Spatially Structured Recurrent Modules

Nasim Rahaman¹² Anirudh Goyal² Muhammad Waleed Gondal¹ Manuel Wuthrich¹ Stefan Bauer¹ Yash Sharma³ Yoshua Bengio² Bernhard Schölkopf¹

Abstract

Capturing the structure of a data-generating process by means of appropriate inductive biases can help in learning models that generalize well and are robust to changes in the input distribution. While methods that harness spatial and temporal structures find broad application, recent work (Goyal et al., 2019) has demonstrated the potential of models that leverage sparse and modular structure using an ensemble of sparingly interacting modules. In this work, we take a step towards dynamic models that are capable of simultaneously exploiting both modular and spatiotemporal structures. We accomplish this by abstracting the modeled dynamical system as a collection of autonomous but sparsely interacting sub-systems. The sub-systems interact according to a topology that is learned, but also informed by the spatial structure of the underlying real-world system. This results in a class of models that are well suited for modeling the dynamics of systems that only offer local views into their state, along with corresponding spatial locations of those views. On the tasks of video prediction from cropped frames and multi-agent world modeling from partial observations in the challenging Starcraft2 domain, we find our models to be more robust to the number of available views and better capable of generalization to novel tasks without additional training, even when compared against strong baselines that perform equally well or better on the training distribution.

1. Introduction

Many spatiotemporal complex systems can be abstracted as a collection of autonomous but sparsely interacting subsystems, where sub-systems tend to interact if they are in each others' *local vicinity* in some sense. As an illustrative example, consider a grid of traffic intersections, wherein traffic flows from a given intersection to the adjacent ones, and the actions taken by some "agent", say an autonomous vehicle, may at first only affect its immediate surroundings. Now suppose we want to forecast the future state of the traffic grid (say for the purpose of avoiding traffic jams).

There is a spectrum of possible strategies to model the system at hand. On one end of it lies the most general strategy: namely, one that calls for considering the entirety of all intersections simultaneously to predict the next state of the grid (Figure 1c). The resulting model class can in principle account for interactions between any two intersections, irrespective of their spatial distance. However, the number of interactions such models must consider does not scale well with the size of the grid, and the strategy might be rendered infeasible for large grids with hundreds of intersections.

On the other end of the spectrum is a specialized strategy that involves abstracting the dynamics of each intersection as an autonomous sub-system, and having each sub-system interact only with other sub-systems associated with the four (or more) neighboring intersections (Figure 1c). The interactions may manifest as messages that one sub-system passes to another and possibly contain information about how many vehicles are headed towards which direction, resulting in a collection of message passing entities (i.e., subsystems) that collectively model the entire grid. By adopting this strategy, one assumes that the immediate future of any given intersection is affected only by the present states of the neighboring intersections, and not some intersection at the opposite end of the grid. The resulting class of models scales well with the size of the grid, but is possibly unable to model certain long-range interactions that could be leveraged to efficiently distribute traffic flow.

The spectrum above parameterizes the extent to which the spatial structure of the underlying system being modeled is incorporated into the design of the model. The former extreme ignores spatial structure altogether, resulting in a class of models that can be expressive but whose sample and computational complexity do not scale well with the size of the system. The latter extreme results in a class of

^{*}Equal contribution ¹Max Planck Institute for Intelligent Systems, Tübingen ²Mila, Québec ³Bethgelab, University of Tübingen. Correspondence to: Nasim Rahaman <nasim.rahaman@tuebingen.mpg.de>.

S2RM: Spatially Structured Recurrent Modules



Figure 1. A schematic representation of the spectrum of modeling strategies. Solid arrows with speech bubbles denote (dynamic) messages being passed between *sub-systems* (dotted arrows denote the lack thereof). **Gist:** on the one end of the spectrum, (Figure 1a), we have the strategy of abstracting each intersection as a sub-system that interact with neighboring sub-systems. On the other end of the spectrum (Figure 1c) we have the strategy of modeling the entire grid with one monolithic system. The middle ground (Figure 1b) we explore involves letting the model *develop* a notion of locality by (say) abstracting entire avenues with a single sub-system.

models that can scale well, but its adequacy (in terms of expressivity) is contingent on a predefined notion of locality (in the example above: the immediate four-neighborhood of an intersection). In this work, we aim to explore a middleground between the two extremes: namely, by proposing a class of models that does leverage the spatial structure, but by developing a notion of locality instead of relying on a predefined one (Figure 1b). Reconsidering the traffic grid example: the proposed strategy results in a model that can potentially learn to abstract (say) entire avenues with a single sub-system. The interactions between intersections are therefore replaced by those between avenues, resulting in a scheme where a single sub-system might account for events that are spatially distant (such as those in the opposite ends of an avenue), but two events that are spatially closer together (such as those on two adjacent avenues of the same street, where streets run perpendicular to avenues) might be accounted for by different sub-systems.

To implement this scheme, we will model the sub-systems as independent recurrent neural networks (RNNs) that interact sparsely via a bottleneck of attention (Goyal et al., 2019), but extend this idea along two salient dimensions. First, we relax the assumption that the interaction topology between sub-systems (i.e., RNNs) is all-to-all, in the sense that all sub-systems are allowed to interact with all other sub-systems. We achieve this by learning to embed each sub-system in an embedding space endowed with a metric, and attenuate the interaction between two given subsystems by their distance in this space (i.e., sub-systems too far away from each other in this space are not allowed to interact). Second, instead of assuming that the entire system is perceived simultaneously, we only assume access to local (partial) observations alongside with the associated spatial locations, resulting in a setting that partially resembles that of Eslami et al. (2018). Expressed in the language of the example above: we do not expect a birds eye view of the traffic grid, but only (say) LIDAR observations from

autonomous vehicles at known GPS coordinates, or video streams from traffic cameras at known locations. The spatial location associated with an observation plays a crucial role in the proposed architecture in that we map it to the embedding space of sub-systems and *address* the corresponding observation only to sub-systems whose embeddings lie in close vicinity. Likewise, to predict future observations at a queried spatial location, we again map said location to the embedding space and poll the states of sub-systems situated nearby. The result is a model that can *learn* which spatial locations are to be associated with each other and be accounted for by the same sub-system. As an added plus, the parameterization we obtain is not only agnostic to the number of available observations and query locations, but also to the number of sub-systems.

To evaluate the proposed model, we choose a problem setting where (a) the task is composed of different sub-systems or processes that locally interact both spatially and temporally, and (b) the environment offers local views into its state paired with their corresponding spatial locations. The challenge here lies in building and maintaining a consistent representation of the global state of the system given only a set of partial observations. To succeed, a model must learn to efficiently capture the available observations and place them in appropriate spatial context. The first problem we consider is that of video prediction from crops, analogous to that faced by visual systems of many animals: given a set of small crops of the video frames centered around stochastically sampled pixels (corresponding to where the fovea is focused), the task is to predict the content of a crop around any queried pixel position at a future time. The second problem is that of multi-agent world modeling from partial observations in spatial domains, such as the challenging Starcraft2 domain (Samvelyan et al., 2019; Vinyals et al., 2017). The task here is to model the dynamics of the global state of the environment given local observations made by cooperating agents and their corresponding actions. Importantly and unlike prior work (Sun et al., 2019), our parameterization is agnostic to the number of agents in the environment, which can be flexibly adjusted on the fly as new agents become available or existing agents retire. This is beneficial for generalization in settings where the number of agents during training and testing are different.

Contributions. (a) We propose a new class of models, which we call Spatially Structured Recurrent Modules or S2RMs, which perform attention-driven spatially local modular computations. (b) We evaluate S2RMs (along with several strong baselines) on a selection of challenging problems to find that S2RMs are robust to the number of available observations and can generalize to novel tasks.

2. Problem Statement

In this section, we build on the intuition from the previous section to formally specify the problem we aim to approach with the methods described in the later sections.

Let \mathcal{X} be a metric space, \mathcal{O} some set of possible *observations*, and $\mathcal{D}_{\mathcal{X}}$ a set of mappings $\mathcal{X} \to \mathcal{O}$. Now, consider the *evolution function* of a discrete-time dynamical system:

$$\phi : \mathbb{Z} \times \mathfrak{O}_{\mathcal{X}} \to \mathfrak{O}_{\mathcal{X}} \text{ satisfying}$$
(1)
$$\phi(0, \mathfrak{o}) = \mathfrak{o} \text{ where } \mathfrak{o} \in \mathfrak{O}_{\mathcal{X}} \text{ and}$$

$$\phi(t_2, \phi(t_1, \mathfrak{o})) = \phi(t_1 + t_2, \mathfrak{o}) \text{ for } t_1, t_2 \in \mathbb{Z}$$

Informally, \mathfrak{o} can be interpreted as the *world state* of the system; together with a spatial *location* $\mathbf{x} \in \mathcal{X}$, it gives the *local* observation $\mathbf{O} = \mathfrak{o}(\mathbf{x}) \in \mathcal{O}$. Given an initial world state \mathfrak{o} , the mapping $\phi(t, \mathfrak{o})$ yields the world state at some (future) time *t*, thereby characterizing the dynamics of the system (which might be stochastic).

While the above class of dynamical systems is fairly general, we now place a crucial restriction: namely, that for any pair of space-time *events* $(t_1, \mathbf{x}_1) \in (\mathbb{Z} \times \mathcal{X})$ and (t_2, \mathbf{x}_2) and given any initial state $\mathfrak{o}_0 \in \mathfrak{O}_{\mathcal{X}}$, there exists a finite C > 0such that the observation $\phi(t_1, \mathfrak{o}_0)(\mathbf{x}_1)$ can influence the observation $\phi(t_2, \mathfrak{o}_0)(\mathbf{x}_2)$ only if $t_2 \ge t_1$ and $d_{\mathcal{X}}(\mathbf{x}_1, \mathbf{x}_2) \le$ $C \cdot (t_2 - t_1)$, where $d_{\mathcal{X}}$ is the metric on \mathcal{X} . This assumption induces a notion of spatio-temporal locality by imposing that the effect of any given event can only propagate at a finite speed, where the latter is upper bounded by C.

In this work, we are concerned with modelling systems that are subject to the above restriction. Assuming the system satisfies said restriction, we have the following

Problem: At every time step t = 0, ..., T, we are given a set of positions $\{\mathbf{x}_t^a\}_{a=1}^A$ and the corresponding observations $\{\mathbf{O}_t^a\}_{a=1}^A$, where $\mathbf{O}_t^a := \phi(t, \mathfrak{o}_0)(\mathbf{x}^a)$ for some initial world state \mathfrak{o}_0 . The task is to infer the world state $\phi(t', \mathfrak{o}_0)$ at some future time-step t' > T.

In the traffic grid example of Section 1, one could imagine

a as indexing traffic cameras or autonomous vehicles (i.e., *observers*), \mathbf{x}_t^a as the GPS coordinates of observer *a*, and \mathbf{O}_t^a as the corresponding sensor feed (e.g. LIDAR observations or video streams from vehicles or traffic cameras).

3. Modelling Assumptions

Given the problem in Section 2, we now constrain it by placing certain structural assumptions. These assumptions will ultimately inform the inductive biases we select for the model (proposed in Section 4); nevertheless, we remark beforehand that as with any inductive bias, their applicability is subject to the properties of the system being modeled and the objectives¹ being optimized.

Recurrent Dynamics Modeling. While there exist multiple ways of modeling dynamical systems, we shall focus on recurrent neural networks (RNNs). Typically, RNN-based dynamics models are expressed as functions of the form:

$$\mathbf{h}_{t+1} = F(\mathbf{O}_t, \mathbf{h}_t) \qquad \mathbf{O}_t = D(\mathbf{h}_t) \tag{2}$$

where O_t is the observation at time $t \in \mathbb{Z}$, and h_{t+1} is the hidden state of the model. *F* can be thought of as the parameterized *forward-evolution function* the hidden state h conditioned on the observation O, whereas *D* is a *decoder* that maps the hidden state to observations. Here, the evolution function of the modelled dynamical system (as defined in Equation 1) can be obtained by rolling out the forward-evolution function in time.

Decomposition into Sub-systems. Without loss of generality, one may assume that the dynamical system ϕ defined in Equation 1 can be decomposed into constituent systems $(\phi_1, \phi_2, ..., \phi_M)$, such that the interaction between all pairs of sub-systems (ϕ_i, ϕ_j) satisfy some criterion. Now, the strength of this assumed criterion lies on a spectrum. On one end of the spectrum is the case where such a criterion is non-existent, i.e., no such decomposition is assumed and full generality is restored; this is the modeling assumption made when using conventional recurrent models like GRUs (Cho et al., 2014), LSTMs (Hochreiter & Schmidhuber, 1997) and vanilla RNNs. On the other end of the spectrum lies a setting where the decomposition is required to be such that sub-systems do not interact, i.e., ϕ_i and ϕ_i have independent dynamics (Li et al., 2018). Goyal et al. (2019) explore a middle ground, where the interaction between sub-systems (ϕ_i, ϕ_i) are possible but constrained. In particular, they investigate a setting where the sub-systems are assumed to interact sparsely, and the interaction pattern (i.e., which sub-systems interact with which others) is dynamic and may depend on the world state o. In this work, we adopt the assumption of sparsely interacting sub-systems, but subject the interaction pattern to an additional spatial constraint.

¹E.g. generalization, sample complexity, robustness, etc.



Figure 2. Schematic representation of the proposed architecture.

Local Interactions Between Sub-systems. In addition to assuming dynamic sparse interactions between sub-systems, we also assume that a given sub-system ϕ_j may *preferentially* interact with another given sub-system ϕ_i . Intuitively, one may think of ϕ_j as lying in *vicinity* of ϕ_i . This naturally leads us to a notion of topology over sub-systems, one where sub-systems situated in each other's local neighborhood are less constrained in their interactions. In the next section, we will discuss how we model this topology by associating each sub-system ϕ_i with a learned embedding \mathbf{p}_i in an existing metric space, which we will call S. Subsequently, the affinity of sub-system ϕ_i to interact with another sub-system ϕ_j will be quantified by a *similarity measure* Z, such that $Z(\mathbf{p}_i, \mathbf{p}_j)$ is large if ϕ_i and ϕ_j prefer to interact.

Locality of Observations. Recall from Section 2 that the observations available to the model respect a notion of spatio-temporal locality. However, this notion of locality is distinct from the one between sub-systems (induced via Z), and one important modeling decision is how the two should interact. We propose to embed the position $\mathbf{x} \in \mathcal{X}$ associated with an observation O to the metric space of sub-systems S via a continuous and one-to-one mapping $P: \mathcal{X} \to \mathcal{S}$, which allows us to *match* the observation **O** to all sub-systems ϕ_m in the vicinity of $P(\mathbf{x}) \in \mathcal{S}$, i.e., where $Z(P(\mathbf{x}), \mathbf{p}_m)$ is sufficiently large. Likewise, the same subsystems ϕ_m are polled if the model is queried for a prediction at x. On a high level, this results in a scheme where each subsystem ϕ_m can account for observations made at a set of positions $\mathcal{X}_m \subset \mathcal{X}$, which we call its *enclave*. In particular, the enclaves \mathcal{X}_i and \mathcal{X}_i corresponding to subsystems ϕ_i and ϕ_j may overlap, and we do not constrain the distance between two given points in \mathcal{X}_m to be small.

4. Proposed Model

Informed² by the model assumptions detailed in the previous section, we now proceed to describe the proposed model – Spatially Structured Recurrent Modules or **S2RM** – which comprise the following components (Figure 2):

Model Inputs. Recall from Section 2 that we have for every time step t = 0, ..., T a set of tuples of positions and

observations $\{(\mathbf{x}_t^a, \mathbf{O}_t^a)\}_{a=1}^A$ where $\mathbf{x}_t^a \in \mathcal{X}$ and $\mathbf{O}_t^a \in \mathcal{O}$ for all t and a. To simplify, we assume that $\mathcal{X} \subset \mathbb{R}^n$, and denote by x_i the *i*-th component of the vector $\mathbf{x} \in \mathcal{X}$. **Encoder.** The encoder E is a parameterized function mapping observations \mathbf{O} to a corresponding vector representation $\mathbf{e} = E(\mathbf{O})$. Here, E processes all observations in parallel across t and a to yield representations \mathbf{e}_t^a .

Positional Embedding. The positional embedding P is a fixed mapping from \mathcal{X} to \mathcal{S} . We choose \mathcal{S} to be the unit sphere in *d*-dimensions, *d* being a multiple of 2n, and the positional encoder as the following function:

$$P(\mathbf{x}) = \mathbf{s}/\|\mathbf{s}\| \in \mathcal{S}$$
 where (3)

$$s_{2i+m} = \sin\left(\frac{x_m}{10000^i}\right) \ s_{2i+1+m} = \cos\left(\frac{x_m}{10000^i}\right) \ (4)$$

with m = 1, ..., n and i = 1, ..., d/2n. While the above function is commonly used (Vaswani et al., 2017), other choices might also be viable. Accommodating a slight abuse of notation, we will refer to $P(\mathbf{x})$ as \mathbf{s} and $P(\mathbf{x}_t^a)$ as \mathbf{s}_t^a .

Set of Interacting RNNs. To model the dynamics of the world state, we use a set of M independent RNN modules (like in Goyal et al. (2019)), which we denote as $\{F_m\}_{m=1}^M$. To each F_m , we associate an embedding vector $\mathbf{p}^m \in S$, where all $\{\mathbf{p}^m\}_{m=1}^M$ are learnable parameters. On a high level, RNNs F_m interact with each other via an *inter-cell attention*, and with the input representations \mathbf{e}_t^a via *input attention*. More precisely, at a given time step t, each F_m expects an input \mathbf{u}_t^m , together with an *aggregated hidden state* $\bar{\mathbf{h}}_t^m$ and optionally a memory state \mathbf{c}_t^m to yield the hidden and memory states at the next time step:

$$(\mathbf{h}_{t+1}^m, \mathbf{c}_{t+1}^m) = F_m(\mathbf{u}_t^m, \bar{\mathbf{h}}_t^m, \mathbf{c}_t^m)$$
(5)

where the input \mathbf{u}_t^m results from the input attention and $\bar{\mathbf{h}}_t^m$ from the inter-cell attention (both described below). If available, the memory state \mathbf{c}_t^m resembles the cell state in an LSTM (Hochreiter & Schmidhuber, 1997).

Input Attention. Similar to MHDPA (multi-head dotproduct attention, Vaswani et al. (2017)), the input attention mechanism is a mapping between sets: namely, from that of observation encodings $\{\mathbf{e}_t^a\}_{a=1}^A$ to that of RNN inputs $\{\mathbf{u}_t^m\}_{m=1}^M$. In what follows, we use the einsum notation³ to succintly describe the exact mechanism. But before that, we define the truncated spherical Gaussian kernel (Fasshauer, 2011) to quantify the similarity between two points $\mathbf{p}, \mathbf{s} \in S$:

$$Z(\mathbf{p}, \mathbf{s}) = \begin{cases} \exp\left[-2\epsilon(1 - \mathbf{p} \cdot \mathbf{s})\right], & \text{if } \mathbf{p} \cdot \mathbf{s} \ge \tau \\ 0, & \text{otherwise} \end{cases}$$
(6)

where $\epsilon \in \mathbb{R}^+$ and $\tau \in [-1, 1)$ are hyper-parameters (kernel bandwidth and truncation parameter, respectively), and

 $^{^{2}}$ In doing so, we use the assumptions merely as guiding principles; we do not claim that we *infer* e.g. the true decomposition of the ground-truth system, even if all assumptions are satisfied.

³Indices not appearing on both sides of an equation are summed over; this is implemented as einsum in most DL frameworks.

 $0 \le Z \le 1$ since **p** and **s** are unit vectors. Now, we use k to index the attention heads, d to index the dimension of the key and query vectors, and denote with e_{ai} the *i*-th component of \mathbf{e}_t^a and with h_{mj} the *j*-th component of \mathbf{h}_t^m . Given learnable parameters $\Theta^{(K)}, \Theta^{(Q)}, \Theta^{(V)}$, we obtain:

$$Q_{akd} = e_{ai}\Theta_{ikd}^{(Q)} \qquad K_{mkd} = h_{mj}\Theta_{jkd}^{(K)} \tag{7}$$

$$V_{akv} = e_{ai}\Theta_{ikv}^{(V)} \qquad \tilde{W}_{mak} = Q_{akd}K_{mkd} \tag{8}$$

$$\bar{W}_{mak} = \operatorname{sm}_a(\tilde{W}_{mak}) \qquad W_{ma}^{(L)} = Z(\mathbf{p}^m, \mathbf{s}^a) \tag{9}$$

$$W_{mak} = W_{ma}^{(L)} \bar{W}_{mak} \qquad \tilde{u}_{m(kv)} = W_{mak} V_{akv} \quad (10)$$

where: sm_a denotes softmax along the *a*-dimension, $W^{(L)}$ is what we will call the *local weights*, we omit the time subscript in \mathbf{s}^a for notational clarity, and $\tilde{u}_{m(kv)}$ is the (kv)-th component of a vector $\tilde{\mathbf{u}}^m$. Finally, we obtain the components u_{mi} of RNN inputs \mathbf{u}_t^m via a gating operation:

$$u_{mi} = G_m^{(\text{inp})} \cdot b_{mi} + (1 - G_m^{(\text{inp})}) \cdot \tilde{u}_{mi}$$
(11)

where the gating weight $G_m^{(inp)} \in (0, 1)$ is obtained by passing \tilde{u}_{mi} and $b_{mi} = W_{ma}^{(L)} e_{ai}$ through a two-layer MLP with sigmoidal output (in parallel across m). Now, observe that by weighting the MHDPA attention outputs (\bar{W} in Equation 10) by the kernel Z (via $W^{(L)}$), we construct a scheme where the interaction between input \mathbf{O}_t^a and RNN F_m is allowed only if the embedding \mathbf{s}_t^a of the corresponding position \mathbf{x}_t^a has a large enough cosine similarity ($\geq \tau$) to the embedding \mathbf{p}_m of F_m . This partially implements the assumption of *Locality of Observation* detailed in Section 3.

Inter-cell Attention. The inter-cell attention maps the hidden states of each RNN $\{\mathbf{h}_{t}^{m}\}_{m=1}^{M}$ to the set of *aggregated* hidden states $\{\bar{\mathbf{h}}_{t}^{m}\}_{m=1}^{M}$, thereby enabling interaction between the RNNs F_{m} . While its mechanism is identical to that of the input attention, we formulate it below for completeness. To proceed, we denote with h_{li} the *i*-th component of \mathbf{h}_{t}^{l} (in addition to the notation introduced before Equation 7), and take $\Phi^{(Q)}$, $\Phi^{(K)}$ and $\Phi^{(V)}$ to be learnable parameters. We have:

$$Q_{mkd} = h_{mj} \Phi_{jkd}^{(Q)} \qquad K_{lkd} = h_{li} \Phi_{ikd}^{(K)}$$
(12)

$$V_{lkv} = h_{li}\Phi_{ikv}^{(V)} \qquad \tilde{W}_{mlk} = Q_{mkd}K_{lkd}$$
(13)

$$W_{mlk} = \operatorname{sm}_l(W_{mlk}) \qquad W_{ml}^{(L)} = Z(\mathbf{p}^m, \mathbf{p}^l) \qquad (14)$$

$$W_{mlk} = \bar{W}_{mlk} W_{ml}^{(L)} \qquad \tilde{h}_{m(kv)} = W_{mlk} V_{lkv}$$
(15)

where $\bar{h}_{m(kv)}$ is the (kv)-th component of a vector $\tilde{\mathbf{h}}^m$. Finally, the *j*-th component \bar{h}_{mj} of the aggregated hidden state $\bar{\mathbf{h}}_t^m$ in Equation 5 is given by a gating operation:

$$\bar{h}_{mj} = G_m^{(\mathrm{ic})} \cdot c_{mj} + (1 - G_m^{(\mathrm{ic})}) \cdot \tilde{h}_{mj}$$
(16)

where the gating weight $G_m^{(ic)} \in (0,1)$ is obtained by passing \tilde{h}_{mj} and $c_{mj} = W_{ml}^{(L)} h_{lj}$ through a two-layer MLP with sigmoid output (in parallel across m). The weighting by Z (in Equation 15, left) ensures that the interaction is constrained to be only between RNNs whose embeddings in S are similar enough, thereby implementing the assumption of *Local Interactions between Sub-systems* in Section 3.

Output Attention. The output attention mechanism together with the decoder (described below) serve as an apparatus to evaluate the world state modeled (implicitly) by the set of RNNs ($\{F_m\}_{m=1}^M$) at time t + 1 (for one-step forward models). Given a query location $\mathbf{x}^q \in \mathcal{X}$ and its corresponding embedding $\mathbf{s}^q = P(\mathbf{x}^q) \in \mathcal{S}$, the output attention mechanism polls the RNNs F_m whose embeddings \mathbf{p}^m are similar enough to \mathbf{s}^q , as measured by the kernel Z. Denoting h_{mj} the j-th component of \mathbf{h}_{t+1}^m , we have:

$$d_j^q = Z(\mathbf{s}^q, \mathbf{p}^m) h_{mj} \tag{17}$$

where d_j^q can be interpreted as the *j*-th component of the vector \mathbf{d}_{t+1}^q associated with the query location \mathbf{x}^q .

Decoder. The decoder *D* is a parameterized function that predicts the observation $\hat{\mathbf{O}}_{t+1}^q \in \mathcal{O}$ at \mathbf{x}^q given the representation \mathbf{d}_{t+1}^q from the output attention.

This concludes the description of the generic architecture, which allows for flexibility in the choice of the RNN architecture (i.e., the internal architecture of F_m). In practice, we find Gated Recurrent Units (GRUs) (Cho et al., 2014) to work well, and call the resulting model Spatially Structured GRU or **S2GRU**. Moreover, Relational Memory Cores (RMCs) (Santoro et al., 2018) also profit from our architecture (with a minor modification detailed in Appendix B.3), and we refer to the resulting model as **S2RMC**.

5. Related Work

Problem Setting. Recall that the problem setting we consider is one where the environment offers local (partial) views into its global state paired with the corresponding spatial locations. With Generative Query Networks (GQNs), Eslami et al. (2018) investigate a similar setting where the 2D images of 3D scenes are paired with the corresponding *viewpoint* (camera position, yaw, pitch and roll). Given that GQNs are feedforward models, they do not consider the dynamics of the underyling scene and as such cannot be expected to be consistent over time (Kumar et al., 2018). Singh et al. (2019) and Kumar et al. (2018) propose variants that are temporally consistent, but unlike us, they do not focus on the problem of modeling the forward dynamics.

Modularity. Modularity has been a recurring topic in the context of meta-learning (Alet et al., 2018; Bengio et al., 2019; Ke et al., 2019), sequence modeling (Henaff et al., 2016; Goyal et al., 2019; Li et al., 2018) and beyond (Jacobs et al., 1991; Shazeer et al., 2017; Parascandolo et al., 2017). However, unlike prior work, we integrate modularity and

spatio-temporal structure in a unified framework.

Spatial Attention. Mechanisms for spatial attention have been well studied (Jaderberg et al., 2015; Wang et al., 2017; Zhang et al., 2018; Parmar et al., 2018), but they typically operate on image pixels. Our setting is more general in the sense that we do not necessarily require that the world state of the underlying system be represented by images.

Attention Mechanisms and Information Flow. Attention mechanisms have been used to attenuate the flow of information between components of the network, e.g. in NTMs (Graves et al., 2014), DNCs (Graves et al., 2016), RMCs (Santoro et al., 2018), SAB (Ke et al., 2018) and Graph Attention Networks (Veličković et al., 2017; Battaglia et al., 2018). Our work contributes to this body of literature.

6. Experiments

In this section, we present a selection of experiments to empirically evaluate S2RMs and gauge their performance against strong baselines on two data domains. We proceed as follows: in Section 6.1 we introduce the baselines, followed by experimental results on a video prediction task (Section 6.2) and on the multi-agent world modeling task in the challenging Starcraft2 domain (Section 6.3). Additional results and supporting plots can be found in Appendix C.

6.1. Baseline Methods

To draw fair comparisons, we require a baseline architecture that is agnostic to the number of observations A, is invariant to the ordering of $\{(\mathbf{x}_t^a, \mathbf{O}_t^a)\}_a^A$ with respect to aand features a querying mechanism to extract a predicted observation $\mathbf{O}_{t'}^q$ at a given query location \mathbf{x}^q in a future time-step t' > t. Fortunately, it is possible to obtain a performant class of models fulfilling our requirements by extending prior work on Generative Query Networks or GQNs (Eslami et al., 2018). The resulting model has three components:

Encoder. At a given timestep t, the encoder E jointly maps the embedding $\mathbf{s}_t^a \in S$ of the position $\mathbf{x}_t^a \in \mathcal{X}$ and the corresponding observations \mathbf{O}_t^a to encodings \mathbf{e}_t^a , which are then summed over a to obtain an *aggregated representation*:

$$\mathbf{r}_t = \sum_{a=1}^{A} E(\mathbf{O}_t^a, \mathbf{s}_t^a) \tag{18}$$

The additive aggregation scheme we use is well known from prior work (Santoro et al., 2017; Eslami et al., 2018; Garnelo et al., 2018) and makes the model agnostic to Aand to permutations of $(\mathbf{x}_t^a, \mathbf{O}_t^a)$ over a. The encoder E is a seven-layer CNN with residual layers, and the positional embedding \mathbf{s}_t^a is injected after the second convolutional layer via concatenation with the feature tensor. The exact architectures can be found in Appendices B.1 and B.2.

RNN. The aggregated representation \mathbf{r}_t is used as an input to a RNN model F as following:

$$\mathbf{h}_{t+1}, \mathbf{c}_{t+1} = F(\mathbf{r}_t, \mathbf{h}_t, \mathbf{c}_t)$$
(19)

where \mathbf{h}_t and \mathbf{c}_t are hidden and memory states of the RNN F respectively. We experiment with various RNN models, including LSTMs (Hochreiter & Schmidhuber, 1997), RMCs (Santoro et al., 2018) and Recurrent Independent Mechanisms (RIMs) (Goyal et al., 2019).

As a sanity check, we also show results with a *Time Travelling Oracle* (TTO), which has access to r_{t+1} (but at time step *t*), and produces $\mathbf{h}_{t+1} = F_{TTO}(\mathbf{r}_{t+1})$ with a two layer MLP F_{TTO} . TTO therefore does not model the dynamics, but merely verifies that the additive aggregation scheme (Equation 18) and the querying mechanism (Equation 20) are sufficient for the task at hand.

Decoder. Given the embedding \mathbf{s}^q of the query position \mathbf{x}^q , the decoder D predicts the corresponding observation $\hat{\mathbf{O}}_{t+1}^q$:

$$\hat{\mathbf{O}}_{t+1}^q = D(\mathbf{h}_{t+1}, \mathbf{s}^q) \tag{20}$$

We parameterize D with a deconvolutional network with residual layers, and inject the positional embedding of the query s^q after a single convolutional layer by concatenating with the layer features (see Appendices B.1 and B.2).

The architecture described above therefore extends the framework of GQNs by predicting the forward dynamics of the aggregated representation; nevertheless, we do not consider it a novel contribution of this work.

6.2. Video Prediction from Crops

Task Description. We consider the problem of predicting the future frames of simulated videos of balls bouncing in a closed box (Miladinović et al., 2019), given only crops from the past video frames which are centered at known pixel positions. Using the notation introduced in Section 2: at every time step t, we sample A = 10 pixel positions $\{\mathbf{x}_t^a\}_{a=1}^{10}$ from the t-th full video frame \mathbf{o}_t of size 48×48 . Around the sampled central pixel positions \mathbf{x}_t^a , we extract 11×11 crops, which we use as the *local* observations \mathbf{O}_t^a . The task now is to predict 11×11 crops $\mathbf{O}_{t'}^q$ corresponding to query central pixel positions $\mathbf{x}_{t'}^q$ at a future time-step t' > t. Observe that at any given time-step t, the model has access to at most 52% of the global video frame assuming that the crops never overlap (which is rather unlikely).

Dataset. We train all models on a training dataset of 20K video sequences with 100 frames of 3 balls bouncing in an arena of size 48×48 . We also include an additional fixed ball in the center to make the task more challenging. We use another 1K video sequences of the same length and the same



Figure 3. Rollouts (OOD) with 5 bouncing balls, from top to bottom: ground-truth, S2GRU, RIMs, RMC, LSTM. Note that all models were trained on sequences with 3 bouncing balls, and the global state was reconstructed by stitching together 11×11 patches from the models (queried on a 4×4 grid). **Gist:** S2GRU succeeds at keeping track of all bouncing balls over long rollout horizons (25 frames).



Figure 4. Performance metrics on OOD one-step forward prediction task. **Gist:** S2GRU outperforms all RNN baselines OOD.



Figure 5. Visualization of the spatial locations each module is responsible for modeling (i.e. the *enclaves* \mathcal{X}_m , defined in Section 3). The central ball does not bounce, i.e. it is stationary in all sequences. **Gist:** the modules *focus attention* on challenging regions, e.g. the corners of the arena and the surface of the fixed ball.

number of balls as a held-out validation set. In addition, we also have 5 out-of-distribution (OOD) test sets with various number of bouncing balls (ranging from 1 to 6) and each containing 1K sequences of length 100.

Training. We train all models until the validation loss is saturated, and select the best of three runs (more details in Appendix B.4.3). During training, we automatically decay the learning rate by a factor of 2 if the validation loss does not decrease by at least 0.01% for five consecutive epochs.

Evaluation Criteria. After having trained on the training dataset with 3 bouncing balls, we evaluate the performance on all test datasets with 1 to 6 bouncing balls. In Figure 4, we report the balanced accuracy (i.e. arithmetic mean of recall and specificity) and F1-scores (i.e. harmonic mean of precision and recall) to account for class-imbalance. Additionally, in Figure 3, we qualitatively show reconstructions from 25 step rollouts on the OOD dataset with 5 balls (see Appendix C.1). Finally in Figure 5, we show for each module its corresponding enclave, which is the spatial region that it is *responsible* for modelling, i.e. for pixels at position **x**, we plot $\{Z(P(\mathbf{x}), \mathbf{p}^m)\}_{m=1}^{10}$ (cf. Section 4).

	UT-F1	FM-F1	NMSE	LL
(1s2z)				
LSTM	0.6267	0.8464	-0.0040	-0.0382
RMC	0.6839	0.8597	-0.0033	-0.0334
S2GRU	0.7488	0.8627	-0.0023	-0.0233
S2RMC	0.7317	0.8563	-0.0026	-0.0261
TTO	0.7518	0.8883	-0.0025	-0.0259
(5s3z)				
LSTM	0.4975	0.7123	-0.0134	-0.1251
RMC	0.5414	0.7486	-0.0132	-0.1167
S2GRU	0.5310	0.7058	-0.0119	-0.1108
S2RMC	0.5114	0.6945	-0.0124	-0.1205
TTO	0.6115	0.7872	-0.0107	-0.0940

Table 1. Performance metrics on OOD scenarios 1s2z and 5s3z (larger numbers are better). The metrics are: unit-type macro F1 score (UT-F1), friendly-marker F1 score (FM-F1), HECS Negative Mean Squared Error (NMSE) and Log Likelihood (LL).

Results. In Figure 4, we see that S2GRUs out-perform all non-oracle baselines OOD on the one-step forward prediction task and strike a good balance in regard to indistribution and OOD performance. Note, however, that the additive aggregation scheme and querying mechanism (Equations 18 and 20) can indeed generalize, as shown by the good performance of the oracle (TTO). Figure 5 shows how the modules *share responsibility* of modelling the entire spatial domain, whereas Figure 3 shows that S2GRUs can perform OOD rollouts over long horizons (25 frames) without losing track of balls. Additional results in Appendix C.1.

6.3. Multi-Agent World Modeling on Starcraft2

Task Description. In Section 2, we formulated the problem of modeling what we called the *world state* \mathfrak{o} of a dynamical system ϕ given *local observations* $\{(\mathbf{x}_t^a, \mathbf{O}_t^a)\}_{a=1}^A$ where $\mathbf{O}_t^a = \phi(t, \mathfrak{o})(\mathbf{x}_t^a)$. Under certain restrictions, this problem can be mapped to that of multi-agent world modeling from partial and local observations, allowing us to evaluate the proposed model in a rich and challenging setting. In particular, we consider environments that are (**a**) *spatial*, i.e.



Figure 6. Performance metrics (larger the better) as a function of the probability that an agent will not supply information to the world model but still query it. **Gist:** while all models lose performance as fewer agents share observations, we find S2RMs to be most robust.

all agents *a* in it have a well-defined and known location \mathbf{x}_t^a (at time *t*), (**b**) the agents' actions \mathbf{u}_t^a are local, in that their effects propagate away (from the agent) only at a finite speed, (**c**) the observations are local and centered around agents, in the sense that the agent only observes the events in its local vicinity, i.e., \mathbf{O}_t^a . Observe that we do not fix the number of agents in the environment, and allow for agents to dynamically enter or exit the environment. Now, the task is: given observations \mathbf{O}_t^a from a team of (cooperating) agents at position \mathbf{x}_t^a and their corresponding actions \mathbf{u}_t^a , predict the observation $\mathbf{O}_{t'}^q$ that would be made by an agent at time t' = t + 1 if it were at position \mathbf{x}_t^a and \mathbf{x}_{t+1}^a are no longer independent and depend on the agents' behaviour.

The SC2 Domain. Starcraft2 unit-micromanagement (Samvelyan et al., 2019) is a multi-agent reinforcement learning benchmark, wherein teams of heterogeneously typed units must defeat a team of opponents in melee and ranged combat. Each unit type has its own characteristics, e.g. maximum *health*, *shields*, weapon abilities (*cool-down*, *damage per second*, *splash damage*, etc), and strengths (vulnerabilities) against (towards) other unit types, making the world-modeling task all the more rich and challenging.

Dataset. The observations O_t^a and actions u_t^a are both multi-channel images represented in polar coordinates centered around the agent position \mathbf{x}_a^t . The field of view (FOV) of each agent is therefore a circle of fixed radius centered around it. The channels of the image correspond to (a) a binary indicator marking whether a position in FOV is occupied by a living friendly agent (friendly marker), (b) a categorical indicator marking the type of living units at a given position in FOV (unit-type marker), and (c) four channels marking the health, energy, weapon-cooldown and shields (HECS markers) of all agents in FOV. With a heuristic, we gather a total of 9K trajectories $(\{\mathbf{x}_t^a, \mathbf{O}_t^a, \mathbf{u}_t^a\}_{a=1}^A)_{t=1}^{100}$ spread over three training scenarios, corresponding to $1c3s5z^4$, 3s5z and 2s5z in Samvelyan et al. (2019). In addition, we also sample 1K trajectories (each) from two OOD scenarios 1s2z and 5s3z. Details in Appendix A.

Training. While adopting the training protocol detailed in Appendix B.4.3, we adapt the encoder and decoder architecture to match the state representation by including circular convolutions (cf. Appendix B.2). Now, recall that predicting the next state entails predicting images of binary friendly markers, categorical unit type markers and real valued HECS markers. Accordingly, the loss function is a sum of a binary cross-entropy term (on friendly markers), a categorical cross-entropy term (on unit-type markers) and a mean squared error term (on HECS markers).

Evaluation Criteria. After having trained all models on scenarios 1c3s5z, 3s5z and 2s5z, we test their robustness to dropped agents (Figure 6) and their performance on OOD scenarios (Table 1). We only show baselines that achieve similar or better validation scores than S2RMs, and report the F1 scores for binary friendly markers, multiclass (macro) F1 score for unit-type markers, negative mean squared error for HECS markers (tables in Appendix C.2).

Results. Figure 6 shows that S2RMs remain robust when fewer agents supply their observations to the world model, whereas Table 1 shows that S2GRU outperforms the baselines in the OOD scenario 1s2z but is matched by RMCs in 5s3z (see Appendix C.2 for details).

Conclusions and Outlook

We proposed Spatially Structured Recurrent Modules, a new class of models constructed to jointly leverage both spatial and modular structure in data, and explored its potential in the challenging problem setting of predicting the forward dynamics from partial observations at known spatial locations. In the tasks of video prediction from crops and multi-agent world modeling in the Starcraft2 domain, we found that it compares favorably against strong baselines in terms of out-of-distribution generalization and robustness to the number of available observations. Future work may focus on exploring efficient implementations using block-sparse methods (Gray et al., 2017), which could potentially unlock applications to significantly larger scale spatial problems encountered in domains such as humanitarian aid and climate change research (Rolnick et al., 2019).

⁴Here, the code 1c3s5z refers to a scenario where each team comprises 1 *colossus* (1c), 3 *stalkers* (3s), and 5 *zealots* (5z).

References

- Alet, F., Lozano-Pérez, T., and Kaelbling, L. P. Modular meta-learning. arXiv preprint arXiv:1806.10166, 2018.
- Battaglia, P. W., Hamrick, J. B., Bapst, V., Sanchez-Gonzalez, A., Zambaldi, V., Malinowski, M., Tacchetti, A., Raposo, D., Santoro, A., Faulkner, R., et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- Bengio, Y., Deleu, T., Rahaman, N., Ke, R., Lachapelle, S., Bilaniuk, O., Goyal, A., and Pal, C. A meta-transfer objective for learning to disentangle causal mechanisms. *arXiv preprint arXiv:1901.10912*, 2019.
- Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- Eslami, S. A., Rezende, D. J., Besse, F., Viola, F., Morcos, A. S., Garnelo, M., Ruderman, A., Rusu, A. A., Danihelka, I., Gregor, K., et al. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018.
- Fasshauer, G. E. Positive definite kernels: past, present and future. 2011.
- Garnelo, M., Schwarz, J., Rosenbaum, D., Viola, F., Rezende, D. J., Eslami, S., and Teh, Y. W. Neural processes. arXiv preprint arXiv:1807.01622, 2018.
- Goyal, A., Lamb, A., Hoffmann, J., Sodhani, S., Levine, S., Bengio, Y., and Schölkopf, B. Recurrent independent mechanisms. arXiv preprint arXiv:1909.10893, 2019.
- Graves, A., Wayne, G., and Danihelka, I. Neural turing machines. *arXiv preprint arXiv:1410.5401*, 2014.
- Graves, A., Wayne, G., Reynolds, M., Harley, T., Danihelka, I., Grabska-Barwińska, A., Colmenarejo, S. G., Grefenstette, E., Ramalho, T., Agapiou, J., et al. Hybrid computing using a neural network with dynamic external memory. *Nature*, 538(7626):471–476, 2016.
- Gray, S., Radford, A., and Kingma, D. P. Gpu kernels for block-sparse weights. arXiv preprint arXiv:1711.09224, 3, 2017.
- Henaff, M., Weston, J., Szlam, A., Bordes, A., and LeCun, Y. Tracking the world state with recurrent entity networks, 2016.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., and Hinton, G. E. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.
- Jaderberg, M., Simonyan, K., Zisserman, A., and Kavukcuoglu, K. Spatial transformer networks, 2015.
- Ke, N. R., GOYAL, A. G. A. P., Bilaniuk, O., Binas, J., Mozer, M. C., Pal, C., and Bengio, Y. Sparse attentive backtracking: Temporal credit assignment through reminding. In *Advances in neural information processing systems*, pp. 7640–7651, 2018.
- Ke, N. R., Bilaniuk, O., Goyal, A., Bauer, S., Larochelle, H., Pal, C., and Bengio, Y. Learning neural causal models from unknown interventions. *arXiv preprint arXiv:1910.01075*, 2019.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- Kumar, A., Eslami, S., Rezende, D. J., Garnelo, M., Viola, F., Lockhart, E., and Shanahan, M. Consistent generative query networks. arXiv preprint arXiv:1807.02033, 2018.
- Li, S., Li, W., Cook, C., Zhu, C., and Gao, Y. Independently recurrent neural network (indrnn): Building a longer and deeper rnn. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5457–5466, 2018.
- Miladinović, ., Gondal, M. W., Schölkopf, B., Buhmann, J. M., and Bauer, S. Disentangled state space representations. arXiv preprint arXiv:1906.03255, 2019.
- Parascandolo, G., Kilbertus, N., Rojas-Carulla, M., and Schölkopf, B. Learning independent causal mechanisms. arXiv preprint arXiv:1712.00961, 2017.
- Parmar, N., Vaswani, A., Uszkoreit, J., Łukasz Kaiser, Shazeer, N., Ku, A., and Tran, D. Image transformer, 2018.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. Pytorch: An imperative style, high-performance deep learning library. In Wallach, H., Larochelle, H., Beygelzimer, A., dAlché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019.
- Rolnick, D., Donti, P. L., Kaack, L. H., Kochanski, K., Lacoste, A., Sankaran, K., Ross, A. S., Milojevic-Dupont, N., Jaques, N., Waldman-Brown, A., et al. Tackling climate change with machine learning. *arXiv preprint arXiv:1906.05433*, 2019.

- Samvelyan, M., Rashid, T., de Witt, C. S., Farquhar, G., Nardelli, N., Rudner, T. G. J., Hung, C.-M., Torr, P. H. S., Foerster, J., and Whiteson, S. The starcraft multi-agent challenge, 2019.
- Santoro, A., Raposo, D., Barrett, D. G., Malinowski, M., Pascanu, R., Battaglia, P., and Lillicrap, T. A simple neural network module for relational reasoning. In *Advances in neural information processing systems*, pp. 4967–4976, 2017.
- Santoro, A., Faulkner, R., Raposo, D., Rae, J., Chrzanowski, M., Weber, T., Wierstra, D., Vinyals, O., Pascanu, R., and Lillicrap, T. Relational recurrent neural networks. In Advances in Neural Information Processing Systems, pp. 7299–7310, 2018.
- Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., and Dean, J. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017.
- Singh, G., Yoon, J., Son, Y., and Ahn, S. Sequential neural processes, 2019.
- Sun, C., Karlsson, P., Wu, J., Tenenbaum, J. B., and Murphy, K. Predicting the present and future states of multiagent systems from partially-observed visual data. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum? id=rlxdH3CcKX.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in neural information* processing systems, pp. 5998–6008, 2017.
- Veličković, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., and Bengio, Y. Graph attention networks. arXiv preprint arXiv:1710.10903, 2017.
- Vinyals, O., Ewalds, T., Bartunov, S., Georgiev, P., Vezhnevets, A. S., Yeo, M., Makhzani, A., Küttler, H., Agapiou, J., Schrittwieser, J., et al. Starcraft ii: A new challenge for reinforcement learning. *arXiv preprint arXiv:1708.04782*, 2017.
- Wang, X., Girshick, R., Gupta, A., and He, K. Non-local neural networks, 2017.
- Zhang, H., Goodfellow, I., Metaxas, D., and Odena, A. Selfattention generative adversarial networks. arXiv preprint arXiv:1805.08318, 2018.

A. Starcraft2

The Starcraft2 Environment we use is a modified version of the SMAC-Env proposed in Samvelyan et al. (2019) and built on PySC2 wrapper around Blizzard SC2 API (Vinyals et al., 2017). Starcraft2 is a real-time-strategy (RTS) game where players are tasked with manufacturing and controlling armies of *units* (airborne or land-based) to defeat the opponent's army (where the opponent can be an AI or another human). The players must choose their *alien race*⁵ before starting the game; available options are *Protoss*, *Terran* and *Zerg*. All unit types (of all races) have their strengths and weaknesses against other unit types, be it in terms of maximum health, shields (Protoss), energy (Terran), DPS (damage per second, related to weapon cooldown), splash damage, or manufacturing costs (measured in *minerals* and *vespene gas*, which must be mined).

The key engineering contribution of Samvelyan et al. (2019) is to repurpose the RTS game as a multi-agent environment, where the individual units in the army become individual agents⁶. The result is a rich and challenging environment where heterogeneous teams of agents must defeat each other in melee and ranged combat. The composition of teams vary between *scenarios*, of which Samvelyan et al. (2019) provide a selection. Further, new scenarios can be easily created with the SC2MapEditor, which allows for practically endlessly many possibilities.

We build on Samvelyan et al. (2019) by modifying their environment to better expose the transfer and out-of-distribution aspects of the domain by (a) standardizing the state and action space across a large class of scenarios and (b) standardizing the unit stats to better reflect the game-defined notion of hit-points.

A.1. Standardized State Space for All Scenarios

In the environment provided by Samvelyan et al. (2019), the dimensionality of the vector state space varies with the number of friendly and enemy agents, which in turn varies with the scenario. While this is not an issue in the typical use case of training MARL agents in a fixed scenario, it is not convenient for designing models that seamlessly handle multiple scenarios. In the following, we propose an alternate state representation that preserves the spatial structure and is consistent across multiple scenarios.

Instead of representing the state of an agent a with a vector of variable dimension, we represent it with a multi-channel *polar image* \mathfrak{I}^a of shape $C \times I \times J$, where C is the number of channels and (I, J) is the image size. Given the *radial* and *angular* resolutions ρ and φ (respectively), the pixel coordinate i = 0, ..., I - 1, j = 0, ..., J - 1 corresponds to coordinates $(i \cdot \rho, j \cdot \varphi)$ with respect to a polar coordinate system centered on the agent a, where the positive x-axis (j = 0) points towards the east. Further, the field of view (FOV) of an agent is characterized by a circle of radius $I \cdot \rho$ centered on the agent at 2D game-coordinates $\mathbf{x}^a = (x_1^a, x_2^a)$, to which the Starcraft2 API (Vinyals et al., 2017) provides raw access.

The polar image \mathfrak{I}^a therefore provides an agent-centric view of the environment, where pixel coordinates i, j in \mathfrak{I}^a can be mapped to global game coordinates $\mathbf{x} = (x_1, x_2)$ in FOV via:

$$x_1 = i \cdot \rho \cos\left[j \cdot \varphi\right] + x_1^a \tag{21}$$

$$x_2 = i \cdot \rho \sin\left[j \cdot \varphi\right] + x_2^a \tag{22}$$

In what follows, we denote this transformation with T_a , as in $T_a(i, j) = (x_1, x_2)$.

Now, the channels in the polar image can encode various aspects of the observation; in our case: friendly markers (one channel), unit-type markers (nine channels, one-hot), health-energy-cooldown-shields (HECS, four channels) and terrain height (one channel). As an example, let us consider the friendly markers, which is a binary indicator marking units that are friendly. If we have an agent at game position (x_1, x_2) that is friendly to agent a, then we would expect the pixel coordinate $(i, j) = T_a^{-1}(x_1, x_2)$ of the corresponding channel in the polar image \Im^a to be 1, but 0 otherwise. Likewise, the value of \Im at the channels corresponding to HECS at pixel position *i*, *j* gives the HECS of the corresponding unit⁷ at $T_a(i, j)$. This representation has the following advantages: (a) it does not depend on the number of units in the field of view, (b) it exposes the spatial structure in the arrangement of units which can naturally processed by convolutional neural networks (e.g. with circular convolutions).

Nevertheless, it has the disadvantage that the positions are *quantized* to pixels, but the euclidean distance between the locations represented by pixels (i, j) and (i, j+1) increases with increasing *i*. Consequently, this representation may not remain suitable for larger FOVs.

Further, this representation is also appropriate for the action space. Given an agent, we represent the one-hot categorical actions of all friendly agents in FOV as a multi-channel polar image. In this representation, the pixel position i, j gives the action taken by an agent at at position $T_a(i, j)$. Unfriendly agents get assigned an "unknown action", whereas positions not occupied by a living agent are assigned a "no-op" action.

⁵Please note that this is a game-specific notion.

⁶Note that this is rather unconventional, since each player usually controls entire armies and must switch between macro- and micro-management of units or unit-groups.

⁷If health drops to zero, the unit is considered dead and the representation does not differentiate between dead and absent units.



(a) 1s2z (1 Stalker and 2 Zealots per team).



(c) 2s3z (2 Stalkers and 3 Zealots per team).



(b) 5s3z (5 Stalkers and 3 Zealots per team).



(d) 3s5z (3 Stalkers and 5 Zealots per team).



(e) 1c3s5z (1 Colossus, 3 Stalkers and 5 Zealots per team).

Figure 7. Human readable illustrations of the Starcraft2 (SMAC) scenarios we consider in this work. Figures 7a and 7b show the OOD scenarios, whereas Figures 7c, 7d and 7e show the training scenarios (provided by Samvelyan et al. (2019)).

A.2. Standardized Unit Stats

At any given point in time, an active unit in Starcraft2 has certain *stats*, e.g. its health, energy (Terran), shields (Protoss) and weapon-cooldown (for armed units). A large and expensive unit-type like the Colossus has more max-health (*hit-points*) than smaller units like Stalkers and Marines⁸. Likewise, unit-types differ in the rate at which they deal damage (measured in damage-per-second or DPS, excluding splash damage), which in turn depends on the cooldown duration of the active weapon.

Now, the environment provided by Samvelyan et al. (2019) normalizes the stats by their respective maximum value, resulting in values between 0 and 1. However, given that different units may have different normalization, the stats are rendered incomparable between unit types (without additionally accounting the unit-type). We address this by stan-

dardizing stats (instead of normalizing) by dividing them by a fixed value. In this scheme, the stats are scaled uniformly across all unit-types, enabling models to directly rely on them instead of having to account for the respective unit-types.

B. Hyperparameters and Architectures

B.1. Encoder and Decoder for Bouncing Balls

The architectures of image encoder and decoder was fixed for all models after initial experimentation. We converged to the following architectures.

B.1.1. S2RMs

The encoder (decoder) is a (de)convolutional network with residual connections (Figure 9).

⁸These stats may change with game-versions, and are catalogued here: https://liquipedia.net/starcraft2/ Units_(StarCraft).





(a) Encoder.

(b) Decoder.

Observation

(IC)

ŧ

Convolution (128 → IC) size 3, padding 1, stride 1

Deconvolution

Skip-Add

f

Convolution

 $(128 \rightarrow 128)$ nel size 3, padding 1, str

Deconvolution

(128 → 128) size 3, padding 0,

Skip-Add

t

Convolution

Deconvolution

(128 → 128) kernel size 3, padding 0, stride 1)

Ť

Deconvolution

Representation (RC)

(RC \rightarrow 128) size 5, padding 0, stride 1)

(128 → 128) size 3, padding 1, stride 1)

 $(128 \rightarrow 128)$ size 3, padding 0, stride 1)

Figure 8. Baseline encoder and decoder architectures for the Bouncing Ball task.

Figure 9. S2RM encoder and decoder architectures for the Bouncing Ball task.

B.1.2. BASELINES

Like in the case of S2RMs, the encoder (decoder) is a (de)convolutional network with residual connections (Figure 8), but with the positional embeddings injected after the second convolutional layer. This is loosely inspired by the encoders used in Eslami et al. (2018).

B.2. Encoder and Decoder for Starcraft2

B.2.1. S2RMs

Recall from Appendix A that the states are polar images. We therefore use *polar convolutions*, which entails zero-padding the input image along the first (*radial*) dimension but circular padding along the second (*angular*) dimension. The encoder and decoder architectures can be found in Figure 11.

B.2.2. BASELINES

Like for S2RMs, we use polar convolutions while injecting the positional embeddings further downstream in the network. The corresponding encoder and decoder architectures are illustrated in Figure 10.

B.3. Spatially Structured Relational Memory Cores (S2RMCs)

Embedding Relational Memory Cores (Santoro et al., 2018) naïvely in the S2RM architecture did not result in a working model. We therefore had to adapt it by first projecting the memory matrix (M in Santoro et al. (2018)) of the m-th RMC to a message \mathbf{h}_t^m . This message is then processed by the intercell attention to obtain $\bar{\mathbf{h}}_t^m$, which is finally concatenated with the memory matrix and current input before applying the attention mechanism (i.e. in Equation 2 of Santoro et al. (2018), we replace [M; x] with $[M; x, \bar{\mathbf{h}}_t^m]$).

B.4. Hyperparameters

B.4.1. BOUNCING BALL MODELS

The hyperparameters we used can be found in Table 2. Further, note that in Equation 6, we pass the gradients through the constant region of the kernel as if the kernel had not been truncated.

B.4.2. STARCRAFT2 MODELS

The hyperparameters we used can be found in Table 3. Note that we only report models that attained a validation loss similar to or better than S2RMs.



(a) Encoder.

(b) Decoder.

Figure 10. Baseline encoder and decoder architectures for the Starcraft2 task.



(a) Encoder.

(b) Decoder.

Figure 11. S2RM encoder and decoder architectures for the Starcraft2 task.

Model	
Hyperparameter	Value
S2GRU	
Number of modules (M)	10
GRU: hidden size per module	128
Module embedding size (d)	16
Kernel bandwidth (ϵ)	1
Kernel truncation (τ)	0.6
shape $\Theta^{(Q/K)}$	(128, 2, 016)
shape $\Theta^{(V)}$	(128, 2, 128)
shape $\Phi^{(Q/K)}$	(128, 4, 016)
shape $\Phi^{(V)}$	(128, 4, 128)
RMC (Santoro et al., 2018)	
Number of attention heads	4
Size of attention head	128
Number of memory slots	1
Key size	128
LSTM (Hochreiter & Schmidhuber, 1997)	
Hidden size	512
RIMs (Goyal et al., 2019)	
Number of RIMs (k_T)	6
Update Top-k (k_A)	5
Hidden size (h_{size})	510
Input key size	32
Input value size	400
ТТО	
MLP hidden size	512

Table 2. Hyperparameters used for various models on the Bouncing Ball task. Hyperparameters not listed here were left at their respective default values.

Model						
Hyperparameter	Value					
S2GRUs						
Number of modules (M)	10					
GRU: hidden size per module	128					
Module embedding size (d)	8					
Kernel bandwidth (ϵ)	1					
Kernel truncation (τ)	0.5					
shape $\Theta^{(Q/K)}$	(128, 2, 016)					
shape $\Theta^{(V)}$	(128, 2, 128)					
shape $\Phi^{(Q/K)}$	(128, 4, 016)					
shape $\Phi^{(V)}$	(128, 4, 128)					
S2RMC						
Number of modules (M)	10					
RMC: number of attention heads	4					
RMC: size of attention head	64					
RMC: number of memory slots	4					
RMC: key size	64					
Module embedding size (d)	8					
Kernel bandwidth (ϵ)	1					
Kernel truncation (τ)	0.5					
shape $\Theta^{(Q/K)}$	(128, 2, 016)					
shape $\Theta^{(V)}$	(128, 2, 128)					
shape $\Phi^{(Q/K)}$	(128, 4, 016)					
shape $\Phi^{(V)}$	(128, 4, 128)					
RMC (Santoro et al., 2018)	i					
Number of attention heads	4					
Size of attention head	128					
Number of memory slots	1					
Key size	16					
LSTM (Hochreiter & Schmidhuber 1997)						
Hidden size	2048					
тто						
MLP hidden size	512					

Table 3. Hyperparameters used for various models on the Starcraft2 task. Hyperparameters not listed here were left at their respective default values.

B.4.3. TRAINING

All models were trained using Adam (Kingma & Ba, 2014) with an initial learning rate 0.0003⁹. We use Pytorch's (Paszke et al., 2019) ReduceLROnPlateau learning rate scheduler to decay the learning rate by a factor of 2 if the validation loss does not improve by at least 0.01% over the span of 5 epochs. We initially train all models for 100 epochs, select the best of three successful runs, fine-tune it for another 100 epochs, and finally select the checkpoint with the lowest validation loss (i.e. we early stop). We train all models with batch-size 8 (Starcraft2) or 32 (Bouncing Balls) on a single V100-32GB GPU (each).

C. Additional Results

C.1. Bouncing Balls

C.1.1. ROLLOUTS

To obtain the rollouts in Figure 3, we adopt the following strategy. For the first 20 *prompt-steps*, we present all models with exactly the same 11×11 crops around randomly sampled pixel positions for 20 time-steps. For the next 25 steps, all models are queried at random pixel positions¹⁰, and the resulting predictions (on crops) are thresholded at 0.5 and fed back in to the model for the next step (at known pixel positions from the previous step).

Also at every time-step, the models are queried for their predictions on 16 pixel locations placed on a 4×4 grid. The resulting predictions are stitched together and shown in Figures 12, 13, 14, 15, 3 and 16.

C.1.2. ROBUSTNESS TO DROPPED VIEWS

In this section, we evaluate the robustness of all models to dropped crops on in-distribution and OOD data. We measure the performance metrics on one-step forward prediction task on all datasets (with 1-6 balls), albeit by dropping a given fraction of the available input observations.

Figure 17 and 18 visualize the performance of all evaluated models. We find that S2GRU maintains performance on OOD data even with fewer views (or crops) than it was trained on. Interestingly, we find that the time-travelling or-acle (TTO), while robust OOD, is adversely affected by the number of available views. This could be because unlike the other models, it cannot leverage the temporal information to compensate for the missing observations.

⁹https://twitter.com/karpathy/status/ 801621764144971776?s=20

¹⁰These random pixel positions are the same for all models, but change between time-steps



Figure 12. Rollouts (OOD) with 1 bouncing ball, from top to bottom: ground-truth, S2GRU, RIMs, RMC, LSTM. Note that all models were trained on sequences with 3 bouncing balls, and the global state was reconstructed by stitching together 11×11 patches from the models (queried on a 4×4 grid).



Figure 13. Rollouts (OOD) with 2 bouncing balls, from top to bottom: ground-truth, S2GRU, RIMs, RMC, LSTM. Note that all models were trained on sequences with 3 bouncing balls, and the global state was reconstructed by stitching together 11×11 patches from the models (queried on a 4×4 grid).



Figure 14. Rollouts (ID) with 3 bouncing balls, from top to bottom: ground-truth, S2GRU, RIMs, RMC, LSTM. Note that all models were trained on sequences with 3 bouncing balls, and the global state was reconstructed by stitching together 11×11 patches from the models (queried on a 4×4 grid).



Figure 15. Rollouts (OOD) with 4 bouncing balls, from top to bottom: ground-truth, S2GRU, RIMs, RMC, LSTM. Note that all models were trained on sequences with 3 bouncing balls, and the global state was reconstructed by stitching together 11×11 patches from the models (queried on a 4×4 grid).



Figure 16. Rollouts (OOD) with 6 bouncing balls, from top to bottom: ground-truth, S2GRU, RIMs, RMC, LSTM. Note that all models were trained on sequences with 3 bouncing balls, and the global state was reconstructed by stitching together 11×11 patches from the models (queried on a 4×4 grid).

C.2. Starcraft2

C.2.1. TABULAR RESULTS

The results used to plot Figure 6 can be found tabulated in Tables 4, 5, 6 and 7.





Figure 17. Balanced accuracy (arithmetic mean of recall and specificity) achieved by all evaluated models for one-step forward prediction task with various number of balls and fractions of available views. All models were trained on video sequences with 3 balls and a constant number of crops / views (10 views, corresponding to the right-most columns labelled 1.0). The color map is consistent across all plots.



Figure 18. F1-Score (harmonic mean of precision and recall) achieved by all evaluated models for one-step forward prediction task with various number of balls and fractions of available views. All models were trained on video sequences with 3 balls and a constant number of crops / views (10 views, corresponding to the right-most columns labelled 1.0). The color map is consistent across all plots.

Model	LSTM	RMC	S2GRU	S2RMC	TTO
% of Active Agents					
20%	0.570565	0.586541	0.642292	0.637618	0.550806
30%	0.599391	0.606114	0.660127	0.653950	0.578965
40%	0.630606	0.640435	0.678752	0.671476	0.605867
50%	0.638374	0.657472	0.688528	0.685988	0.627444
60%	0.681040	0.704552	0.713851	0.708786	0.671961
70%	0.709861	0.737436	0.734256	0.727980	0.723238
80%	0.721041	0.748138	0.738611	0.732114	0.740936
90%	0.750449	0.778647	0.755476	0.747613	0.786931
100%	0.765592	0.795049	0.763126	0.754637	0.813504

Table 4. Friendly marker F1 scores on the validation set of the training distribution. Larger numbers are better.

Model	LSTM	RMC	S2GRU	S2RMC	TTO
70 OI ACTIVE Agents					
20%	0.323482	0.326685	0.435318	0.377538	0.297192
30%	0.345108	0.350621	0.491934	0.433945	0.323736
40%	0.373612	0.387733	0.540163	0.485278	0.350733
50%	0.385550	0.406048	0.552589	0.510371	0.371088
60%	0.430793	0.481986	0.599470	0.566149	0.435724
70%	0.497964	0.590214	0.635928	0.606039	0.539652
80%	0.579952	0.649277	0.650682	0.623040	0.617973
90%	0.657643	0.694158	0.675294	0.655581	0.699008
100%	0.677952	0.715929	0.689669	0.672186	0.737745

Table 5. Unit-type marker (macro averaged) F1 scores on the validation set of the training distribution. Larger numbers are better.

Model	LSTM	RMC	S2GRU	S2RMC	TTO
% of Active Agents					
20%	-0.014035	-0.013569	-0.011491	-0.011921	-0.014174
30%	-0.013355	-0.012747	-0.010631	-0.011101	-0.013539
40%	-0.012567	-0.011808	-0.009906	-0.010367	-0.012916
50%	-0.012220	-0.011305	-0.009637	-0.009887	-0.012481
60%	-0.010888	-0.009799	-0.008751	-0.009034	-0.010929
70%	-0.009738	-0.008469	-0.008068	-0.008359	-0.009184
80%	-0.009081	-0.008027	-0.007873	-0.008162	-0.008466
90%	-0.007970	-0.007180	-0.007347	-0.007615	-0.007038
100%	-0.007638	-0.006823	-0.007103	-0.007362	-0.006401

Table 6. HECS Negative MSE on the validation set of the training distribution. Larger numbers are better.

Model	LSTM	RMC	S2GRU	S2RMC	TTO
% of Active Agents					
20%	-0.303051	-0.300892	-0.141989	-0.146553	-0.434099
30%	-0.258878	-0.256878	-0.126037	-0.137025	-0.347899
40%	-0.216924	-0.211048	-0.113317	-0.126882	-0.276596
50%	-0.206582	-0.191644	-0.108643	-0.113293	-0.245019
60%	-0.158170	-0.142643	-0.094380	-0.099989	-0.175233
70%	-0.126446	-0.109634	-0.084129	-0.089527	-0.120694
80%	-0.111735	-0.099229	-0.081624	-0.086723	-0.104135
90%	-0.082463	-0.074518	-0.073439	-0.078197	-0.071243
100%	-0.070488	-0.063183	-0.069856	-0.074041	-0.057276

Table 7. Log Likelihood (negative loss) on the validation set of the training distribution. Larger numbers are better.