

STRESZCZENIE

Słowa kluczowe:

Dziedziny nauki i techniki zgodne z wymogami OECD:

ABSTRACT

Keywords:

OECD field of science and technology:

CONTENTS

List of abbreviations and nomenclature	7
1. Introduction	8
2. Theoretical background	10
2.1. Partially Observable Markov Decision Processes	10
2.2. Reinforcement Learning	11
2.3. Deep Learning	11
2.4. Variational Inference	14
3. Related work	15
3.1. Model learning	15
3.1.1. World Models	15
4. Planning with learned model	18
4.1. HumbleRL framework	18
4.1.1. Architecture	18
4.2. World Models with the AlphaZero planner	20
4.2.1. Data collection	21
4.2.2. Preprocessing	21
4.2.3. World Models architecture: Vision and Memory	21
4.2.4. AlphaZero architecture: Controller	22
4.3. PlaNet with the CEM planner	23
4.3.1. High-level idea and argument “why?!”	24
4.3.2. Data collection i.e. an iterative approach	24
4.3.3. Preprocessing	24
4.3.4. RSSM architecture	24
4.3.5. CEM planner	25
5. Experiments	26
5.1. Benchmarks	26
5.1.1. Arcade Learning Environment	26
5.1.2. Sokoban	28
5.2. Hardware	29
5.3. World Models for Sokoban	29
5.3.1. Train the unchanged World Models implementation in the Sokoban environment	29
5.3.2. Train World Models in Sokoban environment on 10x10 grid world states	31
5.3.3. Train World Model in Sokoban with auxiliary tasks	32
5.4. World Models for Atari	33
5.4.1. Train unchanged World Models in the Boxing environment	33
5.4.2. Train World Models with AlphaZero planner in the Boxing environment	35
5.5. PlaNet for Sokoban	36

5.5.1. Train unchanged PlaNet in the Sokoban environment.....	36
5.6. PlaNet for Atari	36
5.6.1. Train unchanged PlaNet in the Boxing environment.....	36
5.6.2. Train PlaNet in the Boxing environment with the increased action repeat.....	37
5.6.3. Train PlaNet in the Boxing environment with the lowered decoder variance	37
5.6.4. Train PlaNet in the Boxing environment with increased free nats.....	38
5.6.5. Train tuned PlaNet in the Freeway environment	38
5.6.6. Train tuned PlaNet in the Freeway environment with a longer planning horizon	39
6. Conclusion	40
References	41
List of figures	43
List of tables	44

LIST OF ABBREVIATIONS AND NOMENCLATURE

1. INTRODUCTION

[What is RL and model-free RL...]* Reinforcement learning, a subfield of artificial intelligence (AI), formalise rather the most obvious and common among animals learning strategy. It is learning how to achieve predefined goals through interaction with an environment [32]. Progress has been made in developing capable agents for numerous domains using deep neural networks in conjunction with model-free reinforcement learning [14][22][28], where raw observations directly map to agent's actions. However, current state-of-the-art approaches are very sample inefficient, they sometimes require tens or even hundreds of millions of interactions with the environment [21], and lack the behavioural flexibility of human intelligence, hence the resulting policies poorly generalize to novel tasks in the same environment.

[Model-based RL with its benefits...]* The other branch of reinforcement learning algorithms, called model-based reinforcement learning, aims to address these shortcomings by endowing agents with a model of the world. There are many ways of using the model: one can use the model for data augmentation for model-free methods [9], some methods use the model as the imagined environment to learn model-free policy in it [12], other methods focus on simulation-based search using the model [30] and there are even methods that integrate model-free and model-based approaches [35]. The model allows the agent to simulate an outcome of an action taken in a given state. The main upside of this is that it allows the agent to plan by thinking ahead, seeing what would happen for a range of possible choices, and explicitly deciding between possible options without the risk of the adverse consequences of trial-and-error in the real environment - including making poor, irreversible decision. Agents can then distill the results from planning ahead into a policy. Even if the model needs to be synthesized from past real experience first it can exploit additional unsupervised learning signals, like rewards function modeling or future observations prediction [17], thus it results in a substantial improvement in sample efficiency over model-free methods. Furthermore, the same model can be used by the agent to complete other tasks in the same environment [35]. It gives AI hint of human intelligence flexibility and versatility.

[Planning vs. learning differences/similarities...]* Learning is different from planning. In the former the agent samples episodes of real experience through interaction with an environment and updates its policy or states' value estimates based on them. In the latter the agent also updates its policy or states' value estimates, but this time based on simulated experience gained through evaluation of a model. It is worth noting the symmetry which yields one important implication: algorithms for reinforcement learning can also become algorithms for planning, simply by substituting simulated experience in place of real experience. Moreover, planning carries the promise of increasing performance just by increasing the computational budget for searching for good actions [31]. Model-free methods to improve their performance need more interactions with a real environment and hence scale in amount of data not amount of computation, which very often is much cheaper than collecting more data.

[Really briefly about what is problem of learning world dynamics of POMDP...]* Model-free

methods are more popular and have been more extensively developed and tested than model-based methods. While model-free methods forego the potential gains in sample efficiency from using a model, they tend to be easier to implement and tune. The main downside of model-based reinforcement learning is that a ground-truth model of the environment is usually not available to an agent. If the agent wants to use a model in this case, it has to learn the model from experience, which creates several challenges. One of them is bias in the model that can be exploited by the agent [12], resulting in an agent which performs well with respect to the learned model, but behaves sub-optimally in the real environment. Different challenge also comes from fundamental downside of function approximation, it result in models that are inherently imperfect. The performance of agents employing common planning methods usually suffer from seemingly minor model errors [33]. Those errors compound during planning, causing more and more inaccurate predictions the further horizon of a plan. [34]

[Long-term ambitious goal...]* There are many real-world problems that could benefit from

application of general planning AI system. Company called DeepMind, driven by their experience from creating winning Go search algorithm AlphaZero [30], published AlphaFold [8], a system that predicts protein structure. The 3D models of proteins that AlphaFold generates are far more accurate than any that have come before making significant progress on one of the core challenges in biology. Real-world applications of AI algorithms like this are often limited by the problem of sample inefficiency. In a setting with e.g. a physical robot the AI agent can not afford much trial-and-error behaviour, that could cause damage to the robot, and do this over hundreds of millions of time steps, for each task separately, in order to build a sufficiently large training dataset. Those machines work in the real world, not accelerated computer simulation, and often need a human assistance. To apply sample efficient model-based systems, that can generalize their knowledge, accurate learned models and robust planning algorithms are needed.

[Explain your topic...]* The aim of this work is to derive from previous work on model-

learning in complex high-dimensional decision making problems [4][20][3][13] and apply them to plan in Atari 2600 games, a platform for evaluating general competency in artificial intelligence [21]. Those methods proved to train accurate models, at least in short horizon, and should open a path for application of simulation-based search planning algorithms like TD-Search [29] or AlphaZero [30] to complex planning problems in environments with discrete state-action spaces and without access to a ground-truth model. This should allow for improvement in data efficiency and generalization *[We don't test generalization, maybe it should be omitted then?]* without loss in performance. This work focuses on three benchmarks: an arcade game with dense rewards Boxing, a challenging environment with sparse rewards Freeway and a complex puzzle environment MsPacman.

2. THEORETICAL BACKGROUND

2.1. *Partially Observable Markov Decision Processes*

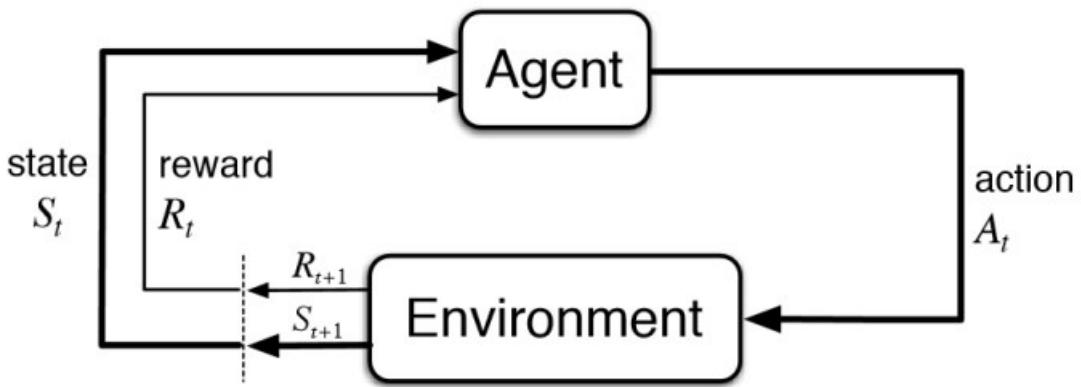
Markov Decision Processes (MDPs) are a classical formalization of sequential decision making, where actions influence not just immediate rewards, but also subsequent situations (states) and through those future rewards. Therefore, MDPs involve delayed rewards and the need to trade-off these with an immediate reward. Partially observable Markov Decision Processes (POMDPs), which describe a more general class of problems, have one major difference, the full state of an environment is unknown. The environment is perceived through observations that provide only partial information about the state. A good example are Atari games. Individual frames often doesn't provide full information about the game's state which is held in the game's RAM.

In this work following definition of POMDP is used: it consists of a set of hidden states S , a set of observations O and a set of actions A . The dynamics of the MDP, from any state $s \in S$ and for any action $a \in A$, are determined by transition function, $P_{ss'}^a = p(S_{t+1} = s' | S_t = s, A_t = a)$, specifying the distribution over the next state $s' \in S$. A reward function, $R_{ss'}^a = p(R_{t+1} | S_t = s, A_t = a, S_{t+1} = s')$, specifies the distribution over rewards for a given state transition. Finally, as mentioned earlier, POMDP is perceived through partial observations specified via probability distribution $P_s = p(O_t | S_t = s)$. A fixed initial state s_0 is assumed. In episodic POMDPs, which this work considers, an environment terminates with probability 1 in one of distinguished terminal states, $s_T \in S$, after finite number of transitions T . A return $G_t = \sum_{k=t+1}^T \gamma^{k-t-1} r_k$ is the total reward accumulated in that episode from time t until reaching the terminal state at time T . $0 \leq \gamma \leq 1$ is a discount factor that trade-offs short-term rewards with long-term rewards. A policy, $\pi(s, a) = p(A_T = a | S_t = s)$, maps a state s to a probability distribution over actions. A value function, $V_\pi(s) = \mathbb{E}_\pi[G_t | s_t = s]$, is the expected return from state s when following policy π where the expectation is over the distributions of the environment and the policy. An action value function, $Q_\pi(s, a) = \mathbb{E}_\pi[G_t | s_t = s, a_t = a]$, is the expected return after selecting action a in state s and then following policy π where, again, the expectation is over the distributions of the environment and the policy. An optimal value function is the unique value function that maximises the value of every state, $V^*(s) = \max_\pi V_\pi(s), \forall s \in S$ and $Q^*(s, a) = \max_\pi Q_\pi(s, a), \forall s \in S, a \in A$. An optimal policy $\pi^*(s, a)$ is a policy that maximises the action value function for every state in the POMDP, $\pi^*(s, a) = \underset{a}{\operatorname{argmax}} Q^*(s, a)$.

POMDPs are a mathematically idealized form of the reinforcement learning problem for which precise theoretical statements can be made. In reinforcement learning the dynamics, the observations distribution and the reward function are hidden behind an environment. Consequently, we can not directly use them for planning, but we can learn them through interaction with the environment.

2.2. Reinforcement Learning

Reinforcement learning (RL) is learning what to do, how to map situations to actions, so as to maximize a numerical reward signal. [32] This mapping is called a policy π . RL consists of an agent that, in order to learn a good policy, acts in an environment. The environment provides a response to each agent's action a that is interpreted and fed back to the agent. Reward r is used as a reinforcing signal and state s is used to condition agent's decisions. Fig. 21 explains it in the diagram. Each action-response-interpretation sequence is called a step or a transition. Multiple steps form an episode. The episode finishes in a terminal state s_T and the environment is reset in order to start the next episode from scratch. Very often, RL agents need dozens and dozens of episodes to gather enough experience to learn the (near) optimal policy.



Rys. 21. Reinforcement Learning [32]

[TODO: Describe here General Policy Iteration, Monte-Carlo Control, TD-Learning, ... what else? Rather than guessing what needs to be described, continue work and see what needs more attention.]

2.3. Deep Learning

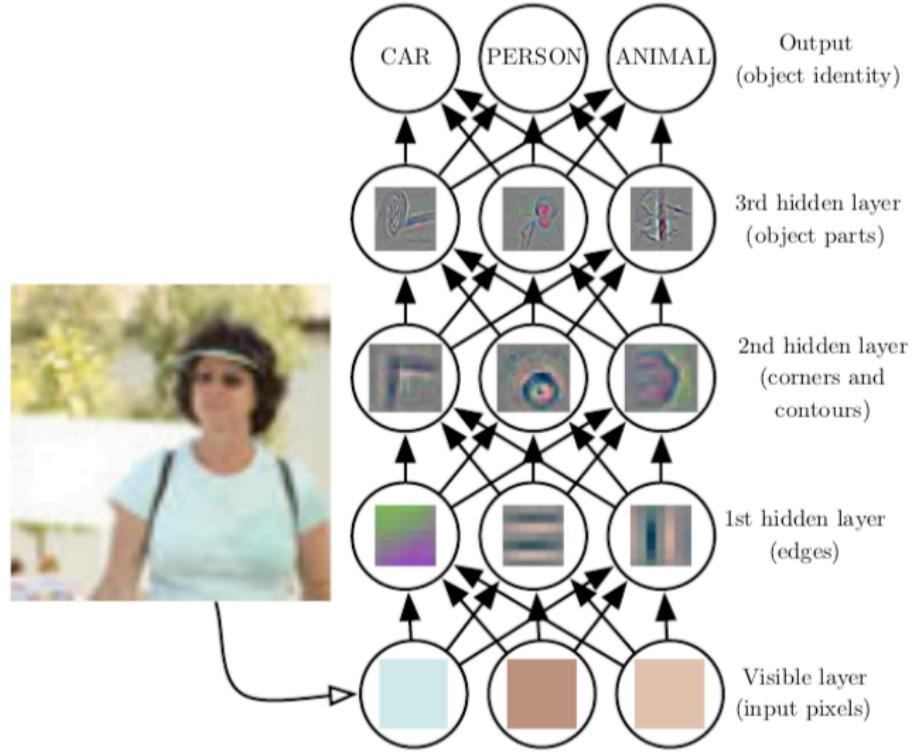
Machine learning gives AI systems the ability to acquire their own knowledge, by extracting patterns from raw data. It stands in opposition to classical computer programs which execute explicit instructions hand-coded by a programmer. One example of machine learning algorithm is logistic regression. It can determine whether to recommend cesarean delivery or not [23]. Another widely used machine learning algorithm called naive Bayes can distinguish between legitimate and spam e-mail.

The performance of these machine learning algorithms depends heavily on the representation of the problem they are given. For example, when logistic regression is used to recommend cesarean delivery, the AI system does not examine the patient's MRI scan directly. It would not be able to make useful predictions as individual pixels in an MRI scan have negligible correlation with any complications that might occur during delivery. It, instead, gets several pieces of relevant information, such as the presence or absence of a uterine scar, from the doctor. Each piece of

information included in the representation of the data is known as a feature. Logistic regression learns the relation between those features and various outcomes, such as a recommendation of cesarean delivery. The algorithm does not influence the way that the features are defined in any way.

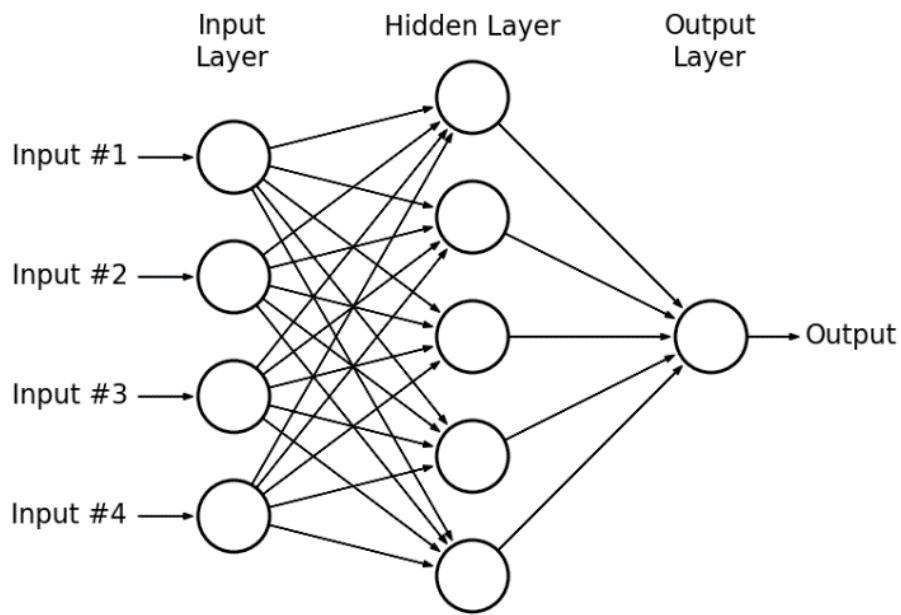
Sometimes it can be hard to hand-craft a good problem's representation. For example, suppose that we would like to write a program to detect cats in photographs. We know that cats are furry and have whiskers, so we might like to use the presence of a fur and whiskers as features. Unfortunately, it is difficult to describe exactly what a fur or a whisker looks like in terms of pixel values. This gets even more complicated when we take into account e.g. shadows falling on the cat or an object in the foreground obscuring part of the animal. One solution to this problem is to use machine learning to discover not only the mapping from representation to output but also the representation itself. This approach is known as representation learning. Learned representations often result in much better performance of machine learning algorithms than can be obtained with hand-designed representations. They also allow AI systems to rapidly adapt to new tasks with minimal human intervention.

Deep learning is a particular kind of machine learning that achieves great power and flexibility by learning to represent the world as a nested hierarchy of concepts, with each concept defined in relation to simpler concepts. Deep learning solves representation learning problem by introducing representations that are expressed in terms of other, simpler representations. Fig. 22 shows how a deep learning system can represent the concept of an image of a person by combining simpler concepts, such as corners and contours, which are in turn defined in terms of edges.



Rys. 22. Deep Learning [10]

The fundamental example of a deep learning model is a feedforward deep network or multilayer perceptron (MLP). A multilayer perceptron is just a mathematical function mapping some set of input values to output values. The function is formed by composing many simpler functions, called perceptrons, each providing a new representation of its input to next functions. Fig. 23 shows example MLP and dependencies between perceptrons. The input is presented at the input layer. Then a hidden layer (or series of them) extracts increasingly abstract features from the image. These layers are called “hidden” because their values are not given in the data. Instead, the model must determine which concepts are useful for explaining the relationships in the observed data. Finally, this description of the input in terms of the features can be used to produce the output at the output layer.



Rys. 23. Multilayer perceptron

[TODO: One word about RNNs, Conv. and other used layers here.] [TODO: Also write here about backprop.]

2.4. Variational Inference

[TODO: Describe VAEs in structured POMDPs: what are zero step, one step and open loop predictions, what is latent state/code, etc.]

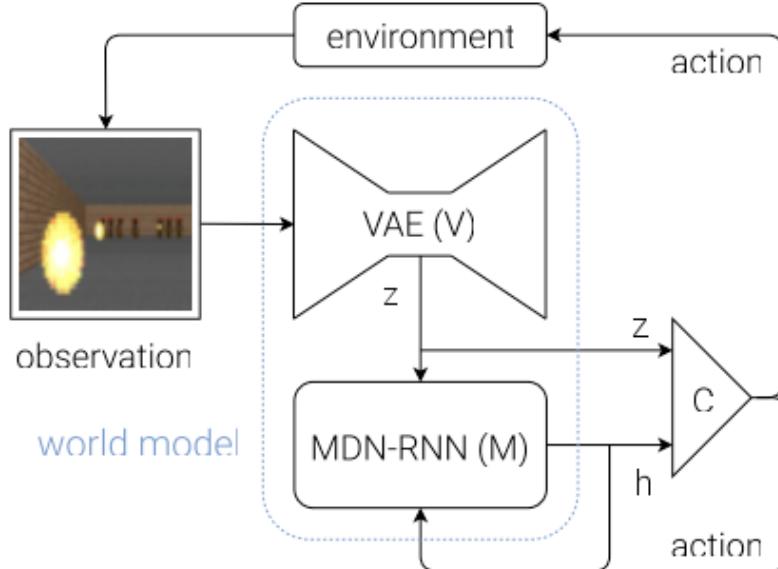
3. RELATED WORK

3.1. Model learning

3.1.1. World Models

In World Models[12] paper, the authors explore the idea of using large and highly expressive neural networks, that can learn rich spatial and temporal representation of data, and applying them to reinforcement learning. The RL algorithm is often bottlenecked by the credit assignment problem, which makes it hard for traditional RL algorithms to learn millions of weights of a large model. To accomplish their goal, they decompose the problem of an agent training into two stages: they first train a generative neural network to learn a model of the agent's world in an unsupervised manner. Thereafter, by using a compressed spatial and temporal representation of the environment extracted from the world model as inputs to the agent, they train a linear model to learn to perform a task in the environment. The small linear model lets the training algorithm focus on the credit assignment problem on a small search space, while not sacrificing capacity and expressiveness via the larger world model.

Their solution consists of three components: Vision (V) for encoding the spatial information, Memory (M) for encoding the temporal information and Controller (C) which represents the agent's policy. Fig. 31 depicts a flow diagram of the agent's model.



Rys. 31. Flow diagram of the agent's model[12]. The raw observation is first processed by V at each time step t to produce z_t . The input into C is this latent vector z_t concatenated with M's hidden state h_t at each time step. C will then output an action vector a_t and will affect the environment. M will then take the current z_t and action a_t as an input to update its own hidden state to produce h_{t+1} to be used at time $t + 1$.

The environment provides the agent with high dimensional visual observation at each time step. The essential task of the Vision model is to encode this high dimensional observation into a low dimensional latent state. To do this, Vision is implemented as Variational Autoencoder[19]. It is trained in an unsupervised manner on randomly generated experience from the environment.

Since many complex environments are partially observable, the visual observation at each time step, and hence the latent state, doesn't include full information about the current situation in the environment. To acquire full knowledge, the agent needs to encode what happens over time. This is the role of the Memory model. It is implemented as Recurrent Neural Network[16] (RNN) and trained on the same data as Vision to predict the future latent state that Vision is expected to produce. Because many environments are stochastic in nature, the RNN is trained to output a probability density of the next latent state approximated as a mixture of Gaussian distribution - in literature, this approach is known as Mixture Density Network combined with a RNN (MDN-RNN)[11]. More specifically, the RNN will model $P(z_{t+1}|a_t, z_t, h_t)$, where z_{t+1} is the output next latent state, a_t is the action taken at time t , z_t is the latent state of the current time step t and h_t is the hidden state of the RNN that encodes past information made available to the agent from the beginning of the episode until the time step t .

The Controller model represents the agent's policy. It is responsible for determining course of actions to take in order to solve a given task. Controller is a simple linear model that maps the concatenated latent state z_t and hidden state h_t at the time step t directly to the action a_t at that time step: $a_t = W[z_t h_t] + b$, where W and b are the weight matrix and bias vector of that model. The authors deliberately made Controller as simple as possible, and trained it separately from Vision and Memory, so that most of the agent's complexity resides in the world model (V and M). The latter can take advantage of current advances in deep learning that provide tools to train large models efficiently when well-behaved and differentiable loss function can be defined. Shift in the agent's complexity towards the world model allows the Controller model to stay small and focus its training on tackling the credit assignment problem in challenging RL tasks. It is trained using evolution strategy, which is rather an unconventional choice that only currently have been considered as a viable alternative to popular RL techniques[26].

Their solution was able to solve an OpenAI Gym's CarRacing environment, which is the continuous-control, top-down racing task. It is the first known solution to achieve the score required to solve this task. In the process, the Memory model have learned to simulate the original environment of CarRacing, that is simulated by Memory. In the second experiment, they show that the agent is able to learn from imagined experience, produced by its Memory, and successfully transfer this policy back to the actual environment of VizDoom (see fig. 32). This result indicates that the world model is able to model complex environments from visual observations and it can be used for planning. Therefore, it may prove useful for the topic of this thesis. *[Should we expand on that in here? Or rather place for that is in "Planning with learned model" chapter where I describe design of my solution?]*



Rys. 32. VizDoom: the agent must learn to avoid fireballs shot by monsters from the other side of the room with the sole intent of killing the agent[12].

4. PLANNING WITH LEARNED MODEL

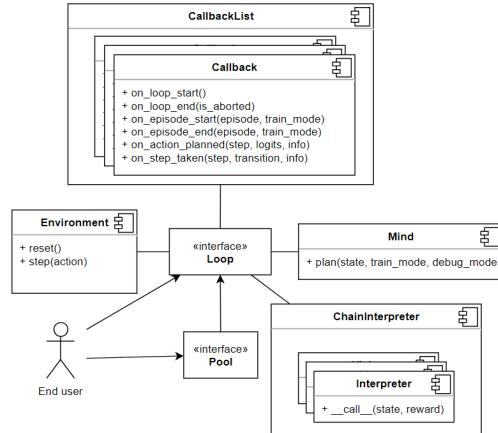
In this section two architectures are described. Both involve a similar model learning approach, but differ substantially in technical details and planning algorithms. Nevertheless, the goal stays the same: train a sufficient environment model, or such that accurately predicts future latent states of an environment to predefined cut-off point in time, and use it to plan and solve the environment. The more accurate the model to the cut-off point, which might be environment dependent, the better model learning algorithm. Before all of that, the code architecture and the framework that was created to accelerate this research are described.

4.1. *HumbleRL framework*

Reinforcement Learning scientists tend to write the entire code from scratch by themselves, instead of using existing RL frameworks. This is justified by the fact, that the commonly available frameworks are not flexible enough for intended experiments or require a specific backend like TensorFlow, which might be disfavored. HumbleRL [18] was created with this problem in mind. Its simple API allows to perform a variety of RL experiments without any restrictions on the algorithms used. Since the backend is not tied to any specific technology, it is possible to mix different neural network frameworks or not use any at all. HumbleRL provides the boilerplate code of RL loop in fig.21 and determines the common interface, the rest is done by the user.

4.1.1. *Architecture*

Framework architecture is depicted in fig. 41. An agent is represented by the Mind class. Mind encapsulates action planning logic and provides it via the plan method. In order to learn, the agent acts in the world represented by the Environment class. The Environment class provides methods for resetting, taking steps, rendering and getting information about the world. The agent is not usually presented with raw environment observations. Instead, it looks at states preprocessed by the Interpreter. Different interpreters can be joined together with the ChainInterpreter class. It acts as a preprocessing pipeline, with each subsequent interpreter using the output of a previous one as an input.

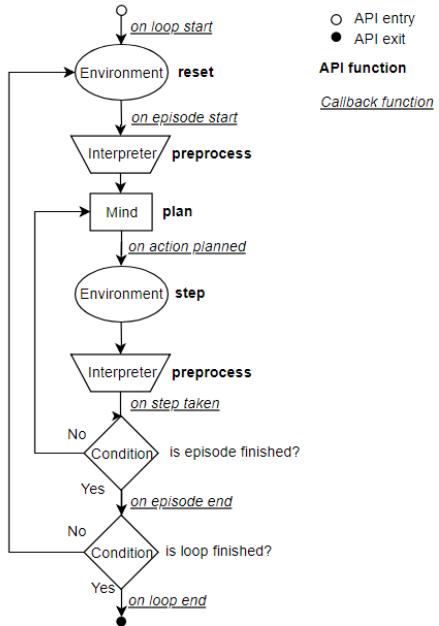


Rys. 41. HumbleRL architecture

Framework user does not need to call all of those methods directly, those are utilized by the loop function. This function gets an action from the Mind, executes it in the Environment and then next observation is preprocessed with the Interpreter in preparation for the next step. To extend basic loop functionality, user can define callbacks that implement the Callback interface. Callbacks can react to events:

- at the beginning and ending of the loop,
- at the beginning and ending of each episode,
- after action is planned by the Mind,
- after step is taken in the Environment.

Callbacks are accumulated in the CallbackList. The entire loop function logic is shown in fig. 42. Parallel version of loop function is available as the pool function. It uses predefined number of workers to execute a pool of Minds in their own Environments in parallel.



Rys. 42. HumbleRL loop function overview

World Models with the AlphaZero planner uses this framework.

4.2. *World Models with the AlphaZero planner*

1. High-level idea and argument “why?!”:

World Models shown it can plan (offline planning, training a policy) in imagination and AlphaZero is incredibly powerful search algorithm.

2. Data collection i.e. a random agent.

3. Preprocessing.

4. World Models architecture: Vision and Memory.

(a) Which does what (brief reminder) e.g. Vision encode a current observation and Memory encodes history and predicts future.

(b) How Memory is used to predict future: MDN and LMH.

(c) Training procedure of each module i.e. in separation exactly like described in the related work (don't describe loss etc. only high-level).

5. AlphaZero architecture: Controller.

(a) What it does i.e. uses the memory module to plan in imagination.

(b) How next action is planned i.e. MCTS.

(c) Training procedure of Value and Policy networks i.e. policy iteration after Vision and Memory are already trained (don't go into AZ details, it's already in the related work, here how you modify it or apply to your case).

World Models' agent, as shown in the paper [12], is able to learn from simulated experience. It is an example of successful planning using learned model. This section describes attempt to utilize the world model part of the agent in the AlphaZero search algorithm. Moreover, the Vision module encodes environment observations into low level representation. The latent state of the world models encodes abstract information about the environment and allow the planning controller to work fast, because there is no need to generate high-dimensional observations, only low-dimensional latent states are generated by the model and processed by the controller.

4.2.1. *Data collection*

To train Vision and Memory modules first collection of 10,000 random rollouts of the environment are gathered to create a dataset. An agent is acting randomly to explore the environment multiple times and records of the random actions taken and the resulting observations from the environment are saved. This dataset is used to train the Vision module to learn a latent state of each frame observed. Next, it is used to preprocess each frame into its latent state to prepare a dataset for the Memory module training.

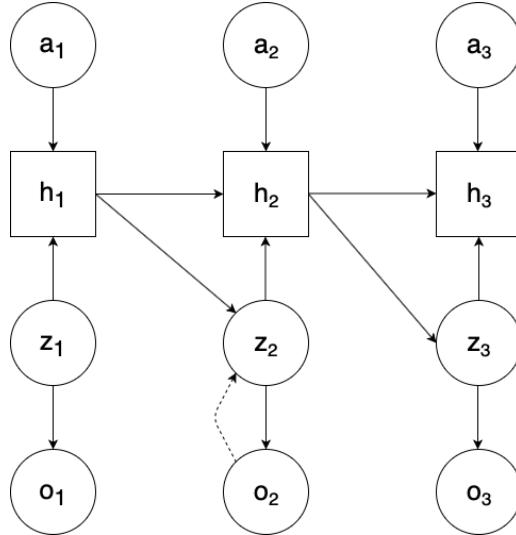
4.2.2. *Preprocessing*

Each frame, before it is used for any training, is central cropped if a frame from an environment includes e.g. some border. This operation depends on a specific environment. It is then resized to 64 x 64 pixels for all environments. All three colour channels are preserved. Actions get one-hot encoded.

4.2.3. *World Models architecture: Vision and Memory*

HumbleRL was used in a few stages of training a model. First of all, OpenAI Gym was used behind the Environment interface from the framework. Loop together with an agent performing random actions (Mind) and a callback was used to gathers transitions and save them to an external storage. The framework allows to focus strictly on collecting trajectories and not worry about agent-environment interactions. Transitions are used to train the Vision and Memory components. For Vision, which function is to encode observations into latent state-space, Keras [5] framework was used and for Memory, which function is to encode history, PyTorch [25] was used, since it was easier to use for this case than Keras. HumbleRL is not constricted to work with any particular deep learning library, so it's not a problem to mix the solutions, as long as trained models are wrapped in proper interfaces. As described in the related work chapter, the Vision module was trained to encode each frame into low dimensional latent vector by minimizing the difference between a given frame and the reconstructed version of the frame produced by the decoder. After data processing using the trained Vision model, the Memory module can be trained to model environment's dynamics, predicting a future latent state from history of previous latent states, as a mixture of Gaussians. Figure 43 depicts how the world model works. Solid arrows describe the predictive model and dotted arrow describes the inference model *[NOTE: It needs to be explained]*

[in Theoretical Background what is the difference.] This model uses stochastic nodes (circles) to model uncertainty in the environment, but also create more robust environment's representation. Uncertainty can originate not only from fundamental stochastic nature of the environment, but also partial observability. Squares depict deterministic nodes. The Memory module model is implemented as the recurrent neural network.



Rys. 43. Stochastic World Models probabilistic graphical model diagram

[TODO: Describe this diagram in probabilistic notation (see note).]

Because benchmarks include deterministic and fully observable case, World Models was tested also with the Memory module without Mixture Density Network on top of the recurrent neural network. Instead a linear model was used to output the next latent state.

4.2.4. AlphaZero architecture: Controller

Once the world model is ready it is used to train the final piece - the planner, which uses the AlphaZero algorithm. The Vision module is used as the Interpreter which encodes incoming observations into latent space. The Memory module gets wrapped in the MDP interface, it is used for AlphaZero simulations. The Mind interface is implemented by AlphaZero algorithm, it implements described in the related work search strategy. It uses a neural network to evaluate each node and guide search direction on each edge. Returned actions' scores are visit counts from the root node, which are then used to choose an action by a policy. Details are the same as in original AlphaZero described in the related work. Pseudo-code of the planning algorithm in Python is shown below.

```

1 def plan(self, state):
2     # Get/create root node
3     root = self.query_tree(state)
4
5     # Perform simulations
6     simulations = 0
7     start_time = time()
8     while time() < start_time + self.timeout and simulations < self.max_simulations:
9         # Simulate
10        simulations += 1
11        leaf, path = self.simulate(root)
  
```

```

12     # Expand and evaluate
13     value = self.evaluate(leaf)
14
15     # Backup value
16     self.backup(path, value)
17
18     # Get actions' visit counts
19     actions = np.zeros(self.model.action_space.num)
20     for action, edge in root.edges.items():
21         actions[action] = edge.num_visits
22
23     return actions

```

Agent's experience and score statistics used for training are gathered using callbacks during self-play phase. The neural network training phase takes place after a predefined number of self-play games. The training phase is performed using the Keras [5] framework. Next, the self-play phase takes place once again and the two further interchange. The self-play phase uses the loop function to effortlessly run AlphaZero for given number of games (episodes).

4.3. PlaNet with the CEM planner

1. High-level idea and argument “why?!”:

PlaNet shown it can successfully plan (online planning, evolutional strategy) in imagination for complex continuous control task with iterative data collection and in a clean algorithm.

2. Data collection i.e. iterative approach.

3. Preprocessing.

4. RSSM architecture.

(a) What it does (brief reminder)? It predicts future, observations and rewards, where the latter is more important for planner which uses the model to evaluate plans.

(b) How it's used. Paper details are already described in the related work (like overshooting etc.), here write how you use it e.g. you encode actions and RSSM is used to sample future latent states that are then used to predict rewards and it receives current observation to update its belief state etc.

(c) Training procedure i.e. interchanged training, test and collection phases (like in paper, nothing changed).

5. CEM planner.

(a) What it does i.e. it's like optimization of actions scores based on model evaluations.

(b) How next action is planned i.e. argmax.

(c) Planning procedure i.e. evolutional strategy.

4.3.1. High-level idea and argument “why?!”

They show working example on continuous control tasks of online planning (wheres World Models was offline planning). Recurrent state space model: We design a latent dynamics model with both deterministic and stochastic components (Buesing et al., 2018; Chung et al., 2015). Our experiments indicate having both components to be crucial for high planning performance. Latent overshooting: We generalize the standard variational bound to include multi-step predictions. Using only terms in latent space results in a fast regularizer that can improve long-term predictions and is compatible with any latent sequence model.

4.3.2. Data collection i.e. an iterative approach

Since the agent may not initially visit all parts of the environment, we need to iteratively collect new experience and refine the dynamics model. We do so by planning with the partially trained model, as shown in Algorithm 1. Starting from a small amount of S seed episodes collected under random actions, we train the model and add one additional episode to the data set every C update steps. When collecting episodes for the data set, we add small Gaussian exploration noise to the action. To reduce the planning horizon and provide a clearer learning signal to the model, we repeat each action R times, as common in reinforcement learning (Mnih et al., 2015; 2016).

4.3.3. Preprocessing

1. action repeat 4
2. action one-hot
3. minimum duration 50
4. maximum duration 2000/4
5. resize (64 x 64 x 3)
6. cast to uint8
7. 5 bit quantisation + random noise data augmentation
8. Cast to from -0.5 to 0.5.

4.3.4. RSSM architecture

It uses code from official implementation.

[*TODO: Recreate prob. model diagram from your WM vs. PlaNet comparison notes in Draw.io.]* [*TODO: Describe this diagram in probabilistic notation (see note).*]

[...] we name recurrent state-space model (RSSM), where $f(ht-1, st-1, at-1)$, deterministic state model, is implemented as a recurrent neural network (RNN). Intuitively, we can understand this model as splitting the state into a stochastic part st and a deterministic part ht , which depend

on the stochastic and deterministic parts at the previous time step through the RNN. We use the encoder $q(s_1:T | o_1:T, a_1:T) = \prod_{t=1}^T q(s_t | h_t, o_t)$ to parameterize the approximate state posteriors. Importantly, all information about the observations must pass through the sampling step of the encoder to avoid a deterministic shortcut from inputs to reconstructions.

4.3.5. CEM planner

[TODO: Cut it down and focus on your modification with cast to one-hot max action.] We use the cross entropy method (CEM; Rubinstein, 1997; Chua et al., 2018) to search for the best action sequence under the model, as outlined in Algorithm 2. We decided on this algorithm because of its robustness and because it solved all considered tasks when given the true dynamics for planning. CEM is a population-based optimization algorithm that infers a distribution over action sequences that maximize the objective. As detailed in Algorithm 2 in the appendix, we initialize a time-dependent diagonal Gaussian belief over optimal action sequences at $t:H \sim \text{Normal}(\mu_t:t+H, \sigma^2 I_{t:t+H})$, where t is the current time step of the agent and H is the length of the planning horizon. Starting from zero mean and unit variance, we repeatedly sample J candidate action sequences, evaluate them under the model, and re-fit the belief to the top K action sequences. After I iterations, the planner returns the mean of the belief for the current time step, μ_t . Importantly, after receiving the next observation, the belief over action sequences starts from zero mean and unit variance again to avoid local optima. To evaluate a candidate action sequence under the learned model, we sample a state trajectory starting from the current state belief, and sum the mean rewards predicted along the sequence. Since we use a population-based optimizer, we found it sufficient to consider a single trajectory per action sequence and thus focus the computational budget on evaluating a larger number of different sequences. Because the reward is modeled as a function of the latent state, the planner can operate purely in latent space without generating images, which allows for fast evaluation of large batches of action sequences. The next section introduces the latent dynamics model that the planner uses.

[TODO: Add pseudo-code of planning procedure.]

5. EXPERIMENTS

5.1. Benchmarks

5.1.1. Arcade Learning Environment

The Arcade Learning Environment (ALE) has became a platform for evaluating artificial intelligence agents. Originally proposed by Bellemare et. al. [1], the ALE makes available dozens of Atari 2600 games for an agent training and evaluation. The agent is expected to do well in as many games as possible without game-specific information, generally perceiving the world through a video stream. Atari 2600 games are excellent environments for evaluating AI agents for three main reasons: they are varied enough to provide multiple different tasks, requiring general competence, they are interesting and challenging for humans and they are free of experimenter's bias, having been developed by an independent party.

In the context of the ALE, a discrete action is a number in range from 0 to 17 inclusive which encodes the composition of a joystick direction and an optional button press. The agent observes a reward signal, which is typically the change in the player's score (the difference in score between the previous time step and the current time step), and an observation $o_t \in O$ of the environment. This observation can take form of a single 210×160 image and/or the current 1024-bit RAM state. Because a single image typically does not satisfy the Markov property the ALE is formalised as POMDP. Observations and the environment state are distinguished, with the RAM data being the real state of the emulator. A frame (as a unit of time) corresponds to 1/60th of a second, the time interval between two consecutive images rendered to the television screen. The ALE is deterministic, which means that given a particular emulator state s and a action a there is a unique next state s' , that is, $P_{ss'}^a = p(s'|s, a) = 1$.

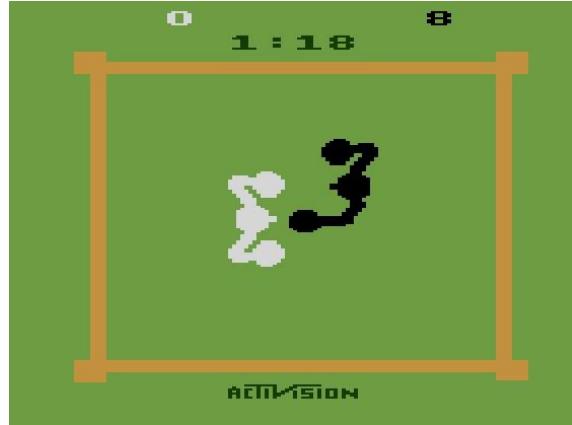
Agents interact with the ALE in an episodic fashion. An episode begins by resetting the environment to its initial configuration, s_0 , and ends at a given endpoint depending on a game. The primary measure of an agent's performance is the score achieved during an episode, namely the undiscounted sum of rewards for that episode. While this performance measure is quite natural, it is important to realize that score is not necessarily an indicator of AI progress. In some games, agents can exploit the game's mechanics to maximize sum of rewards, but not complete the game's goal in human's understanding. [6]

Preprocessing include frame skipping [24] which restricts the agent's decision points by repeating a selected action for 4 consecutive frames. Frame skipping results in a simpler reinforcement learning problem and speeds up execution. *[TODO: Describe other preprocessing techniques used here.]*

This work uses ALE through OpenAI Gym API [2], specifically two games are used as benchmarks: Boxing and Freeway.

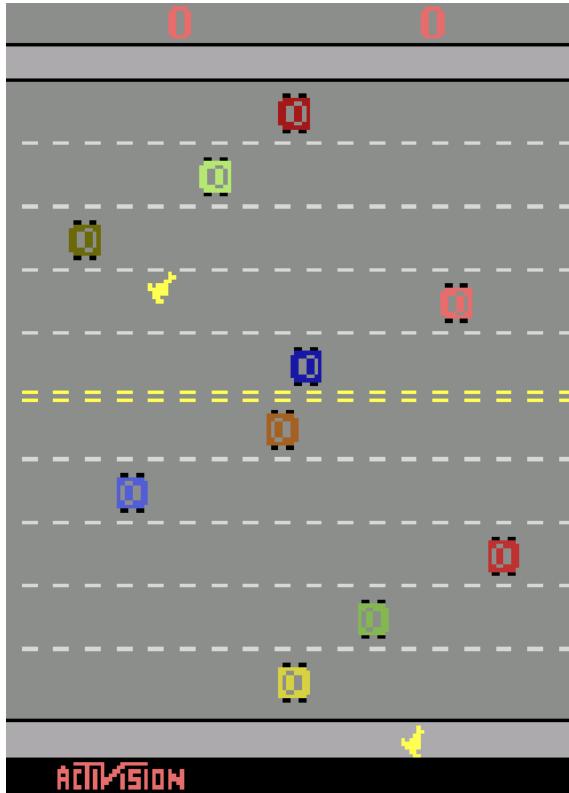
Boxing is a video game based on the sport of boxing. Boxing shows a top-down view of two boxers, one white and one black. When close enough, a boxer can hit his opponent with a punch. This causes his opponent to reel back slightly and the boxer scores a point, a reward

of 1. In the other situation, when the boxer gets hit, he gets a negative reward of -1. There are no knockdowns or rounds. A match is completed either when one player lands 100 punches (a 'knockout') or two minutes have elapsed. In the case of a decision, the player with the most landed punches is the winner. Ties are possible. While the gameplay is simple, there are subtleties, such as getting an opponent on the 'ropes' and 'juggling' him back and forth between alternate punches.



Rys. 51. Example of Boxing level

In Freeway an agent controls a chicken who can be made to run across a ten lane highway filled with traffic in an effort to "get to the other side." Every time a chicken gets across a reward of 1 is earned by the agent. If hit by a car, then a chicken is forced back slightly. The goal is to score as much points as possible in the two minutes. The chicken is only allowed to move up or down. The major challenge in this environment are sparse rewards. The agent scores only when successfully crosses the highway, which is not a trivial task.



Rys. 52. Example of Freeway level

[TODO: Add more games if needed.]

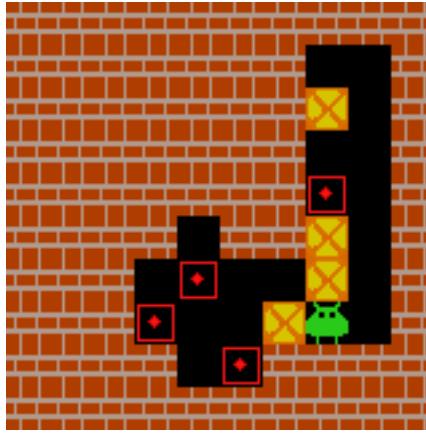
5.1.2. Sokoban

Sokoban is a classic planning problem. It is a challenging one-player puzzle game in which the goal is to navigate a grid world maze and push boxes onto target tiles. A Sokoban puzzle is considered solved when all boxes are positioned on top of target locations. The player can move in all 4 cardinal directions and only push boxes into an empty space (as opposed to pulling). For this reason many moves are irreversible and mistakes can render the puzzle unsolvable. A human player is thus forced to plan moves ahead of time. Artificial agents should similarly benefit from a learned model and simulation.

Despite its simple rule set, Sokoban is an incredibly complex game for which no general solver exists. It can be shown that Sokoban is NP-Hard and PSPACE-complete [7]. Sokoban has an enormous state space that makes it inassailable to exhaustive search methods. An efficient automated solver for Sokoban must have strong heuristics, just as humans utilize their strong intuition, so that it is not overwhelmed by the number of possible game states.

The implementation of Sokoban[27] used for those experiments procedurally generates a new level each episode. This means an agent cannot memorize specific puzzles. Together with the planning aspect, this makes for a very challenging environment. While the underlying game logic operates in a 10×10 grid world, agents were trained directly on RGB sprite graphics. Fig. 53 shows an example of Sokoban level with 4 boxes.

[TODO: Go into deeper details about e.g. how rewards are obtained etc.]



Rys. 53. Example of Sokoban level

5.2. Hardware

[TODO: Add HW specification.]

5.3. World Models for Sokoban

This section focuses on a problem of model learning in Sokoban that is used to solve the environment. The goal was to train a model to generate sharp and accurate open loop predictions of observations and obtain high score using it. In theory, how probable are sequences of ground truth observations is measured using log probability. In practice, though, it is far more useful to compare those generated sequences of future observations with ground truth sequences with an eye of the researcher. What the researcher looks for are sharp generated images which accurately resemble frames from the game. Moreover, the sequence needs to simulate subsequent actions properly, otherwise it is told that the sequence is noisy or does not model actions issued well.

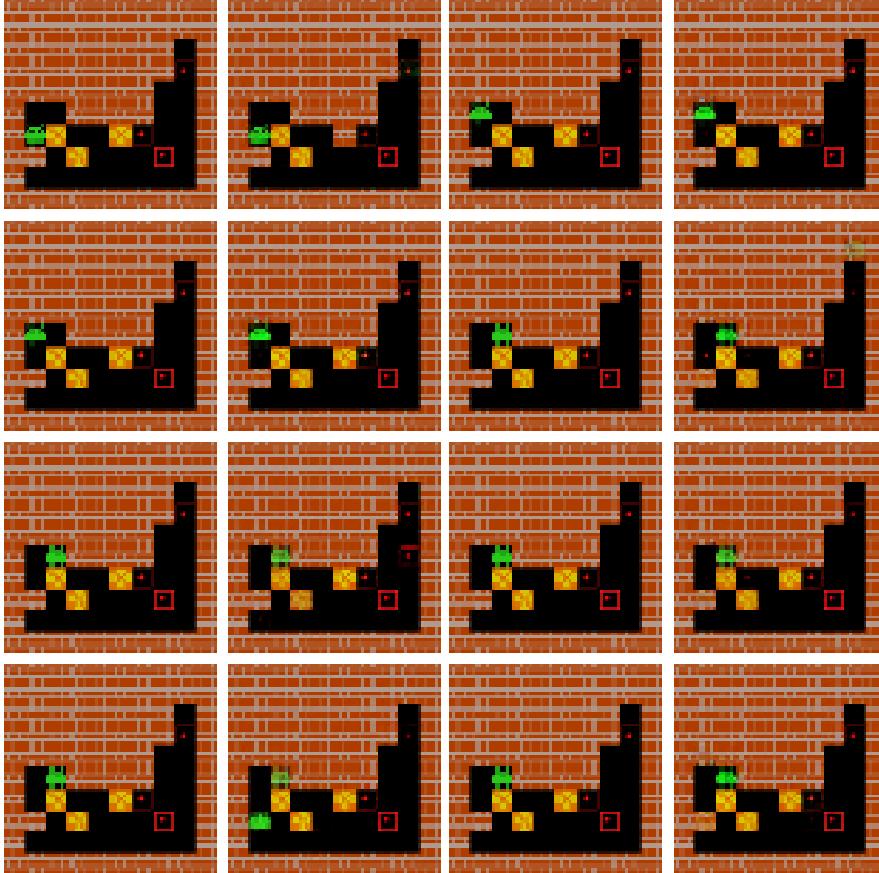
To generate an open loop prediction 60 consecutive frames and actions on average are fed into the model to initialize its hidden state first. Then, the model is ready to generate an open loop future prediction with its memory module by feeding it with subsequent actions and its own predictions as a contemporary state at following time steps. Reader should notice that what gets generated are not frames, but latent state's vectors. Those latent states can then be decoded into images for inspection.

5.3.1. Train the unchanged World Models implementation in the Sokoban environment

In this experiment, the original World Models was trained in the Sokoban environment. No modification to the original method described in the related work chapter was made, beyond addition of the deterministic variant of memory module described in the previous chapter.

The vision model successfully learned to encode high dimensional observations into low dimensional latent states. Fig. 54 shows original observations (first and third columns) side by side with reconstructed observations from their encodings (second and fourth columns). These

are zero step predictions, no future is predicted only encoding to latent space and decoding to image space again is done.



Rys. 54. Qualitative result of the vision model training in Sokoban. First and third columns include original observations. Second and fourth columns include reconstructions. Each reconstruction was obtained by first encoding the original observation and then decoding it, using VAE encoder and decoder respectively.

The stochastic and deterministic memory models were not able to learn Sokoban's dynamics. Fig. 55 shows that the stochastic model very often can not determine the agents position. The agent disappears and blocks change their types. The eighth row shows that pushing mechanics are not modeled, the agent passes through boxes. The deterministic model does not do better. The controller model failed to learn how to solve any level, it behaved comparably to a random play. We suspect that VAE is unable to generate usable abstract Sokoban representation and the shallow memory and controller models can not grasp complex dynamics of Sokoban using this poor representation. This idea is further developed in the next experiment.



Rys. 55. Qualitative result of the memory model training in Sokoban. Each row depicts the memory model rollout in one episode. The first column include original observations from the evaluation dataset from which the rollouts start. The RNN's hidden state was initialized on preceding transitions in each episode. Each subsequent reconstruction was obtained by first predicting the next latent state by the memory model and then decoding it using the VAE decoder.

5.3.2. Train World Models in Sokoban environment on 10x10 grid world states

The latent state vector size is set to 64. This means that in theory this vector can accommodate full information about an observation. As noted before, Sokoban underlying game logic operates in a 10×10 grid world, where far edges of a level are always walls. This means that the level is described by 64 block types organized in an 8×8 grid. In this experiment, this domain knowledge is exploited and the agent uses those 64 block types as an input vector to the memory module, bypassing the vision model. It is worth noting, that the vision model should learn this representation as it is the optimal encoding when the objective is to compress a pixel image into

a 64-dimensional vector and then reconstruct the original observation from it. However, despite use of the optimal encoding, the results have not been improved.

The proposed input format is optimal encoding if one wants to compress a pixel image and then reconstruct it. However, it is really poor representation of a current state of the environment if one wants to use linear combination of those features (block types in each position) to infer optimal next action and this is exactly what the controller model is trying to do. Modeling a value function could have more sense e.g. the value function could learn that a box on a target position yields higher value, but even it would have a hard time modeling more complex relations between entities in the environment. More useful for the controller would be e.g. representation that includes information about distance between the box and each target position. Nevertheless, this could be not enough too. The box on the target position would get discounted for not being on some other target positions. Hence, there is need for feature saying “the box X placed on the target position Y”. In the end, the linear combination of the proposed latent features can not model useful policy.

On the other hand, this representation includes, not well represented, but perfect information about an environment state. The memory model creates its own environment representation encoded in its hidden state and then uses this representation to predict the next latent state. This memory’s hidden state is also utilized by the controller. Still, it does not seem to encode useful enough information for the two, memory and controller, to do well on their tasks. One way to improve the hidden state representation is explored in the next experiment.

5.3.3. *Train World Model in Sokoban with auxiliary tasks*

Auxiliary tasks[17] have proved to help create more informative representation of an environment. In this experiment, reward and value prediction tasks are added to the memory model. In short, two additional linear models are added on top of the RNN to predict the next reward in the environment and model a value function. In theory, it should help form a more informative hidden state of the memory model. Consequently, it should help learn Sokoban’s dynamics, but also generate representation on a higher level of abstraction that could prove useful for the controller. Moreover, a reward prediction will be needed in further work on planning with learned model.

For all that, the memory model have not been able to learn to predict the rewards and values. Also, there was no improvement in memory’s and controller’s performance. It is suspected, that the main cause of this failure are sparse rewards in the training dataset. A random agent used to generate the dataset does not receive many positive rewards. Effectively, most of the episodes do not have any positive reward. Hence, the memory model soon overfit on more or less constant reward and value. This yields insight that the data generation procedure does not cover state-space well. Iterative approach to gathering data, from a better and better agent, could solve this problem.

It is not without significance that Sokoban has enormous state-space. Because each episode, or level, is randomly generated it is much different from the others - it is nearly impossible for an agent to see a similar state in a different episode. Hence, Sokoban requires strong

generalization capability from the memory module. Simple RNN can lack capacity to create good representation and in turn achieve good prediction performance. For instance I2A[35] uses deep neural network architecture to handle Sokoban complexity. A more flexible memory model with larger capacity could manage this complexity and need for generalization. The two insights are explored in the PlaNet experiments, which have larger model and uses iterative training procedure. However before that, World Models are tested in an easier environment, Boxing from Atari.

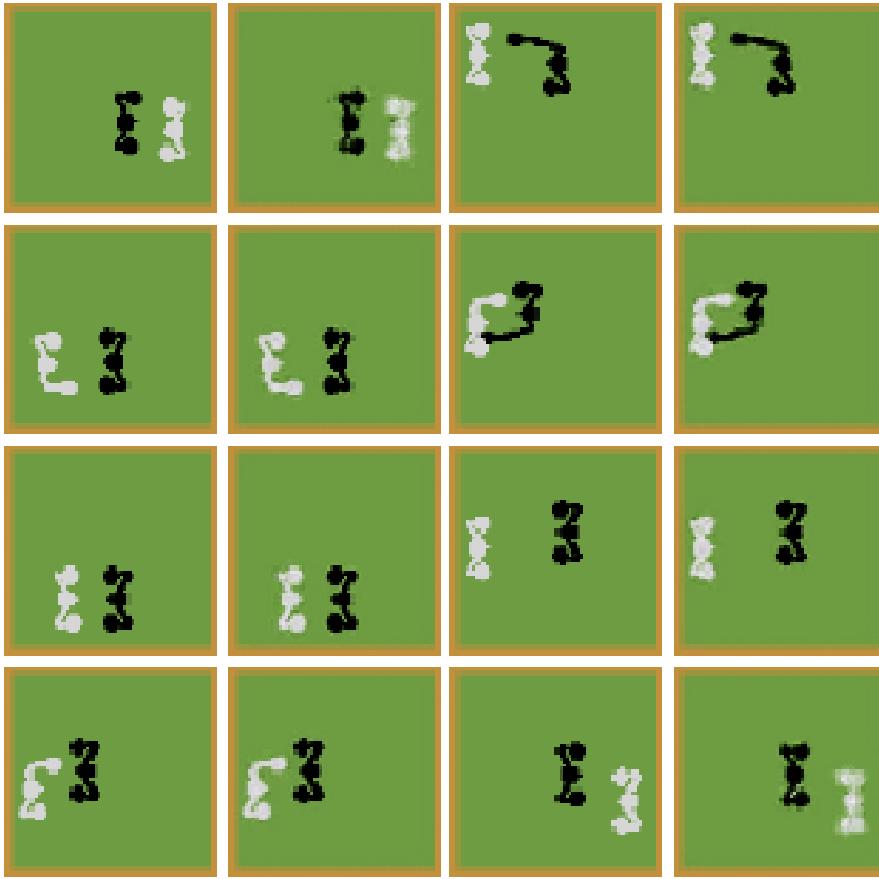
5.4. **World Models for Atari**

In two experiments below ideas from previous section were put into the test. Firstly, the original World Models is trained for the Boxing environment, which has dense rewards and data collection using a random agent cover most of the state-space. Then, World Models is coupled with AlphaZero planner and both are trained jointly.

5.4.1. *Train unchanged World Models in the Boxing environment*

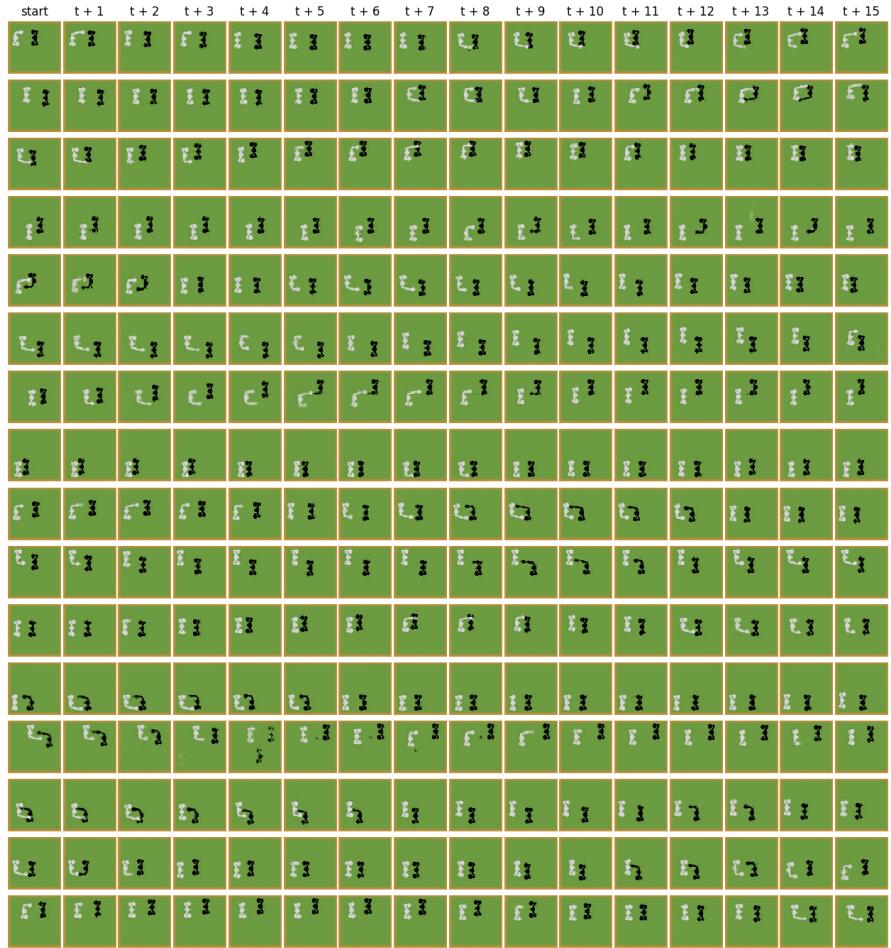
In this experiment, the original World Models was trained in the Boxing environment. No modification to the original method described in the related work chapter was made. In this experiment discrete memory was not tested as original stochastic one did well.

The vision model successfully learned to encode high dimensional observations into low dimensional latent states. Fig. 56 shows original observations (first and third columns) side by side with reconstructed observations from their encodings (second and fourth columns). These are zero step predictions, no future is predicted only encoding to latent space and decoding to image space again is done.



Rys. 56. Qualitative result of the vision model training in Boxing. First and third columns include original observations. Second and fourth columns include reconstructions. Each reconstruction was obtained by first encoding the original observation and then decoding it, using VAE encoder and decoder respectively.

The stochastic memory models was able to learn Boxing's dynamics. Fig. 57 shows that the stochastic model generates very sharp and accurate predictions that model agents movement and punches really well. The agents does not disappear like in Sokoban and actions are smooth. The controller model successfully learned how to solve the game scoring above 18 points on average across 5 runs. We suspect that World Models with latent state of size 16 was forced to encode two characters positions and hands states which are useful high-level features when deciding on the next action. It is worth pointing out here that similar experiment with such a small latent space did not yield improvement in Sokoban.



Rys. 57. Qualitative result of the memory model training in Boxing. Each row depicts the memory model rollout in one episode. The first column include original observations from the evaluation dataset from which the rollouts start. The RNN's hidden state was initialized on preceding transitions in each episode. Each subsequent reconstruction was obtained by first predicting the next latent state by the memory model and then decoding it using the VAE decoder.

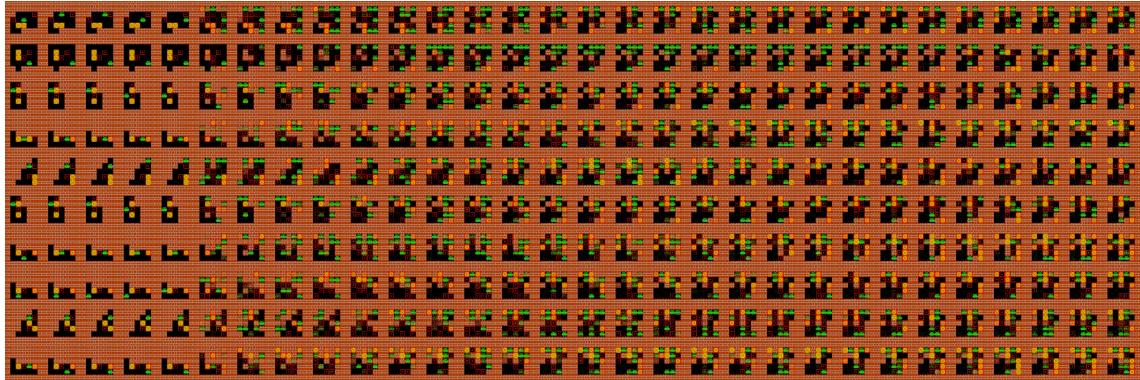
5.4.2. Train World Models with AlphaZero planner in the Boxing environment

Despite getting really accurate future predictions in the previous experiment and hyper-parameter tuning of this architecture the AlphaZero planner training was unstable. It did not train to properly plan in latent space and play the game. Decision was to abandon this solution and move to PlaNet which shown that, in deed, it is possible to plan in continuous control tasks.

5.5. *PlaNet* for Sokoban

5.5.1. Train unchanged *PlaNet* in the Sokoban environment

PlaNet did not capture Sokoban dynamics too. In the figure below (Fig. 58) future predictions are blurred, multiple agents appear and other artifacts, like changing block type, are present. Similarly like in World Models case decision was to move to Atari games as easier environments to start with.



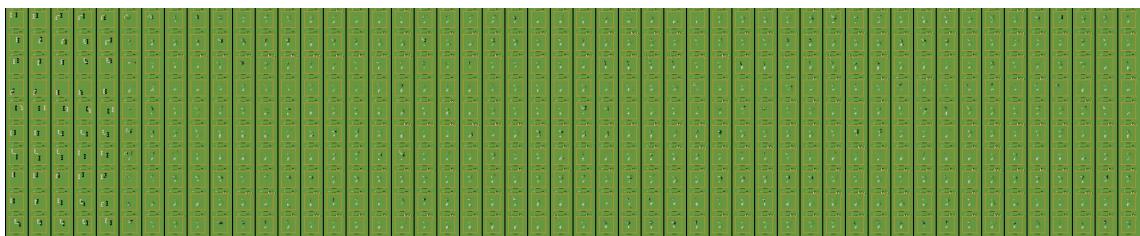
Rys. 58. Qualitative result of the model training in Sokoban. Each row depicts the model rollout in one episode. The first five columns include original consecutive observations from the evaluation dataset from which the rollouts start. The model hidden state was initialized on these transitions. Each subsequent reconstruction was obtained by first predicting the next latent state by the model and then decoding it using the decoder.

5.6. *PlaNet* for Atari

In this section experiments that lead to first successful case are described. The next section will focus on tuning this method to yield high scores, comparable to model-free methods, with the smallest amount of data possible.

5.6.1. Train unchanged *PlaNet* in the Boxing environment

It did not start to work out of the box of course. Fig. 59 shows that future predictions turn into a blurry blob, where it is not possible to distinguish one player from another.



Rys. 59. Qualitative result of the model training in Boxing. Each row depicts the model rollout in one episode. The first five columns include original consecutive observations from the evaluation dataset from which the rollouts start. The model hidden state was initialized on these transitions. Each subsequent reconstruction was obtained by first predicting the next latent state by the model and then decoding it using the decoder.

By default a decoder variance of 1 is used, which means the model explains a lot of variation in the image as random noise. While this leads to more robust representations, it also leads to more blurry images. If the changes in consecutive frames are minor, then the posterior collapses because the model explains everything as observations noise. There are two possible solutions to this issue: one is to increase an action repeat and the other is to try to reduce the decoder variance. These are examined next.

5.6.2. Train PlaNet in the Boxing environment with the increased action repeat

The action repeat will result in a bigger difference between consecutive frames and thus more signal for the model to learn from, that cannot be easily modeled as noise. In practice though, it did not help and even made the agent play worse than a random agent. The random agent is taking moves at random.

5.6.3. Train PlaNet in the Boxing environment with the lowered decoder variance

The predictions are more blurry with a higher variance because the decoder generate more observations that differ slightly from the same latent code. This leads to the posterior explaining more similar observations with the same code. If consecutive frames are very similar, then the posterior collapses and explain them with one code. By lowering the variance it becomes more sensitive to small changes in observations. *[TODO: Proofread this explanation once more and/or ask Danijar if it is right.]*

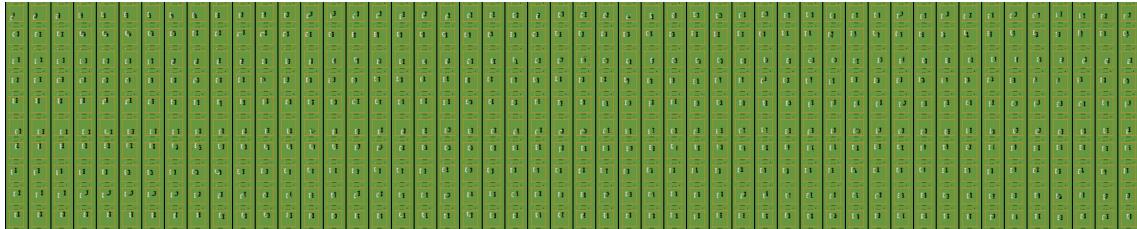
Lowering the decoder variance is equivalent to lowering a KL divergence scale in the PlaNet loss. It can be seen by writing the ELBO for a Gaussian decoder in the standard form $E_q(z)[lnp(x|z)] - KL[q(z)||p(z)]$. The log-likelihood terms is $lnp(x|z) = -0.5(x - f(z))/\sigma^2 - lnZ$. Multiplying the ELBO by σ^2 removes it from the log-probability term and puts it in front of the KL term as in beta-VAE [15]. The objectives have different values because of the Gaussian normalizer Z but they share the same gradient since the normalizer is a constant. Other reason that lowering the divergence scale can help with collapsing posterior is that it allows the model to absorb more information from its observations by loosening the information bottleneck.

On the other hand it is recommended to keep the divergence scale as high as possible while still allowing for good performance. For example, when the divergence scale is set to zero it could learn to become a deterministic autoencoder which reconstruct observations well but is less likely to generalize to state in latent space that the decoder hasn't seen during training.

Random search resulted in the best divergence scale being around 0.03. It was tuned jointly with a free nats parameter which is described in the next section. *[QUESTION: I should add diagram with random search results, but how to do this if those are evaluated with a researcher eye?] [TODO: You should better describe this random search experiment. What parameters where tuned, which turned out to be the most important, for how long and how much runs you were running etc.]*

5.6.4. Train PlaNet in the Boxing environment with increased free nats

Free nats technique is often used for static Variational Autoencoders. The model is allowed to use given amount of nats without KL penalty in variational objective. It helps the model focus on smaller details which do not contribute much to improving the reconstruction loss. Intuitively to this threshold of KL divergence (between a prior and a posterior) reconstruction loss is favoured. In case of Boxing, it helped to model boxers moves and actions more accurately. The best free nats turned out to be 12. Fig. 510 shows final result.

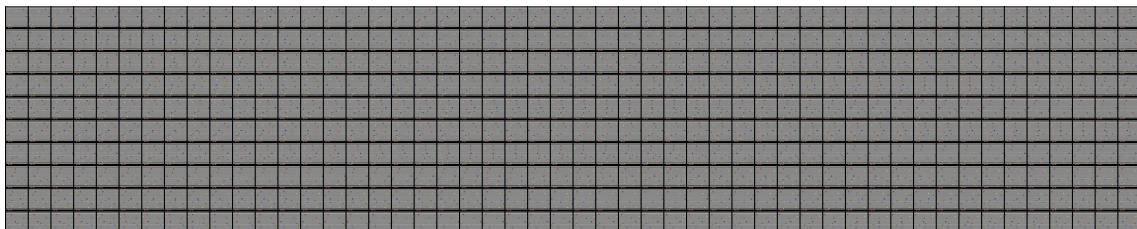


Rys. 510. Qualitative result of the model training in Boxing. Each row depicts the model rollout in one episode. The first five columns include original consecutive observations from the evaluation dataset from which the rollouts start. The model hidden state was initialized on these transitions. Each subsequent reconstruction was obtained by first predicting the next latent state by the model and then decoding it using the decoder.

It achieved final score around 30. *[TODO: Download .csv with results from five Boxing training and plot nice curve.]*

5.6.5. Train tuned PlaNet in the Freeway environment

The same random search procedure was applied to the Freeway environment described earlier. Fig. 511 shows final result.



Rys. 511. Qualitative result of the model training in Freeway. Each row depicts the model rollout in one episode. The first five columns include original consecutive observations from the evaluation dataset from which the rollouts start. The model hidden state was initialized on these transitions. Each subsequent reconstruction was obtained by first predicting the next latent state by the model and then decoding it using the decoder.

Despite really good future observations prediction the agent failed to solve the task. *[TODO: In benchmark description you should write what you consider as solved task in each environment.]* Possibly planner horizon is to short to cover a plan which ends with positive reward on the other side of the road *[TODO: In nomenclature or somewhere else you should clearly describe what is “a plan”.]* This is explored in the next experiment.

5.6.6. *Train tuned PlaNet in the Freeway environment with a longer planning horizon*

[NOTE: This is the last experiment in progress, but it does not seem promising.]

6. CONCLUSION

REFERENCES

- [1] Marc G. Bellemare et al. “The Arcade Learning Environment: An Evaluation Platform for General Agents”. In: *arXiv e-prints* (2012). arXiv: 1207.4708.
- [2] Greg Brockman et al. *OpenAI Gym*. 2016. arXiv: 1606.01540.
- [3] Lars Buesing et al. “Learning and Querying Fast Generative Models for Reinforcement Learning”. In: *arXiv e-prints* (2018). arXiv: 1802.03006.
- [4] Silvia Chiappa et al. “Recurrent Environment Simulators”. In: *arXiv e-prints* (2017). arXiv: 1704.02254.
- [5] François Chollet et al. *Keras*. <https://keras.io>. 2015.
- [6] Jack Clark and Dario Amodei. *Faulty Reward Functions in the Wild*. <https://openai.com/blog/faulty-reward-functions/>. Accessed: 2019-05-25.
- [7] Dorit Dor and Uri Zwick. “SOKOBAN and other motion planning problems”. In: *Computational Geometry* (1999). ISSN: 0925-7721. DOI: [https://doi.org/10.1016/S0925-7721\(99\)00017-6](https://doi.org/10.1016/S0925-7721(99)00017-6). URL: <http://www.sciencedirect.com/science/article/pii/S0925772199000176>.
- [8] Richard Evans et al. “De novo structure prediction with deep-learning based scoring”. Dec. 2018. URL: <https://deepmind.com/blog/alphafold/>.
- [9] Vladimir Feinberg et al. “Model-Based Value Estimation for Efficient Model-Free Reinforcement Learning”. In: *arXiv e-prints* (2018). arXiv: 1803.00101.
- [10] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016.
- [11] Alex Graves. “Generating Sequences With Recurrent Neural Networks”. In: *arXiv e-prints* (2013). arXiv: 1308.0850.
- [12] David Ha and Jürgen Schmidhuber. “World Models”. In: *arXiv e-prints* (2018). arXiv: 1803.10122.
- [13] Danijar Hafner et al. “Learning Latent Dynamics for Planning from Pixels”. In: *arXiv e-prints* (2018). arXiv: 1811.04551.
- [14] Matteo Hessel et al. “Rainbow: Combining Improvements in Deep Reinforcement Learning”. In: *arXiv e-prints* (2017). arXiv: 1710.02298.
- [15] Irina Higgins et al. “beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework”. In: *ICLR* (2017).
- [16] Sepp Hochreiter and Jürgen Schmidhuber. “Long Short-Term Memory”. In: *Neural Comput.* (1997). ISSN: 0899-7667. DOI: 10.1162/neco.1997.9.8.1735. URL: <http://dx.doi.org/10.1162/neco.1997.9.8.1735>.
- [17] Max Jaderberg et al. “Reinforcement Learning with Unsupervised Auxiliary Tasks”. In: *arXiv e-prints* (2016). arXiv: 1611.05397.
- [18] Piotr Januszewski, Grzegorz Beringer, and Mateusz Jablonski. *HumbleRL - Straightforward reinforcement learning Python framework*. 2019. URL: <https://github.com/piojanu/humblerl>.
- [19] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. “Stochastic Backpropagation and Approximate Inference in Deep Generative Models”. In: *arXiv e-prints* (2014). arXiv: 1401.4082.
- [20] Felix Leibfried, Nate Kushman, and Katja Hofmann. “A Deep Learning Approach for Joint Video Frame and Reward Prediction in Atari Games”. In: *arXiv e-prints* (2016). arXiv: 1611.07078.
- [21] Marlos C. Machado et al. “Revisiting the Arcade Learning Environment: Evaluation Protocols and Open Problems for General Agents”. In: *arXiv e-prints* (2017). arXiv: 1709.06009.
- [22] Volodymyr Mnih et al. “Asynchronous Methods for Deep Reinforcement Learning”. In: *arXiv e-prints* (2016). arXiv: 1602.01783.
- [23] S Mor-Yosef et al. “Ranking the risk factors for cesarean: Logistic regression analysis of a nationwide study”. In: *Obstetrics and gynecology* 75 (July 1990), pp. 944–7.

- [24] Y. Naddaf and University of Alberta. Department of Computing Science. *Game-independent AI Agents for Playing Atari 2600 Console Games*. University of Alberta, 2010. URL: <https://books.google.pl/books?id=c85vnQAACAAJ>.
- [25] Adam Paszke et al. “Automatic Differentiation in PyTorch”. In: *NIPS Autodiff Workshop*. 2017.
- [26] Tim Salimans et al. “Evolution Strategies as a Scalable Alternative to Reinforcement Learning”. In: *arXiv e-prints* (2017). arXiv: 1703.03864.
- [27] Max-Philipp B. Schrader. *gym-sokoban*. <https://github.com/mpSchrader/gym-sokoban>. 2018.
- [28] John Schulman et al. “Proximal Policy Optimization Algorithms”. In: *arXiv e-prints* (2017). arXiv: 1707.06347.
- [29] David Silver, Richard S. Sutton, and Martin Müller. “Temporal-difference search in computer Go”. In: *Machine Learning* 87.2 (2012), pp. 183–219. ISSN: 1573-0565. DOI: 10.1007/s10994-012-5280-0. URL: <https://doi.org/10.1007/s10994-012-5280-0>.
- [30] David Silver et al. “Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm”. In: *arXiv e-prints* (2017). arXiv: 1712.01815.
- [31] David Silver et al. “Mastering the game of Go without human knowledge | Nature”. In: *Nature* 550 (2017). URL: <http://dx.doi.org/10.1038/nature24270>.
- [32] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction* 2nd. The MIT Press, 2018. ISBN: 9780262039246.
- [33] Erik Talvitie. “Agnostic System Identification for Monte Carlo Planning”. In: AAAI’15 (2015), pp. 2986–2992. URL: <http://dl.acm.org/citation.cfm?id=2888116.2888132>.
- [34] Erik Talvitie. “Model Regularization for Stable Sample Rollouts”. In: UAI’14 (2014), pp. 780–789. URL: <http://dl.acm.org/citation.cfm?id=3020751.3020832>.
- [35] Théophane Weber et al. “Imagination-Augmented Agents for Deep Reinforcement Learning”. In: *arXiv e-prints* (2017). arXiv: 1707.06203.

LIST OF FIGURES

21.	Reinforcement Learning	11
22.	Deep Learning	13
23.	Multilayer perceptron	14
31.	Flow diagram of the World Models agent's model	15
32.	VizDoom	17
41.	HumbleRL architecture	19
42.	HumbleRL loop function overview	20
43.	Stochastic World Models probabilistic graphical model diagram	22
51.	Boxing	27
52.	Freeway	28
53.	Sokoban.....	29
54.	Qualitative result of the World Models vision model training in Sokoban	30
55.	Qualitative result of the World Models memory model training in Sokoban	31
56.	Qualitative result of the World Models vision model training in Boxing	34
57.	Qualitative result of the World Models memory model training in Boxing.....	35
58.	Qualitative result of the PlaNet model training in Sokoban.....	36
59.	Qualitative result of the original PlaNet model training in Boxing	36
510.	Qualitative result of the PlaNet model training with a lower divergence scale in Boxing	38
511.	Qualitative result of the PlaNet model training with a lower divergence scale in Freeway	38

LIST OF TABLES