

Towards real-world navigation with deep differentiable planners

Shu Ishida João F. Henriques
 Visual Geometry Group
 University of Oxford
 {ishida, joao}@robots.ox.ac.uk

Abstract

We train embodied neural networks to plan and navigate unseen complex 3D environments, emphasising real-world deployment. Rather than requiring prior knowledge of the agent or environment, the planner learns to model the state transitions and rewards. To avoid the potentially hazardous trial-and-error of reinforcement learning, we focus on differentiable planners such as Value Iteration Networks (VIN), which are trained offline from safe expert demonstrations. Although they work well in small simulations, we address two major limitations that hinder their deployment. First, we observed that current differentiable planners struggle to plan long-term in environments with a high branching complexity. While they should ideally learn to assign low rewards to obstacles to avoid collisions, these penalties are not strong enough to guarantee collision-free operation. We thus impose a structural constraint on the value iteration, which explicitly learns to model impossible actions and noisy motion. Secondly, we extend the model to plan exploration with a limited perspective camera under translation and fine rotations, which is crucial for real robot deployment. Our proposals significantly improve semantic navigation and exploration on several 2D and 3D environments, succeeding in settings that are otherwise challenging for differentiable planners. As far as we know, we are the first to successfully apply them to the difficult Active Vision Dataset, consisting of real images captured from a robot.¹

1. Introduction

Continuous advances in robotics have enabled robots to be deployed to a wide range of scenarios, from manufacturing in factories and cleaning in households, to the emerging applications of autonomous vehicles and delivery drones [2]. Improving their autonomy is met with many challenges, due to the difficulty of planning from uncertain sensory data. In classical robotics, the study of planning has a long tradi-

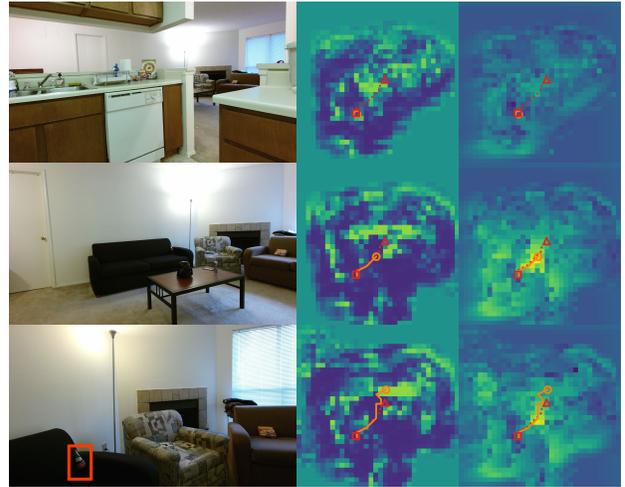


Figure 1. (1st column) Input images seen during a run of our method on AVD (Sec. 5.2.2). This embodied neural network has learned to efficiently explore and navigate unseen indoor environments, to seek objects of a given class (highlighted in the last image). (2nd-3rd columns) Predicted rewards and values (resp.), for each spatial location (higher for brighter values). The unknown optimal trajectory is dashed, while the robot’s trajectory is solid.

tion [46], using detailed knowledge of a robot’s configuration and sensors, with little emphasis on learning from data. An almost orthogonal approach is to use deep learning, an intensely data-driven, non-parametric approach [12]. Modern deep neural networks excel at pattern recognition [42], although they do not offer a direct path to planning applications. While one approach would be to parse a scene into pre-defined elements (*e.g.* object classes and their poses) to be passed to a more classical planner, an end-to-end approach where all modules are learnable has the chance to improve with data, and be adaptable to novel settings with no manual tuning. Because of the data-driven setup, a deep network has the potential to learn behaviour that leverages the biases of the environment, such as likely locations for certain types of rooms. Value Iteration Networks (VINs) [45] emerged as an elegant way to merge classical planning and data-driven deep networks, by defining a *differentiable* plan-

¹Code available: <https://github.com/shuishida/calvin>

ner. Being sub-differentiable, like all other elements of a deep network, allows the planner to include learnable elements, trained end-to-end from example data. For example, it can learn to identify and avoid obstacles, and to recognise and seek classes of target objects, without explicitly labelled examples. However, there are gaps between VIN’s idealised formulation and realistic robotics scenarios, which some works address [13, 23]. The CNN-based VIN [45] considers that the full environment is visible and expressible as a 2D grid. As such, it does not account for embodied (first-person) observations in 3D spaces, unexplored and partially-visible environments, or the mismatch between egocentric observations and idealised world-space discrete states.

In this paper we address these challenges, and close the gap between current differentiable planners and realistic robot navigation applications. Our contributions are:

1. A constrained transition model for value iteration, following a rigorous probabilistic formulation, which explicitly models illegal actions and task termination (Sec. 4.1). This is our main contribution.
2. A 3D state-space for embodied planning through the robot’s translation and rotation (Sec. 4.2.1). Planning through *fine-grained* rotations is often overlooked (Sec. 2), requiring better priors for transition modelling (Sec. 4.1.2).
3. A trajectory reweighting scheme that addresses the unbalanced nature of navigation training distributions (Sec. 4.1.5).
4. We demonstrate for the first time that differentiable planners can learn to navigate in both complex 3D mazes, and the challenging Active Vision Dataset [1], with images from a real robot platform (Sec. 5.2.2). Thus our method can be trained with limited data collected offline, as opposed to limitless data from a simulator as in prior work.

Sec. 2 discusses related work, while Sec. 3 gives a short introduction to differentiable planning. Sec. 4 presents our technical proposals, and Sec. 5 evaluates them.

2. Related work

Planning and reinforcement learning. Planning in fully-known deterministic state spaces was partially solved by graph-search algorithms [10, 16, 24–26, 43]. A Markov Decision Process (MDP) [3, 44] considers probabilistic transitions and rewards, enabling planning on stochastic models and noisy environments. With known MDPs, Value Iteration achieves the optimal solution [3, 44]. Reinforcement Learning (RL) focuses on unknown MDPs [31, 31, 38, 47, 48]. Model-free RL systems are reactive; a policy (*i.e.* a generic plan) is built over a long period of trial and error, and will be specific to the training environment (there is no online planning). This makes them less ideal for robotics, where risky failures must be avoided, and new plans (policies) must be created on-the-fly in new environments. Model-based RL

methods attempt to learn an MDP and often use it for planning online [9, 14, 15, 21, 22]. The environment is assumed to provide a reward signal at every step in both cases.

Deep networks for navigation. Several works have made advances into training deep networks for navigation. The Neural Map [36] is an A2C [31] agent (thus reactive) that reads and writes to a differentiable memory (*i.e.* a map). Similarly, Mirowski *et al.* [30] train reactive policies which are specific to an environment. The Value Prediction Network [35] learns an MDP and its state-space from data, and plans with a single roll-out over a short horizon. Hausknecht *et al.* [18] add recurrence to deep Q-learning to address partially-observable environments, but considering only single-frame occlusions. Several works treat planning as a non-differentiable module, and focus on training neural networks for other aspects of the navigation system. Savinov *et al.* [39] do this by composing siamese networks and value estimators trained on proxy tasks, and use an initial map built from footage of a walk through the environment. Active Neural SLAM [6] trains a localisation and mapping component (outputting free space and obstacles), a policy network to select a target for the non-differentiable planner, and another to perform low-level control. Subsequent work [5] complemented this map with semantic segmentation. This map is different from ours, which consists of embeddings learned by a differentiable planner, and so are not constrained to correspond to annotated labels. The later two works navigate successfully in large simulations of 3D-scanned environments, and were transferred to real robots.

Imitation learning. To avoid learning by trial-and-error, Imitation Learning instead uses expert trajectories. Inverse Reinforcement Learning (IRL) [20, 33, 50, 51] achieves this by learning a reward function that best explains the expert trajectories. However, this generally leads to a difficult nested optimisation and is an ill-posed problem, since many reward functions can explain the same trajectory.

Differentiable planning. Tamar *et al.* [45] introduced the Value Iteration Network (VIN) for Imitation Learning, which generates a plan online (with value iteration) and back-propagates errors through the plan to train estimators of the rewards (Sec. 3.1). VINs were also applied to localisation from partial observations by Karkus *et al.* [23], but assuming a full map of the environment. Lee *et al.* [27] replaced the maximum over actions in the VI formula with a LSTM. The evaluation included an extension to rotation, but assuming a fully-known four directional view for every gridded state, which is often not available in practice. A subtle difference is that they handle rotation by applying a linear policy layer to the hidden channels of a 2D state-space grid, which is less interpretable than our 3D translation-rotation grid.

Most previous work on deploying VIN to robots [34, 40] assume an occupancy map and goal map rather than learning the map embedding and goal themselves from data. To

the best of our knowledge, Gupta *et al.*'s Cognitive Mapper and Planner (CMP) [13] is the only work which evaluated a differentiable planner on real-world data using map embeddings learned end-to-end. CMP uses a hierarchical VIN [45] to plan on larger environments, and updates an egocentric map. It handles rotation as a warping operation external to the VIN (i.e. it plans in a 2D translation state-space, not 3D translation-rotation space). Warping an egocentric map per update achieves scalability at the cost of progressive blurring of the embeddings. CMP was evaluated for rotations of 90° and deterministic motion, rather than noisy translation and rotation. Our 3D translation-rotation space and learned transition model allows smooth trajectories with finer rotations. CMP also requires gathering new trajectories online during training (with DAGger [37]), while ours is trained solely with a fixed set of training trajectories, foregoing the need for a simulator and the associated domain gap.

3. Background

A Markov Decision Process (MDP) [3, 44] formalises sequential decision-making. It consists of states $s \in \mathcal{S}$ (e.g. locations), actions available at each state $a \in \mathcal{A}(s)$, a reward function $R(s, a, s')$ to be maximised (e.g. reaching the target), and a transition probability $P(s'|s, a)$ (the probability of the next state s' given the current state s and action a). The objective of an agent is to learn a policy $\pi(a|s)$, specifying the probability of choosing action a for any state s , chosen to maximise the expected return G_t at every time step t . A return is a sum of discounted rewards, $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$, where a discount factor $\gamma \in (0, 1)$ prevents divergences of the infinite sum. The value function $V_{\pi}(s) = \mathbb{E}_{\pi}[G_t|s_t = s]$ evaluates future returns from a state, while the action-value function $Q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t|s_t = s, a_t = a]$ considers both a state and the action taken. An optimal policy π_* should then maximise the expected return for all states, i.e. $\forall s \in \mathcal{S}, v_*(s) = \max_{\pi} v_{\pi}(s)$. Value Iteration (VI) [3, 44] is an algorithm to obtain an optimal policy, by alternating a refinement of both value (V) and action-value (Q) function estimates in each iteration k . When s and a are discrete, $Q^{(k)}$ and $V^{(k)}$ can be implemented as simple tables (tensors). In particular, we will consider as states the cells of a 2D grid, corresponding to discretised locations in an environment, i.e. $s = (i, j) \in \mathcal{S} = \{1, \dots, N\}^2$. Furthermore, transitions are local (only to coordinates offset by $\delta \in K = \{-1, 0, 1\}^2$):

$$Q_{a,s}^{(k)} = \sum_{\delta \in K} P_{a,\delta,s} \left(R_{a,\delta,s} + \gamma V_{s+\delta}^{(k-1)} \right)$$

$$\forall a \in \mathcal{A}(s), \quad \text{with} \quad V_s^{(k)} = \max_{a \in \mathcal{A}(s)} Q_{a,s}^{(k)}, \quad (1)$$

Note that to avoid repetitive notation s and δ are 2D indices, so V is a 2D matrix and both R and P are 5D tensors. The policy $\pi_s^{(k)} = \operatorname{argmax}_{a \in \mathcal{A}} Q_{a,s}^{(k)}$ simply chooses the action

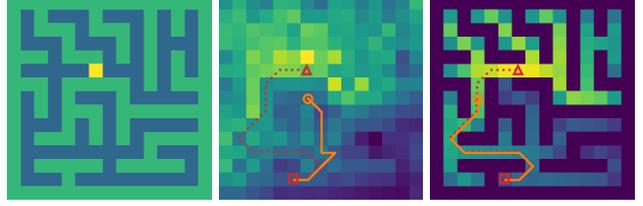


Figure 2. (left) A 2D maze, with the target in yellow. (middle) Values produced by the VIN for each 2D state (actions are taken towards the highest value). Higher values are brighter. The correct trajectory is dashed, the current one is solid. The agent (orange circle) is stuck due to the local maximum below it. (right) Same values for our method, CALVIN. There are no spurious maxima, and the values of walls are correctly considered low (dark).

with the highest action-state value. While the sum in Eq. 1 resembles a convolution, the filters (P) are space-varying (depend on $s = (i, j)$), so it is not directly expressible as such. Eq. 1 represents the “ideal” VI for local motion on a 2D grid, without further assumptions.

3.1. Value Iteration Network

While VI guarantees the optimal policy, it requires that the functions for transition probability P and reward R are known (e.g. defined by hand). Tamar *et al.* [45] pointed out that all VI operations are (sub-)differentiable, and as such a model of P and R can be trained from data by back-propagation. For the case of planning on a 2D grid (navigation), they related Eq. 1 to a CNN, as:

$$Q_{a,s}^{(k)} = \sum_{\delta \in K} \left(P_{a,\delta}^R \hat{R}_{s+\delta} + P_{a,\delta}^V V_{s+\delta}^{(k-1)} \right),$$

$$\forall a \in \mathcal{A}, \quad \text{with} \quad V_s^{(k)} = \max_{a \in \mathcal{A}} Q_{a,s}^{(k)}, \quad (2)$$

where $P^R, P^V \in \mathbb{R}^{A \times |K|}$ are two learned convolutional filters that represent the transitions (P in Eq. 1), and \hat{R} is a predicted 2D reward map. Note that \mathcal{A} is independent of s , i.e. all actions are allowed in all states. This turns out to be detrimental (see next paragraph). The reward map \hat{R} is predicted by a CNN, from an input of the same size that represents the available observations. In Tamar *et al.*'s experiments [45], the observations were a fully-visible overhead image of the environment, from which negative rewards such as obstacles and positive rewards such as navigation targets can be located. Each action channel in \mathcal{A} corresponds to a move in the 2D grid, typically 8-directional or 4-directional. Eq. 2 is attractive, because it can be implemented as a CNN consisting of alternating convolutional layers (Q) and max-pooling along the actions (channels) dimension (V).

Motivating experiment. Since the VIN allows all actions at all states (\mathcal{A} does not depend on s), collisions must be modelled as states with low rewards. In practice, the VIN does not learn to forbid such actions completely, resulting in

propagation of values from states which cannot be reached directly due to collision along the way. We verified this experimentally, by training a VIN according to Tamar *et al.* [45] on 4K mazes (see Sec. 5 for details). We then measured, for each state, whether the predicted scores for all valid actions are larger than those for all invalid actions (*i.e.* collisions). Intuitively, this means the network always prefers free space over collisions. Surprisingly, we found that this was not true for 24.6% of the states. For a real-world robot to work reliably, this is an unacceptably high chance of collisions. As a comparison, for our method (Sec. 4) this rate is only 1.6%. In the same experiment, the VIN often gets trapped in local minima of the value function and does not move (Fig. 2), which is another failure mode. We aim to fix these issues, and push VINs towards realistic scenarios. An alternative would be to employ online retraining (DAGger) [37], which cannot use solely a fixed set of offline trajectories.

4. Proposed method

We propose a transition model that accounts for illegal actions and termination. We then extend it to embodied planning (rigid 3D motion and partially-observed environments).

4.1. Augmented navigation state-action space

In this section we will derive a probabilistic transition model from first principles, with only two assumptions and no extra hyper-parameters. The first assumption is locality and translation invariance of the agent motion, which was introduced in the VIN to allow efficient learning with shared parameters. Unlike the VIN, we will decompose the transition model $P(s'|s, a)$ into two components: the agent motion model $\hat{P}(s' - s|a)$, which is translation invariant and shared across states (depending only on the spatial difference between states $s' - s$); and an observation-dependent predictor $\hat{A}(s, a) \in [0, 1]$ which evaluates whether action a is available from state s , to disqualify illegal actions.

In robotics, it is essential that the agent understands that the current task has been completed to move on to the next one. Since in small environments there is a high chance that a randomly-acting agent will stumble upon the target, Anderson *et al.* [2] suggested that an explicit termination action must be taken at the target to finish successfully. Therefore, in addition to positional states, we assume a success state W (“win”) that is reached only by triggering a termination action D (“done”), and a failure state F (“fail”) that is reached upon triggering an incorrect action. We denote the reward for reaching F as R_F , and the translation-invariant rewards as $\hat{R}(a, s' - s)$. For simplicity, we consider the reward for success a special case of $\hat{R}(a, s' - s)$ where $a = D$ and $s' = W$. With these assumptions, the reward function $R(s, a, s')$ is:

$$R(s, a, s') = \begin{cases} \hat{R}_F, & s' = F \\ \hat{R}(a, s' - s), & s' \neq F. \end{cases} \quad (3)$$

Together with the agent motion model $\hat{P}(s' - s|a)$, action validity predictor $\hat{A}(s, a)$ and the definition of a failure state F , the transition model $P(s'|s, a)$ can be derived as Eq. 4:

$$P(s'|s, a) = \begin{cases} 1 - \hat{A}(s, a), & s' = F \\ \hat{A}(s, a)\hat{P}(s' - s|a), & s' \neq F. \end{cases} \quad (4)$$

From the above, we calculate the reward $R(s, a)$ by marginalising over the neighbour states s' :

$$\begin{aligned} R(s, a) &= \sum_{s'} P(s'|s, a)R(s, a, s') \\ &= \hat{R}_F(1 - \hat{A}(s, a)) + \hat{A}(s, a) \sum_{s'} \hat{P}(s' - s|a)\hat{R}(a, s' - s) \end{aligned} \quad (5)$$

Finally, Eq. 5 can be plugged into Eq. 1 to obtain our proposed value iteration’s $Q(s, a)$:

$$\begin{aligned} Q(s, a) &= R(s, a) + \gamma \mathbb{I}_{a \in \mathcal{A}} \sum_{s'} P(s'|s, a)V(s') \\ &= R(s, a) + \gamma \hat{A}(s, a) \mathbb{I}_{a \neq D} \sum_{s'} \hat{P}(s' - s|a)V(s') \end{aligned} \quad (6)$$

where \mathbb{I} is an indicator function. Eq. 5 and 6 essentially express a constrained VI, which models the case of an MDP on a grid with unknown illegal states and termination at a goal state. The inputs to this model are three learnable functions — the motion model $\hat{P}(s' - s|a)$, the action validity $\hat{A}(s, a)$ (*i.e.* obstacle predictions), and the rewards $\hat{R}(a, s' - s)$ and \hat{R}_F . These are implemented as CNNs with the observations as inputs (\hat{R}_F is a single learned scalar). All the constraints follow from a well-defined world-model, with very interpretable predicted quantities, unlike previous proposals [27, 45]. We named this method Collision Avoidance Long-term Value Iteration Network (CALVIN).

4.1.1 Training

Similarly to Tamar *et al.* [45], we train our method with a softmax cross-entropy loss L , comparing an example trajectory $\{(s_t, a_t^*) : t \in T\}$ with predicted action scores $Q(s, a)$:

$$\min_{\hat{P}, \hat{A}, \hat{R}} \frac{1}{|T|} \sum_{t \in T} w_t L(Q(s_t), a_t^*), \quad (7)$$

where $Q(s)$ is a vector with one element per action ($Q(s, a)$), and w_t is an optional weight that can be used to bias the loss if $w_t \neq 1$ (Sec. 4.1.5). Example trajectories are shortest paths computed from random starting points to the target (also chosen randomly). Note that the learned functions can be conditioned on input observations. These are 2D grids of features, with the same size as the considered state-space (*i.e.* a map of observations), which is convenient since it enables \hat{P} , \hat{A} and \hat{R} to be implemented as CNNs.

4.1.2 Transition modelling

Similarly to the VIN (Eq. 2), the motion model $\widehat{P}(s' - s|a)$ is implemented as a convolutional filter $\widehat{P} \in \mathbb{R}^{|\mathcal{A}| \times |\mathcal{K}|}$, so it only depends on the relative spatial displacement $s' - s$ between the two states s and s' . We can use the transitions already observed in the example trajectory to constrain the model, by adding a cross-entropy loss term $L(\widehat{P}(a_t^*), s_{t+1} - s_t)$ for each step in the example trajectory. After training, the filter $\widehat{P}(s' - s|a)$ will consist of a distribution over possible state transitions for each action (visualised in Fig. 5).

4.1.3 Action availability

Although the available actions \widehat{A} could, in theory, be learned completely end-to-end, in practice we found that additional regularisation is necessary. If we had a reliable log-probability of each action being taken, $\widehat{A}_{\text{logit}}(s, a)$, then by thresholding it at some point $\widehat{A}_{\text{thresh}}(s)$ we would distinguish between available and unavailable actions. Using a sigmoid function σ as a soft threshold, we can write this as:

$$\widehat{A}(s, a) = \sigma(\widehat{A}_{\text{logit}}(s, a) - \widehat{A}_{\text{thresh}}(s)). \quad (8)$$

Both $\widehat{A}_{\text{logit}}(s, a)$ and $\widehat{A}_{\text{thresh}}(s)$ are predicted by the network, given the observations at each time step. In order to ground the probabilities of each action being taken, we encourage the actions $\widehat{A}_{\text{logit}}(s_t)$ to match the example trajectory action a_t^* for all steps t , with an additional cross-entropy loss $L(\widehat{A}_{\text{logit}}(s_t), a_t^*)$. Note that *there is no additional ground truth supervision* – we use strictly the same data as a VIN.

4.1.4 Fully vs. partially observable environments

Some previous works [27, 45] assume that the entire environment is static and fully observable, which is often unrealistic. Partially observable environments, involving unknown scenes, require exploring to gather information and are thus more challenging. We account for this with a simple but significant modification. Note that $Q(s_t)$ in Eq. 7 depends on the learned functions (CNNs) \widehat{P} , \widehat{A} and \widehat{R} through Eq. 6, and these in turn are computed from the observations. We extend the VIN framework to the case of partial observations by ensuring that $Q(s_t|O_{\leq t})$ depends *only* on the observations up to time t , $O_{\leq t}$. This means that the VIN recomputes a plan at each step t (since the observations are different), instead of once for a whole trajectory. Unobserved locations have their features set to zero, so in practice knowledge of the environment is slowly built up during an expert demonstration, which enables exploration behaviour to be learned.

4.1.5 Trajectory reweighting

Exploration provides observations $O_{\leq t}$ of the same locations at different times (Sec. 4.1.4). We can thus augment Eq. 7

with these partial observations as extra samples. The sum in Eq. 7 becomes $\sum_{t' \in T} \sum_{t \in T_{1:t'}} w_t L(Q(s_t|O_{\leq t'}), a_t^*)$. For long trajectories, the training data then becomes severely unbalanced between exploration (target not visible) and exploitation (target visible), with a large proportion of the former (90% of the data in Sec. 5.1.2). We addressed this by reweighting the samples. We set $w_t = \beta^{d_t} / \max_{j \in T} \beta^{d_j}$ in Eq. 7, *i.e.* a geometric decay scaling with the topological distance to the target d_t . Since expert trajectories are shortest paths, this simplifies to $w_t = \beta^{|T|-t}$.

4.2. Embodied navigation in 3D environments

Embodied agents such as robots have a pose in 3D space, which for non-holonomic robots limits their available actions (*e.g.* moving forward or rotating), and their observations.

4.2.1 Embodied pose states (position and orientation)

To address the first limitation, we augment the 2D state-space with an extra dimension, which corresponds to Θ discretised orientations: $\mathcal{S} = \{1, \dots, N\}^2 \times \{1, \dots, \Theta\}$. This can be achieved by directly adding one extra dimension to each spatial tensor in Eq. 5 and 6. A table of tensor sizes is in Appendix A. Note that the much larger state-space makes long-term planning more difficult to learn. We observed that when naively augmenting the state-space in this way, the models fail to learn correct motion kernels $\widehat{P}(s' - s|a)$. This further reinforces the need for an auxiliary motion loss (Sec. 4.1.2), which overcomes this obstacle, and may explain why prior works did not plan through fine-grained rotations.

4.2.2 3D embeddings for geometric reasoning

To aggregate image information (CNN embeddings) on to the 2D grid (map) where VIN performs planning, we follow the same strategy as MapNet [19]. Each embedding is associated with a 3D point in world-space via projective geometry, assuming that the camera poses and depths are known (as assumed in prior work [4, 13, 27, 28, 45], which can be estimated from monocular vision [32]). The embeddings of 3D points *from all past frames* that fall into each cell of the world-space 2D grid are aggregated with mean-pooling. Due to the use of PointNet [7] aggregation on cells of a lattice, we named this Lattice PointNet (LPN). A self-contained description is in Appendix A. The LPN has some appealing properties in our context of navigation: 1) it allows reasoning about far away, observed but yet unvisited locations; 2) it fuses multiple observations of the same location, whether from different points-of-view or different times.

Memory-efficient mapping. Temporal aggregation during rollout can be computed recursively as $e_{t,i,j,k}/n_{t,i,j,k}$ for $e_{t,i,j,k} = e_{t-1,i,j,k} + e'_{t,i,j,k}$ and $n_{t,i,j,k} = n_{t-1,i,j,k} + n'_{t,i,j,k}$, where $e'_{t,i,j,k}$ is the summed embedding for the

Table 1. Navigation success rate (fraction of trajectories that reach the target) on unseen 2D mazes. Partial observations (exploring an environment gradually) and embodied navigation (translation-rotation state-space) are important yet challenging steps towards full 3D environments.

Env.	Standard loss			Reweighted loss (ours)		
	VIN	GPPN	CALVIN	VIN	GPPN	CALVIN
Full obs.	75.6 \pm 20.6	91.3 \pm 8.1	99.0 \pm 1.0	77.5 \pm 26.6	96.6 \pm 4.0	99.7 \pm 0.5
Partial	3.6 \pm 0.6	8.5 \pm 3.5	48.0 \pm 5.2	1.7 \pm 1.7	11.25 \pm 3.7	92.2 \pm 1.3
Embod.	11.0 \pm 1.0	14.5 \pm 2.1	90.0 \pm 7.9	15.2 \pm 3.6	28.5 \pm 3.5	93.7 \pm 6.2

points in cell (i, j, k) at time t , and $n'_{t,i,j,k}$ is the number of points per cell. Only the previous map $e_{t-1,i,j,k}$ and previous counts $n_{t-1,i,j,k}$ must be kept, not all past observations. Thus at run-time the memory cost is *constant* over time, allowing unbounded operation (unlike methods that do not have an explicit map [39]).

4.3. Limitations

Our main contribution is to improve the VIN algorithm itself (sec. 3.1) with correct termination, transition and availability probabilistic models, which is orthogonal to works which build on top of VIN. We assume the agent’s pose, depth image and the camera parameters to be known. Other dynamic objects in the environment are not modelled.

5. Experiments

In this section we will gradually build up the capabilities of our method with the proposals from Sec. 4, comparing it to various baselines on increasingly challenging environments, and leading up to unseen real-world 3D environments.

5.1. 2D grid environments

We start with 2D environments, where the observations are top-down views of the whole scene, and which thus do not require dealing with perspective images (Sec. 4.2.2). Since Tamar *et al.* [45] obtain near-perfect performance in their 2D environments, after reproducing their results we focused on 2D mazes, which are much more challenging since they require frequent backtracking to navigate if exploration is required. The 15×15 mazes are generated using Wilson’s algorithm [49], and an example can be visualised in Fig. 3. The allowed moves \mathcal{A} are to any of the 8 neighbours of a cell. As discussed in Sec. 4, a termination action D must be triggered at the target to successfully complete the task. The target is placed in a free cell chosen uniformly at random, with a minimum topological distance from the (random) start location equal to the environment size to avoid trivial tasks.

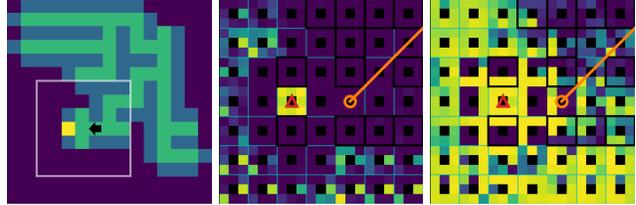


Figure 3. Our method on 2D mazes (Sec. 5.1.3). (left) Input visualisation: unexplored cells are dark, the target is yellow (just found by the agent), and a black arrow shows the agent’s position and orientation. (middle) Close-up of predicted rewards (higher values are brighter) inside the white rectangle of the left panel. The 3D state-space (position/orientation) is shown, with rewards for the 8 orientations in a radial pattern within each cell (position). Explored cells have low rewards, with the highest reward at the target. (right) Close-up of predicted values. They are higher facing the direction of the target. Obstacles (black border) have low values.

5.1.1 Fully-known environment with positional states

Baselines and training. For the first experiment, we compare our method (CALVIN) with other differentiable planners: the VIN [45] and the more recent GPPN [27], on fully-observed environments. Other than using mazes instead of convex obstacles, this setting is close to Tamar *et al.*’s [45]. The VIN, GPPN and CALVIN all use 2-layer CNNs to predict their inputs (details in Appendix A). All networks are trained with $4K$ example trajectories in an equal number of different mazes, using the Adam optimiser with the optimal learning rate chosen from $\{0.01, 0.005, 0.001\}$, until convergence (up to 30 epochs). Navigation success rates (fraction of trajectories that reach the target) for epochs with minimum validation loss are reported. Our reweighted loss (Sec. 4.1.5) is equally applicable to all differentiable planners, so we report results both with and without it.

Results. Table 1 (first row) shows the navigation success rate, averaged over 3 random seeds (and the standard deviation). The VIN has a low success rate, showing that it does not scale to large mazes. GPPN achieves a high success rate, and CALVIN performs near-perfectly. This may be explained by the GPPN’s higher capacity, as it contains a LSTM with more parameters. Nevertheless, CALVIN has a more constrained architecture, so its higher performance hints at a better inductive bias for navigation. It is interesting to note that the proposed reweighted loss has a beneficial effect on all methods, not just CALVIN. With the correct data distribution, any method with sufficient capacity can fit the objective. This shows that addressing the unbalanced nature of the data is an important, complementary factor.

5.1.2 Partially observable environment

Setting. Next, we compare the same methods in unknown environments, where the observation maps only contain ob-

Table 2. Navigation success rate on unseen 3D mazes (MiniWorld). Note that the baselines do not generalise to larger mazes.

Size	CNN backbone		LPN backbone (ours)		
	A2C	PPO	VIN	GPPN	CALVIN (ours)
3 × 3	98.7 \pm 1.9	81.0 \pm 26.9	90.3 \pm 3.1	91.3 \pm 4.7	97.7 \pm 1.7
8 × 8	23.6 \pm 4.9	14.7 \pm 6.2	41.2 \pm 9.5	33.3 \pm 8.6	69.2 \pm 5.3

served features up to the current time step (sec. 4.1.4). To simulate local observations, we perform ray-casting to identify cells that are visible from the current position, up to 2 cells away. Example observations are in Appendix A.

Results. In this case, the agent has to take significantly more steps to explore compared to a direct route to the target. From Tab. 1 (2nd row) we can see that partial observability causes most methods to fail catastrophically. The sole exception is CALVIN with the reweighted loss (proposed), which performs well. Note that to succeed, an agent must acquire several complex behaviours: directing exploration to large unseen areas; backtracking from dead ends; and seeking the target when seen. Our method displays all of these behaviours (see Appendix A for visualisations). While the model initially assigns high values to all unexplored states, as soon as the target is in view, the model assigns a high value to the target state and its neighbours. Since only a combination of CALVIN and a reweighted loss works at all, we infer that a correct inductive bias and a balanced data distribution are both necessary for success.

5.1.3 Embodied navigation with orientation

Setting. Now we consider embodied navigation, where transitions depend on the agent’s orientation (sec. 4.2.1). We augment the state-space of all methods (Sec. 4.2.1) with 8 orientations at 45° intervals, and allow 4 move actions A : forward, backward, and rotating in either direction.

Results. In Tab. 1 (3rd row) we observe that VIN and GPPN perform slightly better, but still have a low chance of success. CALVIN outperforms them by a large margin. We also visualise a typical run in Fig. 3 (refer to the caption for more detailed analysis). One advantage of CALVIN displayed in Figs. 2 and 3 is that values and rewards are fully interpretable and play the expected roles in value iteration (Eq. 1). Less constrained architectures [27, 45] insert operators that deviate from the value iteration formulation, and thus lose their interpretability as rewards and values (*cf.* Appendix A).

5.2. 3D environments

Having validated embodied navigation and exploration, we now integrate the Lattice PointNet (LPN, Sec. 4.2.2) to handle first-person views of 3D environments.

5.2.1 Synthetically-rendered environments

Dataset. We used the MiniWorld simulator [8], which allowed us to easily generate 3D maze environments with arbitrary layouts. Only a monocular camera is considered (not 360° views [27]). The training trajectories now consist of first-person videos of the shortest path to the target, visualised in Fig. 4. We randomly generate 1K trajectories in mazes on either small (3 × 3) or large (8 × 8) grids by adding or removing walls at the boundaries of this grid’s cells. Note that the maze’s layout and the agent’s location do not necessarily align with the 2D grids used by the planners (as in [27]). Thus, planning happens on a fine discretisation of the state-space (30 × 30 for small mazes and 80 × 80 for large ones, with 8 orientations). This allows smooth motions and no privileged information about the environment. Both translation and rotation are perturbed by Gaussian noise, forcing all agents to model uncertain dynamics.

Baselines and training. We compare several baselines: two popular RL methods, A2C [31] and PPO [41], as well as the VIN, GPPN and the proposed CALVIN. Since A2C and PPO are difficult to train if triggering the “done” action is strictly required, we relaxed the assumption, allowing the agent to terminate once it is in close proximity to the target. All methods use as a first stage a simple 2-layer CNN (details in Appendix A). Since this CNN was not enough to get the VIN, GPPN and CALVIN to work well (see Appendix A), they all use our proposed LPN backbone. We could not include GPPN’s strategy of taking as input views at all possible states [27], due to the high memory requirements and the environment not being fully visible (only local views are available). Other training details are identical to Sec. 5.1.

Results. In Tab. 2 we observe that, although the RL methods are successful in small mazes, their reactive policies cannot scale to large mazes. CALVIN succeeds reliably even for longer trajectories, outperforming the others. It is interesting to note that the proposed LPN backbone is important for all differentiable planners, and it allows them to achieve very high success rates for small environments (though our CALVIN method performs slightly better). We show an example run of our method in Fig. 4, where it found the target after an efficient exploration period, despite not knowing its location and never encountering this maze before.

Figure 5 visualises the learnt state transitions $\hat{P}(s' - s|a)$ for the *move forward* action in CALVIN. We observe that the learning mechanism outlined in Sec. 4.1.2 works even for transitions with added noise, which is the case for MiniWorld experiments. The network learns to propagate values probabilistically from the possible next states.

5.2.2 Indoor images from a real-world robot

Finally, we tested our method on real images obtained with a robotic platform.

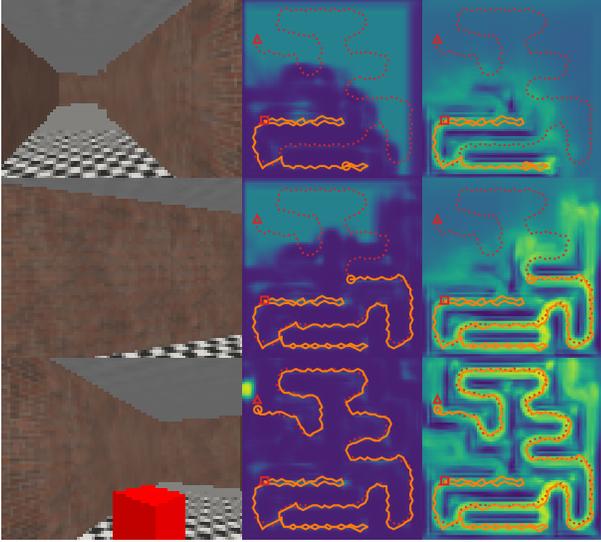


Figure 4. Example results on MiniWorld (Sec. 5.2.1). Left to right: input images, predicted rewards and values. The format is as in Fig. 1. Notice the high reward on unexplored regions, replaced with a single peak around the target when it is seen (last row).

Dataset. We used the Active Vision Dataset [1] (AVD), which allows interactive navigation with real image streams, without synthetic rendering (as opposed to [29]). This is achieved using monocular RGBD images from 19 indoor environments, densely collected by a robot on a 30cm grid and at 30° rotations. The set of over 30K images can be composed to simulate any trajectory, up to the some spatial granularity. There are also bounding box annotations of object instances, which we use to evaluate semantic navigation. We use the last four environments as a validation set, and the rest for training (by sampling 1K shortest paths to the target). A visualisation is shown in Fig. 1.

Tasks and training. We considered a semantic navigation task of seeking an object of a learnt class. We chose the most common class (“soda bottle”) as a target object. Training follows Sec. 5.2.1.

Results. We report the performances of VIN, GPPN and CALVIN after 8 epochs of training in Tab. 3. As in Sec. 5.2.1, they also fail without the LPN, so all results are with the LPN backbone. A similar conclusion to that for synthetic environments can be drawn: proper spatio-temporal aggregation of local observations is essential for differentiable planners to scale realistically. CALVIN achieves a signifi-



Figure 5. Transition model learnt from MiniWorld trajectories for the *move forward* action at each discretised orientation, at 45° intervals. Higher values are brighter, and lower values are darker (purple for a probability of 0).

Table 3. Navigation success rate on the Active Vision Dataset, with real robot images taken in indoor spaces. The task is to navigate to an object of a learned class. All methods use the proposed LPN backbone, as they fail without it.

Subset	VIN	GPPN	CALVIN (ours)
Training	61.6±4.5	50.6±9.2	70.3±4.9
Validation	45.0±1.0	44.0±3.5	47.6±6.0

cantly higher success rate on the training set than the other methods. On the other hand, while it has higher mean validation success, the high variance of this estimate does not allow the result to be as conclusive as for training. We can attribute this to the small size of AVD in general, and of the validation set in particular, which contains only 3 different indoor scenes. Nevertheless, this shows that CALVIN learns effective generic strategies to seek a specific object, and that VIN and GPPN equipped with a similar backbone can achieve partial success in several training environments.

6. Discussion

We analysed several shortcomings of current differentiable planners, with the goal of deploying to real robot platforms. We found that they can be addressed by highly complementary solutions: a constrained transition model (CALVIN) correctly incorporates illegal actions (as opposed to simply discouraging them via rewards) and task termination; a 3D state-space accounts for robot orientation; a LPN backbone efficiently fuses spatio-temporal information; and trajectory reweighting addresses unbalanced training data. We provide empirical evidence in several settings, including on images of unseen environments from a real robot.

Although deep networks hold promise for making navigation more robust under uncertainty, their uninterpretable failure modes mean that they are not yet mature enough for safety-critical applications, and more research to close this gap is still needed. Vehicle applications are especially important, due to the high masses and velocities involved. We propose a more interpretable MDP structure for VINs, and train our model solely with safe offline demonstrations. However, their reliability is far from guaranteed, and complementary safety systems in hardware must be considered in any deployment. In addition to further improving robustness to failures, future work may investigate more complex tasks, such as natural language commands, and study the effect of sensor drift on navigation performance.

Acknowledgements. The authors would like to acknowledge the generous support of the Royal Academy of Engineering (RF\201819\18\163) and the Ezoë Memorial Recruit Foundation. They would also like to thank the authors of MiniWorld [8] and the Active Vision Dataset [1] for making their source code and data available.

References

- [1] Phil Ammirato, Patrick Poirson, Eunbyung Park, Jana Košecká, and Alexander C. Berg. A dataset for developing and benchmarking active vision. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1378–1385, 2017. [2](#), [8](#)
- [2] Peter Anderson, Angel Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, et al. On evaluation of embodied navigation agents. *arXiv preprint arXiv:1807.06757*, 2018. [1](#), [4](#)
- [3] Richard Bellman. A Markovian Decision Process. *Indiana University Mathematics Journal*, 6(4):679–684, 1957. [2](#), [3](#)
- [4] Vincent Cartillier, Zhile Ren, Neha Jain, Stefan Lee, Irfan Essa, and Dhruv Batra. Semantic mapnet: Building allocentric semantic maps and representations from egocentric views, 2021. [5](#), [12](#)
- [5] Devendra Singh Chaplot, Dhiraj Gandhi, Abhinav Gupta, and Ruslan Salakhutdinov. Object goal navigation using goal-oriented semantic exploration. In *In Neural Information Processing Systems*, 2020. [2](#)
- [6] Devendra Singh Chaplot, Dhiraj Gandhi, Saurabh Gupta, Abhinav Gupta, and Ruslan Salakhutdinov. Learning to Explore using Active Neural SLAM. In *International Conference on Learning Representations*, page 18, 2020. [2](#)
- [7] R. Qi Charles, Hao Su, Mo Kaichun, and Leonidas J. Guibas. PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 77–85, Honolulu, HI, July 2017. IEEE. [5](#), [11](#)
- [8] Maxime Chevalier-Boisvert. Gym-MiniWorld environment for OpenAI Gym. <https://github.com/maximecb/gym-miniworld>, 2018. [7](#), [8](#)
- [9] Ignasi Clavera, Jonas Rothfuss, John Schulman, Yasuhiro Fujita, Tamim Asfour, and Pieter Abbeel. Model-based reinforcement learning via meta-policy optimization. In *Conference on Robot Learning*, pages 617–629. PMLR, 2018. [2](#)
- [10] E. W. Dijkstra. A note on two problems in connexion with graphs. *Numerische Mathematik*, 1959. [2](#)
- [11] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2002–2011, 2018. [11](#)
- [12] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT press Cambridge, 2016. [1](#)
- [13] Saurabh Gupta, Varun Tolani, James Davidson, Sergey Levine, Rahul Sukthankar, and Jitendra Malik. Cognitive Mapping and Planning for Visual Navigation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. arXiv: 1702.03920. [2](#), [3](#), [5](#), [11](#), [12](#)
- [14] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019. [2](#)
- [15] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In *International Conference on Machine Learning*, pages 2555–2565. PMLR, 2019. [2](#)
- [16] P. E. Hart, N. J. Nilsson, and B. Raphael. A formal basis for the heuristic determination of minimum cost paths. *IEEE Transactions on Systems Science and Cybernetics*, 1968. [2](#), [12](#)
- [17] R Hartley and A Zisserman. Multiple view geometry in computer. *Vision, 2nd ed., New York: Cambridge*, 2003. [11](#)
- [18] Matthew Hausknecht and Peter Stone. Deep recurrent q-learning for partially observable mdps. *arXiv preprint arXiv:1507.06527*, 2015. [2](#)
- [19] Joao F. Henriques and Andrea Vedaldi. MapNet: An Allocentric Spatial Memory for Mapping Environments. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8476–8484, June 2018. ISSN: 1063-6919. [5](#), [12](#)
- [20] Jonathan Ho and Stefano Ermon. Generative Adversarial Imitation Learning. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 4565–4573, 2016. [2](#)
- [21] Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: Model-based policy optimization. *arXiv preprint arXiv:1906.08253*, 2019. [2](#)
- [22] Lukasz Kaiser, Mohammad Babaeizadeh, Piotr Milos, Blazej Osinski, Roy H Campbell, Konrad Czechowski, Dumitru Erhan, Chelsea Finn, Piotr Kozakowski, Sergey Levine, et al. Model-based reinforcement learning for atari. *arXiv preprint arXiv:1903.00374*, 2019. [2](#)
- [23] Peter Karkus, David Hsu, and Wee Sun Lee. QMDP-Net: Deep Learning for Planning under Partial Observability. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4694–4704, 2017. [2](#)
- [24] Lydia Kavvaki, Petr Svestka, J.C. Latombe, and M.H. Overmars. Probabilistic roadmaps for path planning in high-dimensional configuration spaces. *IEEE Transactions on Robotics and Automation*, 1996. [2](#)
- [25] Steven M. Lavalle. Rapidly-exploring random trees: A new tool for path planning. *Computer Science Department, Iowa State University*, 1998. [2](#)
- [26] Steven M. Lavalle and Jr. James J. Kuffner. Rapidly-exploring random trees - progress and prospects. *Algorithmic and Computational Robotics: New Directions*, 2000. [2](#)
- [27] Lisa Lee, Emilio Parisotto, Devendra Singh Chaplot, Eric Xing, and Ruslan Salakhutdinov. Gated Path Planning Networks. *arXiv:1806.06408 [cs, stat]*, June 2018. arXiv: 1806.06408. [2](#), [4](#), [5](#), [6](#), [7](#), [11](#), [12](#)
- [28] Daniel James Lenton, Stephen James, Ronald Clark, and Andrew Davison. End-to-end egospheric spatial memory. In *International Conference on Learning Representations*, 2021. [5](#), [12](#)
- [29] Manolis Savva*, Abhishek Kadian*, Oleksandr Maksymets*, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia

- Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A Platform for Embodied AI Research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 8
- [30] Piotr Mirowski, Razvan Pascanu, Fabio Viola, Hubert Soyer, Andrew J. Ballard, Andrea Banino, Misha Denil, Ross Goroshin, Laurent Sifre, Koray Kavukcuoglu, Dharshan Kumaran, and Raia Hadsell. Learning to Navigate in Complex Environments. *arXiv:1611.03673 [cs]*, Jan. 2017. arXiv: 1611.03673. 2
- [31] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous Methods for Deep Reinforcement Learning. In *Proceedings of the International Conference on Machine Learning*, pages 1928–1937, June 2016. ISSN: 1938-7228 Section: Machine Learning. 2, 7
- [32] Raul Mur-Artal and Juan D. Tardos. ORB-SLAM2: an Open-Source SLAM System for Monocular, Stereo and RGB-D Cameras. *IEEE Transactions on Robotics*, 33(5):1255–1262, Oct. 2017. arXiv: 1610.06475. 5, 11
- [33] Andrew Y. Ng and Stuart J. Russell. Algorithms for Inverse Reinforcement Learning. In *Proceedings of the International Conference on Machine Learning, ICML '00*, pages 663–670, San Francisco, CA, USA, June 2000. 2
- [34] Buqing Nie, Yue Gao, Yidong Mei, and Feng Gao. Capability iteration network for robot path planning. *International Journal of Robotics and Automation*, 36(0), 2021. 2
- [35] Junhyuk Oh, Satinder Singh, and Honglak Lee. Value prediction network. In *Advances in Neural Information Processing Systems*, pages 6118–6128, 2017. 2
- [36] Emilio Parisotto and Ruslan Salakhutdinov. Neural map: Structured memory for deep reinforcement learning. *arXiv preprint arXiv:1702.08360*, 2017. 2
- [37] Stephane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 627–635, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR. 3, 4
- [38] G A Rummery and M Niranjan. Online Q-learning using Connectionist Systems. *Department of Engineering, University of Cambridge*, page 21, 1994. 2
- [39] Nikolay Savinov, Alexey Dosovitskiy, and Vladlen Koltun. Semi-parametric Topological Memory for Navigation. *International Conference on Learning Representations (ICLR)*, Mar. 2018. 2, 6
- [40] Daniel Schleich, Tobias Klamt, and Sven Behnke. Value iteration networks on multiple levels of abstraction. *Robotics: Science and Systems XV*, Jun 2019. 2
- [41] John Schulman, F. Wolski, Prafulla Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *ArXiv, abs/1707.06347*, 2017. 7
- [42] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 806–813, 2014. 1
- [43] Anthony Stentz. Optimal and Efficient Path Planning for Partially-Known Environments. In *Proceedings of the International Conference on Robotics and Automation*, page 8, 1994. 2
- [44] Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2nd edition, 2018. 2, 3
- [45] Aviv Tamar, YI WU, Garrett Thomas, Sergey Levine, and Pieter Abbeel. Value Iteration Networks. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 2154–2162, 2016. 1, 2, 3, 4, 5, 6, 7, 11, 12, 14
- [46] Sebastian Thrun. Probabilistic robotics. *Communications of the ACM*, 45(3):52–57, 2002. 1
- [47] Christopher Watkins. *Learning From Delayed Rewards*. PhD thesis, University of Cambridge, Jan. 1989. 2
- [48] Ronald J. Williams. Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning. In *Machine Learning*, pages 229–256, 1992. 2
- [49] David Bruce Wilson. Generating random spanning trees more quickly than the cover time. In *Proceedings of the ACM symposium on Theory of Computing, STOC '96*, pages 296–303, Philadelphia, Pennsylvania, USA, July 1996. 6
- [50] Brian D Ziebart, J Andrew Bagnell, and Anind K Dey. Modeling Interaction via the Principle of Maximum Causal Entropy. In *27th International Conference on Machine Learning*, page 8, 2010. 2
- [51] Brian D Ziebart, Andrew Maas, J Andrew Bagnell, and Anind K Dey. Maximum Entropy Inverse Reinforcement Learning. *AAAI Conference on Artificial Intelligence*, page 6, 2008. 2

Appendix A.

A.1. Implicit assumptions of the VIN architecture

Differences between VIN and the idealised VI on a grid.

Comparing eq. 1 in sec. 3 to eq. 2 in sec. 3.1, we note several differences:

1. The VIN value estimate V takes the maximum action-value Q across *all possible actions* \mathcal{A} , even illegal ones (e.g. moving into an obstacle). Similarly, the Q estimate is also updated even for illegal actions. In contrast, VI only considers legal actions for each state (cell), i.e. $\mathcal{A}(i, j)$.
2. The VIN reward \hat{R} is assumed independent of the action. This means that, for example, a transition between two states cannot be penalised directly; a penalty must be assigned to one of the states (regardless of the action taken to enter it).
3. The VIN transition probability is expanded into 2 terms, P^R which affects the reward and P^V which affects the estimated values. This decoupling means that they do not enjoy the physical interpretability of VI's P (i.e. probability of state transitions), and rewards and values can undergo very different transition dynamics.
4. Unlike the VI, the VIN considers the state transition translation-invariant. This means that it cannot model obstacles (illegal transitions) using P , and must rely on assigning a high penalty to the reward \hat{R} of those states instead.

A.2. Implementation details

A.2.1 Embodied pose states

One of our contributions is extending the VIN framework to accommodate embodied pose states, i.e. states which encode both position and orientation. We achieve this by augmenting the 2D state-space with an extra dimension for orientations. Table 4 shows correspondences between tensor dimensions of the positional method and the embodied method for each component of the architecture. X and Y are the size of the internal spatial discretisation of the environment, M is the internal discretisation of the orientation, A is the number of actions, and K is the kernel dimension for spatial locality.

Note that the value iteration step in CALVIN performs a 2D convolution of \hat{P} over a 2D value map in the case of positional states and a 3D convolution over a 3D value map with orientation in the case of embodied pose states. In the embodied case, the second dimension of \hat{P} corresponds to the orientation of the current state, and the third dimension corresponds to that of the next state.

A.2.2 3D embeddings for geometric reasoning

Since the learnable functions (\hat{P} , \hat{A} and \hat{R}) in our proposed method (and other VIN-based methods) are 2D CNNs, their natural input is a 2D grid of m -dimensional embeddings, denoted $e_{tij} \in \mathbb{R}^m$, for time t and discrete world-space coordinates (i, j) . This can be interpreted as a spatio-temporal

Table 4. Comparison of individual components in the implementation of CALVIN for positional states and for embodied pose states.

	Positional	Embodied
State s	(x, y)	(θ, x, y)
VI step	Conv2d	Conv3d
$V(s)$	$X \times Y$	$M \times X \times Y$
$Q(s, a)$	$A \times X \times Y$	$A \times M \times X \times Y$
$\hat{A}(s, a)$	$A \times X \times Y$	$A \times M \times X \times Y$
\hat{P}	$A \times K \times K$	$A \times M \times M \times K \times K$
\hat{R}	$A \times K \times K$	$A \times M \times M \times K \times K$

map tensor. We then wish to project and aggregate useful semantic information from an image I_t , extracted by a CNN ϕ , into this tensor. This requires both knowledge of the camera position c_t and rotation matrix R_t , which we assume following previous work [13, 27, 45] (and which can be estimated from monocular vision [32]). Spatial projection also requires knowing (or estimating) the depths $d_t(p)$ of each pixel p in I_t (either with a RGBD camera as in our experiments, or monocular depth estimation [11]). We can then write the homogenous 3D coordinates of each pixel p in the absolute reference frame using projective geometry [17]:

$$[x_t(p), y_t(p), z_t(p), 1] = c_t + R_t K [p_1, p_2, d_t(p), 1]^T, \quad (9)$$

where K is the camera's intrinsics matrix. Given these absolute coordinates of pixel p , we can calculate the closest map embedding e_{tij} to it, and thus aggregate the CNN embeddings $\phi(I_t)$ associated with all pixels close to a map cell. Inspired by PointNet [7], we choose mean-pooling for aggregation. Since we have spatial aggregation, we can easily extend this framework to work spatio-temporally, aggregating information from past frames $t' \leq t$. More formally:

$$e_{tij} = \text{avg}_{t' \leq t} \{ \phi_p(I_{t'}) : \tau i \leq x_{t'}(p) < \tau(i+1), \\ \tau j \leq y_{t'}(p) < \tau(j+1), \\ \tau k \leq z_{t'}(p) < \tau(k+1), p \in I_{t'} \} \quad (10)$$

where τ is the absolute size of each square grid cell, avg averages the elements of a set, and $\phi_p(I_{t'})$ retrieves the CNN embedding of image $I_{t'}$ for pixel p . Due to the similarity between Eq. 10 and a PointNet embedded on a 2D lattice, we named it Lattice PointNet (LPN). Other than the lattice embedding, there are other major differences from the PointNet: we apply it spatio-temporally with a causal constraint ($t' \leq t$), and the downstream predictors that take it as input ($\hat{P}(e_t)$, $\hat{A}(e_t)$ and $\hat{R}(e_t)$) are 2D CNNs that can reason spatially in the lattice, as opposed to the PointNet's unstructured multi-layer perceptrons [7]. A related proposal for SLAM used spatial max-pooling but more complex LSTMs/GRUs

for temporal aggregation [4, 19]. Another related work on end-to-end trainable spatial embeddings uses egospherical memory [28].

A.2.3 Architectural design of Lattice PointNet

The Lattice PointNet described in Appendix A.2.2 consists of three stages: a CNN that extracts embeddings from observations in image-space (image encoder), a spatial aggregation step (eq. 10 in sec. 4.2.2) that performs mean pooling of embeddings for each map cell, and another CNN that refines the map embedding (map encoder). The image encoder consists of two CNN blocks, each consisting of the following layers in order: optional group normalisation, 2D convolution, dropout, ReLU and 2D max pooling. The map encoder consists of 2D convolution, dropout, ReLU, optional group normalisation, and finally, another 2D convolution. The number of channels of each convolutional layer are (80, 80, 80, 40) for MiniWorld and (40, 40, 40, 20) for AVD respectively. The point clouds can consume a significant amount of memory for long trajectories. Hence, we use the most recent 40 frames for the 8×8 MiniWorld maze.

The input to the LPN is a 3-channel RGB image for the MiniWorld experiment, and a 128-channel embedding extracted using the first 2 blocks of ResNet18 pre-trained on ImageNet for the AVD experiment.

A.2.4 Architectural design of the CNN backbone

This CNN backbone is used in a control experiment in Appendix A.4.2 to show the effectiveness of the LPN backbone. In contrast to LPN which performs spatial aggregation of embeddings, the CNN backbone is a direct application of an encoder-decoder architecture that transforms image-space observations into map-space embeddings. Gupta *et al.* [13] employed a similar architecture to obtain their map embeddings. While they use ResNet50 as the encoder network, we used a simple CNN for the MiniWorld experiment to match the result obtained with LPN.

The CNN backbone consists of three stages: a CNN encoder, two fully-connected layers with ReLU to transform embeddings from image-space to map-space, and a CNN decoder. The encoder consists of 3 blocks of batch normalisation, 2D convolution, dropout, ReLU and 2D max pooling, and a final block with just batch normalisation and 2D convolution. The number of channels of each convolutional layer are (64, 128, 128, 128), respectively.

The fully-connected layers take in an input size of $128 \times 5 \times 7$, reduces it to a hidden size of 128, and outputs either $128 \times 5 \times 5$ for the smaller maze or $128 \times 4 \times 4$ for the larger maze, which is then passed to the decoder.

The decoder consists of 3 blocks of batch normalisation, 2D deconvolution, dropout and ReLU, and a final block with just 2D deconvolution. The number of channels

of each deconvolution layers are (128, 128, 64, 20), respectively. The output size of the decoder depends on the map resolution, hence we chose appropriate strides, kernel sizes and paddings in the decoder network to match the output sizes of 30×30 and 80×80 . This approach is not scalable to maps with high resolution or with arbitrary size, which is one of the drawbacks of this approach.

A.3. Experiment setup

A.3.1 Expert trajectory generation

Expert trajectories are generated by running an A* [16] planner from the start state to the target state. We assigned Euclidean costs to every transition in the 2D grid environments, and a cost of 1 per move for the MiniWorld and AVD environments. In the case of MiniWorld, an additional cost is assigned to locations near obstacles to ensure that the trajectories are not in close proximity to the walls.

A.3.2 Hyperparameter choices

Similarly to VIN [45] which uses a 2-layer CNN to predict the reward map, and GPPN [27], which uses a 2-layer CNN to produce inputs to the LSTM, CALVIN uses a 2-layer CNN as an available actions predictor $\hat{A}(s, a)$. For each experiment, we chose the size of the hidden layer from $\{40, 80, 150\}$. 150 was used for all the grid environments, 80 for MiniWorld and 40 for AVD, partially due to memory constraints.

VIN has an additional hyperparameter for the number of hidden action channels, which we set to 40, which is sufficiently bigger than the number of actual actions in all of our experiments. While the kernel size K for VIN and CALVIN were set to 3 for experiments in the grid environment, it was noted in [27] that GPPN works better with larger kernel size. Therefore, we chose the best kernel size out of $\{3, 5, 7, 9, 11\}$ for GPPN. For experiments on MiniWorld and AVD, there are state transitions with step size of 2, hence we chose $K = 5$ for VIN and CALVIN.

The number of value iteration steps k was chosen from $\{20, 40, 60, 80, 100\}$. For trajectory reweighting, β was chosen from $\{0.1, 0.25, 0.5, 0.75, 1.0\}$.

A.3.3 Rollout at test time

We test the performance of the model by running navigation trials (rollouts) on a randomly generated environment. At every time step, the model is queried the set of Q -values $\{Q(s, a) : a \in \mathcal{A}\}$ for the current state s , and an action which gives the maximum predicted Q value is taken.

While VIN is trained with $V^{(0)}$ initialised with zeros, in a true Value Iteration algorithm, the value function must converge for an optimal policy to be obtained. To help the value function converge faster under a time and compute budget,

we initialise the value function with predicted values from the previous time step at test time with online navigation.

We set a limit to the maximum number of steps taken by the agent, which were 200 for the fully-known 15×15 grid, 500 for the partially known grid, 300 for MiniWorld (3×3), 1000 for MiniWorld (8×8), and 100 for AVD.

A.4. Additional experiments

A.4.1 Ablation study of removing loss components

CALVIN is trained on three additive loss components: a loss term for the predicted Q-values L_Q (sec.4.1.1), a loss term for the transition models L_P (sec.4.1.2), and a loss term for the action availability L_A (sec.4.1.3). We assessed the contribution of each loss component to the overall performance.

We conducted the experiments on the partially observable grid environment (sec. 5.1.2). The results in Tab. 5 indicate that all loss components, in particular the transition model loss, contributes to the robust performance of the network.

Table 5. Navigation success rate of CALVIN in partially observable 2D mazes with loss components removed.

Loss	$L_Q + L_P + L_A$	$L_Q + L_P$	$L_Q + L_A$
Success rate	92.2	84.1	8.3

A.4.2 Comparison of LPN against CNN backbone

We compared our proposed LPN backbone against a typical encoder-decoder CNN backbone as a component that maps observations to map embeddings. We evaluated the performance of the two methods for VIN, GPPN and CALVIN. In Tab. 6, we observe that LPN backbone is highly effective, especially for larger environments where long-term planning based on spatially aggregated embeddings is necessary.

Table 6. Navigation success rate on unseen 3D mazes (MiniWorld). Most methods do not generalise to larger mazes. The proposed LPN demonstrates robust performance in larger unseen mazes.

Size	CNN backbone			LPN backbone (ours)		
	VIN	GPPN	CALVIN	VIN	GPPN	CALVIN
3×3	89.4	73.1	75.2	90.3	91.3	97.7
8×8	0.6	18.3	8.6	41.2	33.3	69.2

A.5. Example rollout in a partially observable maze

We present an example of a trajectory taken by CALVIN at runtime, with corresponding observation maps and predicted values in Figure 6. At each rollout step, CALVIN performs inference on the best action to take based on its current observation map. No information about the location of the target is given until it is within view of the agent. This

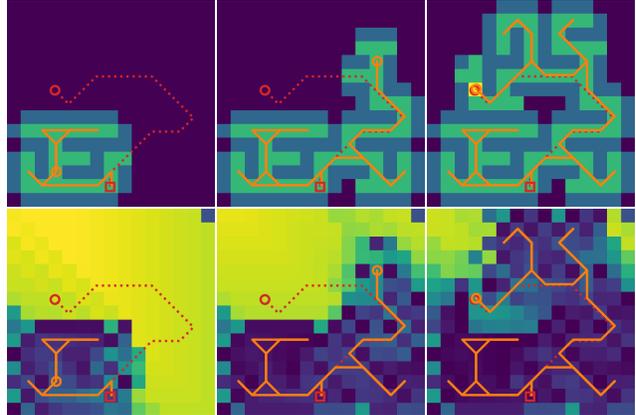


Figure 6. Example rollout of CALVIN after 21 steps (left column), 43 steps (middle column) and 65 steps (right column). CALVIN successfully terminated at 65 steps. **(top row)** Input visualisation: unexplored cells are dark, the discovered target is yellow. The correct trajectory is dashed, the current one is solid. The orange circle shows the position of the agent. **(bottom row)** Predicted values (higher values are brighter). Explored cells have low values, while unexplored cells and the discovered target are assigned high values.

makes the problem challenging, since the agent may have to take significantly more steps compared to an optimal route to reach the target. In this example, the agent managed to backtrack every time it encountered a dead end, successfully reaching the target after 65 steps. The model initially assigns high values to all unexplored states. When the target comes into view, the model assigns a high probability to the availability of the “done” action at the corresponding state. The agent learns a sufficiently high reward for a successful termination so that the “done” action is triggered at the target.

A.6. Comparison of embodied navigation

For visual comparison of CALVIN, VIN and GPPN, we generated a maze and performed rollouts using each of the algorithms, assuming partial observability and embodied navigation.

A.6.1 Rollout of CALVIN

Figure 7 shows an example of a trajectory taken by CALVIN at runtime, with corresponding observation maps, predicted values and predicted rewards for taking the “done” action. Similarly to Appendix A.5, the agent manages to explore unvisited cells and backtrack upon a dead end until the target is discovered. One key difference is that now the agent learns to predict rewards and values for every discretised orientation as well as the discretised location. Upon closer inspection, we observe that the predicted values are higher facing the direction of unexplored cells and towards the discovered target.

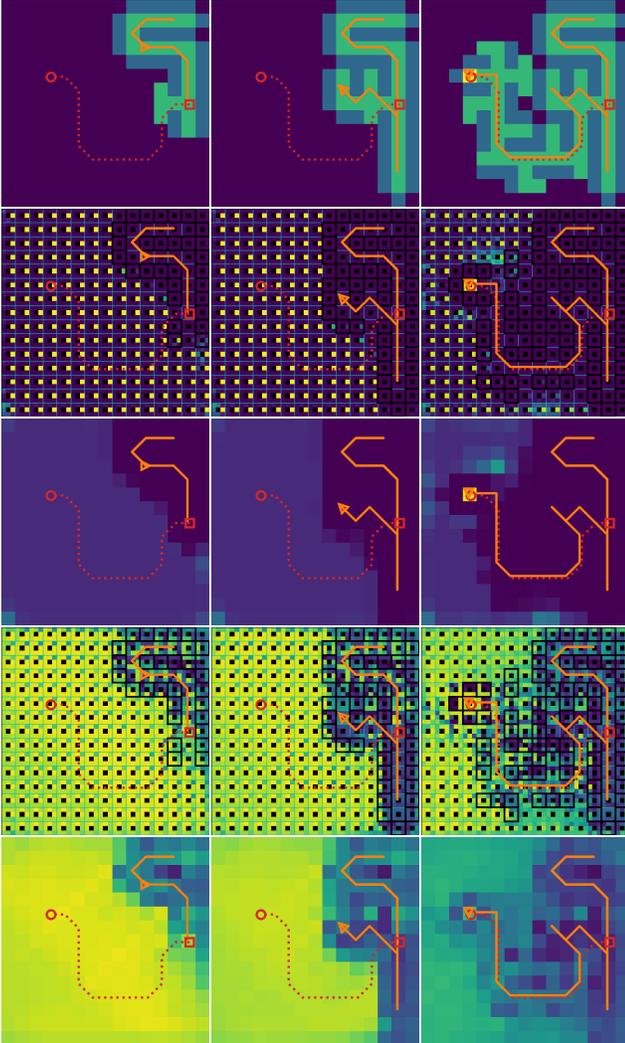


Figure 7. Example rollout of embodied CALVIN after 30 steps (left column), 60 steps (middle column) and 90 steps (right column). CALVIN successfully terminated at 91 steps. **(first row)** Input visualisation: unexplored cells are dark, the discovered target is yellow. The correct trajectory is dashed, the current one is solid. The orange triangle shows the position and the orientation of the agent. **(second row)** Predicted rewards (higher values are brighter). The 3D state-space (position/orientation) is shown, with rewards for the 8 orientations in a radial pattern within each cell (position). Explored cells have low rewards, while unexplored cells and the discovered target are assigned high rewards. **(third row)** Predicted rewards averaged over the 8 orientations. **(fourth row)** Predicted values following the same convention. Values are higher facing the direction of unexplored cells and the target (if discovered). **(fifth row)** Predicted values averaged over the 8 orientations.

Since rotation is a relatively low cost operation, in this training example, the network seems to have learnt to assign high rewards to a particular orientation at unexplored cells, from which high values propagate. Rewards and values averaged over orientations yield a more intuitive visualisation.

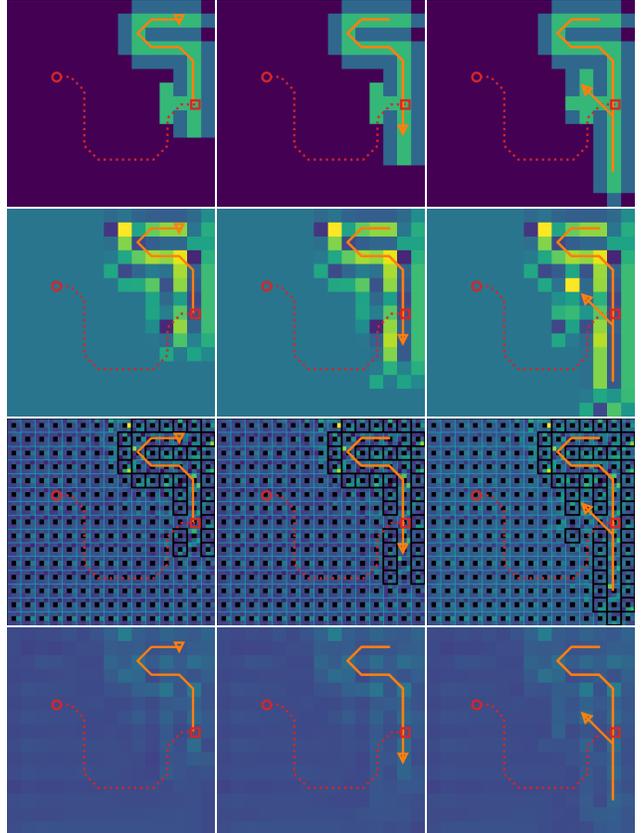


Figure 8. Example rollout of embodied VIN after 20 steps (left column), 40 steps (middle column) and 60 steps (right column). VIN kept oscillating between the same two states after 57 steps. The convention is the same as for Fig. 6, except that a single reward map is shared across all orientations. **(first row)** Input visualisation. **(second row)** Predicted rewards. **(third row)** Predicted rewards averaged over the 8 orientations. **(fourth row)** Predicted values.

A.6.2 Rollout of VIN

A corresponding visualisation for VIN is shown in Fig. 8. Unlike CALVIN’s implementation of rewards (eq. 5 in sec. 4.1) as a function of discretised states and actions, the “reward map” produced by the VIN does not offer a direct interpretation, as it is shared across all actions as implemented by Tamar *et al.* [45], and is also shared across all orientations in the case of embodied navigation. The values are also not well learnt, with some of the higher values appearing in obstacle cells. The unexplored cells are not assigned sufficiently high values to incentivise exploration by the agent. In this example, the agent gets stuck and starts oscillating between two orientations after 57 steps.

A.6.3 Rollout of GPPN

Finally, a visualisation for GPPN is shown in Fig. 9. Unlike VIN and CALVIN, GPPN does not have an explicit

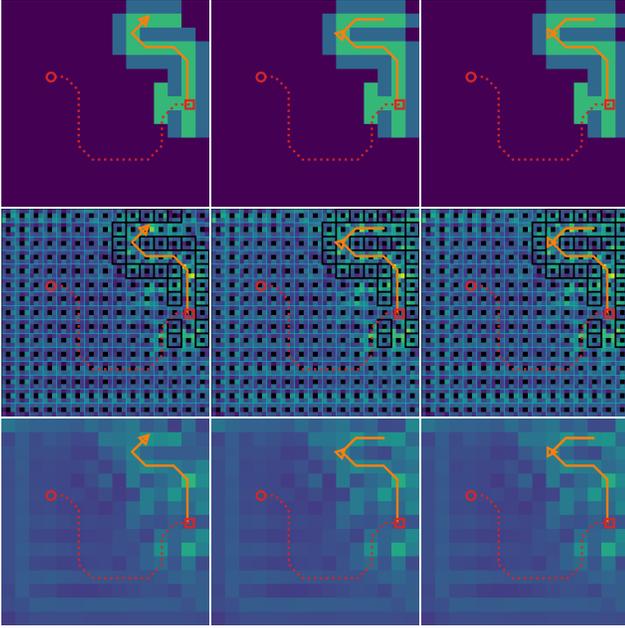


Figure 9. Example rollout of embodied GPPN after 15 steps (left column), 30 steps (middle column) and 45 steps (right column). GPPN revisits the same sequences of states leading to a dead end after 45 steps. The convention is the same as for Fig. 6. **(first row)** Input visualisation. **(second row)** Predicted rewards. **(third row)** Predicted rewards averaged over the 8 orientations.

reward map predictor, but performs value propagation using an LSTM before outputting a final Q-value prediction. Similarly to Appendix A.6.2, the values predicted is not very interpretable, and does not incentivise exploration or avoidance of dead ends. In this example, the agent keeps revisiting a dead end that has already been explored in the first 20 steps.