# Learning models for visual 3D localization with implicit mapping

**Dan Rosenbaum, Frederic Besse, Fabio Viola, Danilo J. Rezende, S. M. Ali Eslami**
DeepMind, London, UK
{danro,fbesse,fviola,danilor,aeslami}@google.com

## Abstract

We propose a formulation of visual localization that does not require construction of explicit maps in the form of point clouds or voxels. The goal is to learn an implicit representation of the environment at a higher, more abstract level, for instance that of objects. To study this approach we consider procedurally generated Minecraft worlds, for which we can generate visually rich images along with camera pose coordinates. We first show that Generative Query Networks (GQNs) enhanced with a novel attention mechanism can capture the visual structure of 3D scenes in Minecraft, as evidenced by their samples. We then apply the models to the localization problem, investigating both generative and discriminative approaches, and compare the different ways in which they each capture task uncertainty. Our results show that models with implicit mapping are able to capture the underlying 3D structure of visually complex scenes, and use this to accurately localize new observations, paving the way towards future applications in sequential localization. Supplementary video available at `https://youtu.be/iHEXX5wXbCI`.

## 1 Introduction

The problem of identifying the position of the camera that captured an image of a scene has been studied extensively for decades [17, 2]. With applications in domains such as robotics and autonomous driving, a considerable amount of engineering effort has been dedicated to developing systems for different versions of the problem. One formulation, often referred to as simply 'localization', assumes that a map of the 3D scene is provided in advance and the goal is to localize any new image of the scene relative to this map. A second formulation, commonly referred to as 'Simultaneous Localization and Mapping' (SLAM), assumes that there is no prespecified map of the scene, and that it should be estimated concurrently with the locations of each observed image. Research on this topic has focused on different aspects and challenges including: estimating the displacement between frames in small time scales, correcting accumulated drifts in large time scales (also known as 'loop closure'), extracting and tracking features from the observed images [12], reducing computational costs of inference (graph-based SLAM [5]) and more.

Although this field has seen huge progress in recent years, performance is still limited by a number of factors [2]. One key limitation stems from reliance on hand-engineered representations of various components of the problem (e.g. key-point descriptors as representations of images and occupancy grids as a representations of the map), and it has been argued that moving towards systems that operate using more abstract representations is likely to be beneficial [16, 10]. Consider the map, which is typically either part of the input provided to the system (in localization) or part of the output expected of the system (in SLAM). This forces algorithm designers to define its structure explicitly in advance, e.g. either as 3D positions of keypoints, or an occupancy grid of a pre-specified resolution. It is not always clear what the optimal representation is for a given domain, and the need to pre-define this structure is restricting and can lead to sub-optimal performance. Even though machine learning
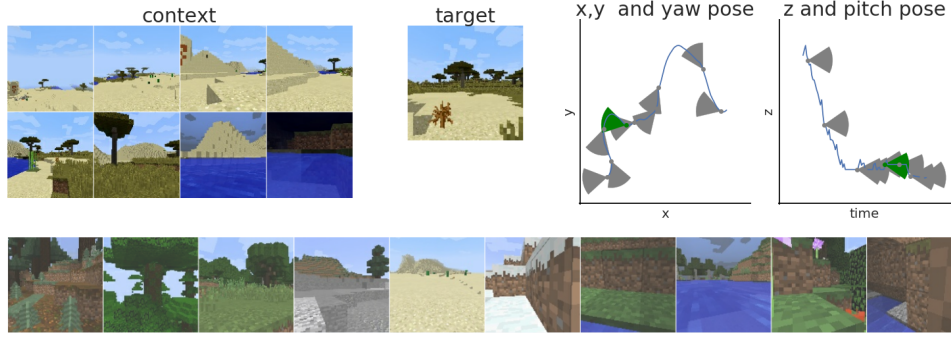
Figure 1: The Minecraft random walk dataset for localization in 3D scenes. We generate random trajectories in the Minecraft environment, and collect images along the trajectory labelled by the camera pose coordinates (x,y,z yaw and pitch). Bottom: Images from random scenes. Top: The localization problem setting - for a new trajectory in a new scene, given a set of images along the trajectory and their corresponding camera poses (the 'context'), predict the camera pose of an additional observed image in the trajectory (the 'target', shown in green).

has seen increased use for localization and SLAM problems in recent years [22, 6, 13], most methods still rely on pre-specified map representations. Agents trained with reinforcement learning have been demonstrated to solve navigation tasks that implicitly require localization and mapping [11, 3, 1, 20], suggesting that it is possible to learn these abilities without such pre-specification. Other methods where these abilities emerge in more specific settings include [8, 21, 23, 18].

In this work, we consider the problem of localization with *implicit* mapping. Like SLAM, we assume that a map is not available in advance, but unlike SLAM, we do not expect the system to produce an explicit map either. Given a collection of 'context' images with known camera poses, the task is defined as finding the relative camera pose of a new image which we call the 'target'. This task can also be viewed as doing loop closure without an explicit map. We consider deep models that learn implicit representations of the map, and can therefore capture abstract descriptions of the environment and exploit abstract cues for the localization problem.

A recent generative model that has shown promise in learning representations for 3D scene structure is the Generative Query Network (GQN) [4]. The GQN is conditioned on different views of a 3D scene, and trained to generate images from new views of the same scene. Although the model demonstrates generalization of the representation and rendering capabilities to held-out 3D scenes, the experiments are in simple environments that only exhibit relatively well-defined, small-scale structure, e.g. a few objects in a room, or small mazes with a few rooms.

In this paper, we ask two questions: 1) Does the GQN scale to more complex, visually rich environments? 2) Can the model be used for localization? To this end, we propose a dataset of random walks in the game Minecraft (figure 1), using the Malmo platform [7]. We compare a generative and discriminative approach, and show that coupled with a novel attention mechanism, the GQN can capture the 3D structure of complex Minecraft scenes, and can use this for accurate localization.

In summary, our contributions are the following: 1) We investigate the task of localization with implicit mapping and propose a dataset for it. 2) We enhance the GQN using a novel sequential attention mechanism, and show that it can capture the complexity of 3D scenes in Minecraft evidenced by generated samples. 3) We show that the GQN can be used for localization in both generative and discriminative directions, and compare the way each of the approaches captures the uncertainty.

## 2 Data for localization with implicit mapping

In order to study machine learning approaches to visual localization with implicit mapping, we create a dataset of random walks in a visually rich simulated environment. We use the Minecraft environment through the open source Malmo platform [7]. Minecraft is a popular computer game where a player can wander in a virtual world, and the Malmo platform allows us to connect to the game engine, control the player and record data. Since the environment is procedurally generated,
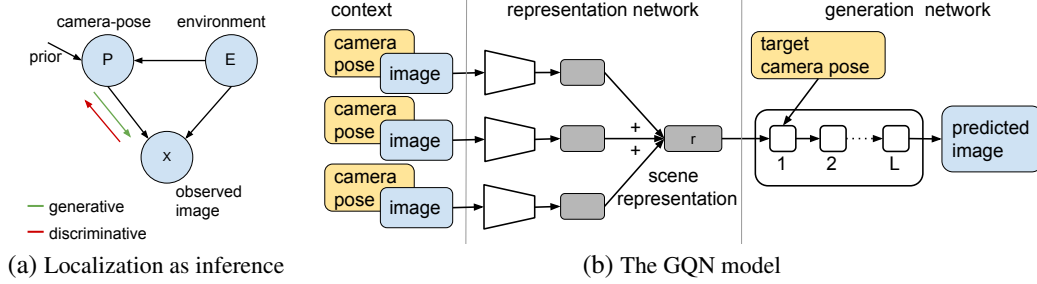
(a) Localization as inference           (b) The GQN model

Figure 2: Localization as probabilistic inference (a). The observed images $X$ depend on the environment $E$ and the camera pose $P$. To predict $Pr(P|X)$, in the generative approach, we use a model of $Pr(X|P)$ (green) and apply Bayes' rule. In the discriminative approach we directly train a model of $Pr(P|X)$ (red). In both cases the environment is implicitly modeled given a context of (image, camera pose) pairs $C = \{x_i, p_i\}$. In the GQN model (b), the context is processed by a representation network and the resulting scene representation is fed to a recurrent generation network with L layers, predicting an image given a camera pose.

we can generate sequences of images in an unlimited number of different scenes. Figure 1 shows the visual richness of the environment and the diversity of the scene types including forests, lakes, deserts, mountains etc. The images are not realistic but nevertheless contain many details such as trees, leaves, grass, clouds, and also different lighting conditions.

Our dataset consists of 6000 sequences of 100 images each generated by a blind exploration policy based on simple heuristics (see appendix A for more details). We record images at a resolution of $128 \times 128$ although for all the experiments in this paper we downscale to $32 \times 32$. Each image is recorded along with a 5 dimensional vector of camera pose coordinates consisting of the $x, y, z$ position and yaw and pitch angles (omitting the roll as it is constant). Figure 1 also demonstrates the localization task: For every sequence we are given a context of images along with their camera poses, and an additional target image with unknown pose. The task is to find the camera pose of the target image.

## 3 Model

The localization problem can be cast as an inference task in the probabilistic graphical model presented in figure 2a. Given an unobserved environment $E$ (which can also be thought of as the map), any observed image $X$, depends on the environment $E$ and on the pose $P$ of the camera that captured the image $Pr(X|P, E)$. The camera pose $P$ depends on the environment and perhaps some prior, e.g. a noisy odometry sensor. In this framing, localization is an inference task involving the posterior probability of the camera pose which can be computed using Bayes' rule:

$$Pr(P|X, E) = \frac{1}{Z} Pr(X|P, E) Pr(P|E) \tag{1}$$

In our case, the environment is only implicitly observed through the context, a set of image, and camera pose pairs, which we denote by $C = \{x_i, p_i\}$. We can use the context to estimate the environment $\hat{E}(C)$ and use the estimate for inference:

$$Pr(P|X, C) = \frac{1}{Z} Pr(X|P, \hat{E}(C)) Pr(P|\hat{E}(C)) \tag{2}$$

Relying on the context to estimate the environment $E$ can be seen as an empirical Bayes approach, where the prior for the data is estimated (in a point-wise manner) at test time from a set of observations. In order to model the empirical Bayes likelihood function $Pr(X|P, \hat{E}(C))$ we consider the Generative Query Network (GQN) model [4] (figure 2b). We use the same model as described in the GQN paper, which has two components: a representation network, and a generation network. The representation network processes each context image along with its corresponding camera pose coordinates using a 6-layer convolutional neural network. The camera pose coordinates are combined to each image by concatenating a 7 dimensional vector of $x, y, z, sin(yaw), cos(yaw), sin(pitch), cos(pitch)$ to the

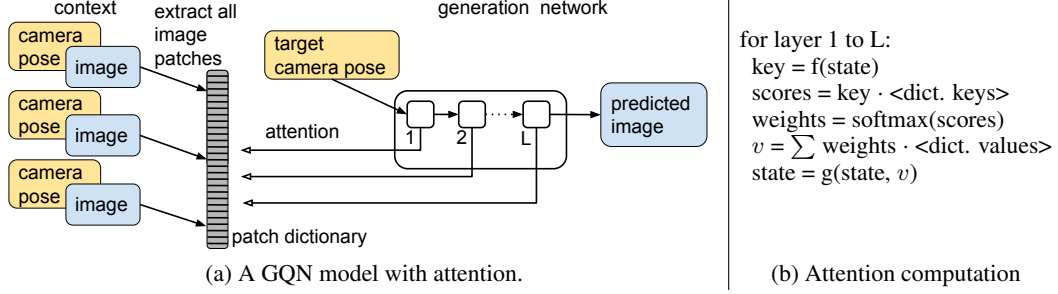| | |
|---|---|
| (a) A GQN model with attention. | (b) Attention computation |

Figure 3: (a) Instead of conditioning on a parametric representation of the scene, all patches from context images are stored in a dictionary, and each layer of the generation network can access them through an attention mechanism. In each layer, the key is computed as a function of the generation network's recurrent state, and the result $v$ is fed back to the layer (b).

features of the neural network in the middle layer. The output of the network for each image is then added up resulting in a single scene representation. Conditioned on this scene representation and the camera pose of a new target view, the generation network which consists of a conditional latent-variable model DRAW [4] with 8 recurrent layers, generates a probability distribution of the target image $Pr_{GQN}(X|P,C)$. We use this to model $Pr(X|P, \hat{E}(C))$, where the scene representation $r$ serves as the point-wise estimate of the environment $\hat{E}(C)$.

Given a pre-trained GQN model as a likelihood function, and a prior over the camera pose $Pr(P|C)$, localization can be done by maximizing the posterior probability of equation 2 (MAP inference):

$$\arg\max_P Pr(P|X,C) = \arg\max_P \log Pr_{GQN}(X|P,C) + \log P(P|C) \qquad (3)$$

With no prior, we can simply resort to maximum likelihood, omitting the second term above. The optimization problem can be a hard problem in itself, and although we show some results for simple cases, this is usually the main limitation of this generative approach.

## 3.1 A discriminative approach

While the generative approach is to learn models in the $P \rightarrow X$ direction and invert them using Bayes' rule as described above, a discriminative approach would be to directly learn models in the $X \rightarrow P$ direction. In order to model $Pr(P|X,C)$ directly, we construct a 'reversed-GQN' model, by using the same GQN model described above, but 'reversing' the generation network, i.e. using a model that is conditioned on the target image, and outputs a distribution over the camera pose. The output of the model is a set of categorical distributions over the camera pose space, allowing multi-modal distributions to be captured. In order to keep the output dimension manageable, we split the camera pose coordinate space to four components: 1) the $x, y$ positions, 2) the $z$ position, 3) the yaw angle, and 4) the pitch angle. To implement this we process the target image using a convolutional neural network, where the scene representation is concatenated to the middle layer, and compute each probability map using an MLP. For more details on the model see appendix B.

Since it outputs probability maps with a certain structure, the discriminative model is already less 'implicit' than the generative model, making it closer to previous methods which use pre-defined map structure and discriminatively train models to predict location [13, 1]. Although the discriminative approach is a more direct way to solve the problem, the generative approach has some aspects that can prove advantageous: it learns a model in the causal direction, which might be easier to capture; and it learns a more general problem not specific to any particular task and free from any prior on the solution. In this way it can be used for tasks which it was not trained for (for example predict the x,y location in a different spatial resolution).

## 3.2 Attention

One limitation of the GQN and reversed-GQN models is the fact that the representation network compresses all the information from the context images and poses to a single global representation vector. Although the GQN has shown a remarkable ability to capture the 3D structure of scenes
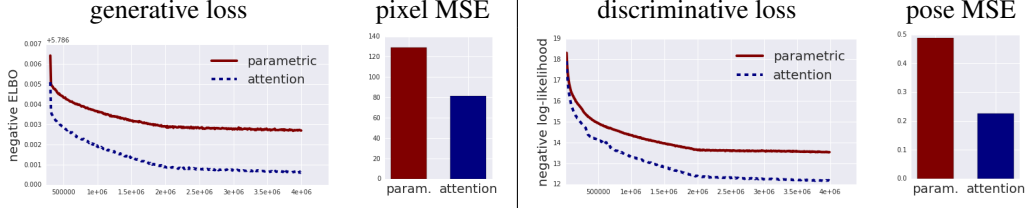
4

Figure 4: The loss and predictive MSE for both the generative direction and the discriminative direction. The attention model results in lower loss and lower MSE in both directions.

and generate samples of new views that are almost indistinguishable from ground truth, it has done so in fairly simple scenes with only a few objects in each scene and limited variability. Since we are interested here in more complex scenes with much more visual variability, it is unlikely for a simple representation network with a fixed size scene representation to be able to retain all the important information of the scene. Following previous work on using attention for making conditional predictions [19, 14, 13], we propose another variation of the model that relies on attention rather than a parametric representation of the scene.

We enhance the GQN model using an attention mechanism on patches inspired by [14] (figure 3), where instead of passing the context images through a representation network, we extract from each of the images all the 8 by 8 patches and concatenate to each the camera poses and the 2D coordinates of the patch within the image. For each patch we also compute a key using a convolutional neural network, and place all patches and their corresponding keys in a large dictionary. The GQN's and reversed-GQN's decoders use soft attention to access the patch dictionary multiple times, and condition their prediction on the attention results.

For the GQN, in the generation network based on the DRAW model with 8 recurrent layers, a key is computed at every layer from the current state, and used to perform soft attention on the patch dictionary. The attention is performed using a dot-product with all the patches' keys, normalizing the results to one using a softmax, and using them as the weights in a weighted sum of all patches. The resulting vector is then used by the generation network in the same way as the global representation vector is used in the parametric GQN model. In the reversed GQN decoder, we implement a similar sequential attention mechanism with 10 layers. In both models the attention is sequential and the key in each layer depends on the result of the attention at the previous layer. This allows the model to use a more sophisticated attention strategy, which depends both on the initial query and on patches and coordinates that previous layers attended to. See appendix B for more details.

## 4 Training results

We train the models on the Minecraft random walk data where each sample consists of 21 images divided to 20 context images and 1 target image. The images are drawn randomly from the sequence to reduce the effect of the prior on the target image camera pose due to the sequentiality of the original captured images. The loss we optimize in the generative direction is the negative variational lower bound on the log-likelihood over the target image, since the GQN contains a DRAW model with latent variables. For the discriminative direction with the reversed GQN which is fully deterministic, we minimize the negative log-likelihood over the target camera pose.

Figure 4 shows the training curves for the generative and discriminative directions. In both cases we compare between the standard parametric model and the attention model. We also compare the MSE of the trained models, computed in the generative direction by the L2 distance between the predicted image and the ground-truth image, and in the discriminative direction by the L2 distance between the most likely camera pose and the ground truth ones. All results are computed on a held out test set containing scenes that were not used in training, and are similar to results on the training set. A notable result is that the attention models improve the performance significantly. This is true for both the loss and the predictive MSE, and for both the discriminative and generative directions.

Figure 5 shows generated image samples from the generative model, and the predicted distribution over the camera pose space from the reversed GQN model. All results are from attention based models, and each image and camera pose map comes from a different scene from a held out test set,

ground truth

samples
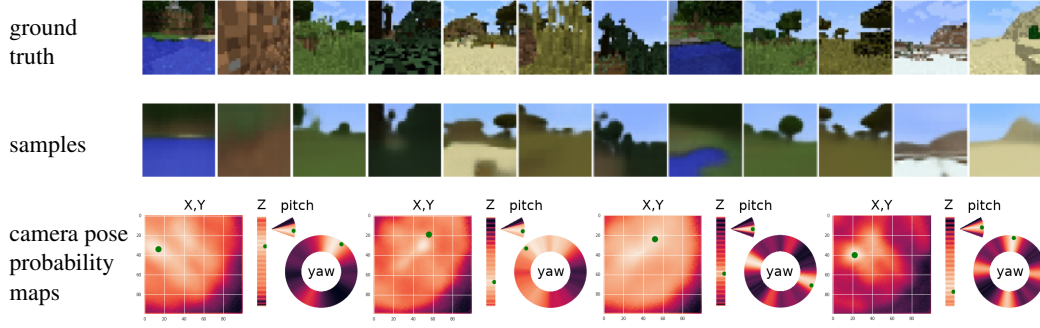
camera pose probability maps

Figure 5: Generated samples from the generative model (middle), and the whole output distribution for the discriminative model (bottom). Both were computed using the attention models, and each image and pose map is from a different scene conditioned on 20 context images. The samples capture much of the structure of the scenes including the shape of mountains, the location of lakes, the presence of large trees etc. (see supplementary video for more samples). The distribution of camera pose computed with the reversed GQN shows the model's uncertainty, and the distribution's mode is usually close to the ground truth (green dot).

conditioned on 20 context images. Image samples are blurrier than the ground truth but they capture the underlying structure of the scene, including the shape of mountains, the location of lakes, the presence of big trees etc (see supplementary video for more samples). For the reversed GQN model, the predicted distribution shows the uncertainty of the model, however the mode of the distribution is usually close to the ground truth camera pose. This is by itself an interesting result suggesting an easy way to perform localization even in complex scenes such as Minecraft. The results show that some aspects of the camera pose are easier than others. Namely, predicting the height $z$ and the pitch angle are much easier than the $x, y$ location and yaw angle. At least for the pitch angle this can be explained by the fact that it can be estimated from the target image alone, not relying on the context. As for the yaw angle, we see a frequent pattern of high probability areas around multiples of 90 degrees. This is due to the cubic nature of the Minecraft graphics, allowing the estimate of the yaw angle up to any 90 degrees rotation given the target image only. This last phenomena is a typical one when comparing discriminative methods to generative ones, where the former is capable of exploiting shortcuts in the data.

# 5 Attention

One interesting feature of models with attention, is that they can be analyzed by observing where they are attending, shedding light on some aspects of the way they work. Figure 6 shows the attention in three different scenes. Each row shows 1) the target image, 2) the x,y trajectory containing the poses of all context images (in black) and the pose of the target image (in green), and 3) the context images with a red overlay showing the patches with high attention weights. In each row we show the context images twice, once with the attention weights of the generative model, and once with the attention weights of the discriminative model. In both cases we show the total attention weights summed over all recurrent layers.

A first point to note is that the attention weights are sparse, focusing on a small number of patches in a small number of images. The attention focuses on patches with high contrast like the edges between the sky and the ground. In these aspects, the model has learned to behave similarly to the typical components of hand crafted localization methods - detecting informative feature points (e.g. Harris corner detection), and constructing a sparse computation graph by pruning out the irrelevant images (e.g. graph-SLAM). A second point to note is the difference between the generative and discriminative attention, where the first concentrates on fewer images, covering more space within them, and the second is distributed between more images. Since the generative model is queried using a camera pose, and the discriminative model using an image, it is perhaps not surprising that the resulting attention strategies tend to be position-based and appearance-based respectively. See supplementary video for more examples.
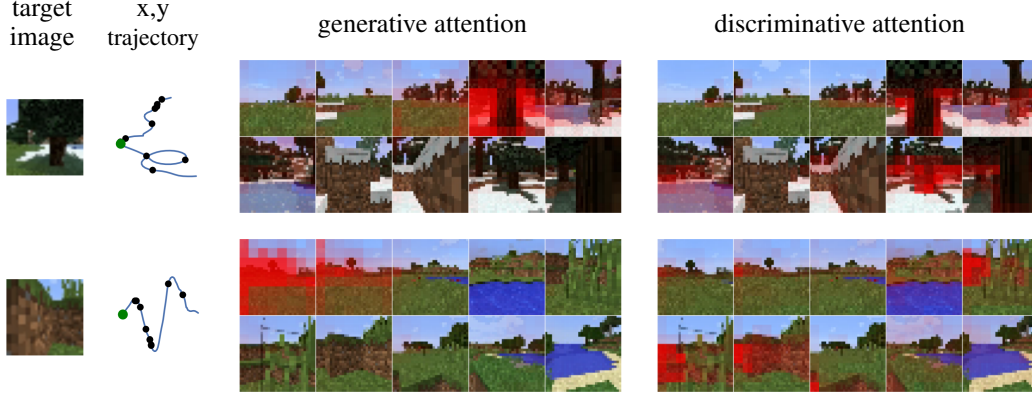
Figure 6: Attention over the context images in the generative and discriminative directions. The total attention weights are shown as a red overlay, using the same context images for both directions. Similar to hand-crafted feature point extraction and graph-SLAM, the learned attention is sparse and prunes out irrelevant images. While the generative attention is mainly position based, focusing on one or two of the nearest context images, the discriminative attention is based more on appearance, searching all context images for patches resembling the target. See also supplementary video.

# 6   Localization

We compare the generative and discriminative models ability to perform localization, i.e. find the camera pose of a target image, given a context of image and camera pose pairs. For the discriminative models we simply run a forward pass and get the probability map of the target's camera pose. For the generative model, we fix the target image and optimize the likelihood by searching over the camera pose used to query the model. We do this for both the x,y position, and the yaw, using grid search with the same grid values as the probability maps that are predicted by the discriminative model. The optimization of the x,y values, and the yaw values is done separately while fixing all other dimensions of the camera pose to the ground truth values. We do this without using a prior on the camera pose, essentially implementing maximum likelihood inference rather than MAP.

Figure 7 shows the results for a few held-out scenes comparing the output of the discriminative model, with the probability map computed with the generative model. For both directions we use the best trained models which are the ones with attention. We see that in most cases the maximum likelihood estimate (magenta star) is a good predictor of the ground truth value (green circle). However it is interesting to see the difference in the probability maps of the generative and discriminative directions. One aspect of this difference is that while in the discriminative direction the x,y values near the context points tend to have high probability, for the generative model it is the opposite. This can be explained by the fact that the discriminative model captures the prior of the data, consisting of the assumption that the target image was taken near context images. The generative model cannot capture this prior, sometimes resulting in an opposite effect where a new image with is more likely to be taken from an unexplored location. See supplementary document for a figure with more examples.

Table 1 shows a quantitative comparison of models with and without attention, with varying number of context points. The first table compares the MSE of the maximum likelihood estimate of the x,y position and the yaw angle, by computing the square distance from the ground truth values (again, for the generative model we fix all other dimensions of the pose to the ground truth values). The table shows that as expected, more context images result in better estimates. We also see that the generative models estimates have a lower MSE for all context sizes. This is true even though the generative model does not capture the prior of the data and can therefore make very big mistakes by estimating the target pose far from the context poses. The second table compares the log probability of the ground truth location for the x,y position and yaw angles (by using the value in the probability map cell which is closest to the ground truth). Here we see that the discriminative model results in higher probability. This can be explained again by the fact that the generative model does not capture the prior of the data and therefore distributes a lot of probability mass across all of the grid. We validate this claim by computing the log probability in the vicinity of all context poses (summing the

Table 1: Comparing different models for localization with a varying number of context images. Using attention ('+att') and increasing the number of context points lead to improved results. Even without a prior on the data, the generative model with attention results in better MSE scores. However the discriminative model with attention assigns higher probability to the ground-truth location since the prior allows it to concentrate more mass in the vicinity of the context poses.

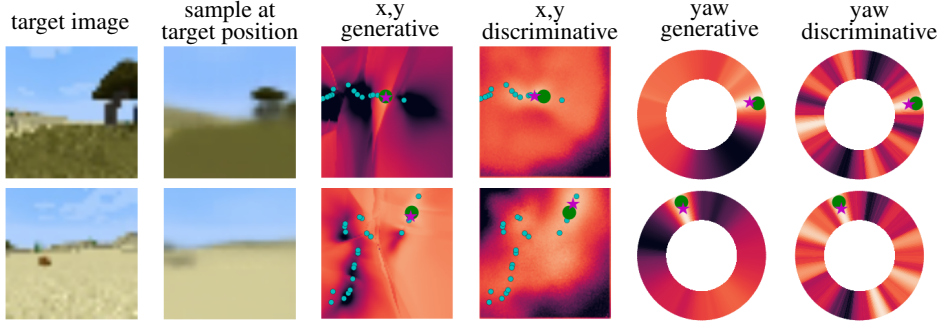| | MSE | | | | | | log-probability | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | x,y position | | | yaw angle | | | x,y position | | | yaw angle | | |
| context | 5 | 10 | 20 | 5 | 10 | 20 | 5 | 10 | 20 | 5 | 10 | 20 |
| disc | 0.27 | 0.25 | 0.27 | 0.25 | 0.21 | 0.17 | -8.14 | -7.24 | -7.08 | -3.27 | -3.08 | -3.02 |
| gen | 0.41 | 0.36 | 0.35 | 0.16 | 0.14 | 0.15 | -10.19 | -9.08 | -8.55 | -5.95 | -5.16 | -4.74 |
| disc+att | 0.17 | 0.12 | 0.09 | 0.21 | 0.15 | 0.11 | -7.65 | **-6.47** | **-5.74** | **-3.56** | **-3.29** | **-3.08** |
| gen+att | **0.11** | **0.08** | **0.07** | **0.06** | **0.06** | **0.05** | **-7.56** | -6.62 | -6.00 | -4.48 | -4.06 | -3.66 |



Figure 7: Localization with the discriminative and generative models - target image, a sample drawn using the ground-truth pose, and probability maps for x,y position and yaw angle (bright=high prob.). The generative maps are computed by querying the model with all possible pose coordinates, and the discriminative maps are simply the model's output. The poses of the context images are shown in cyan, target image in green and maximum likelihood estimates in magenta. The generative maps are free from the prior on poses in the training data giving higher probability to unexplored areas.

probability of a 3 by 3 grid around each context point) and indeed get a log probability of -4.1 for the generative direction and -1.9 for the discriminative (using attention and 20 context images).

# 7 Discussion

We have proposed a formulation of the localization problem that does not involve an explicit map, where we can learn models with implicit mapping in order to capture higher level abstractions. We have enhanced the GQN model with a novel attention mechanism, and showed its ability to capture the underlying structure of visually complex scenes demonstrated by generating samples of new views of the scene. We showed that GQN-like models can be used for localization in a generative and discriminative approach, and described the advantages and disadvantages of each.

When comparing the approaches for localization we have seen that the generative model is better in capturing the underlying uncertainty of the problem and is free from the prior on the pose that is induced by the training data. However, the clear disadvantage of the generative approach is that it requires performing optimization at test time, while the discriminative model can be used by running one forward pass. We believe that future research should focus on combining the approaches which can be done in many different ways, e.g. using the discriminative model as a proposal distribution for importance sampling, or accelerating the optimization using amortized inference by training an inference model on top of the generative model [15].

While localization relative to a context of images as formulated here can already be used as a 'loop closure' component in a SLAM system, we think that further research along the lines we presented, can lead to improved methods for full SLAM in a sequential setting, where the context of images is constructed online. Coupled with an increasing availability of real data, this can lead to efficient methods for real-world applications.

# References

[1] A. Banino, C. Barry, B. Uria, C. Blundell, T. Lillicrap, P. Mirowski, A. Pritzel, M. J. Chadwick, T. Degris, J. Modayil, et al. Vector-based navigation using grid-like representations in artificial agents. *Nature*, page 1, 2018.

[2] 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*, 32(6):1309–1332, 2016.

[3] C. J. Cueva and X.-X. Wei. Emergence of grid-like representations by training recurrent neural networks to perform spatial localization. *arXiv preprint arXiv:1803.07770*, 2018.

[4] S. M. A. Eslami, D. Jimenez Rezende, F. Besse, F. Viola, A. S. Morcos, M. Garnelo, A. Ruderman, A. A. Rusu, I. Danihelka, K. Gregor, D. P. Reichert, L. Buesing, T. Weber, O. Vinyals, D. Rosenbaum, N. Rabinowitz, H. King, C. Hillier, M. Botvinick, D. Wierstra, K. Kavukcuoglu, and D. Hassabis. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018.

[5] G. Grisetti, R. Kummerle, C. Stachniss, and W. Burgard. A tutorial on graph-based slam. *IEEE Intelligent Transportation Systems Magazine*, 2(4):31–43, 2010.

[6] S. Gupta, D. Fouhey, S. Levine, and J. Malik. Unifying map and landmark based representations for visual navigation. *arXiv preprint arXiv:1712.08125*, 2017.

[7] M. Johnson, K. Hofmann, T. Hutton, and D. Bignell. The malmo platform for artificial intelligence experimentation.

[8] A. Kendall, M. Grimes, and R. Cipolla. Posenet: A convolutional network for real-time 6-dof camera relocalization. In *Computer Vision (ICCV), 2015 IEEE International Conference on*, pages 2938–2946. IEEE, 2015.

[9] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[10] J. McCormac, A. Handa, A. Davison, and S. Leutenegger. Semanticfusion: Dense 3d semantic mapping with convolutional neural networks. In *Robotics and Automation (ICRA), 2017 IEEE International Conference on*, pages 4628–4635. IEEE, 2017.

[11] P. Mirowski, R. Pascanu, F. Viola, H. Soyer, A. J. Ballard, A. Banino, M. Denil, R. Goroshin, L. Sifre, K. Kavukcuoglu, et al. Learning to navigate in complex environments. *arXiv preprint arXiv:1611.03673*, 2016.

[12] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos. Orb-slam: a versatile and accurate monocular slam system. *IEEE Transactions on Robotics*, 31(5):1147–1163, 2015.

[13] E. Parisotto, D. S. Chaplot, J. Zhang, and R. Salakhutdinov. Global pose estimation with an attention-based recurrent network. *arXiv preprint arXiv:1802.06857*, 2018.

[14] S. Reed, Y. Chen, T. Paine, A. v. d. Oord, S. Eslami, D. J. Rezende, O. Vinyals, and N. de Freitas. Few-shot autoregressive density estimation: Towards learning to learn distributions. 2017.

[15] D. Rosenbaum and Y. Weiss. The return of the gating network: combining generative models and discriminative training in natural image priors. In *Advances in Neural Information Processing Systems*, pages 2683–2691, 2015.

[16] R. F. Salas-Moreno, R. A. Newcombe, H. Strasdat, P. H. Kelly, and A. J. Davison. Slam++: Simultaneous localisation and mapping at the level of objects. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 1352–1359. IEEE, 2013.

[17] S. Thrun, W. Burgard, and D. Fox. *Probabilistic robotics*. MIT press, 2005.

[18] S. Tulsiani, A. A. Efros, and J. Malik. Multi-view consistency as supervisory signal for learning shape and pose prediction. *arXiv preprint arXiv:1801.03910*, 2018.

[19] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, et al. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems*, pages 3630–3638, 2016.

[20] G. Wayne, C.-C. Hung, D. Amos, M. Mirza, A. Ahuja, A. Grabska-Barwinska, J. Rae, P. Mirowski, J. Z. Leibo, A. Santoro, et al. Unsupervised predictive memory in a goal-directed agent. *arXiv preprint arXiv:1803.10760*, 2018.

[21] A. R. Zamir, T. Wekel, P. Agrawal, C. Wei, J. Malik, and S. Savarese. Generic 3d representation via pose estimation and matching. In *European Conference on Computer Vision*, pages 535–553. Springer, 2016.

[22] J. Zhang, L. Tai, J. Boedecker, W. Burgard, and M. Liu. Neural slam. *arXiv preprint arXiv:1706.09520*, 2017.

[23] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, volume 2, page 7, 2017.

# A The Minecraft random walk dataset for localization

In order to generate a dataset, we chose the Minecraft environment since it allows generating unlimited labeled data, and also because it is an interesting environment by itself, containing visually rich scenes, developed not for the purpose of research but rather as a game with a goal of recreating the experience of exploration in new worlds, where navigation and localization are key ingredients.

To generate the data we use the Malmo platform [7], an open source project allowing access to the Minecraft engine. Our dataset consists of 6000 sequences of random walks consisting of 100 images each. We use 1000 sequences as a held-out test set. The sequences are generated using a simple heuristic-based blind policy as follows:

We generate a new world and a random initial position, and wait for the agent to drop to the ground.

For 100 steps:

1. Make a small random rotation and walk forward.
2. With some small probability make a bigger rotation.
3. If no motion is detected jump and walk forward, or make a bigger rotation
4. Record image and camera pose.

We prune out sequences where there was no large displacement or where all images look the same (e.g. when starting in the middle of the ocean, or when quickly getting stuck in a hole), however in some cases our exploration policy results in close up images of walls, or even underwater images which might make the localization challenging.

We use 5 dimensional camera poses consisting of $x, y, z$ position and yaw and pitch angles (omitting the roll as it is constant). We record images in a resolution of $128 \times 128$ although for all the experiments in this paper we downscale to $32 \times 32$. We normalize the x,y,z position such that most scenes contain values between -1 and 1.

# B Model details

The basic model that we use is the Generative Query Network (GQN) as described in [4]. In the representation network, every context image is processed by a 6 layer convolutional neural network (CNN) outputing a representation with spatial dimension of $8 \times 8$ and $64$ channels. We add the camera pose information of each image after 3 layers of the CNN by broadcasting the values of $x, y, z, sin(yaw), cos(yaw), sin(pitch), cos(pitch)$ to the whole spatial dimension and concatenate as additional 7 channels. The output of each image is added up to a single scene representation of $8 \times 8 \times 64$. The architecture of the CNN we use is: [k=2,s=2,c=32] $\rightarrow$ [k=3,s=1,c=32] $\rightarrow$ [k=2,s=2,c=64] $\rightarrow$ [k=3,s=1,c=32] $\rightarrow$ [k=3,s=1,c=32] $\rightarrow$ [k=3,s=1,c=64], where for each layer 'k' stands for the kernel size, 's' for the stride, and 'c' for the number of output channels.

In the generation network we use the recurrent DRAW model as described in [4] with $8$ recurrent layers and representation dimension of $8 \times 8 \times 128$ (used as both the recurrent state and canvas dimensions). The model is conditioned on both the scene representation and the camera pose query by injecting them (using addition) to the biases of the LSTM's in every layer. The output of the generation network is a normal distribution with a fixed standard deviation of $0.3$.

In order to train the model, we optimize the negative variational lower bound on the log-likelihood, using Adam [9] with a batch size of 36, where each example comes from a random scene in Minecraft, and contains 20 context images and 1 target image. We train the model for 400M iterations, and anneal the output standard deviation starting from $1.5$ in the first iteration, down to $0.3$ in the 300K'th iteration, keeping it constant in further iterations.

**Reversed GQN**

The reversed GQN model we use as a discriminative model is based on the GQN described above. We use the same representation network, and a new decoder that we call the localization network, which is queried using the target image, and outputs a distribution of camera poses (see figure 8). In order to make the output space manageable, we divide it to 4 log-probability maps:
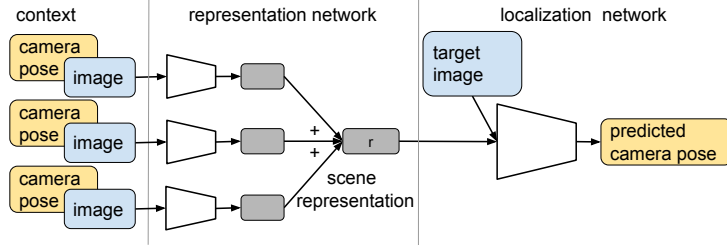
Figure 8: The reversed GQN model. A similar architecture to GQN, where the generation network is 'reversed' to form a localization network, queried using a target image, predicting its camera pose.

1. A matrix over $x, y$ values between -1 and 1, quantized in a $0.02$ resolution.
2. A vector over $z$ values between -1 and 1, quantized in a $0.02$ resolution.
3. A vector over yaw values between -180 and 180 degrees, quantized in a 1 degree resolution.
4. A vector over pitch values between -20 and 30 degrees, quantized in a 1 degree resolution.

The localization network first processes the image query using a CNN with the following specifications: [k=3,s=2,c=32] $\rightarrow$ [k=3,s=2,c=64] $\rightarrow$ [k=3,s=1,c=64] $\rightarrow$ [k=1,s=1,c=64] $\rightarrow$ [k=1,s=1,c=64], resulting in an intermediate representation of $8 \times 8 \times 64$. Then it concatenates the scene representation computed by the representation network and processes the result using another CNN with: [k=3,s=1,c=64] $\rightarrow$ [k=3,s=1,c=64] $\rightarrow$ [k=3,s=1,c=64] $\rightarrow$ [k=5,s=1,c=4]. Each of the 4 channels of the last layer is used to generate one of the 4 log-probability maps using an MLP with 3 layers, and a log-softmax normalizer. Figure 5 in the main paper shows examples of the resulting maps.

**Attention GQN**

In the attention GQN model, instead of a scene encoder as described above, all $8 \times 8 \times 3$ patches are extracted from the context images with an overlap of 4 pixels and placed in a patch dictionary, along with the camera pose coordinates, the 2D patch coordinates within the image (the x,y, position in image space), and a corresponding key. The keys are computed by running a CNN on each image resulting in a $8 \times 8 \times 64$ feature map, and extracting the (channel-wise) vector corresponding to each pixel. The architecture of the CNN is [k=2,s=2,c=32] $\rightarrow$ [k=3,s=1,c=32] $\rightarrow$ [k=2,s=2,c=64] $\rightarrow$ [k=1,s=1,c=32] $\rightarrow$ [k=1,s=1,c=32] $\rightarrow$ [k=1,s=1,c=64], such that every pixel in the output feature map corresponds to a $4 \times 4$ shift in the original image. The keys are also concatenated to the patches.

Given the patch dictionary, we use an attention mechanism within the recurrent layers of the generation network based on DRAW as follows. In each layer we use the recurrent state to compute a key using a CNN of [k=1,s=1,c=64] $\rightarrow$ [k=1,s=1,c=64] followed by spatial average pooling. The key is used to query the dictionary by computing the dot product with all dictionary keys, normalizing the results using a softmax. The normalized dot products, are used as weights in a weighted sum over all dictionary patches and keys (using the keys as additional values). See pseudo-code in figure 3b. The result is injected to the computation of each layer in the same way as the scene representation is injected in the standard GQN. Since the state in each layer of the recurrent DRAW depends on the computation of the previous layer, the attention becomes sequential, with multiple attention keys where each key depends on the result of the attention with the previous key.

In order to implement attention in the discriminative model, and make it comparable to the sequential attention in DRAW, we implement a similar recurrent network used only for the attention. This is done by adding 10 recurrent layers between layer 5 and 6 of the localization network's CNN. Like in DRAW, a key is computed using a 2 layer CNN of [k=1,s=1,c=64] $\rightarrow$ [k=1,s=1,c=64] followed by spatial average pooling. The key is used for attention in the same way as in DRAW, and the result is processed through a 3 layer MLP with 64 channels and concatenated to the recurrent state. This allows the discriminative model to also use a sequential attention strategy where keys depend on the attention in previous layers.
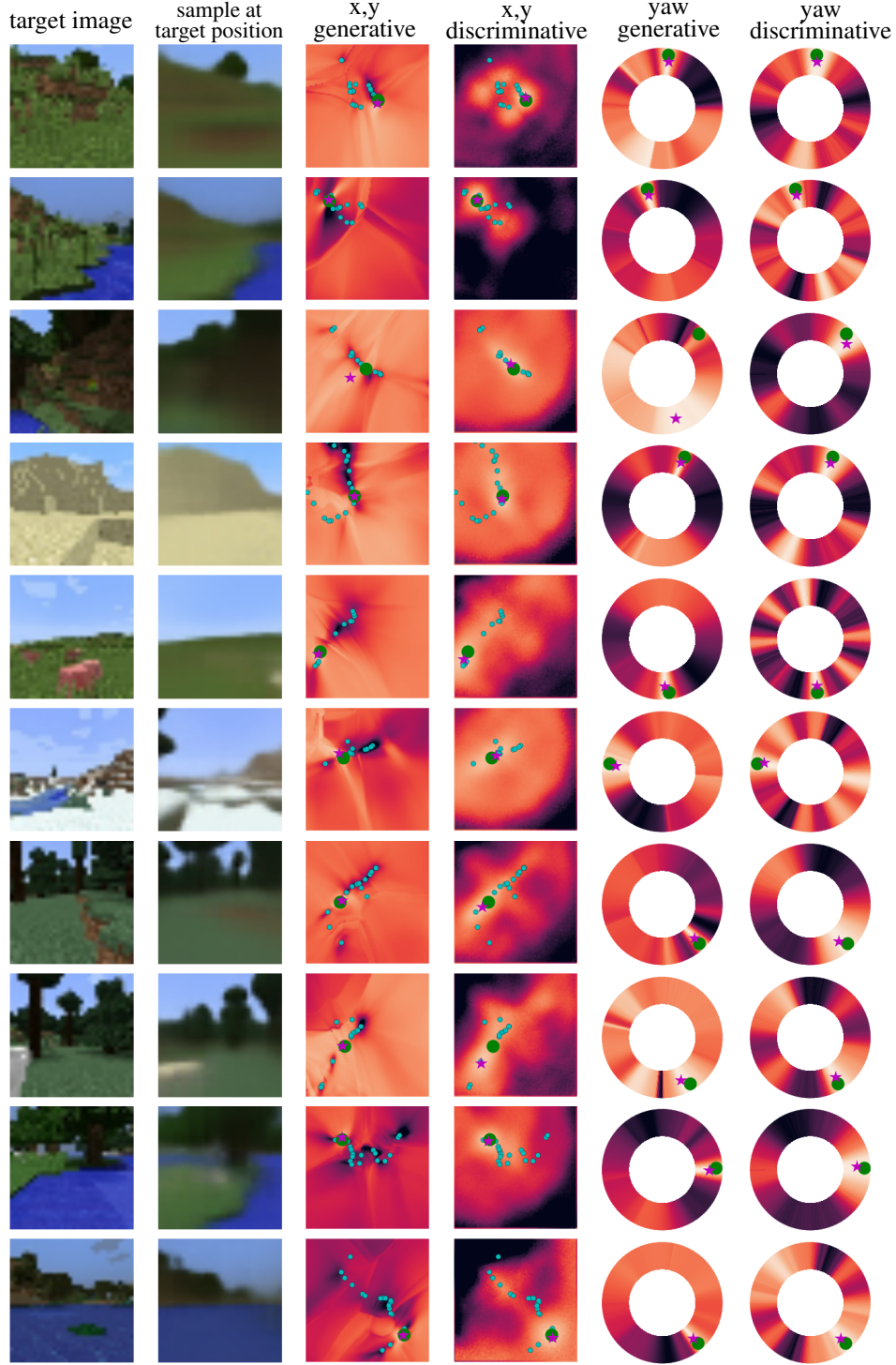
Figure 9: Localization with the discriminative and generative models - target image, a sample drawn using the ground-truth pose, and probability maps for the x,y position and yaw angles. The generative maps are computed by querying the model with all possible pose coordinates, and the discriminative maps are simply the model's output. The poses of the context images are shown in cyan, target image in green and maximum likelihood estimates in magenta. The generative maps are free from the prior in the training data (that target poses are close to the context poses) giving higher probability to unexplored areas and lower probabilities to areas with context images which are dissimilar to the target.

13