

# Reinforcement Learning Soccer Teams with Incomplete World Models

MARCO WIERING, RAFAL SALUSTOWICZ, JÜRGEN SCHMIDHUBER \*

marco@idsia.ch, rafal@idsia.ch, juergen@idsia.ch

*IDSIA, Corso Elvezia 36, 6900 Lugano, Switzerland*

*Received 28 October, 1998; Revised*

Editors:

**Abstract.** We use reinforcement learning (RL) to compute strategies for multiagent soccer teams. RL may profit significantly from world models (WMs) estimating state transition probabilities and rewards. In high-dimensional, continuous input spaces, however, learning accurate WMs is intractable. Here we show that incomplete WMs can help to quickly find good action selection policies. Our approach is based on a novel combination of CMACs and prioritized sweeping-like algorithms. Variants thereof outperform both  $Q(\lambda)$ -learning with CMACs and the evolutionary method Probabilistic Incremental Program Evolution (PIPE) which performed best in previous comparisons.

**Keywords:** reinforcement learning, CMAC, world models, simulated soccer,  $Q(\lambda)$ , evolutionary computation, PIPE

## 1. Introduction

Our goal is to build teams of autonomous agents that learn to play soccer from very sparse reinforcement signals: only scoring a goal yields reward for the successful team. Team members try to maximize reward by improving their adaptive policy mapping sensory inputs to actions. In principle there are at least two types of learning algorithms applicable to such problems: reinforcement learning (RL), e.g., [23, 29, 33, 31], and evolutionary approaches, e.g., [9, 17, 7, 19]. Here we describe a novel RL method and compare its results to those obtained by previous RL methods and an evolutionary approach.

Most existing RL algorithms are based on function approximators (FAs) learning value functions

(VFs) that map state/action pairs to the expected outcome (reward) of a trial [5, 33]. In realistic, partially observable, multiagent environments, learning value functions is hard though. This makes evolutionary methods a promising alternative. For instance, in previous work on learning soccer strategies [22] we found that *Probabilistic Incremental Program Evolution* (PIPE) [19], a novel evolutionary approach to searching program space, outperforms  $Q(\lambda)$  [16, 33, 37] combined with FAs based on linear neural networks [21] or neural gas [20].

We identified several reasons for PIPE's superiority: (1) In complex environments such as ours RL methods tend to be brittle — once discovered, good policies do not stabilize but tend to get destroyed by subsequent “unlucky” experiences. PIPE is less affected by this problem because good policies have a large probability of surviving. (2)

\*This research is supported by SNF grant 2100-49'144.96 “Long Short-Term Memory”

PIPE learns faster by isolating important features in the sensory input, combining them in programs of initially low algorithmic complexity, and subsequently refining the programs. This motivates our present approach: VF-based RL should also be able to (a) stabilize or improve fine policies (as opposed to unlearning them), (b) profit from the existence of low-complexity solutions, and (c) use incremental search to find more complex solutions where simple ones do not work.

**Incomplete world models.** *Direct* RL methods [5, 33] use temporal differences (TD) [29] for training FAs to approximate the VF from simulated trajectories through state/action space. *Indirect* RL, however, learns a world model (WM) [14, 35] estimating the reward function and the transition probabilities between states, then uses dynamic programming (DP) [4] or similar, faster algorithms such as prioritized sweeping (PS — which we will use in the paper) [14] for computing the VF. This can significantly improve learning performance in discrete state/action spaces [14]<sup>1</sup>. In case of continuous spaces, WMs are most effectively combined with *local* FAs transforming the input space into a set of discrete regions and then learning the VF. Similarly, continuous action spaces can be transformed in a set of discrete actions. Previous work has already demonstrated the effectiveness of learning discrete world models for robotic localization and navigation tasks, e.g., [32]. Learning accurate WMs in high-dimensional, continuous, partially observable environments is hard, however, and this motivates our novel approach to learning useful but incomplete models instead.

**CMAC models.** We will present a novel combination of CMACs and world models. CMACs [1] use *filters* mapping sensor-based inputs to a set of activated cells. Each filter partitions the input space into subsections in a prewired way such that each (possibly multi-dimensional) subsection is represented by exactly one discrete cell of the filter. For example, a filter might consist of a finite number of cells representing an infinite set of colors represented by cubes with 3 dimensions red, blue and yellow, and activate the cell which encloses the current color input component.

In an RL context each cell has a Q-value for each action. The Q-values of currently active cells are averaged to compute the overall Q-values re-

quired for action selection. Previous work already combined CMACs with Q-learning [33] and Q( $\lambda$ ) methods [30, 24]. Here we combine CMACs with WMs by learning an independent model for each filter. These models are then exploited by a version of prioritized sweeping (PS) [14, 36] for computing the Q-functions. Later we will find that CMAC models can quickly learn to play a good soccer game and surpass the performance of PIPE and an approach combining CMACs and Q( $\lambda$ ).

**Outline.** Section 2 describes our soccer environment. Section 3 presents CMACs and describes how they can be combined with model-based learning. Section 4 describes experimental results. Section 5 concludes.

## 2. The Soccer Simulator

Our soccer simulator [22] runs discrete-time simulations involving two teams consisting of either 1 or 3 players per team. A game lasts from time  $t = 0$  to time  $t_{end} = 5000$ . The field is represented by a two-dimensional continuous Cartesian coordinate system. As in indoor soccer the field is surrounded by impassable walls except for the two goals centered in the east and west walls. There are fixed initial positions for all players and the ball (see Figure 1).

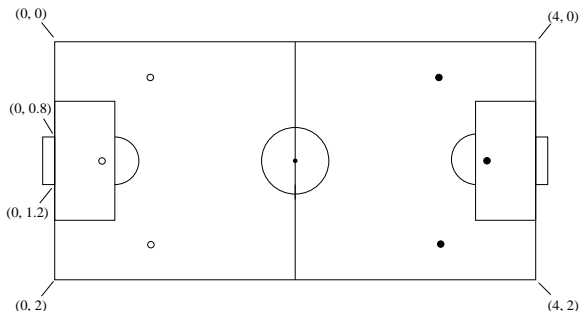


Fig. 1. Players and ball (center) in initial positions. Players of a 1 player team are those furthest in the back.

**Players/Ball.** Each player and the ball are represented by a solid circle and a variable real-valued position and orientation. A player whose circle intersects the ball picks it up and then owns it. The ball owner can move or shoot the ball. A shot is in the direction of the player’s orientation. When shot, the ball’s initial speed is 0.12 units per time step. Each following time step the ball slows down due to friction by 0.005 units per time step

(unless it is picked up by a player) - the ball can travel freely at most 1.5 units. At each discrete time step each player selects one of the following actions:

- *go\_forward*: move 0.025 units in current direction.
- *turn\_to\_ball*: point player’s orientation towards ball.
- *turn\_to\_goal*: point player’s orientation towards opponent’s goal.
- *shoot*: if the player owns the ball then change player’s orientation by a random angle from the interval  $[-5^\circ, 5^\circ]$  (to allow for noisy shots), and shoot ball in the corresponding direction.

A player that makes a step forward such that its circle intersects another player’s circle bounces back to its original position. If one of them owns the ball prior to collision then it will lose it to the collision partner.

**Action framework.** During each time step all players execute one action each, in randomly chosen order. Then the ball moves according to its current speed and direction. If a team scores or  $t = t_{end}$  then all players and ball will be reset to their initial positions.

**Sensory input.** At any given time a player’s input vector  $\vec{x}$  consists of 16 (1 player) or 24 (3 players) components:

- Three Boolean input components that tell whether the player/a team member/opponent team owns the ball.
- Polar coordinates (distance, angle) of both goals and the ball with respect to the player’s orientation and position.
- Polar coordinates of both goals relative to the ball’s orientation and position.
- Ball speed.
- Polar coordinates of all other players w.r.t. the player are ordered by (a) teams and (b) distances to the player.

**Policy-sharing.** All players share the same Q-functions or PIPE-programs. Still their behaviors differ due to different, situation-specific inputs. Policy-sharing has the advantage of greatly reducing the number of adaptive free parameters, which tends to reduce the number of required training

examples (learning time) and increase generalization performance, e.g., [15]. A potential disadvantage of policy sharing, however, is that different players cannot develop truly different strategies to be combined in fruitful ways.

### 3. CMAC Models

CMACs [1] use multiple, *a priori* designed filters to quantize the input space. Each filter consists of several cells with associated Q-values. Applying the filters to the current input yields a set of activated cells (a discrete distributed representation of the input). Their Q-values are averaged to compute the overall Q-value.

**Filter design.** In principle the filters may yield arbitrary divisions of the input space, such as hypercubes. To avoid the curse of dimensionality one may use hashing to group a random set of inputs into an equivalence class, or use hyperslices omitting certain dimensions in particular filters [30]. Although hashing techniques may help to overcome storage problems, we do not believe that random grouping is the best we can do. Since our soccer simulation involves a fair number of input dimensions (16 or 24), we use hyperslices to reduce the number of adjustable parameters. Our filters divide the state-space by splitting it along single input dimensions into a fixed number of cells — input components are treated in a mutually independent way. Multiple filters are applied to the same input component to allow for smoother generalization.

**Partitioning the input space.** We use two filters for each input component, both splitting the same component. Input components representing Boolean values, distances (or speeds), and angles, are split in various ways (see Figure 2): (1) Filters associated with a *Boolean* input component just return its value. (2) *Distance* or *ball-speed* input components are rescaled to values between 0 and 1. Then the filters partition the components into  $n_c$  or  $n_c + 1$  quanta. (3) *Angle* input components are partitioned in  $n_c$  equal quanta in a circular (and thus natural) way — one filter groups the angles  $359^\circ$  and  $0^\circ$  to the same cell, the other separates them by a cell boundary.

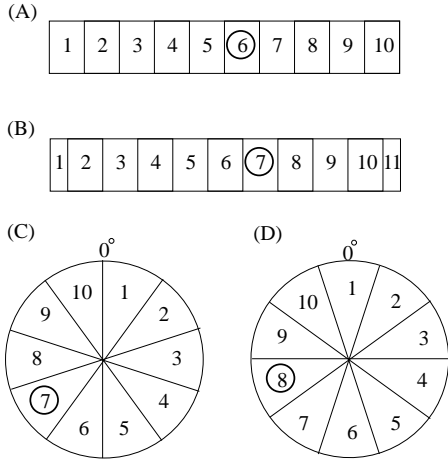


Fig. 2. We use two filters for each input component, resulting in a total of 32 (1 player) or 48 (3 players) filters. Filters of a Boolean input component just return the Boolean value as cell number. The figure (activated cells are marked) illustrates decompositions of (A) a continuous distance input component into 10 discrete cells, (B) the same component into 11 cells, (C) a continuous angle component into 10 cells, (D) the same component into 10 different cells.

**Selecting an action.** Applying all filters on a player’s current input vector at time  $t$  returns the active cells  $\{f_1^t, \dots, f_z^t\}$ , where  $z$  is the number of filters. The Q-value of selecting action  $a$  given input  $\vec{x}$  is calculated by averaging all Q-values of the active cells:

$$Q(\vec{x}, a) = \sum_{k=1}^z Q_k(f_k^t, a)/z,$$

where  $Q_k$  is the Q-function of filter  $k$ . After computing the Q-values of all actions we select an action according to the Max-random exploration rule: select the action with maximal Q-value with probability  $P_{max}$ , and a uniformly random action otherwise.

**Learning with WMs.** Learning accurate models for high-dimensional input spaces is hard. Usually there are so many possible successor states that storing all of them for each different input would be infeasible and updates would cost a lot of time. Instead we introduce a novel combination of model-based RL and CMACs. We use a set of independent models to estimate the dynamics of each filter. To estimate the transition model for filter  $k$ , we count the transitions from activated cell  $f_k^t$  to activated cell  $f_k^{t+1}$  at the next time-

step, given the selected action. These counters are used to estimate the transition probabilities  $P_k(c_j|c_i, a) = P(f_k^{t+1} = c_j | f_k^t = c_i, a)$ , where  $c_j$  and  $c_i$  are cells, and  $a$  is an action. For each transition we also compute the average reward  $R_k(c_i, a, c_j)$  by summing the immediate reinforcements, given that we make a step from active cell  $c_i$  to cell  $c_j$  by selecting action  $a$ .

**Prioritized sweeping (PS).** We could immediately apply dynamic programming (DP) [4] to the estimated models. Online learning with DP, however, is computationally expensive. But fortunately there are more efficient update management methods. We will use a method similar to prioritized sweeping (PS) [14] which may be the most efficient available update mechanism. PS updates the Q-value of the filter/cell/action triple with the largest update size before updating others. Each update is made via the usual Bellman backup [4]:

$$Q_k(c_i, a) \leftarrow \sum_j P_k(c_j|c_i, a)(\gamma V_k(c_j) + R_k(c_i, a, c_j))$$

where  $V_k(c_i) = \max_a Q_k(c_i, a)$  and  $\gamma \in [0, 1]$  is the discount factor. After each player action we update all filter models and use PS to compute the new Q-functions. PS uses a parameter to set the maximum number of updates per time step and a cutoff parameter  $\epsilon$  preventing tiny updates. Note that PS may use different numbers of updates for different filters, since some filters tend to make larger updates than others and the total number of updates per time step is limited. The complete PS algorithm is given in the Appendix.

**Non-pessimistic value functions.** Policy sharing requires the fusion of experimental data from different players into a single representation. This data, however, is generated by different player histories. In fact, certain experiences of certain players will probably never occur to others — there is no obvious and straightforward way of data fusing. For instance, the unlucky experience of one particular player may cause the VF approximation to assign low values to certain actions for all players. After having identified this problem, we tried a heuristic solution to overcome it. We compute *non-pessimistic* value functions: we decrease the probability of the worst transition from each cell/action and renormalize the other proba-

bilities. Then we apply PS to the adjusted probabilities (details of the algorithm are given in the Appendix). The effect is that only frequently occurring bad experiences have high impact on the Q-function. Experiments showed small but significant improvements over the basic algorithm. The method is quite similar to Model-Based Interval Estimation [36], an exploration algorithm extending Interval Estimation [10] by computing optimistic value functions for action selection.

**Multiple restarts.** The method sometimes may get stuck with continually losing policies which hardly ever score and fail to prevent (many) opponent goals (also observed with our previous simulations based on linear networks and neural gas). We could not overcome this problem by adding standard exploration techniques (evaluating alternative actions of losing policies is hard, since the perturbed policy will usually still lead to negative rewards). Instead we reset Q-functions and WMs once the team has not scored for 5 successive games but the opponent scored during the most recent game (we check these conditions every 5 games). After each restart, the team will gather different experiences affecting policy quality. We found that multiple restarts can significantly increase the probability of finding good policies.

We use  $P_{max} = 1.0$  in the Max-random exploration rule, since that worked best. The reason multiple restarts works better without exploration is that it makes the detection of losing policies easier. Hopeless greedy policies will lose something, whereas with exploration our agents may still score although they remain unable to improve their policy from the generated experiences. Thus, using greedy policies we may use a simpler rule for restarting.

**Learning with Q( $\lambda$ ).** Possibly the most widely used RL algorithm is Q-learning [33], which tries out sequences of actions through state/action space according to its policy and uses environmental rewards to estimate the expected long-term reward for executing specific actions in particular states. Q-learning repeatedly performs a one-step lookahead backup, meaning that the Q-value of the current state/action pair (SAP) becomes more like the immediately received reward plus the estimated value of the next state.

Q( $\lambda$ )-learning [33, 16, 37] combines TD( $\lambda$ ) methods [29] with Q-learning to propagate state/action updates back in time such that multiple SAPs which have occurred in the past are updated based on a single current experience. Q( $\lambda$ )-learning has outperformed Q-learning in a number of experiments [13, 18, 37]. For purposes of comparison we also use online Q( $\lambda$ )-learning for training the CMACs to play soccer. The details of the algorithm are given in the Appendix.

**PIPE.** The other competitor is *Probabilistic Incremental Program Evolution* (PIPE) [19]. PIPE is a novel technique for automatic program synthesis. It combines probability vector coding of program instructions [25, 26, 27], Population-Based Incremental Learning [2], and tree-coded programs like those used in some variants of Genetic Programming (GP) [7, 8, 12]. PIPE iteratively generates successive populations of functional programs according to an adaptive probability distribution over all possible programs. Each iteration it lets all programs play one soccer game; then the best program is used to refine the distribution. Thus PIPE stochastically generates better and better programs. All details can be found in [22].

## 4. Experiments

We compare the CMAC model to CMAC-Q( $\lambda$ ) and PIPE [19], which outperformed Q( $\lambda$ )-learning combined with various FAs in previous comparisons [20, 22].

**Task.** We train and test the learners against handmade programs of different strengths. The opponent programs are mixtures of a program which randomly executes actions (random program) and a (good) program which moves players towards the ball as long as they do not own it, and shoots it straight at the opponent’s goal otherwise. Our five opponent programs, called *Opponent( $P_r$ )*, use the random program to select an action with probability  $P_r \in \{0, \frac{1}{4}, \frac{1}{2}, \frac{3}{4}, 1\}$ , respectively, and the good program otherwise.

**CMAC model set-up.** We play a total of 200 games. Every 10 games we test current performance by playing 20 test games against the opponent and summing the score results. The reward is +1 if the team scores and -1 if the opponent

scores. The discount factor is set to 0.98. After a coarse search through parameter space we chose the following parameters: 2 filters per input component (total of 32 or 48 filters) number of cells  $n_c = 20$  (21 for the second filters of distance/speed input components). Q-values are initially zero. PS uses  $\epsilon = 0.01$  and a maximum of 1000 updates per time step. We only compute non-pessimistic value functions for the 3-player teams for which we use  $z_\alpha = 1.96$ .

**CMAC  $Q(\lambda)$  set-up.** We play a total of 200 games. Every 20 games we test current performance of the policy (during tests we continue selecting actions according to the current exploration scheme) by playing 20 test games against the opponent and summing the score results. The reward is +1 if the team scores and -1 if the opponent scores. The discount factor is set to 0.98. We conducted a coarse search through parameter space to select the best learning parameters. We use online  $Q(\lambda)$  with replacing traces [28] and  $\lambda = 0.8$  for the 1-player case, and  $\lambda = 0.5$  for the 3-player case. The initial learning rate is set to  $\alpha_c = 1.0$ , the learning rate decay rate to  $\beta = 0.3$ .

We use Max-random exploration with  $P_{max}$  linearly increased from 0.7 in the beginning of the simulation to 1.0 at the end. As for CMAC-models we use two filters per input component (total of 32 or 48 filters). The number of cells is set to  $n_c = 10$  (11 for the second filters of distance/speed input components). All Q-values are initially zero. In general, learning performance does not very sensitively depend on the used parameters. E.g., using  $n_c = 20$  results in only slightly worse performance. Small values for  $\lambda$  ( $< 0.3$ ) do make things worse though.

**PIPE set-up.** For PIPE we play a total of 1000 games. Every 50 games we test performance of the best program found during the most recent generation. Parameters for all PIPE runs are the same as in previous experiments [22].

**Results : 1-Player case.** We plot number of points (2 for scoring more goals than the opponent during the 20 test games, 1 for a tie, and 0 for scoring less) against number of training games in Figure 3.

We observe that on average our CMAC model wins against almost all training programs. Only against the best 1-player team ( $P_r = 0$ ) it wins

as often as it loses, and often plays ties (it finds a blocking strategy leading to a 0-0 result). Against the worst two teams, CMAC model always finds winning strategies.

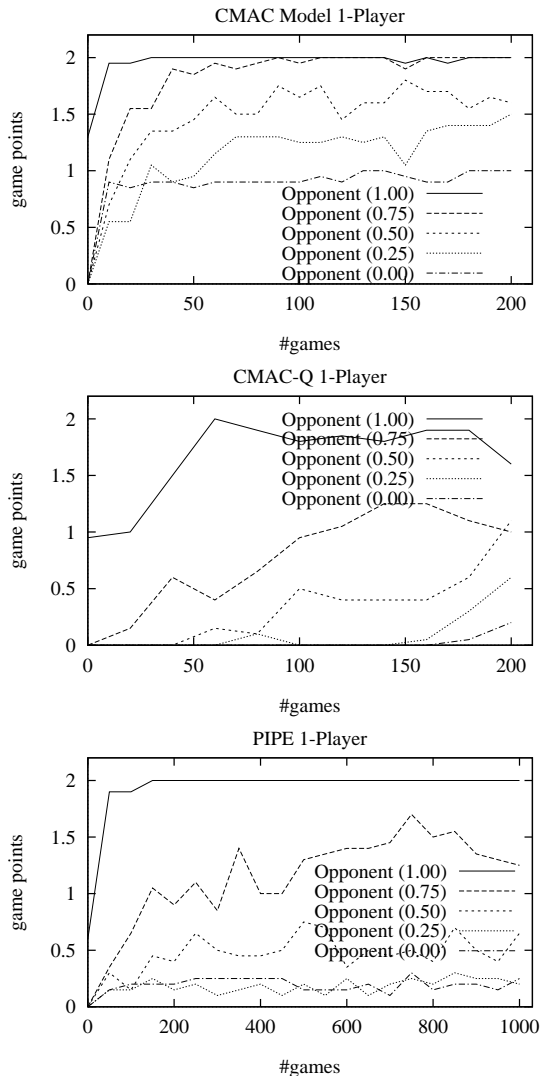


Fig. 3. Number of points (means of 20 simulations) during test phases for teams consisting of 1 player. Note the varying x-axis scalings.

CMAC- $Q(\lambda)$  finds programs that on average win against the random team, although they do not always win. It learns to play about as well as the 75% random and 50% random teams. CMAC- $Q(\lambda)$  is no match against the best opponent, and although it seems that performance jumps up at

the end of the trial, longer trials do not lead to better performances.

PIPE is able to find programs beating the random team and quite often discovers programs that win against 75% random teams. It encounters great difficulties in learning good strategies against the better teams, though: although PIPE may execute more games (1000 vs. 200), the probability of generating programs that perform well against the good opponents is very small. For this reason it tends to learn from the best of the losing programs. This in turn does not greatly facilitate the discovery of winning programs.

**Results : 3-Players case.** We plot number of points (2 for scoring more goals than the opponent during the 20 testgames) against number of training games in Figure 4.

Again, CMAC model always learns winning strategies against the worst 2 opponents. It loses on average against the best 3-player team (with  $P_r = 0.25$ ) though. Note that this strategy mixture works better than always using the deterministic program ( $P_r = 0$ ) against which CMAC model plays ties or even wins. In fact, the deterministic program tends to clutter agents such that they obstruct each other. The deterministic opponent’s behavior also is easier to model. All of this makes the stochastic version a more difficult opponent.

CMAC-Q is clearly worse than CMAC model — it learns to win only against the worst opponent.

PIPE performs well only against random and 75% random opponents. For the better opponents it runs into the same problems as mentioned above.

**Score differences.** We show maximal obtained score differences in Table 1 (1 player) and Table 2 (3 players). Although PIPE performs better against the weakest opponent than CMAC-models or CMAC-Q, PIPE often cannot score against strong opponents. CMAC-models, however, do score against the good opponents, and

are able to find (at least once) winning policies against all opponents.

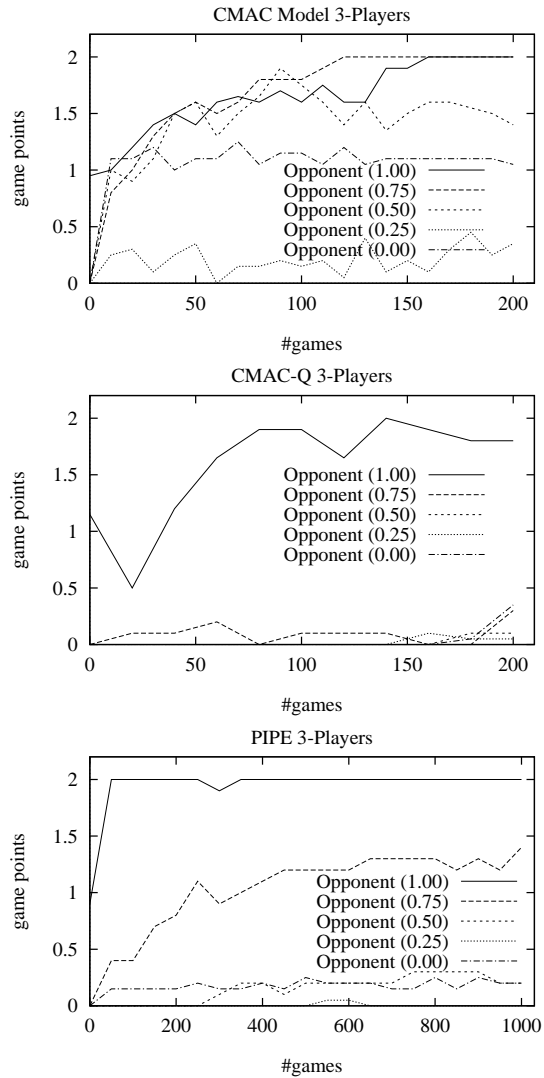


Fig. 4. Number of points (means of 20 simulations) during test phases for teams consisting of 3 players. Note the varying x-axis scalings.

We should keep in mind that score differences may have a large variance. For instance, in some experiments with the 1-player CMAC model, opponent(0.0) may continuously win 760-0 in test matches. This extreme score difference is caused by resetting the (losing) CMAC model just before testing. PIPE has a small advantage here, since it uses the best program of the last generation for

testing and thus almost lets vanish the probability of testing a really bad policy.

*Table 1. Best average score differences for different learning methods against 1-player opponents of varying strengths. \* = Although CMAC-models were sometimes able to score 7 goals, they also sometimes lost 0-760.*

Learning Alg.	1.0	0.75	0.5	0.25	0.0
CMAC model	85-2	68-3	27-1	6-15	1-146*
CMAC-Q	92-6	52-23	10-7	1-4	0-13
PIPE	225-18	127-35	19-13	0-3	0-10

*Table 2. Best average score differences for different learning methods against 3-player opponents of varying strengths.*

Learning Alg.	1.0	0.75	0.5	0.25	0.0
CMAC model	161-31	236-100	84-70	6-20	0.3-0
CMAC-Q	111-26	36-73	13-58	3-23	0-24
PIPE	297-18	163-64	30-31	0-11	0-21

**Discussion.** Despite treating all components independently the CMAC model is able to learn good reactive soccer strategies preferring actions that activate those cells of a filter which promise highest average reward. The use of a model often quickly stabilizes good strategies: given sufficient experiences (5-20 learning games), the policy will hardly change anymore. The reason is that deterministic policies generate similar experiences.

Early experiences with a random initial policy may greatly impact the final policy’s quality. They may result in continually losing policies (especially against better opponents) that are unable to improve, due to the near-impossibility of learning from bad experiences. For such reasons we performed multiple restarts (1 up to more than 10). Since tested strategies often remain either winners or losers, the step-wise improvements shown in the learning curves are mainly due to multiple restarts. The ups and downs in the learn-

ing curves are caused by unstable policies with unstable score results.

CMAC model tends to be quite robust under variations of filter design (e.g., combining multiple input components) and number of cells. Conducting additional experiments with filters combining distance and angle input components, or 10/11 instead of 20/21 cells per filter, we obtained similar levels of performance.

Multiple restarts helped CMAC model to avoid getting stuck with losing policies. When we tried CMAC-Q with multiple restarts and without exploration, it often found blocking strategies (leading to 0-0 results) against all 1-player opponents, but did not learn to win. Learning blocking strategies against multiagent teams, however, is much harder.

All methods perform better in the single agent case. This can probably be explained by the fact that the multiagent case yields more significantly different game configurations so that finding a policy that works fine for all of them is more difficult.

## 5. Conclusion

Model-based RL is a promising method for learning to control autonomous agents. Since learning accurate world models in high dimensional, continuous spaces is difficult, we have focused on learning useful but incomplete models instead. Here we have described a novel combination of CMACs and incomplete world models which allows for discovering successful soccer strategies and tends to outperform both PIPE and a  $Q(\lambda)$ /CMAC combination. Especially against better opponents CMAC models proved superior.

In some environments certain more complex filters grouping multiple context-dependent input components may be necessary. Filters combining many different, mutually dependent input components for a particular task may require a lot of storage space. Many of the possible input combinations, however, will never be experienced. A more space-efficient approach will use decision tree models to keep track of rewards and transition probabilities between leaf nodes defining “interesting” input component combinations. Starting with an initial set of low-complexity decision trees consisting of single root components, new



leaf nodes may be generated online using statistical tests as done in, e.g., the G-algorithm [6].

### Appendix Prioritized Sweeping

An efficient method determining which updates to perform is prioritized sweeping (PS) [14]. PS assigns priorities to updating the Q-values of different states according to a heuristic estimate of the size of the Q-values' updates. The algorithm keeps track of a "backward model" relating states to predecessor state/action pairs. After the update of a state value the state's predecessors are inserted in a priority queue which is then used for updating the Q-values of actions that can be performed in those states which have the highest priority.

**Our Prioritized Sweeping.** Moore and Atkeson's PS (M+A's PS) calculates the priority of some state by checking all transitions to updated successor states and identifying the one whose update contribution is largest. Our variant allows for computing the *exact* size of updates of state values since they have been used for updating the Q-values of their predecessors, and yields more appropriate priorities. Unlike our PS, M+A's PS cannot detect large state-value changes due to many small update steps, and will not process the corresponding states. A complete description of both algorithms is given in [35].

Our implementation for CMAC models uses a set of predecessor lists  $Preds_k(j)$  containing all predecessor cells of cell  $j$  in filter  $k$ . We denote the priority of cell  $i$  of filter  $k$  by  $|\Delta_k(i)|$ , where the value  $\Delta_k(i)$  equals the change of  $V_k(i)$  since the last time it was processed by the priority queue. To calculate it, we constantly update all Q-values of predecessor cells of currently processed cells, and track changes of  $V_k(i)$ .

The model-based update of the Q-value  $Q_k(c, a)$ , **Q-update**( $k, c, a$ ) looks as follows:

$$Q_k(c, a) \leftarrow \sum_j P_{c_j}^k(a)(R_k(c, a, j) + \gamma V_k(j)),$$

where  $P_{c_j}^k(a) = P_k(j|c, a)$ . The details of our PS look as follows:

1. **Our-Prioritized-Sweeping**( $\mathbf{x}$ ):
2.     Compute active cells:  $f_1, \dots, f_z$ ;
3.     For  $k = 1$  to  $z$  do:
4.         Update  $f_k$  —  $\forall a$  do:
5.             **Q-update**( $k, f_k, a$ );
6.             Set  $|\Delta_k(f_k)|$  to  $\infty$ ;
7.             Promote  $(k, f_k)$  to top of queue;
8.             While ( $n < U_{max}$  & queue  $\neq$  nil)
9.                 Remove top  $(k, c)$  from the queue;
10.                  $\Delta_k(c) \leftarrow 0$ ;
11.                  $\forall$  Predecessor cells  $k, i$  of  $k, c$  do:
12.                      $V'_k(i) \leftarrow V_k(i)$ ;
13.                      $\forall a$  do:
14.                         **Q-update**( $k, i, a$ );
15.                          $V_k(i) \leftarrow \max_a Q_k(i, a)$ ;
16.                          $\Delta_k(i) \leftarrow \Delta_k(i) + V_k(i) - V'_k(i)$
17.                         If  $|\Delta_k(i)| > \epsilon$
18.                             Insert  $i$  at priority  $|\Delta_k(i)|$ ;
19.                      $n \leftarrow n + 1$ ;
20.             Empty queue, but keep  $\Delta_k(i)$  values;

---

Here  $U_{max}$  is the maximal number of updates to be performed per update-sweep. The parameter  $\epsilon \in \mathbb{R}^+$  controls update accuracy. Note that another difference to M+A's PS is that we remove all entries from the queue after having processed all updates.

### Appendix Non-Pessimistic Value Functions

To compute non-pessimistic value functions we decrease the probability of the worst transition from each filter/cell/action and then renormalize the other probabilities. Then we use the adjusted probabilities to compute the Q-functions. Thus we substitute the following for **Q-update**( $k, c, a$ ):

1. **Q-update-Non-Pessimistic**( $\mathbf{k}, \mathbf{i}, \mathbf{a}$ ):
2.      $m \leftarrow \arg \min_j \{R_k(i, a, j) + \gamma V_k(j)\}$ ;
3.      $n \leftarrow C_i^k(a)$ ;
4.      $P \leftarrow \hat{P}_{im}^k(a)$ ;
5.      $P_{im}^k(a) \leftarrow \frac{(P - \frac{z^2}{2n} + \frac{z_n}{\sqrt{n}} \sqrt{P(1-P) + \frac{z^2}{4n}})}{1 + \frac{z^2}{n}}$ ;

6.  $\Delta_P \leftarrow P_{im}^k(a) - \hat{P}_{im}^k(a);$
7.  $\forall j \neq m$
8.  $P_{ij}^k(a) \leftarrow \hat{P}_{ij}^k(a) - \frac{\Delta_P C_{ij}^k(a)}{C_i^k(a) - C_{im}^k(a)};$
9. **Q-update**( $k, i, a$ );

Here  $C_{ij}^k(a)$  counts the number of transitions of cell  $i$  to  $j$  in filter  $k$  after selecting action  $a$  and  $C_i^k(a)$  counts the number of times action  $a$  was selected and cell  $i$  of filter  $k$  was activated. We obtain  $\hat{P}_{ij}^k(a)$ , the estimated transition probability, by dividing them.

The variable  $z_\alpha$  determines the step size for decreasing worst transition probabilities. To select the worst transition in step 2, we only compare existing transitions (we check whether  $\hat{P}_{ij}^k(a) > 0$  holds). Note that if there is only one transition for a given filter/cell/action triplet then there will not be any renormalization. Hence the “probabilities” may not sum up to 1. Consequentially, if some filter/cell/action has not occurred frequently then it will contribute just a comparatively small Q-value and thus have less impact on the computation of the overall Q-value.

## Appendix Q( $\lambda$ )-learning

Q-learning [33, 34] enables an agent to learn a policy by repeatedly executing actions given the current state. At each time step the algorithm uses 1-step lookahead to update the currently selected filter/cell/action pairs (FCAPs):

1. **Q-learning**( $\mathbf{k}, \mathbf{c}_t, \mathbf{a}_t, \mathbf{r}_t, \mathbf{c}_{t+1}$ ):
2.  $e'_t \leftarrow (r_t + \gamma V_k(c_{t+1}) - Q_k(c_t, a_t));$
3.  $Q_k(c_t, a_t) \leftarrow Q_k(c_t, a_t) + \alpha_n(k, c_t, a) e'_t;$

Here  $V_k(c) = \max_a Q_k(c, a)$ ,  $\alpha_n(k, c, a)$  is the learning rate for the  $n^{\text{th}}$  update of FCAP ( $k, c, a$ ), and  $e'_t$  is the temporal difference or TD(0)-error, which tends to decrease over time.

The learning rate  $\alpha_n(k, c, a)$  should decrease online, such that it fulfills two conditions for stochastic iterative algorithms [34, 5]. The conditions on the learning rate  $\alpha_n(k, c, a)$  are:

- (1)  $\sum_{n=1}^{\infty} \alpha_n(k, c, a) = \infty$ , and

- (2)  $\sum_{n=1}^{\infty} \alpha_n^2(k, c, a) < \infty$ .

Learning rate adaptations for which the conditions are satisfied may be of the form:  $\alpha_n = \frac{1}{n^\beta}$ , where  $n$  is a variable that counts the number of times an FCAP has been updated.

Q( $\lambda$ )-learning uses eligibility traces  $l_t(k, c, a)$  [3, 29] to allow for updating multiple FCAPs which have occurred in the past. We use the replacing traces algorithm [28]:

$$\begin{aligned} l_{t+1}(k, c, a) &\leftarrow \gamma \lambda l_t(k, c, a) \text{ if } f_t^k \neq c \\ l_{t+1}(k, c, a) &\leftarrow 1 \text{ if } f_t^k = c \text{ and } a_t = a \\ l_{t+1}(k, c, a) &\leftarrow 0 \text{ if } f_t^k = c \text{ and } a_t \neq a \end{aligned}$$

where  $\lambda$  discounts the influence of FCAPs occurring in the distant future relative to immediate FCAPs. After updating the eligibility traces we update the Q-values:  $\forall (k, c, a)$  do :

$$Q_k(c, a) \leftarrow Q_k(c, a) + \alpha [e'_t \eta_k^t(c, a) + e_t l_t(k, c, a)]$$

where  $\eta_k^t(c, a)$  denotes the indicator function which returns 1 if  $(k, c, a)$  occurred at time  $t$ , and 0 otherwise ( $\alpha = \alpha_n(k, c, a)$ ). The TD-error  $e_t$  of the value function is defined as:  $e_t \leftarrow (r_t + \gamma V_k(c_{t+1}) - V_k(c_t))$ .

The procedure described here updates all occurred FCAPs at each time step. This is computationally expensive. We actually used a faster method which allows for updating Q-values in time proportional to  $O(z|A|)$ , the number of filters times actions [37].

## Acknowledgements

We would like to thank Rich Sutton for inspiring suggestions concerning the potential benefits of CMACs.

## Notes

1. A recent theoretical result [11] suggests that computational complexities of certain direct and indirect methods for MDPs are of the same order. This result, however, is irrelevant for most real world RL applications, because its stringent assumptions are violated by FA-based set-ups such as the one studied in this paper.

## References

1. J. S. Albus. A new approach to manipulator control: The cerebellar model articulation controller (CMAC). *Dynamic Systems, Measurement and Control*, pages 220–227, 1975.
2. S. Baluja and R. Caruana. Removing the genetics from the standard genetic algorithm. In A. Prieditis and S. Russell, editors, *Machine Learning: Proceedings of the Twelfth International Conference*, pages 38–46. Morgan Kaufmann Publishers, San Francisco, CA, 1995.
3. A. G. Barto, R. S. Sutton, and C. W. Anderson. Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-13:834–846, 1983.
4. R. Bellman. *Adaptive Control Processes*. Princeton University Press, 1961.
5. D. P. Bertsekas and J. N. Tsitsiklis. *Neuro-dynamic Programming*. Athena Scientific, Belmont, MA, 1996.
6. D. Chapman and L. P. Kaelbling. Input generalization in delayed reinforcement learning. In *Proceedings of the 13th International Joint Conference on Artificial Intelligence (IJCAI)*, volume 2, pages 726–731. Morgan Kaufman, 1991.
7. N. L. Cramer. A representation for the adaptive generation of simple sequential programs. In J. J. Grefenstette, editor, *Proceedings of an International Conference on Genetic Algorithms and Their Applications*, pages 183–187, Hillsdale NJ, 1985. Lawrence Erlbaum Associates.
8. D. Dickmanns, J. Schmidhuber, and A. Winklhofer. Der genetische Algorithmus: Eine Implementierung in Prolog. Fortgeschrittenenpraktikum, Institut für Informatik, Lehrstuhl Prof. Radig, Technische Universität München, 1986.
9. J. H. Holland. *Adaptation in Natural and Artificial Systems*. University of Michigan Press, Ann Arbor, 1975.
10. L. Kaelbling. *Learning in Embedded Systems*. MIT Press, 1993.
11. M. Kearns and S. Singh. Finite-sample convergence rates for Q-learning and indirect algorithms. In M. Kearns, S. A. Solla, and D. Cohn, editors, *Advances in Neural Information Processing Systems 12*. MIT Press, Cambridge MA, 1999.
12. J. R. Koza. Genetic evolution and co-evolution of computer programs. In C. G. Langton, C. Taylor, J. D. Farmer, and S. Rasmussen, editors, *Artificial Life II*, pages 313–324. Addison Wesley Publishing Company, 1992.
13. L.-J. Lin. *Reinforcement Learning for Robots Using Neural Networks*. PhD thesis, Carnegie Mellon University, Pittsburgh, January 1993.
14. A. Moore and C. G. Atkeson. Prioritized sweeping: Reinforcement learning with less data and less time. *Machine Learning*, 13:103–130, 1993.
15. S. J. Nowlan and G. E. Hinton. Simplifying neural networks by soft weight sharing. *Neural Computation*, 4:173–193, 1992.
16. J. Peng and R.J. Williams. Incremental multi-step Q-learning. *Machine Learning*, 22:283–290, 1996.
17. I. Rechenberg. Evolutionsstrategie - Optimierung technischer Systeme nach Prinzipien der biologischen Evolution. Dissertation, 1971. Published 1973 by Fromman-Holzboog.
18. G. A. Rummery and M. Niranjan. On-line Q-learning using connectionist systems. Technical Report CUED/F-INFENG-TR 166, Cambridge University, UK, 1994.
19. R. P. Salustowicz and J. Schmidhuber. Probabilistic incremental program evolution. *Evolutionary Computation*, 5(2):123–141, 1997.
20. R. P. Salustowicz, M. A. Wiering, and J. Schmidhuber. Evolving soccer strategies. In *Proceedings of the Fourth International Conference on Neural Information Processing (ICONIP'97)*, pages 502–506. Springer-Verlag Singapore, 1997.
21. R. P. Salustowicz, M. A. Wiering, and J. Schmidhuber. On learning soccer strategies. In W. Gerstner, A. Germond, M. Hasler, and J.-D. Nicoud, editors, *Proceedings of the Seventh International Conference on Artificial Neural Networks (ICANN'97)*, volume 1327 of *Lecture Notes in Computer Science*, pages 769–774. Springer-Verlag Berlin Heidelberg, 1997.
22. R. P. Salustowicz, M. A. Wiering, and J. Schmidhuber. Learning team strategies: Soccer case studies. *Machine Learning*, 33(2/3):263–282, 1998.
23. A. L. Samuel. Some studies in machine learning using the game of checkers. *IBM Journal on Research and Development*, 3:210–229, 1959.
24. J. C. Santamaria, R. S. Sutton, and A. Ram. Experiments with reinforcement learning in problems with continuous state and action spaces. Technical Report COINS 96-088, Georgia Institute of Technology, Atlanta, 1996.
25. J. Schmidhuber. On learning how to learn learning strategies. Technical Report FKI-198-94, Fakultät für Informatik, Technische Universität München, November 1994. Revised January 1995.
26. J. Schmidhuber, J. Zhao, and N. Schraudolph. Reinforcement learning with self-modifying policies. In S. Thrun and L. Pratt, editors, *Learning to learn*, pages 293–309. Kluwer, 1997.
27. J. Schmidhuber, J. Zhao, and M. Wiering. Shifting inductive bias with success-story algorithm, adaptive Levin search, and incremental self-improvement. *Machine Learning*, 28:105–130, 1997.
28. S. P. Singh and R. S. Sutton. Reinforcement learning with replacing eligibility traces. *Machine Learning*, 22:123–158, 1996.
29. R. S. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3:9–44, 1988.
30. R. S. Sutton. Generalization in reinforcement learning: Successful examples using sparse coarse coding. In D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo, editors, *Advances in Neural Information Processing Systems 8*, pages 1038–1045. MIT Press, Cambridge MA, 1996.
31. R. S. Sutton and A. G. Barto. *Reinforcement learning: an introduction*. MIT Press/Bradford Books, 1988.
32. S. Thrun, D. Fox, and W. Burgard. A probabilistic approach to concurrent mapping and localization for

- mobile robots. *Machine Learning*, (31):29–53, 1998. Also appeared in *Autonomous Robots* 5, 253–271, 1998 as joint issue.
33. C. J. C. H. Watkins. *Learning from Delayed Rewards*. PhD thesis, King’s College, Cambridge, England, 1989.
  34. C. J. C. H. Watkins and P. Dayan. Q-learning. *Machine Learning*, 8:279–292, 1992.
  35. M. A. Wiering. *Explorations in Efficient Reinforcement Learning*. PhD thesis, University of Amsterdam / IDSIA, February 1999.
  36. M. A. Wiering and J. Schmidhuber. Efficient model-based exploration. In J. A. Meyer and S. W. Wilson, editors, *Proceedings of the Sixth International Conference on Simulation of Adaptive Behavior: From Animals to Animats 6*, pages 223–228. MIT Press/Bradford Books, 1998.
  37. M. A. Wiering and J. Schmidhuber. Fast online  $Q(\lambda)$ . *Machine Learning*, 33(1):105–116, 1998.