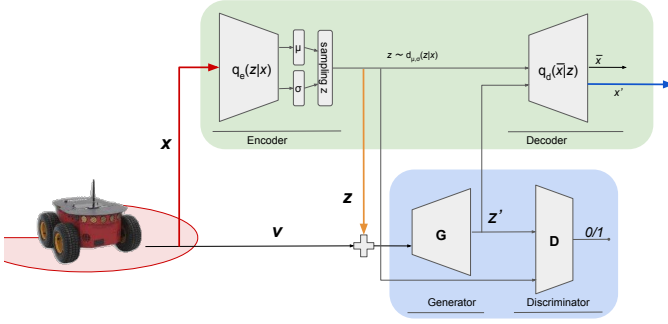


Fig. 1. GUESS system scheme. Data is collected through the robot sensor (red arrow) and given to the system. The VAE is refined to generate latent encodings (orange arrow) paired with velocity commands and collected afterwards. The concatenated feature vector of encodings and velocities  $\{z, v\}$  is then fed to the GAN for refinement and data prediction (blue arrow). The generated latent descriptor  $z'$  is then fed to the decoder to reconstruct a laser scan which is used to extend the robot sensory horizon.

We configure our algorithm to accumulate lasers scans, while navigating the environment and to predict laser scans, expected to be perceived at a given point  $t'$  in the future. The algorithm assumes no prior knowledge about the environment and it continuously performs online aggregation [21] of new data samples to allow the robot to quickly adapt and generalize to new environments.

To infer missing knowledge not (yet) perceivable through robot sensors, we exploit generative adversarial learning. It has been demonstrated that GANs achieve remarkable results in learning data distributions to generate new samples and/or completing missing knowledge [5, 10]. However, since they need to be fed with large datasets and configured with complex networks with high dimensionality inputs and outputs, GANs cannot be easily deployed in dynamic and real-time applications. In robotic settings, for example, large datasets and computational costs are assumptions that cannot always be satisfied. Thus, in order to exploit the potential of GANs in robotics, we need to alleviate dimensionality constraints and network complexity – yet guaranteeing robust performance. To meet such a compelling requirement, we reduce the dimensionality of the generative network input by learning a lower-dimensionality latent representation of it. We exploit an autoencoder [12] to learn the latent – and more compact – representation of the input data. Specifically, a variational autoencoder is used to learn the distribution of input data over the latent space, which allows us to perform sampling directly in the latent space and to generate batches of inputs of lower dimensionality to be fed to the generative network.

Operationally, to enable GUESS to evaluate sequences  $T$  of time-correlated inputs  $x_{t:t+T}$  (e.g. 2D laser scans) and to predict a future observation  $x'$ , we substitute raw sensor data  $x_{t:t+T}$  with sequences of their latent representation  $z_{t:t+T}$  that are iteratively learned through variational autoencoding. This allows for a significant reduction in the dimensionality of the input, yet preserving an informative representation. In the following sections, we describe how the two building blocks

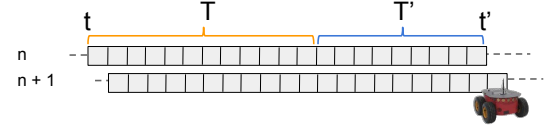


Fig. 2. Buffered time moving-window.

of GUESS are trained and then we describe its algorithmic steps.

### A. Algorithm Overview

In this contribution, we focus on validating the key insights of the GUESS system, and we evaluate its performance by implementing it on a mobile robot equipped with a 2D laser sensor. Fig. 1 summarizes the overall GUESS system. Data is collected through the robot sensors and fed to the VAE for encoding. Then, their encoded representation is (1) collected in a dataset of time-correlated sequences  $x, z$ , (2) paired with velocities commands  $v$ , and (3) provided to the generative network in order to predict  $z'$ . The predicted latent representation is then fed to the VAE decoder which reconstructs the laser scan  $x'$ . It is worth remarking that, in our former formalization [20], the GAN was directly generating the laser scan  $x'$ . However, with such a configuration the prediction is too noisy and scattered which limit a practical use. In this contribution, we improve the GAN predictions by iteratively performing the following steps. GUESS:

- 1) initializes the networks, an empty dataset *Dataset*, and a buffer *Buffer* by collecting sensor readings over  $T + T'$  timesteps.  $T$  and  $T'$  are two metaparameters of the algorithm which determine the length of the sequences of time-correlated samples, and the time interleave for the next scan prediction  $x'$ . Fig. 2 illustrates the moving time window that we use to collect new training samples to train the generator and the discriminator of the GAN. The robot collects sensor readings  $x$  and velocity commands  $v$  over the time interval  $T + T'$ . Then, it (1) pairs the first  $T$  samples  $x_{t:t+T}$  with velocity  $v_{t:t+T}$  commands registered simultaneously, while it uses the latent representation  $z'$  of sample  $x'$  at  $t+T+T'$  to create a reference for the target distribution used to refine the discriminator. At the next iteration the moving window is shifted by one step and a new training sample is collected iteratively. Conversely, we can exploit any sensor reading to train the variational autoencoder. In this implementation, we select the laser scan observed at  $t + T + T'$ .
- 2) iteratively drops the first sample of the buffer and append new sensor readings by shifting the time window;
- 3) randomly selects a batch of  $B$  elements from the dataset. We randomly sample from a dataset since data coming from the sensor stream is sequentially correlated and non-i.i.d. Hence, to properly train the networks, we decorrelate our samples using an experience replay dataset;

- 4) updates the variational autoencoder by minimizing the loss function with the set of random samples  $\mathbf{x}'$  and refining the multivariate Gaussian distribution parameters  $\mu$  and  $\sigma$  characterizing the latent space. Then, GUESS uses the updated parameters to generate latent samples  $\mathbf{z}$  from the time-correlated minibatch  $\mathbf{x}$  by sampling from the latent distribution as in Equation 1;
- 5) exploits the latent samples  $\mathbf{z}$  and the registered velocity commands  $\mathbf{v}$  as condition to the generative network to predict the expected sensor reading latent representation  $\mathbf{z}'_g$  after  $T'$  timesteps. The original latent representation instead, is used as label for the discriminator as a sample drawn from the target distribution.
- 6) as in standard adversarial learning, GUESS updates both the generator and discriminator by applying label smoothing[22].
- 7) lastly, it forwards generated latent representation  $\mathbf{z}'_g$  to the VAE decoder in order to reconstruct the predicted future laser scans  $\mathbf{x}'$ .

GUESS adopts a variational autoencoder to reduce the dimensionality of input laser scans  $x$ . Differently from standard autoencoders, VAEs have the desirable feature of learning a parametric distribution of input data in the latent space. This is a key requirement for generative models, as it allows for random sampling in the latent space. By design, VAEs feature this property by forking the output of a standard encoder network in two separate layers – representing the mean  $\mu$  and standard deviation  $\sigma$  of each feature in the latent space. Then, the encoder uses a sampler over the posterior probability distribution  $q_e(\mathbf{z}|\mathbf{x})$ , where  $\mathbf{x}$  is the raw laser scans retrieved from the mobile platform and  $\mathbf{z}$  their latent representation. Each sample  $z \sim q_e(\mathbf{z}|\mathbf{x})$  from the posterior distribution, is sampled in accordance with Equation 1 by using the re-parameterization trick.

$$\mathbf{z} = \mu + [\exp(\log \sigma/2) \odot \epsilon] \quad (1)$$

where  $\epsilon \sim \mathcal{N}(0, I)$  is drawn from a normal distribution and  $\odot$  is an element-wise product [12].

Then, GUESS exploits generative adversarial learning to learn a target distribution of data samples and to generate missing knowledge. GANs are composed by two components mimicking a zero-sum minimax game – a protagonist and an antagonist – that are designed to learn data distributions and to validate samples drawn from such a learned distribution. GUESS, in particular, exploits two deep networks to implement both the generator  $G$  (the protagonist) and the discriminator  $D$  (the antagonist). In our setting, the GAN is designed to generate a latent observation  $z'$  conditioned over a sequence of time-correlated inputs  $\mathbf{z}$  and velocity commands  $\mathbf{v}$  and it configures the generator to learn the conditional p.d.f.  $z' \sim g(z' | \mathbf{z}, \mathbf{v})$ .

To this end, GUESS makes use of a stochastic gradient descent training that, at each iteration, refines the two networks with a two phase update. More in detail, to train the GAN GUESS iteratively:

- collects  $b$  random samples from a dataset to compose minibatches. In this implementation, each of the  $b$  elements is a tuple of time-correlated latent samples  $\mathbf{z}$  drawn directly from the latent space of the VAE and robot velocity commands  $\mathbf{v}$  registered altogether with laser scans by matching the robot platform timesteps. Such a batch is then used to condition the generator and to output a synthetic sample in the VAE latent space  $z' \sim g(z' | \mathbf{z}, \mathbf{v})$ ;
- updates both the discriminator and the generator to improve the prediction of sample drawn from the learned distribution  $z_g \sim g(z_g | \mathbf{z}, \mathbf{v})$ ;
- returns the generated sample  $z_g \sim g(z)$  and exploits the VAE decode and reconstruct a laser scan  $x'$  which is forwarded to the robot platform.

### III. EXPERIMENTAL EVALUATION

In this section we validate the key insights of GUESS and demonstrate that our approach is feasible and practical in robotics. We validate GUESS in an indoor office environment where a robot is tasked to navigate while collecting data in real-time and refining its predictions iteratively.

#### A. Setup

Experiments have been configured within ROS melodic, running on a single Intel Core i7-5700HQ core, with CPU@2.70GHz and 16GB of RAM. GUESS is configured to evaluate  $T = 16$  time-correlated sequence in batch of  $B = 32$  samples, and to predict sensor readings with a  $T' = 10$  interleave. It is important to highlight that, while GUESS performs inference at each iteration of the systems. The data collection and model refinement is performed with a frequency of  $1Hz$ . Hence, in this setup, each 26 seconds a new training sample is generated, and GUESS predicts laser scans expected to be perceived after 10 seconds of operation. The simulated robot platform is equipped with a laser sensor that emits  $R = 512$  beams and navigates the environment with a maximum cruise speed of  $0.5m/s$ . Moreover, both the VAE and GAN have been configured to satisfy dimensionality constraints, yet showing reliable results.

a) *Variational Autoencoder*: The encoder and the decoder of the VAE have a dual deep structure composed by three layers of 512, 128, and 32 neurons respectively – where  $R=512$  is the dimensionality of the raw laser scan, and  $L=32$  represents the dimension of the latent encodings. In fact, GUESS is able to reduce the dimensionality of the input representation by a factor of 16 (R/L) yet preserving information. Each layer is fully connected, while activation layers are configured with leaky ReLUs among hidden layers, and tanh for the decoder output. The VAE is trained by minimizing the objective function, with an Adam optimizer configured with learning rate  $lr = 1e-4$ .

b) *Generative Adversarial Network*: The generator and discriminator of the GAN architecture have a four layer structure. The generator is configured to accept  $T \times (L + V)$  sequences of time-correlated encodings paired with velocity

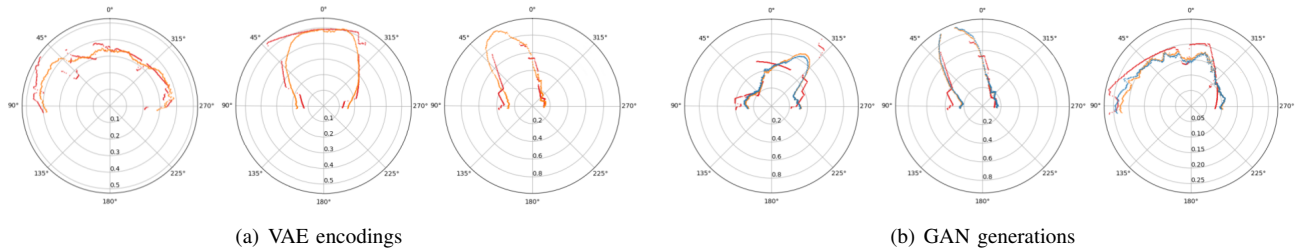


Fig. 3. GUESS training and inference. 3(a) and 3(b) show the performance and generalization capabilities of the two networks after 500 iterations of the algorithm. In 3(a) the original scans (red) and VAE reconstructed (orange) are reported, while in 3(b) in the original scan (red), its encoding (orange), and the generation and decoding (blue).

commands – where  $V$  is a two features vector reporting the linear and angular velocity of the non-holonomic platform. Moreover, the generator is composed by a cascade of three de-convolutional layers with increasing filters depth of 64, 128 and 512; all layers are activated with leaky ReLUs except for the output layer which is characterized by a tanh activation. The discriminator, instead, is composed by four dense fully connected layers of 512, 128, 32, 1 neurons where the first three layers are activated by leaky ReLUs and the output layer by a sigmoid function. Both networks are trained in an adversarial manner in a zero-sum game, with two Adam optimizer with  $lr = 1e-4$  learning rate for the generator and  $lr = 1e-4$  for the discriminator.

### B. Learning Performance

In this experiment a training session is simulated, the system is randomly initialized and ran for 500 iterations. All components of the system are trained simultaneously as previously described. Fig. 3(a) and Fig. 3(b) show the inference of the VAE and GAN over samples collected by the robot during its mission.

Analyzing the performance of the VAE component we observe a reconstruction error that decreases and reaches  $\sim 20$ cm of accuracy already after few minutes of execution. At the 500th iteration, three different samples are selected from the robot experience and reported in Fig. 3(a). From the left to right, samples represent: (1) a corner of a large room, (2) a small office, and (3) a corridor. The red line represents the original laser scan  $x$ , while the orange line represents the VAE reconstructed scan  $\bar{x}$ . It is worth remarking that the latent encodings used in this experiment reduce the dimensionality of the inputted laser scans to  $L=16$  and – as shown in the figure – yet maintain an informative representation.

Also the learning of the GAN evolves as expected and the prediction error decreases as expected reaching  $\sim 20$ cm as in the previous case. Such a result suggests that the GAN predictions are limited by the encoding capability of the VAE. This is also confirmed by the samples reported in Fig. 3(b). From left to right, scans represent (1) a narrow passage, (2) a corridor and (3) a corner of a large room. It is interesting to notice that the GAN perfectly learns to generate samples from the target distribution – which is the latent space of the VAE. Noticeably, the orange and the blue line are very

similar. We did not exclude from the visualization overlapping laser beams in order to display the full generated scans. In spite of that, the planner will eventually exclude that part of the generated scan. In the next step of our development we want to integrate features from the original laser scans as the GAN target distribution. Such a setting will provide more information about sharp edges and corners, thus improving on fidelity of the reconstruction of scans. Overall, it is worth noticing that Fig. 3 confirms that a more compact – but equally informative – representation of the input data, can be used as a surrogate to significantly reduce dimensionality of data and dimension of the networks of the GAN.

### IV. CONCLUSION

In this paper we improved on a novel architecture to achieve online training and inference of missing spatial knowledge for a mobile robot. We configured GUESS to predict 2D laser scans, and we validated the ability of the architecture to exploit a learned latent representation of environment. The experimental evaluation shows the promising prediction errors achieved both in VAE reconstruction and GAN prediction – which confirms that online data aggregation does not invalidate the training. Moreover, the GUESS architecture is agnostic with respect the observation  $x$  – that in this case we instantiate to laser scans. Hence, it is possible to deploy GUESS to learn to different data observation types (e.g. laser scans, RGB images, depth images, point clouds) and have a GUESS instance for each of them. Nevertheless, the performance of the approach is still limited, and it has to be further improved in order to support decision-making. Our ultimate goal is to extend the robot horizon to predict complete spatial entities (e.g. rooms, corridors) and to generalize the architecture to objects in the environment. Finally, several future directions can be pursued to improve the GUESS performance. For example, we are evaluating LSTM [7] recurrent networks to improve predictions of time-correlated samples.

### REFERENCES

- [1] Alireza Ahrary, Amir A.F. Nassiraei, and Masumi Ishikawa. A study of an autonomous mobile robot for a sewer inspection system. *Artificial Life and Robotics*, 11(1):23–27, Jan 2007. ISSN 1614-7456. doi: 10.1007/s10015-006-0392-x.

- [2] S. Bai, J. Wang, F. Chen, and B. Englot. Information-theoretic exploration with bayesian optimization. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1816–1822, Oct 2016. doi: 10.1109/IROS.2016.7759289.
- [3] H. J. Chang, C. S. G. Lee, Y. Lu, and Y. C. Hu. P-slam: Simultaneous localization and mapping with environmental-structure prediction. *IEEE Transactions on Robotics*, 23(2):281–293, April 2007. ISSN 1552-3098. doi: 10.1109/TRO.2007.892230.
- [4] Michael Jae-Yoon Chung\*, Andrzej Pronobis\*, Maya Cakmak, Dieter Fox, and Rajesh P. N. Rao. Autonomous question answering with mobile robots in human-populated environments. In *Proceedings of the 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Daejeon, Korea, October 2016. doi: 10.1109/IROS.2016.7759146.
- [5] Jiankang Deng, Shiyang Cheng, Niannan Xue, Yuxiang Zhou, and Stefanos Zafeiriou. Uv-gan: Adversarial facial uv map completion for pose-invariant face recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [6] Felix Duvallat, Matthew R. Walter, Thomas Howard, Sachithra Hemachandra, Jean Oh, Seth Teller, Nicholas Roy, and Anthony Stentz. *Inferring Maps and Behaviors from Natural Language Instructions*, pages 373–388. Springer International Publishing, 2016. ISBN 978-3-319-23778-7. doi: 10.1007/978-3-319-23778-7\_25.
- [7] D. Eck and J. Schmidhuber. Finding temporal structure in music: blues improvisation with lstm recurrent networks. In *Proceedings of the 12th IEEE Workshop on Neural Networks for Signal Processing*, pages 747–756, Sep. 2002. doi: 10.1109/NNSP.2002.1030094.
- [8] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014.
- [9] G. Grisetti, C. Stachniss, and W. Burgard. Improved techniques for grid mapping with rao-blackwellized particle filters. *IEEE Transactions on Robotics*, 23(1):34–46, Feb 2007. ISSN 1552-3098. doi: 10.1109/TRO.2006.889486.
- [10] Kapil D. Katyal, Katie M. Popek, Chris Paxton, Joseph L. Moore, Kevin C. Wolfe, Philippe Burlina, and Gregory D. Hager. Occupancy map prediction using generative and fully convolutional networks for vehicle navigation. *CoRR*, abs/1803.02007, 2018.
- [11] Mohammad Al Khawaldah and Andreas Nüchter. Enhanced frontier-based exploration for indoor environment with multiple robots. *Advanced Robotics*, 29(10):657–669, 2015. doi: 10.1080/01691864.2015.1015443.
- [12] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2013.
- [13] Stefan Kohlbrecher, Johannes Meyer, Thorsten Graber, Karen Petersen, Uwe Klingauf, and Oskar von Stryk. Hector open source modules for autonomous mapping and navigation with rescue robots. In Sven Behnke, Manuela Veloso, Arnoud Visser, and Rong Xiong, editors, *RoboCup 2013: Robot World Cup XVII*, pages 624–631, Berlin, Heidelberg, 2014. Springer Berlin Heidelberg. ISBN 978-3-662-44468-9.
- [14] Raffaele Limosani, Raffele Esposito, Alessandro Manzi, Giancarlo Teti, Filippo Cavallo, and Paolo Dario. Robotic delivery service in combined outdoor/indoor environments: technical analysis and user evaluation. *Robotics and Autonomous Systems*, 103:56 – 67, 2018. ISSN 0921-8890. doi: <https://doi.org/10.1016/j.robot.2018.02.001>.
- [15] E. Nelson and N. Michael. Information-theoretic occupancy grid compression for high-speed information-based exploration. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4976–4982, Sept 2015. doi: 10.1109/IROS.2015.7354077.
- [16] Tobias Neumann, Alexander Ferrein, Stephan Kallweit, and Ingrid Scholl. Towards a mobile mapping robot for underground mines. In *Proc. of the 2014 PRASA, RobMech and AfLaI Int. Joint Symposium, Cape Town, South Africa*, 2014.
- [17] S. Obwald, M. Bennewitz, W. Burgard, and C. Stachniss. Speeding-up robot exploration by exploiting background information. *IEEE Robotics and Automation Letters*, 1(2):716–723, July 2016. ISSN 2377-3766. doi: 10.1109/LRA.2016.2520560.
- [18] K. Otsu, A. Agha-Mohammadi, and M. Paton. Where to look? predictive perception with applications to planetary exploration. *IEEE Robotics and Automation Letters*, 3(2): 635–642, April 2018. ISSN 2377-3766. doi: 10.1109/LRA.2017.2777526.
- [19] Andrzej Pronobis and Rajesh P. N. Rao. Learning deep generative spatial models for mobile robots. In *Proceedings of the 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Vancouver, BC, Canada, September 2017. doi: 10.1109/IROS.2017.8202235.
- [20] F. Riccio, R. Capobianco, and D. Nardi. Guess: Generative modeling of unknown environments and spatial abstraction for robots. In *International Conference on Autonomous Agents and Multi-Agent Systems 2020, (AAMAS20)*, 2020.
- [21] Stéphane Ross, Geoffrey J Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *International Conference on Artificial Intelligence and Statistics*, pages 627–635, 2011.
- [22] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS’16*, pages 2234–2242, USA, 2016. Curran

Associates Inc. ISBN 978-1-5108-3881-9.

[23] Cyrill Stachniss and Wolfram Burgard. Exploring unknown environments with mobile robots using coverage maps. In *Proceedings of the 18th International Joint Conference on Artificial Intelligence, IJCAI'03*, pages 1127–1132, San Francisco, CA, USA, 2003. Morgan

Kaufmann Publishers Inc.

[24] Joan Vallvé and Juan Andrade-Cetto. Potential information fields for mobile robot exploration. *Robotics and Autonomous Systems*, 69:68 – 79, 2015. ISSN 0921-8890. doi: <https://doi.org/10.1016/j.robot.2014.08.009>. Selected papers from 6th European Conference on Mobile Robots.