
Reasoning-Modulated Representations

Anonymous Author(s)

Anonymous Affiliation

Anonymous Email

Abstract

Neural networks leverage robust internal representations in order to generalise. Learning them is difficult, and often requires a large training set that covers the data distribution densely. We study a common setting where our task is not purely opaque. Indeed, very often we may have access to information about the underlying system (e.g. that observations must obey certain laws of physics) that any “tabula rasa” neural network would need to re-learn from scratch, penalising performance. We incorporate this information into a pre-trained reasoning module, and investigate its role in shaping the discovered representations in diverse self-supervised learning settings from pixels. Our approach paves the way for a new class of representation learning, grounded in algorithmic priors.

1 Introduction

Neural networks are able to learn policies in environments without access to their specifics [1], generate large quantities of text [2], or automatically fold proteins to high accuracy [3]. However, such “tabula rasa” approaches hinge on having access to substantial quantities of data, from which robust representations can be learned. Without a large training set that spans the data distribution, representation learning is difficult [4–6].

Here, we study ways to construct neural networks with representations that are robust, while retaining a data-driven approach. We rely on a simple observation: very often, we have some (partial) knowledge of the underlying dynamics of the data, which could help make stronger predictions from fewer observations. This knowledge, however, usually requires us to be mindful of *abstract* properties of the data—and such properties cannot always be robustly extracted from *natural* observations.

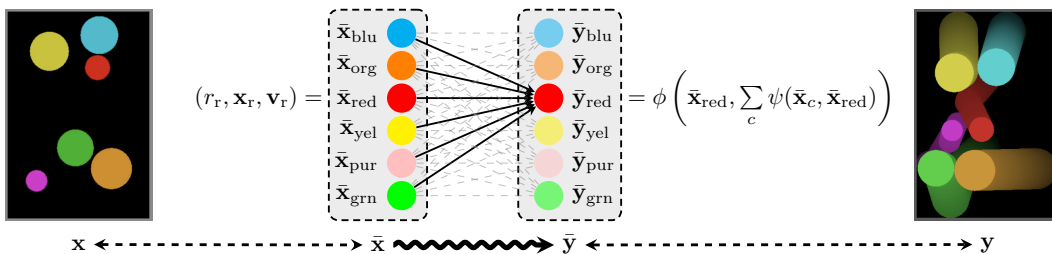


Figure 1: Bouncing balls example (re-printed, with permission, from Battaglia et al. [7]). Natural inputs, x , correspond to pixel observations. Predicting future observations (natural outputs, y), can be simplified as follows: if we are able to extract a set of *abstract* inputs, \bar{x} , (e.g. the radius, position and velocity for each ball), the movements in this abstract space must obey the laws of physics.

Motivation. Consider the task of predicting the future state of a system of n bouncing balls, from a pixel input x (Figure 1). Reliably estimating future pixel observations, y , is a challenging reconstruction task. However, the generative properties of this system are simple. Assuming knowledge of simple abstract inputs (radius, r_c , position, \mathbf{x}_c , and velocity, \mathbf{v}_c) for every ball, \bar{x}_c , the future movements in this abstract space are the result of applying the laws of physics to these

low-dimensional quantities. Hence, future abstract states, \bar{y} , can be computed via a simple algorithm that aggregates pair-wise forces between objects.

While this gives us a potentially simpler path from pixel inputs to pixel outputs, via abstract inputs to abstract outputs ($\mathbf{x} \rightarrow \bar{\mathbf{x}} \rightsquigarrow \bar{\mathbf{y}} \rightarrow \mathbf{y}$), it still places potentially unrealistic demands on our task setup, every step of the way:

$\mathbf{x} \rightarrow \bar{\mathbf{x}}$: Necessitates either upfront knowledge of how to abstract away $\bar{\mathbf{x}}$ from \mathbf{x} , or a massive dataset of *paired* $(\mathbf{x}, \bar{\mathbf{x}})$ to learn such a mapping from;

$\bar{\mathbf{x}} \rightsquigarrow \bar{\mathbf{y}}$: Implies that the algorithm *perfectly* simulates all aspects of the output. In reality, an algorithm may often only give partial context about \mathbf{y} . Further, algorithms often assume that \mathbf{x} is provided without error, exposing an algorithmic bottleneck [8]: if $\bar{\mathbf{x}}$ is incorrectly predicted, this will negatively compound in $\bar{\mathbf{y}}$, hence \mathbf{y} ;

$\bar{\mathbf{y}} \rightarrow \mathbf{y}$: Necessitates a renderer that generates \mathbf{y} from $\bar{\mathbf{y}}$, or a dataset of *paired* $(\bar{\mathbf{y}}, \mathbf{y})$ to learn it.

We will assume a general setting where *none* of the above constraints hold: we know that the mapping $\bar{\mathbf{x}} \rightsquigarrow \bar{\mathbf{y}}$ is likely of use to our predictor, but we do not assume a trivial mapping or a paired dataset which would allow us to convert directly from \mathbf{x} to $\bar{\mathbf{x}}$ or from $\bar{\mathbf{y}}$ to \mathbf{y} . Our only remaining assumption is that the algorithm $\bar{\mathbf{x}} \rightsquigarrow \bar{\mathbf{y}}$ can be efficiently computed, allowing us to generate massive quantities of paired abstract input-output pairs, $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$.

Present work. In this setting, we propose Reasoning-Modulated Representations (RMR), an approach that first learns a latent-space *processor* of abstract data; i.e. a mapping $\bar{\mathbf{x}} \xrightarrow{f} \mathbf{z} \xrightarrow{P} \mathbf{z}' \xrightarrow{g} \bar{\mathbf{y}}$, where $\mathbf{z} \in \mathbb{R}^k$ are high-dimensional latent vectors. f and g are an encoder and decoder, designed to take abstract representations to and from this latent space, and P is a processor network which simulates the algorithm $\bar{\mathbf{x}} \rightsquigarrow \bar{\mathbf{y}}$ in the latent space.

We then observe, in the spirit of neural algorithmic reasoning [9], that such a processor network can be used as a drop-in differentiable component for *any* task where the $\bar{\mathbf{x}} \rightsquigarrow \bar{\mathbf{y}}$ kind of reasoning may be applicable. Hence, we then learn a pipeline $\mathbf{x} \xrightarrow{\tilde{f}} \mathbf{z} \xrightarrow{P} \mathbf{z}' \xrightarrow{\tilde{g}} \mathbf{y}$, which *modulates* the representations \mathbf{z} obtained from \mathbf{x} , forcing them to pass through the pre-trained processor network. By doing so, we have ameliorated the original requirement for a massive *natural* dataset of (\mathbf{x}, \mathbf{y}) pairs. Instead, we inject knowledge from a massive *abstract* dataset of $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ pairs, directly through the pre-trained parameters of P . This has the potential to relieve the pressure on encoders and decoders \tilde{f} and \tilde{g} , which we experimentally validate on several challenging representation learning domains.

Our contributions can be summarised as follows:

- Verifying and extending prior work, we show that meaningful latent-space models can be learned from massive abstract datasets, on physics simulations and Atari 2600 games;
- We then show that these latent-space models can be used as differentiable components within neural pipelines that process raw observations. In doing so, we recover a neural network pipeline that relies solely on the existence of massive abstract datasets (which can often be automatically generated).
- Finally, we demonstrate early signs of processor *reusability*: latent-space abstract models can be used in tasks which do not even directly align with their environment, so long as these tasks can benefit from their underlying reasoning procedure.

2 Related Work

Neural algorithmic reasoning. RMR relies on being able to construct robust latent-space models, akin to world models [10], that imitate abstract reasoning procedures. This makes it well aligned with neural algorithmic reasoning [11], which is concerned with constructing neuralised versions of classical algorithms (typically by learning to execute them in a manner that extrapolates). Leveraging the ideas of algorithmic alignment [12], several known algorithmic primitives have already been successfully neuralised. This includes *iterative computation* [13, 14], *linearithmic algorithms* [15], and *data structures* [16, 17]. Further, the XLVIN model [8] demonstrates how such primitives can be re-used for data-efficient planning, paving the way for a blueprint [9] that we leverage in RMR as well. Our work also relates to transferring algorithmic reasoning knowledge [18], where no evidence

78 of positive transfer was found, in case of cross-algorithm transfer, whereas we did in the case of
79 transfer among aligned algorithms (same underlying algorithm from abstract to raw inputs.)

80 **Physical simulation with neural networks.** Our work also has contact points with prior art in
81 using (graph) neural networks for *physics simulations*. In fact, there is a tight coupling between
82 algorithmic computation and simulations, as the latter are typically realised using the former. Within
83 this space, abstract GNN models of physics have been proposed by works such as interaction
84 networks [7] and NPE [19], and extended to pixel-based inputs by visual interaction networks [20].
85 The generalisation power of these models has increased drastically in recent years, with effective
86 models of systems of particles [21] as well as meshes [22] being proposed. Excitingly, it has also
87 been demonstrated that rudimentary laws of physics can occasionally be recovered from the update
88 rules of these GNNs [23], and that they can be used to uncover new physical knowledge [24].

89 Recent work has also explored placing additional constraints on learning-based physical simulators,
90 for example by using Hamiltonian ODE integrators in conjunction with GNN models [25], or by
91 coupling a (non-neural) differentiable physics engine directly to visual inputs and optimizing its
92 parameters via backprop [26, 27].

93 **Object-centric and modular models for dynamic environments.** RMR with factored latents
94 can be viewed as a form of *object-centric neural network*, in which visual objects in an image or
95 video are represented as separate latent variables in the model and their temporal dynamics and
96 pairwise interactions are modeled via GNNs or self-attention mechanisms. There is a rich literature
97 on discovering objects and learning their dynamics from raw visual data without supervision, with
98 object-centric models such as R-NEM [28], SQAIR [29], OP3 [30], SCALOR [31], G-SWM [32].
99 Recent work has explored using contrastive losses [33] in this context or other losses directly in latent
100 space [34]. Related approaches discover and use keypoints [35] to describe objects and even discover
101 causal relations from visual input [36] using neural relational inference [37] in conjunction with a
102 keypoint discovery method. A related line of works integrate attention-mechanisms in modular and
103 object-centric models to interface latent variables with visual input, including models such as RMC
104 [38], RIM [39], Slot Attention [40], SCOFF [41], and NPS [42].

105 3 RMR architecture

106 Having provided a high-level overview of RMR and surveyed the relevant related work, we proceed
107 to carefully detail the blueprint of RMR’s various components. This will allow us to ground any
108 subsequent RMR experiments on diverse domains directly in our blueprint. Throughout this section,
109 it will be useful to refer to Figure 2 which presents a visual overview of this section.

110 **Preliminaries.** We assume a set of *natural inputs*, \mathcal{X} , and a set of *natural outputs*, \mathcal{Y} . These sets
111 represent the possible inputs and outputs of a target function, $\Phi : \mathcal{X} \rightarrow \mathcal{Y}$, which we would like to
112 learn based on a (potentially small) dataset of input-output pairs, (\mathbf{x}, \mathbf{y}) , where $\mathbf{y} = \Phi(\mathbf{x})$.

113 We further assume that the inner workings of Φ can be related to an *algorithm*, $A : \bar{\mathcal{X}} \rightarrow \bar{\mathcal{Y}}$. The
114 algorithm operates over a set of *abstract inputs*, $\bar{\mathcal{X}}$, and produces outputs from an *abstract output*
115 *set* $\bar{\mathcal{Y}}$. Typically, it will be the case that $\dim \bar{\mathcal{X}} \ll \dim \mathcal{X}$; that is, abstract inputs are assumed
116 substantially lower-dimensional than natural inputs. We do not assume existence of any aligned input
117 pairs $(\mathbf{x}, \bar{\mathbf{x}})$, and we do not assume that A perfectly explains the computations of Φ . What we do
118 assume is that A is either known or can be trivially computed, giving rise to a massive dataset of
119 abstract input-output pairs, $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$, where $\bar{\mathbf{y}} = A(\bar{\mathbf{x}})$.

120 Lastly, we assume a *latent space*, \mathcal{Z} , and that we can construct neural network components to both
121 encode and decode from it. Typically, \mathcal{Z} will be a real-valued vector space ($\mathcal{Z} = \mathbb{R}^k$) which is
122 high-dimensional; that is, $k > \dim \bar{\mathcal{X}}$. This ensures that any neural networks operating over \mathcal{Z} are
123 not vulnerable to bottleneck effects.

124 Note that either the natural or abstract input set may be *factorised*, e.g., into objects; in this case, we
125 can accordingly factorise the latent space, enforcing $\mathcal{Z} = \mathbb{R}^{n \times k}$, where n is the assumed maximal
126 number of objects (typically a hyperparameter of the models if not known upfront).

127 **Abstract pipeline.** RMR training proceeds by first learning a model of the algorithm A , which is
128 bound to pass through a latent-space representation. That is, we learn a neural network approximator

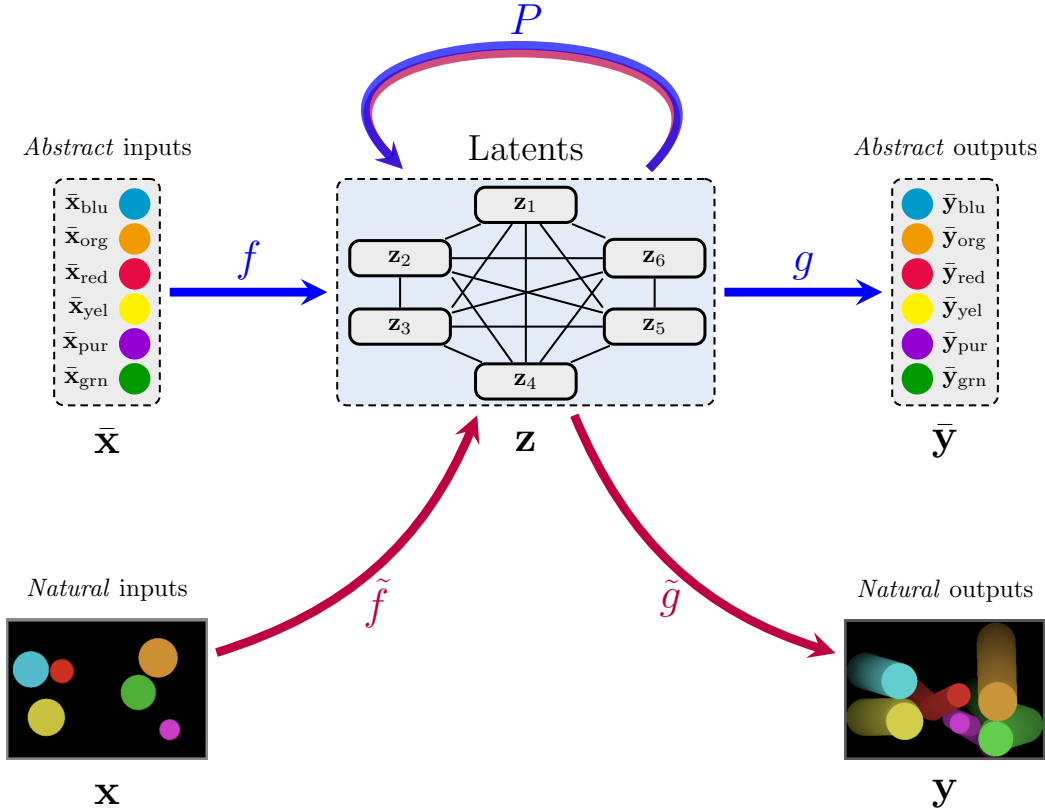


Figure 2: Reasoning-modulated representation learner (RMR).

129 $g(P(f(\bar{\mathbf{x}}))) \approx A(\bar{\mathbf{x}})$, which follows the encode-process-decode paradigm [43]. It consists of the
 130 following three building blocks: Encoder, $f : \bar{\mathcal{X}} \rightarrow \mathcal{Z}$, tasked with projecting the abstract inputs
 131 into the latent space; Processor, $P : \mathcal{Z} \rightarrow \mathcal{Z}$, simulating individual steps of the algorithm in the
 132 latent space; Decoder, $g : \mathcal{Z} \rightarrow \bar{\mathcal{Y}}$, tasked with projecting latents back into the abstract output space.
 133 Such a pipeline is now widely used both in neural algorithmic reasoning [13] and learning physical
 134 simulations [21], and can be trained end-to-end with gradient descent.

135 For reasons that will become apparent, it is favourable for most of the computational effort to be
 136 performed by P . Accordingly, encoders and decoders in the abstract pipeline are often designed to be
 137 simple learnable linear projections as there is usually no need for elaborate encoders/decoders, which,
 138 we will see, is not the case in the natural pipeline, due to nontrivial geometry of inputs and outputs.
 139 The processors tend to be either deep MLPs or graph neural networks—depending on whether the
 140 latent space is factorised into nodes.

141 **Natural pipeline.** Once an appropriate processor, P , has been learned, it may be observed that it
 142 corresponds to a highly favourable component in our setting. Namely, we can relate its operations to
 143 the algorithm A , and since it stays high-dimensional, it is a differentiable component we can easily
 144 plug into other neural networks without incurring any bottleneck effects. This insight was originally
 145 recovered in XLVIN [8], where it yielded a generic implicit planner. As an ablation, we have also
 146 rediscovered the bottleneck effect in our settings; see Appendix C. We now leverage similar insights
 147 for general representation learning tasks.

148 On a high level, what we need to do is simple and elegant: swap out f and g for *natural* encoders
 149 and decoders, $\tilde{f} : \mathcal{X} \rightarrow \mathcal{Z}$ and $\tilde{g} : \mathcal{Z} \rightarrow \mathcal{Y}$, respectively. We are then able to learn a function
 150 $\tilde{g}(P(\tilde{f}(\mathbf{x}))) \approx \Phi(\mathbf{x})$, which is once again to be optimised through gradient descent. We would like P
 151 to retain its semantics during training, and therefore it is typically kept *frozen* in the natural pipeline.
 152 Note that P might not perfectly represent A which in turn might not perfectly represent Φ . While we
 153 rely on a skip connection in our implementation of P , it has no learnable parameters and does not

154 offer the system the ability to learn a correction of P in the natural setting. Our choice is motivated
 155 by the desire to both maintain the semantics and interpretability of P and to force the model to rely
 156 on the processor P , not to simply bypass it. We show empirically that our pipeline is surprisingly
 157 robust to imperfect P models even with weak (linear) encoders/decoders.

158 It is worth noting several potential challenges that may arise while training the natural pipeline,
 159 especially if the training data for it is sparsely available. We also suggest remedies for each:

- 160 • If \mathbf{x} and/or \mathbf{y} exhibit any nontrivial geometry, simple linear projections will rarely suffice for \tilde{f}
 161 and \tilde{g} . For example, our natural inputs will often be pixel-based, necessitating a convolutional
 162 neural network for \tilde{f} .
- 163 • Further, since the parameters of P are kept frozen, \tilde{f} is left with a challenging task of mapping
 164 natural inputs into an appropriate manifold that P can meaningfully operate over. While we
 165 demonstrate clear empirical evidence that such meaningful mappings definitely occur, we remark
 166 that its success may hinge on carefully tuning the hyperparameters of \tilde{f} .
- 167 • A very common setting assumes that the abstract inputs and latents are factorised into objects,
 168 but the natural inputs are not. In this case, \tilde{f} is tasked with predicting appropriate object
 169 representations from the natural inputs. This is known to be a challenging feat [44], but can be
 170 successfully performed. Sometimes arbitrarily factorising the feature maps of a CNN [33] is
 171 sufficient, while at other times, models such as slot attention [40] may be required.
- 172 • One corollary of using automated object extractors for \tilde{f} is that it’s very difficult to enforce their
 173 slot representations to line up in the same way as in the abstract inputs. This implies that P
 174 should be permutation equivariant (and hence motivates using a GNN for it).

175 4 RMR for bouncing balls

176 To evaluate the capability of the RMR pipeline for transfer from the abstract space to the pixel
 177 space, we apply it on the “bouncing balls” problem. The bouncing balls problem is an instance of a
 178 physics simulation problem, where the task is to predict the next state of an environment in which
 179 multiple balls are bouncing between each other and a bounding box. Though this problem had been
 180 studied in the the context of physics simulation from (abstract) trajectories [7] and from (natural)
 181 videos [20, 28, 45], here we focus on the aptitude of RMR to transfer learned representations from
 182 trajectories to videos.

183 Our results affirm that strong abstract models can be trained on such tasks, and that including
 184 them in a video pipeline induces more robust representations. See Appendix A for more details on
 185 hyperparameters and experimental setup.

186 **Preliminaries.** Here, trajectories are represented by 2D coordinates of 10 balls through time,
 187 defining our abstract inputs and outputs $\mathcal{X} = \mathcal{Y} = \mathbb{R}^{10 \times 2}$. We slice these trajectories into a series
 188 of moving windows containing the input, $\bar{\mathbf{x}}^*$, spanning a history of three previous states, and the
 189 target, $\bar{\mathbf{y}}$, representing the next state. We obtain these trajectories from a 3D simulator (MuJoCo [46]),
 190 together with their short-video renderings, which represent our natural input and output space
 191 $\mathcal{X} = \mathcal{Y} = \mathbb{R}^{64 \times 64 \times 3}$. Our goal is to train an RMR abstract model on trajectories and transfer learned
 192 representations to improve a dynamics model trained on these videos.

193 **Abstract pipeline.** So as to model the dynamics of trajectories, we closely follow the RMR
 194 desiderata for the abstract model. We set f to a linear projection over the input concatenation, P
 195 to a Message Passing Neural Network (MPNN), following previous work [7, 21], and g to a linear
 196 projection.

197 Our model learns a transition function $g(P(f(\bar{\mathbf{x}}^*)) \approx \bar{\mathbf{y}}$, supervised using Mean Squared Error
 198 (MSE) over ball positions in the next step. It achieves an MSE of 4.59×10^{-4} , which, evaluated
 199 qualitatively, demonstrates the ability of the model to predict physically realistic behavior when
 200 unrolled for 10 steps (the model is trained on 1-step dynamics only). Next, we take the processor P
 201 from the abstract pipeline and re-use it in the natural pipeline.

202 **Natural pipeline.** Here we evaluate whether the pre-trained RMR processor can be reused for
 203 learning the dynamics of the bouncing balls from videos. The pixel-based encoder \tilde{f} simply runs a

204 Slot Attention model [40] on each input image, concatenates its outputs and passes them through a
 205 linear layer. The pixel-based decoder \bar{g} is a Broadcast Decoder [47].

206 The full model is a transition function $\bar{g}(P(\bar{f}(\mathbf{x}^*)) \approx \mathbf{y}$, supervised by pixel reconstruction loss over
 207 the next step image. We compare the performance of the RMR model with a pre-trained processor
 208 P against a baseline in which P is trained fully end-to-end. The RMR model achieves an MSE of
 209 $7.94 \pm 0.41 (\times 10^{-4})$, whereas the baseline achieves $9.47 \pm 0.24 (\times 10^{-4})$. We take a qualitative look
 210 at the reconstruction rollout of the RMR model in Figure 3, and expose the algorithmic bottleneck
 211 properties for this task in Appendix C.

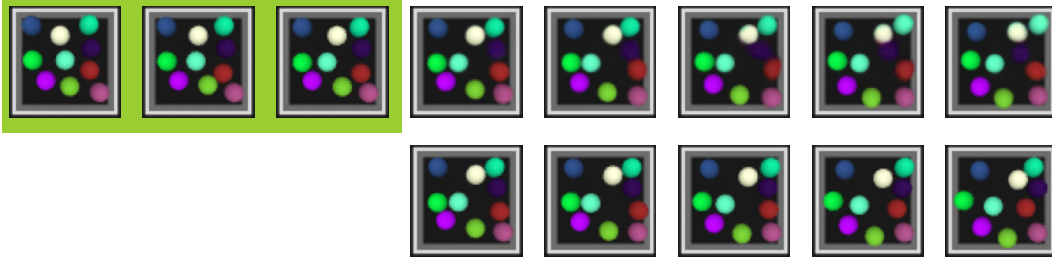


Figure 3: RMR for bouncing balls reconstruction rollout. States marked in green are the natural input, followed by the reconstructed output. The states below the reconstruction are the ground truth.

212 5 Contrastive RMR for Atari

213 We evaluate the potential of our RMR pipeline for state representation learning on the Atari 2600
 214 [48]. We find the RMR applicable here because there is a potential wealth of information that can be
 215 obtained about the Atari’s operation—namely, by inspecting its RAM traces.

216 **Preliminaries.** Accordingly, we will define our set of abstract inputs and outputs as Atari RAM
 217 matrices. Given that the Atari has 128 bytes of memory, $\mathcal{X} = \mathcal{Y} = \mathbb{B}^{128 \times 8}$ (where $\mathbb{B} = \{0, 1\}$ is the
 218 set of bits). We collect data about how the console modifies the RAM by acting in the environment
 219 and recording the trace of RAM arrays we observe. These traces will be of the form $(\bar{\mathbf{x}}, a, \bar{\mathbf{y}})$ which
 220 signify that the agent’s initial RAM state was $\bar{\mathbf{x}}$, and that after performing action $a \in \mathcal{A}$, its RAM
 221 state was updated to $\bar{\mathbf{y}}$. We assume that a is encoded as an 18-way one-hot vector.

222 We would like to leverage any reasoning module obtained over RAM states to support representation
 223 learning from raw pixels. Accordingly, our natural inputs, \mathcal{X} , are pixel arrays representing the Atari’s
 224 framebuffer.

225 Mirroring prior work, we perform contrastive learning directly in the latent space, and set $\mathcal{Y} = \mathcal{Z}$;
 226 that is, our natural outputs correspond to an estimate of the “updated” latents after taking an action.
 227 All our models use latent representations of 64 dimensions per slot, meaning $\mathcal{Z} = \mathbb{R}^{128 \times 64}$.

228 We note that it is important to generate a diverse dataset of experiences in order to train a robust
 229 RAM model. To simulate a dataset which might be gathered by human players of varying skill, we
 230 sample our data using the 32 policy heads of a pre-trained Agent57 [49]. Each policy head collects
 231 data over three episodes in the studied games. Note that this implies a substantially more challenging
 232 dataset than the one reported by [50], wherein data was collected by a purely random policy, which
 233 may well fail to explore many relevant regions of the games.

234 **Abstract pipeline.** Firstly, we set out to verify that it is possible to train nontrivial Atari RAM
 235 transition models. The construction of this abstract experiment follows almost exactly the abstract
 236 RMR setup: f and g are appropriately sized linear projections, while P needs to take into account
 237 which action was taken when updating the latents. To simplify the implementation and allow further
 238 model re-use, we consider the action a part of the P ’s inputs. See Appendix F for detailed equations.

239 This implies that our transition model learns a function $g(P(f(\bar{\mathbf{x}}), a)) \approx \bar{\mathbf{y}}$. We supervise this model
 240 using binary cross-entropy to predict each bit of the resulting RAM state. Since RAM transitions are
 241 assumed deterministic, we assume a fully Markovian setup and learn 1-step dynamics.

242 For brevity purposes, we detail our exact hyperparameters and results per each Atari game considered
 243 in Appendix B. Our results ascertain the message passing neural network (MPNN) [51] as a highly
 244 potent processor network in Atari: it ranked most potent in 17 out of 24 games considered, compared
 245 to MLPs and Deep Sets [52]. Accordingly, we will focus on leveraging pre-trained MPNN processors
 246 for the next phase of the RMR pipeline.

247 **Natural pipeline.** We now set out to evaluate whether our pre-trained RMR processors can be
 248 meaningfully re-used by an encoder in a pixel-based contrastive learning pipeline.

Table 1: Natural modelling results for Atari 2600. Bit-level F_1 reported for slots with high entropy, as in Anand et al. [50]. Results assumed **significant** at $p < 0.05$ (one-sided paired Wilcoxon test).

Game	Agent57	C-SWM	RMR	p -value
Asteroids	0.514±0.001	0.582±0.009	0.593 ±0.004	< 10^{-5}
Battlezone	0.351±0.003	0.592±0.005	0.589±0.007	0.056
Berzerk	0.454±0.084	0.463±0.053	0.470±0.025	0.364
Bowling	0.554±0.004	0.944±0.006	0.946±0.003	0.071
Boxing	0.558±0.002	0.667±0.012	0.669±0.011	0.215
Breakout	0.657±0.001	0.836±0.009	0.852 ±0.008	< 10^{-5}
Demon Attack	0.539±0.004	0.653±0.006	0.658 ±0.004	0.002
Freeway	0.424±0.052	0.912±0.025	0.919 ±0.035	0.032
Frostbite	0.405±0.001	0.580±0.025	0.594 ±0.016	0.035
H.E.R.O.	0.481±0.001	0.729±0.026	0.779 ±0.021	< 10^{-5}
Montezuma’s Revenge	0.743±0.003	0.824±0.012	0.821±0.016	0.156
Ms. Pac-Man	0.506±0.001	0.599±0.004	0.602 ±0.006	0.038
Pitfall!	0.495±0.003	0.626 ±0.015	0.603±0.010	< 10^{-5}
Pong	0.392±0.001	0.750±0.016	0.762 ±0.010	0.001
Private Eye	0.594±0.001	0.863±0.010	0.867 ±0.008	0.045
Q*Bert	0.536±0.010	0.588±0.015	0.590±0.017	0.165
River Raid	0.686±0.001	0.762±0.005	0.764 ±0.007	0.032
Seaquest	0.472±0.007	0.634±0.013	0.653 ±0.008	< 10^{-5}
Skiing	0.599±0.007	0.766±0.028	0.775±0.014	0.174
Space Invaders	0.588±0.002	0.719±0.012	0.761 ±0.006	< 10^{-5}
Tennis	0.533±0.008	0.724±0.007	0.729 ±0.005	0.007
Venture	0.567±0.001	0.632±0.005	0.633±0.004	0.392
Video Pinball	0.375±0.011	0.724±0.009	0.745 ±0.008	< 10^{-5}
Yars’ Revenge	0.608±0.001	0.715±0.008	0.751 ±0.010	< 10^{-5}

249 For our pixel-based encoder \tilde{f} , we use the same CNN trunk as Anand et al. [50]—however, as we
 250 require slot-level rather than flat embeddings, the final layers of our encoder are different. Namely,
 251 we apply a 1×1 convolution computing $128m$ feature maps (where m is the number of feature maps
 252 per-slot). We then flatten the spatial axes, giving every slot $m \times h \times w$ features, which we finally
 253 linearly project to 64-dimensional features per-slot, aligning with our pre-trained P . Note that setting
 254 $m = 1$ recovers exactly the style of object detection employed by C-SWM [33]. Since our desired
 255 outputs are themselves latents, \tilde{g} is a single linear projection to 64 dimensions.

256 Overall, our pixel-based transition model learns a function $\tilde{g}(P(\tilde{f}(\mathbf{x}), a)) \approx \tilde{f}(\mathbf{y})$, where \mathbf{y} is the
 257 next state observed after applying action a in state \mathbf{x} . To optimise it, we re-use exactly the same
 258 TransE-inspired [53] contrastive loss that C-SWM [33] used.

259 Once state representation learning concludes, all components are typically thrown away except for
 260 the encoder, \tilde{f} , which is used for downstream tasks. As a proxy for evaluating the quality of the
 261 encoder, we train linear classifiers on the concatenation of all slot embeddings obtained from \tilde{f} to
 262 predict individual RAM bits, exactly as in the abstract model case. Note that we have *not* violated
 263 our assumption that paired $(\mathbf{x}, \bar{\mathbf{x}})$ samples will not be provided while training the natural model—in
 264 this phase, the encoder \tilde{f} is frozen, and gradients can only flow into the linear probe.

265 Our comparisons for this experiment, evaluating our RMR pipeline against an identical architecture
 266 with an unfrozen P (equivalent to C-SWM [33]) is provided in Table 1. For each of 20 random seeds,

267 we feed identical batches to both models, hence we can perform a paired Wilcoxon test to assess the
 268 statistical significance of any observed differences on validation episodes. The results are in line with
 269 our hypothesis: representations learnt by RMR are significantly better ($p < 0.05$) on 15 out of 24
 270 games, significantly worse only on Pitfall!—and indistinguishable to C-SWM’s on others. This is
 271 despite the fact their architectures are identical, indicating that the pre-trained abstract model induces
 272 stronger representations for predicting the underlying data factors. As was the case for bouncing
 273 balls, we expose the algorithmic bottleneck here too; see Appendix C.

274 As a relevant initial baseline—and to emphasise the difficulty of our task—we also include in Table
 275 1 the performance of the latent embeddings extracted from a pre-trained Agent57 [49]. These
 276 embeddings are, perhaps unsurprisingly, substantially worse than both the RMR and C-SWM models.
 277 As Agent57 was trained on a reward-maximising objective, its embeddings are likely to capture the
 278 controllable aspects of the input, while filtering out the environment’s background.

279 **Abstract model transfer.** While the results above indicate that endowing a self-supervised Atari
 280 representation learner with knowledge of the underlying game’s RAM transitions yields stronger
 281 representations, one may still argue that this constitutes a form of “privileged information”. This is
 282 due to the fact we knew upfront which game we were learning the representations for, and hence
 283 could leverage the specific RAM dynamics of this game.

284 In the final experiment, we study a more general setting: we are given access to RAM traces showing
 285 us how the Atari console manipulates its memory in response to player input, but we *cannot* guarantee
 286 these traces came from the same game that we are performing representation learning over. Can we
 287 still effectively leverage this knowledge?

288 Specifically, we test *abstract model transfer* in the following way: first, we take a pre-trained abstract
 289 processor network P from one game (the “*train game*”) and freeze it. Then, using this processor, we
 290 perform the aforementioned natural pipeline training and testing over frames from another game (the
 291 “*test game*”). We then evaluate whether the recovered performance improves over the C-SWM model
 292 with unfrozen weights—once again using a paired one-sided Wilcoxon test over 20 seeds to ascertain
 293 statistical significance of any differences observed. The final outcome of our experiment is hence
 294 a 24×24 matrix, indicating the quality of abstract model transfer from every game to every other
 295 game. Table 1 corresponds to the “diagonal” entries of this matrix.

296 The results, presented in Figure 4, testify to the performance of RMR. Representations learned by
 297 RMR transfer better in 64.6% of the train/test game pairs, are indistinguishable from C-SWM in
 298 33.7% of the game pairs and perform worse than C-SWM in only 1.7% of game pairs. Therefore, in
 299 plentiful circumstances, the answer to our original question is *positive*: discovering a trace of Atari
 300 RAM transitions of unknown origin can often be of high significance for representation learning from
 301 Atari pixels, regardless of whether the underlying games match.

302 **Qualitative analysis of transfer.** While we find this result interesting in and of itself, it also raises
 303 interesting follow-up questions. Is representation learning on certain games more prone to being
 304 improved just by knowing *anything* about the Atari console? Figure 4 certainly implies so: several
 305 games (such as Seaquest or Space Invaders) have “fully-green” columns, making them “universal
 306 recipients”. Similarly, we may be interested about the “donor” properties of each game – to what
 307 extent are their RAM models useful across a broad range of test games?

308 We study both of the above questions by performing a hierarchical (complete-link) clustering of the
 309 rows and columns of the 24×24 matrix of transfer performances, to identify clusters of related donor
 310 and recipient games. Both clusterings are marked in Figure 4 (on the sides of the matrix).

311 The analysis reveals several well-formed clusters, from which we are able to make some preliminary
 312 observations, based on the properties of the various games. To name a few examples:

- 313 • The strongest recipients (e.g. Yars’ Revenge, H.E.R.O., Seaquest, Space Invaders and Asteroids)
 314 tend to include elements of “shooting” in their gameplay.
- 315 • Conversely, the weakest recipients (River Raid, Berzerk, Private Eye, Ms. Pac-Man and Pitfall!)
 316 are all games in which movement is generally unrestricted across most of the screen, indicating
 317 a larger range of possible coordinate values to model.
- 318 • Pong, Breakout, Battlezone and Skiing cluster closely in terms of donor properties—and they
 319 are all games in which movement is restricted to one axis only.



Figure 4: A hierarchically clustered heatmap depiction of the transfer results, where the *Train Game* abstract models have been trained on the *Test Games*. The results are summarised by a Wilcoxon test, denoting where the performance of RMR is better, does not differ or is worse than C-SWM’s. In each cell, the mean relative % improvement is noted. The significance cutoff is $p < 0.05$.

- 320 • Lastly, strong donors (Montezuma’s Revenge, Yars’ Revenge, Q*Bert and Venture) all feature
 321 massive, abrupt, changes to the game state (as they feature multiple rooms, for example). This
 322 implies that they might generally have seen a more diverse set of transitions. Further, on Yars’
 323 Revenge, a massive laser features, which, conveniently, prints an RGB projection of the RAM
 324 itself on the frame.

325 6 Conclusions

326 We presented Reasoning-Modulated Representations (RMR), a novel approach for leveraging back-
 327 ground algorithmic knowledge within a representation learning pipeline. By encoding the underlying
 328 algorithmic priors as weights of a processor neural network, we alleviate requirements on alignment
 329 between abstract and natural inputs and protect our model against bottlenecks. We believe that RMR
 330 paves the way to a novel class of algorithmic reasoning-inspired representations, with a high potential
 331 for transfer across tasks—a feat that has largely eluded deep reinforcement learning research.

References

- 332
- 333 [1] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Si-
334 mon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering
335 atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
336 1
- 337 [2] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
338 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
339 few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. 1
- 340 [3] Andrew W Senior, Richard Evans, John Jumper, James Kirkpatrick, Laurent Sifre, Tim Green,
341 Chongli Qin, Augustin Žídek, Alexander WR Nelson, Alex Bridgland, et al. Improved protein
342 structure prediction using potentials from deep learning. *Nature*, 577(7792):706–710, 2020. 1
- 343 [4] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-
344 learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy
345 of Sciences*, 116(32):15849–15854, 2019. 1
- 346 [5] Yann LeCun. The power and limits of deep learning. *Research-Technology Management*, 61(6):
347 22–27, 2018.
- 348 [6] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable
349 effectiveness of data in deep learning era. In *Proceedings of the IEEE International Conference
350 on Computer Vision (ICCV)*, Oct 2017. 1
- 351 [7] Peter W Battaglia, Razvan Pascanu, Matthew Lai, Danilo Rezende, and Koray Kavukcuoglu.
352 Interaction networks for learning about objects, relations and physics. *arXiv preprint
353 arXiv:1612.00222*, 2016. 1, 3, 5, 14
- 354 [8] Andreea-Ioana Deac, Petar Veličković, Ognjen Milinkovic, Pierre-Luc Bacon, Jian Tang, and
355 Mladen Nikolich. Neural algorithmic reasoners are implicit planners. *Advances in Neural
356 Information Processing Systems*, 34, 2021. 2, 4, 13, 16
- 357 [9] Petar Veličković and Charles Blundell. Neural algorithmic reasoning. *Patterns*, 2(7):100273,
358 2021. 2
- 359 [10] David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018. 2
- 360 [11] Quentin Cappart, Didier Chételat, Elias Khalil, Andrea Lodi, Christopher Morris, and Petar
361 Veličković. Combinatorial optimization and reasoning with graph neural networks. *arXiv
362 preprint arXiv:2102.09544*, 2021. 2
- 363 [12] Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie
364 Jegelka. What can neural networks reason about? *arXiv preprint arXiv:1905.13211*, 2019. 2
- 365 [13] Petar Veličković, Rex Ying, Matilde Padovano, Raia Hadsell, and Charles Blundell. Neural
366 execution of graph algorithms. *arXiv preprint arXiv:1910.10593*, 2019. 2, 4, 16
- 367 [14] Hao Tang, Zhiao Huang, Jiayuan Gu, Bao-Liang Lu, and Hao Su. Towards scale-invariant
368 graph-related problem solving by iterative homogeneous graph neural networks. *arXiv preprint
369 arXiv:2010.13547*, 2020. 2
- 370 [15] Kārlis Freivalds, Emīls Ozoliņš, and Agris Šostaks. Neural shuffle-exchange networks–sequence
371 processing in $o(n \log n)$ time. *arXiv preprint arXiv:1907.07897*, 2019. 2
- 372 [16] Heiko Strathmann, Mohammadamin Barekatin, Charles Blundell, and Petar Veličković. Persis-
373 tent message passing. *arXiv preprint arXiv:2103.01043*, 2021. 2, 13
- 374 [17] Petar Veličković, Lars Buesing, Matthew C Overlan, Razvan Pascanu, Oriol Vinyals, and
375 Charles Blundell. Pointer graph networks. *arXiv preprint arXiv:2006.06380*, 2020. 2, 13
- 376 [18] Louis-Pascal Xhonneux, Andreea-Ioana Deac, Petar Veličković, and Jian Tang. How to transfer
377 algorithmic reasoning knowledge to learn new algorithms? *Advances in Neural Information
378 Processing Systems*, 34:19500–19512, 2021. 2
- 379 [19] Michael B Chang, Tomer Ullman, Antonio Torralba, and Joshua B Tenenbaum. A compositional
380 object-based approach to learning physical dynamics. *arXiv preprint arXiv:1612.00341*, 2016.
381 3
- 382 [20] Nicholas Watters, Daniel Zoran, Theophane Weber, Peter Battaglia, Razvan Pascanu, and
383 Andrea Tacchetti. Visual interaction networks: Learning a physics simulator from video.
384 *Advances in neural information processing systems*, 30:4539–4547, 2017. 3, 5

- 385 [21] Alvaro Sanchez-Gonzalez, Jonathan Godwin, Tobias Pfaff, Rex Ying, Jure Leskovec, and
386 Peter Battaglia. Learning to simulate complex physics with graph networks. In *International*
387 *Conference on Machine Learning*, pages 8459–8468. PMLR, 2020. 3, 4, 5
- 388 [22] Tobias Pfaff, Meire Fortunato, Alvaro Sanchez-Gonzalez, and Peter W Battaglia. Learning
389 mesh-based simulation with graph networks. *arXiv preprint arXiv:2010.03409*, 2020. 3
- 390 [23] Miles Cranmer, Alvaro Sanchez-Gonzalez, Peter Battaglia, Rui Xu, Kyle Cranmer, David
391 Spergel, and Shirley Ho. Discovering symbolic models from deep learning with inductive
392 biases. *arXiv preprint arXiv:2006.11287*, 2020. 3
- 393 [24] Victor Bapst, Thomas Keck, A Grabska-Barwińska, Craig Donner, Ekin Dogus Cubuk, Samuel S
394 Schoenholz, Annette Obika, Alexander WR Nelson, Trevor Back, Demis Hassabis, et al.
395 Unveiling the predictive power of static structure in glassy systems. *Nature Physics*, 16(4):
396 448–454, 2020. 3
- 397 [25] Alvaro Sanchez-Gonzalez, Victor Bapst, Kyle Cranmer, and Peter Battaglia. Hamiltonian graph
398 networks with ode integrators. *arXiv preprint arXiv:1909.12790*, 2019. 3
- 399 [26] Miguel Jaques, Michael Burke, and Timothy Hospedales. Physics-as-inverse-graphics: Unsu-
400 pervised physical parameter estimation from video. *arXiv preprint arXiv:1905.11169*, 2019.
401 3
- 402 [27] Krishna Murthy Jatavallabhula, Miles Macklin, Florian Golemo, Vikram Voleti, Linda Petrini,
403 Martin Weiss, Breandan Considine, Jerome Parent-Levesque, Kevin Xie, Kenny Erleben, et al.
404 gradsim: Differentiable simulation for system identification and visuomotor control. *arXiv*
405 *preprint arXiv:2104.02646*, 2021. 3
- 406 [28] Sjoerd Van Steenkiste, Michael Chang, Klaus Greff, and Jürgen Schmidhuber. Relational neural
407 expectation maximization: Unsupervised discovery of objects and their interactions. *arXiv*
408 *preprint arXiv:1802.10353*, 2018. 3, 5
- 409 [29] Adam R Kosiorek, Hyunjik Kim, Ingmar Posner, and Yee Whye Teh. Sequential attend, infer,
410 repeat: Generative modelling of moving objects. *arXiv preprint arXiv:1806.01794*, 2018. 3
- 411 [30] Rishi Veerapaneni, John D Co-Reyes, Michael Chang, Michael Janner, Chelsea Finn, Jiajun Wu,
412 Joshua Tenenbaum, and Sergey Levine. Entity abstraction in visual model-based reinforcement
413 learning. In *Conference on Robot Learning*, pages 1439–1456. PMLR, 2020. 3
- 414 [31] Jindong Jiang, Sepehr Janghorbani, Gerard De Melo, and Sungjin Ahn. Scalor: Generative
415 world models with scalable object representations. *arXiv preprint arXiv:1910.02384*, 2019. 3
- 416 [32] Zhixuan Lin, Yi-Fu Wu, Skand Peri, Bofeng Fu, Jindong Jiang, and Sungjin Ahn. Improving
417 generative imagination in object-centric world models. In *International Conference on Machine*
418 *Learning*, pages 6140–6149. PMLR, 2020. 3
- 419 [33] Thomas Kipf, Elise van der Pol, and Max Welling. Contrastive learning of structured world
420 models. *arXiv preprint arXiv:1911.12247*, 2019. 3, 5, 7, 13
- 421 [34] Vincent François-Lavet, Yoshua Bengio, Doina Precup, and Joelle Pineau. Combined reinforce-
422 ment learning via abstract representations. In *Proceedings of the AAAI Conference on Artificial*
423 *Intelligence*, volume 33, pages 3582–3589, 2019. 3
- 424 [35] Tejas D Kulkarni, Ankush Gupta, Catalin Ionescu, Sebastian Borgeaud, Malcolm Reynolds,
425 Andrew Zisserman, and Volodymyr Mnih. Unsupervised learning of object keypoints for
426 perception and control. *Advances in neural information processing systems*, 32:10724–10734,
427 2019. 3
- 428 [36] Yunzhu Li, Antonio Torralba, Animashree Anandkumar, Dieter Fox, and Animesh Garg. Causal
429 discovery in physical systems from videos. *arXiv preprint arXiv:2007.00631*, 2020. 3
- 430 [37] Thomas Kipf, Ethan Fetaya, Kuan-Chieh Wang, Max Welling, and Richard Zemel. Neural
431 relational inference for interacting systems. In *International Conference on Machine Learning*,
432 pages 2688–2697. PMLR, 2018. 3
- 433 [38] Adam Santoro, Ryan Faulkner, David Raposo, Jack Rae, Mike Chrzanowski, Theophane Weber,
434 Daan Wierstra, Oriol Vinyals, Razvan Pascanu, and Timothy Lillicrap. Relational recurrent
435 neural networks. *arXiv preprint arXiv:1806.01822*, 2018. 3
- 436 [39] Anirudh Goyal, Alex Lamb, Jordan Hoffmann, Shagun Sodhani, Sergey Levine, Yoshua Bengio,
437 and Bernhard Schölkopf. Recurrent independent mechanisms. *arXiv preprint arXiv:1909.10893*,
438 2019. 3

- 439 [40] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg
440 Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with
441 slot attention. *arXiv preprint arXiv:2006.15055*, 2020. 3, 5, 6, 13
- 442 [41] Anirudh Goyal, Alex Lamb, Phanideep Gampa, Philippe Beaudoin, Sergey Levine, Charles
443 Blundell, Yoshua Bengio, and Michael Mozer. Factorizing declarative and procedural knowledge
444 in structured, dynamic environments. In *International Conference on Learning Representations*,
445 2021. 3
- 446 [42] Anirudh Goyal, Aniket Didolkar, Nan Rosemary Ke, Charles Blundell, Philippe Beaudoin,
447 Nicolas Heess, Michael Mozer, and Yoshua Bengio. Neural production systems. *arXiv preprint*
448 *arXiv:2103.01937*, 2021. 3
- 449 [43] Jessica B Hamrick, Kelsey R Allen, Victor Bapst, Tina Zhu, Kevin R McKee, Joshua B
450 Tenenbaum, and Peter W Battaglia. Relational inductive bias for physical construction in
451 humans and machines. *arXiv preprint arXiv:1806.01203*, 2018. 4
- 452 [44] Klaus Greff, Sjoerd van Steenkiste, and Jürgen Schmidhuber. On the binding problem in
453 artificial neural networks. *arXiv preprint arXiv:2012.05208*, 2020. 5
- 454 [45] Sindy Löwe, Klaus Greff, Rico Jonschkowski, Alexey Dosovitskiy, and Thomas Kipf. Learning
455 object-centric video models by contrasting sets. *arXiv preprint arXiv:2011.10287*, 2020. 5
- 456 [46] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based
457 control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages
458 5026–5033. IEEE, 2012. 5
- 459 [47] Nicholas Watters, Loic Matthey, Christopher P Burgess, and Alexander Lerchner. Spatial
460 broadcast decoder: A simple architecture for learning disentangled representations in vaes.
461 *arXiv preprint arXiv:1901.07017*, 2019. 6, 13
- 462 [48] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning
463 environment: An evaluation platform for general agents. *Journal of Artificial Intelligence*
464 *Research*, 47:253–279, 2013. 6
- 465 [49] Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvit-
466 skyi, Zhaohan Daniel Guo, and Charles Blundell. Agent57: Outperforming the atari human
467 benchmark. In *International Conference on Machine Learning*, pages 507–517. PMLR, 2020.
468 6, 8
- 469 [50] Ankesh Anand, Evan Racah, Sherjil Ozair, Yoshua Bengio, Marc-Alexandre Côté, and R Devon
470 Hjelm. Unsupervised state representation learning in atari. *arXiv preprint arXiv:1906.08226*,
471 2019. 6, 7, 13, 14
- 472 [51] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
473 message passing for quantum chemistry. In *International Conference on Machine Learning*,
474 pages 1263–1272. PMLR, 2017. 7, 13, 16
- 475 [52] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan Salakhutdinov,
476 and Alexander Smola. Deep sets. *arXiv preprint arXiv:1703.06114*, 2017. 7, 13
- 477 [53] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko.
478 Translating embeddings for modeling multi-relational data. In *Neural Information Processing*
479 *Systems (NIPS)*, pages 1–9, 2013. 7
- 480 [54] Adam Santoro, David Raposo, David GT Barrett, Mateusz Malinowski, Razvan Pascanu, Peter
481 Battaglia, and Timothy Lillicrap. A simple neural network module for relational reasoning.
482 *arXiv preprint arXiv:1706.01427*, 2017. 13
- 483 [55] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
484 *arXiv:1412.6980*, 2014. 13
- 485 [56] Keyulu Xu, Mozhi Zhang, Jingling Li, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie
486 Jegelka. How neural networks extrapolate: From feedforward to graph neural networks. *arXiv*
487 *preprint arXiv:2009.11848*, 2020. 16

488 A Bouncing balls modelling setup

489 **Abstract pipeline.** The f is a linear projection over the concatenation of inputs, outputting a 128-
 490 long representation for each of the 10 balls in the input. The P is a MPNN over the fully connected
 491 graph of the ball representations. It is a 2-pass MPNN with a 3-layered ReLU-activated MLP as a
 492 message function, without the final layer activation, projecting to the same 128-dimensional space,
 493 per object. Finally, g is a linear projection applied on each object representation of the output of P .

494 The model is MSE-supervised with ball position on the next step. It is trained on a 8-core TPU for
 495 10000 epochs, with a batch size of 512, and the Adam optimizer with the initial learning rate of
 496 0.0001.

497 **Natural pipeline.** \bar{f} is a Slot Attention model [40] on each input image, concatenating the images
 498 and passing them through a linear layer, outputting a 128-dimensional vector for each of the objects.
 499 \bar{g} is a Broadcast Decoder [47] containing a sequence of 5 transposed convolutions and a linear layer
 500 mapping before calculating the reconstructions and their masks.

501 The model is MSE-supervised by pixel reconstruction (per-pixel MSE) over the next step image.
 502 Both the RMR and the baseline are trained on a 8-core TPU for 1000 epochs, with a batch size of
 503 512, with the Adam optimiser and the initial learning rate of 0.0001, all over 3 random seeds.

504 B Atari abstract modelling setup and results

505 Our best processor network is a MPNN [51] over a fully connected graph [54] of RAM slots, which
 506 concatenates the action embedding to every node (as done in [33]). It uses three-layer MLPs as
 507 message functions, with the ReLU activation applied after each hidden layer. The entire model
 508 is trained for every game in isolation, over 48 distinct episodes of Agent57 experience. We use
 509 the Adam SGD optimiser [55] with a batch size of 50 and a learning rate of 0.001 across all Atari
 510 experiments. To evaluate the benefits of message passing, we also compare our model to Deep Sets
 511 [52], which is equivalent to our MPNN model—only it passes messages over the identity adjacency
 512 matrix. Lastly, we evaluate the benefits of factorised latents by comparing our methods against a
 513 three-layer MLP applied on the flattened RAM state.

514 One immediate observation is that RAM updates in Atari are extremely sparse, with a *copy baseline*
 515 already being very strong for many games. To prevent the model from having to repeatedly re-learn
 516 identity functions, we also make it predict *masks* of the shape $\mathcal{M} = \mathbb{B}^{128}$, specifying which cells
 517 are to be overwritten by the model at this step. This strategy, coupled with teacher forcing (as
 518 done by [16, 17]) yielded substantially stronger predictors. We also use this observation to prevent
 519 over-inflating our prediction scores: we only display prediction accuracy over RAM slots with label
 520 entropy larger than 0.6 (as done by [50]).

521 The full results of training Atari RAM transition models, for the games studied in [50], are provided
 522 in Table 2. We evaluate both bit-level F_1 scores, as well as slot-level accuracy (for which all 8
 523 bits need to be predicted correctly in order to count), over the remaining 48 Agent57 episodes as
 524 validation. To the best of our knowledge, this is the first comprehensive feasibility study for learning
 525 Atari RAM transition models.

526 C Rediscovering the algorithmic bottleneck

527 As mentioned in the main text body, one of the key reasons in favour of a high-dimensional algorithmic
 528 component is to avoid the *algorithmic bottleneck*, as first exposed by Deac et al. [8].

529 In short, the performance guarantees of running classical algorithms rely on having the exactly correct
 530 inputs for them. If there are any errors in the predictions of these (usually very low-dimensional)
 531 inputs, these errors may propagate to the algorithmic computations and yield suboptimal results.
 532 Further, there is no room for any kind of fallback if such an event occurs.

533 In contrast, the high-dimensional neural processors like the ones we study here are not vulnerable to
 534 bottleneck effects: if any dimensions of the latent state are poorly predicted, the other components of
 535 it could step in and compensate for this. Further, we can easily support skip connections in the case
 536 where the algorithm is not fully descriptive of the problem we’re solving.

Table 2: Abstract modelling results for Atari 2600. Entire-slot accuracies and bit-level F_1 scores are reported only for slots with high entropy, as per [50].

Game	Copy baseline		MLP		Deep Sets		MPNN	
	Slot acc.	Bit F_1	Slot acc.	Bit F_1	Slot acc.	Bit F_1	Slot acc.	Bit F_1
Asteroids	70.65%	0.856	71.28%	0.872	72.84%	0.879	80.69%	0.930
Battlezone	57.09%	0.841	61.19%	0.840	61.71%	0.867	71.06%	0.892
Berzerk	84.32%	0.905	86.17%	0.930	84.16%	0.923	86.67%	0.933
Bowling	93.86%	0.972	97.43%	0.991	90.72%	0.966	98.41%	0.995
Boxing	59.78%	0.848	54.45%	0.834	59.79%	0.890	58.56%	0.877
Breakout	89.80%	0.949	92.77%	0.970	94.34%	0.979	96.45%	0.988
Demon Attack	67.90%	0.850	68.43%	0.864	66.70%	0.877	69.51%	0.879
Freeway	46.65%	0.787	75.93%	0.921	84.68%	0.959	89.17%	0.965
Frostbite	76.83%	0.904	79.09%	0.904	78.25%	0.946	76.52%	0.918
H.E.R.O.	76.71%	0.891	82.96%	0.932	80.17%	0.929	89.07%	0.956
Montezuma’s Revenge	82.58%	0.907	87.30%	0.941	85.90%	0.951	85.44%	0.932
Ms. Pac-Man	83.80%	0.941	80.60%	0.935	81.50%	0.952	85.88%	0.966
Pitfall!	66.60%	0.862	78.28%	0.923	81.92%	0.947	80.40%	0.941
Pong	68.76%	0.873	73.58%	0.911	74.71%	0.920	83.23%	0.952
Private Eye	75.25%	0.889	81.95%	0.932	84.77%	0.954	86.41%	0.955
Q*Bert	83.00%	0.915	90.07%	0.966	87.87%	0.943	89.26%	0.946
River Raid	76.95%	0.895	80.82%	0.927	69.96%	0.865	86.79%	0.954
Seaquest	71.23%	0.859	78.48%	0.898	75.53%	0.906	70.94%	0.798
Skiing	91.02%	0.966	93.42%	0.980	93.51%	0.983	96.37%	0.992
Space Invaders	81.67%	0.942	84.62%	0.957	89.38%	0.974	91.98%	0.985
Tennis	78.13%	0.890	82.13%	0.926	71.60%	0.856	80.20%	0.893
Venture	61.29%	0.858	63.16%	0.863	64.88%	0.886	76.56%	0.935
Video Pinball	76.71%	0.848	85.92%	0.912	78.64%	0.877	86.61%	0.913
Yars’ Revenge	69.25%	0.896	74.87%	0.929	72.69%	0.948	84.11%	0.969

537 In this section, we reaffirm the bottleneck effect for both Atari and bouncing ball experiments by
538 providing additional sets of ablations on the processor network’s latent size.

539 We ablate various RMR processor network architectures, for $\dim \mathbf{z} \in \{2, 4, 8, 16, 32, 64\}$, but leaving
540 all other components and operations unchanged.

541 The results of this ablation for the Atari representation learning setting are provided in Figure 5. It
542 can be clearly observed that, as we reduce the size of the latents, this typically induces a performance
543 regression in the downstream bit F_1 scores. This indicates the algorithmic bottleneck effect.

544 On the bouncing balls task, we observe the algorithmic bottleneck as well, with more pronounced
545 effects. Namely, for all latent sizes 32 and under, the validation MSE shoots up even though the
546 training MSE keeps improving. Ultimately, we also attempted using an Interaction Network [7]
547 as a bottlenecked RMR processor which requires projecting on to $\bar{\mathbf{x}}$ rather than \mathbf{z} , which exhibited
548 exactly the same behaviour, at a larger scale. The IN processor achieved a final validation MSE of
549 $341.40 \pm 5.60 (\times 10^{-4})$. This is roughly $43\times$ worse than the non-bottlenecked RMR.

550 D On the weak acceptor properties of Battlezone

551 Out of all the Atari games studied in our work, Battlezone could be singled out as the least “recep-
552 tive” to RMR processors—with only four games successfully donating their RAM representations to
553 its pixel based encoder (Figure 4).

554 We set out to study why this effect took place. Upon inspection of the game’s dynamics¹, we deter-
555 mined that Battlezone has certain properties that make representation learning uniquely challenging
556 and somewhat decoupled from the underlying console computation.

557 Namely, Battlezone is the only game in our dataset that is played in the “first-person” perspective.
558 Therefore, from the point of view of the pixel inputs, it may seem as if the player is always in the

¹https://www.youtube.com/watch?v=9X4_xy7rC1A

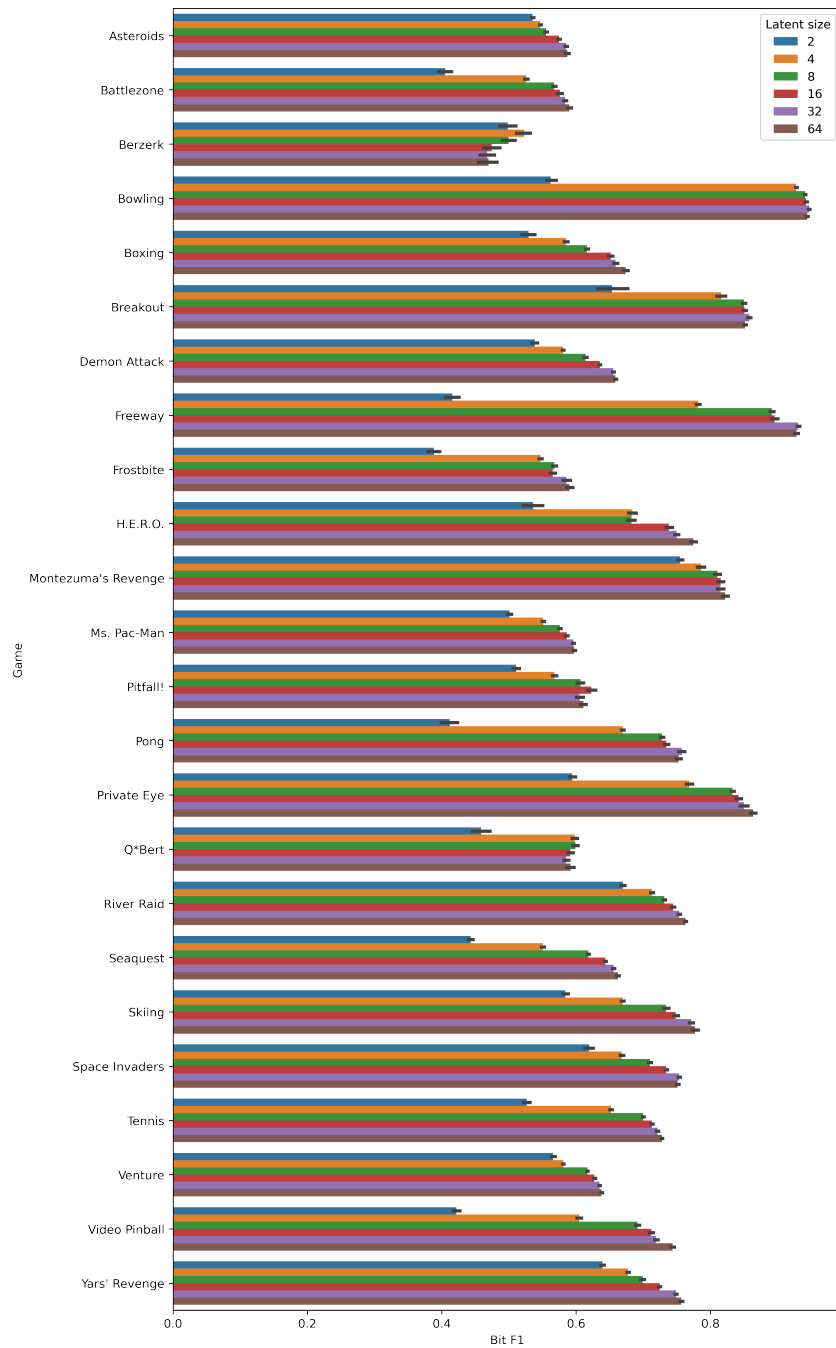


Figure 5: Bit-level F_1 scores for the Atari experiments while varying the latent size. A clear decreasing trend is apparent as $\dim \mathbf{z}$ is reduced, indicating the bottleneck effect.

559 same place, and we would expect the mechanics and the way of thinking about the ‘avatar’ to be
560 substantially different from all other Atari games considered.

561 **E On the data distribution and setup used for training P**

562 In many cases, while an algorithmic prior may be known, generative aspects of the relevant abstract
563 inputs for the natural task could be unknown. For example, it may be known that the natural task
564 could benefit from a sorting subroutine, but it may not be known upfront how many objects will need
565 to be sorted.

566 When such details are unknown, it may still be possible to fall back to abstract inputs sampled from
567 some sensible generic random distribution—so long as the processor network is trained in a way
568 that promotes extrapolation. For a recent theoretical treatment of OOD generalisation in algorithmic
569 reasoners, we refer the reader to Xu et al. [56]. Further, as observed by Deac et al. [8] in the case of
570 implicit planning, even generic random abstract distributions can promote useful transfer to noisy
571 pixel-based acting settings such as Atari.

572 Lastly, it is important to ensure that, in the abstract pipeline, the processor network P carries the
573 brunt of the computational burden; otherwise, the reasoning task may be partially captured in either
574 f or g , and not left to P ’s weights. In our work, we generally promote such behaviour by keeping
575 f and g to be only linear layers; but even in settings where this is not appropriate, we recommend
576 taking care to not overpower the abstract encoder and decoder.

577 **F Atari abstract model equations**

578 In this appendix, we provide a ‘bird’s eye’ view of the modelling steps taken by our abstract Atari
579 pipeline, aiming to support future implementations of the RMR blueprint. We will readily re-use the
580 notation from the main text for the various components.

581 Firstly, the abstract encoder f is applied on the relevant Atari RAM representations, $\bar{\mathbf{x}}$, augmented by
582 a one-hot representation which is aiming to provide a generic embedding of the semantics of each
583 RAM slot. f is implemented as a linear layer, hence:

$$\mathbf{z} = f(\bar{\mathbf{x}}\|\mathbf{o}) = \mathbf{W}^f \bar{\mathbf{x}} + \mathbf{U}^f \mathbf{o} + \mathbf{b}^f \quad (1)$$

584 where \mathbf{W}^f , \mathbf{U}^f , \mathbf{b}^f are learnable weights, and $\mathbf{o} \in \mathbb{B}^{128 \times 128}$ is a one-hot encoding s.t. $\mathbf{o} = \mathbf{I}_{128}$.

585 In a separate pipeline (officially part of the processor network), the performed actions are also encoded
586 using a linear action encoder:

$$\boldsymbol{\alpha} = f_a(\mathbf{a}) = \mathbf{W}^{f_a} \mathbf{a} + \mathbf{b}^{f_a} \quad (2)$$

587 where \mathbf{a} is a one-hot encoded action representation.

588 Then, the processor network GNN is called to update these latents, and a sum-based skip connection
589 is employed to promote a model-free path:

$$\mathbf{z}' = P(\mathbf{z} + \boldsymbol{\alpha}) + \mathbf{z} + \boldsymbol{\alpha} \quad (3)$$

590 Here, the processor network P is implemented as a standard message passing neural network [51]
591 over a fully connected graph. Equations of such a P are commonly exposed in e.g. Veličković et al.
592 [13].

593 Lastly, the relevant decoder networks predict two properties: a mask which signifies which RAM
594 slots have been changed as a result of applying \mathbf{a} , and the updated states $\bar{\mathbf{y}}$.

$$\boldsymbol{\mu} = g_\mu(\mathbf{z}') = \sigma(\mathbf{W}^\mu \mathbf{z}' + \mathbf{b}^\mu) \quad (4)$$

$$\bar{\mathbf{y}} = g(\mathbf{z}') \odot \mathbb{I}_{\boldsymbol{\mu} > 0.5} = (\mathbf{W}^g \mathbf{z}' + \mathbf{b}^g) \odot \mathbb{I}_{\boldsymbol{\mu} > 0.5} \quad (5)$$

595 Here, σ is the logistic sigmoid activation, \odot is the elementwise product, and \mathbb{I} is the indicator function,
596 which thresholds the mask. While this thresholding is non-differentiable, we note that the mask
597 values can be directly supervised from known trajectories, and at training time, teacher forcing can
598 be applied—slotting ground-truth masks in place of the indicator function.