

# Biasing MCTS with Features for General Games

Dennis J. N. J. Soemers, Éric Piette, and Cameron Browne

*Department of Data Science and Knowledge Engineering*

*Maastricht University*

Maastricht, the Netherlands

{dennis.soemers,eric.piette,cameron.browne}@maastrichtuniversity.nl

**Abstract**—This paper proposes using a linear function approximator, rather than a deep neural network (DNN), to bias a Monte Carlo tree search (MCTS) player for general games. This is unlikely to match the potential raw playing strength of DNNs, but has advantages in terms of generality, interpretability and resources (time and hardware) required for training. Features describing local patterns are used as inputs. The features are formulated in such a way that they are easily interpretable and applicable to a wide range of general games, and might encode simple local strategies. We gradually create new features during the same self-play training process used to learn feature weights. We evaluate the playing strength of an MCTS player biased by learnt features against a standard upper confidence bounds for trees (UCT) player in multiple different board games, and demonstrate significantly improved playing strength in the majority of them after a small number of self-play training games.

**Index Terms**—games, features, search, learning

## I. INTRODUCTION

Combinations of search algorithms with learning from self-play have led to strong results in game-playing AI for various board games [1]–[3] and video games [4]. Currently, a common combination is to use *Monte Carlo tree search* (MCTS) [5]–[7] for search and *deep neural networks* (DNNs) for learning.

While DNN-based learning approaches have advanced the state of the art in terms of raw playing strength for game-playing AI, they have a number of disadvantages in comparison to other learning techniques in other aspects. For example, effectively training the large DNNs used to obtain state-of-the-art results in board games [1]–[3] requires significant training time and/or large amounts of hardware. Different games typically require different DNN architectures – the input and output layers in particular are game-specific – and separate training processes starting from scratch per game. Due to the large number of trainable parameters, and the use of low-level inputs (i.e. raw board states), it is often difficult to extract interpretable knowledge such as strategies from a DNN during or after training.

In this paper, we describe and evaluate an approach to simultaneously grow a set of features, and learn weights for a linear policy function using those features, from self-play. Each feature [8] describes a simple local pattern, and is specified in a general manner that is applicable to many different games and easily interpretable.

Funded by a €2m ERC Consolidator Grant (<http://ludeme.eu>).

The focus of this work is not on achieving state-of-the-art game-playing performance. Our main contribution is to demonstrate that, using a simple linear policy function, our learned features and weights can improve the performance of a standard MCTS player across a variety of board games given a set of severe restrictions:

- Features are specified in a general format, compatible with many different games [8].
- We do not manually construct game-specific feature sets.
- We use a low amount of training time; up to 300 games of self-play per game, with 5 seconds of thinking time allowed per move.
- We use relatively little hardware (every sequence of self-play runs sequentially on a single node).
- The learned function must be able to work correctly for any number of legal actions (i.e., we do not provide an upper bound on the size of the action space, as is typically required for the output layer in DNN-based approaches).
- Hyperparameters (for MCTS as well as the self-play learning process) have not been optimised (their values are manually selected), and the same hyperparameter values are used across all games.

Under the restrictions listed above, we find that the performance of MCTS can already be improved relatively easily in multiple different games. We expect that the interpretability of the features, in combination with a simple linear function, could provide insight into what strategies are relevant to the games being modelled, potentially giving some insight into their strategic potential and the relationships between different games in terms of strategy [8].

The generality of the approach may also create opportunities for transferring of strategies between games, which combined with the low computational requirements could allow the approach to be applied to large numbers of games and variants within reasonable times. This is crucial for the Digital Ludeme Project [9], which involves digitally modelling large numbers of ancient games and exploring relationships between them.

Section II provides relevant background information for this paper. The self-play training process used to learn weights for a linear function is described in Section III. Our approach for automatically constructing and growing the set of features per game is explained in Section IV. Section V discusses the experimental evaluation of the approach. The paper is concluded in Section VI.

```

(game "Tic-Tac-Toe"
  (play { (player "P1") (player "P2") })
  (equipment
    { (board "Board" (square 3)) }
    { (disc "Piece") (cross "Cross") }
  )
  (rules
    (moves (to Mover (empty)))
    (end (line length:3) (result Mover win))
  )
)

```

Fig. 1. The game of Tic-Tac-Toe modelled in LUDII.

## II. BACKGROUND

This section provides background information on the LUDII general game system [9], which we use to run different games used for evaluation, and the format in which we specify general game features [8].

### A. The LUDII General Game System

The LUDII *general game system* [9] implements units of game-related information, referred to as *ludemes* [10], in different Java classes. A single ludeme can, for example, be a simple game rule that describes a particular kind of legal move, or it can be a description of a component (a piece or a board) used to play a game. A complete game can be described as a tree of ludemes. For example, Fig. 1 depicts how the game of Tic-Tac-Toe can be modelled in LUDII.

### B. Specification of Features

In *general game playing* (GGP) research based on the Stanford *game definition language* (GDL) [11], various approaches have been proposed for using general game features [12]–[16]. These are all built around the logic-based formalism of GDL, and therefore not directly applicable to general game systems that use a different language (such as LUDII).

Following [8], the features  $x$  used in this paper consist of:

- 1) A *pattern*  $p_x$ , which contains one or more elements that the feature tests for in relative positions. The different types of elements that a feature can test for in different positions are *off-board*, *empty*, *friendly piece*, *enemy piece*, *piece owned by player  $n$*  (for any  $n$ ), and *piece with unique index  $n$*  (in the game’s definition).
- 2) A description of an *action*  $a_x$  which the feature “recommends” playing (in practice a feature may also discourage playing its action if a negative weight is learned for that feature). The action  $a_x$  is specified by two relative positions; a position to move “from” and a position to move “to”. For games like Hex and Go, only the “to” position is relevant. For games such as Draughts and Chess, also the “from” position is relevant.

Relative positions in patterns and action specifications of features are described as sequences of numbers, referred to as *walks*. The length of such a sequence corresponds to the

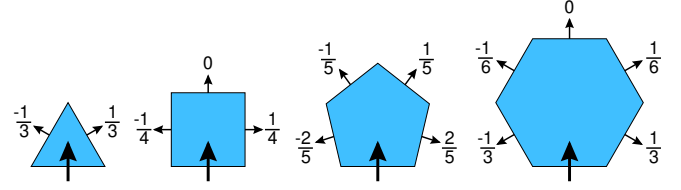


Fig. 2. Fractional steps and their corresponding relative steps.

number of steps to take through a graph representation of the playable area (typically a board), starting from some fixed *anchor position*. Every number in the sequence denotes how far we should rotate (as a fraction of a full 360°, clockwise turn), relative to the “current” direction, before taking the next step. Fig. 2 depicts different fractions and their corresponding movement directions in different types of cells, relative to a “current” direction pointing to the north.

Features defined in this manner are applicable to any game that can be modelled as being played on some graph, where vertices may contain pieces that may be owned by players. This applies to many board games, and potentially also games without an explicit board. The LUDII system uses such a graph to model the playable area of any game. Every vertex contains a list of references to adjacent vertices, sorted in a consistent manner to facilitate indexing based on clockwise turns, with null entries to facilitate off-board checks. Fig. 3 depicts how a  $\{0, 0, \frac{1}{4}\}$  walk can specify relative positions with a similar semantic meaning in two different types of boards.

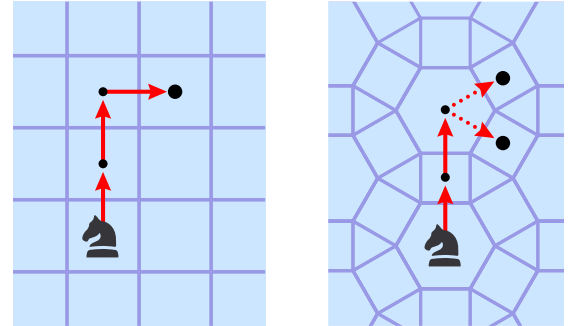


Fig. 3. Relative position(s) reached by a  $\{0, 0, \frac{1}{4}\}$  walk in two different boards. In the board on the right-hand side, the  $\frac{1}{4}$  turn can be rounded to a turn of  $\frac{1}{3}$  or  $\frac{1}{6}$ .

## III. EXPERT ITERATION WITH A LINEAR POLICY

In the *expert iteration* framework [1], [2], an *apprentice* policy and an *expert* policy are used to iteratively improve each other during self-play. The apprentice policy is a trainable, computationally efficient component, such as a DNN. Given any state  $s$ , it can compute a distribution  $p(s)$  over the set of actions  $A(s)$  that are legal in  $s$  in a fixed amount of time. The expert is generally a component that involves more “deliberation” (i.e., search or planning), such as MCTS.

The main idea of expert iteration is to view the expert policy as a *policy improvement operator* for the apprentice.

The apprentice can be used to bias the searching behaviour of the expert, and the expert can use additional computation time to adjust (ideally improve) the distribution computed by the apprentice. The adjusted distribution can subsequently be used as a learning target by the apprentice.

#### A. Formalisation of the Apprentice

The particular combination of a DNN as apprentice, and MCTS as expert, has led to state-of-the-art game-playing performance in various board games [1]–[3], but the DNNs have a number of important drawbacks in aspects other than raw playing strength, as listed in Section I. Therefore we investigate using a linear function rather than a DNN in this paper. We aim to train a function of the form given by (1):

$$f(s, a) = \theta^\top \phi(s, a), \quad (1)$$

where  $\theta$  is a vector of trainable parameters,  $\phi(s, a)$  is a (binary) feature vector for a state-action pair  $(s, a)$ , and  $f(s, a) \in \mathbb{R}$  is a real-valued output for the same state-action pair. Given a set of legal actions  $A(s)$  in a state  $s$ , the complete distribution  $p(s)$  over all actions  $a_i \in A(s)$  is computed by applying the softmax function to a vector of outputs  $f(s, a_i)$ :

$$p_i(s) = \frac{\exp(f(s, a_i))}{\sum_{k=1}^{|A(s)|} \exp(f(s, a_k))}. \quad (2)$$

When DNNs are used as apprentice, it is customary to have an output layer with one output node per unique action that may ever be legal in any given game state. This can easily lead to excessively large numbers of outputs in some games, such as 11,259 outputs in Shogi [3]. It also requires domain knowledge in the form of an accurate upper bound on the number of unique actions, which is a problem in terms of generality. The number of outputs computed in any given state  $s$  by (2) is equal to the number of legal actions  $|A(s)|$  in that state, which is typically multiple orders of magnitude lower in a game like Shogi [17]. Of course, an advantage of DNNs may be that its computational requirements remain constant regardless of  $|A(s)|$ , whereas the computational requirements of (2) scale linearly with the number of legal actions  $|A(s)|$ .

#### B. Formalisation of Feature Vectors

Suppose that we have some set of features  $\mathcal{X}$ , where every feature  $x \in \mathcal{X}$  is specified as described in Subsection II-B. More details on how such a feature set is created will be provided in the next section. We say that a feature  $x$  is *active* for a state-action pair  $(s, a)$  if there exists some *anchor* position in the game's underlying graph such that, after applying any necessary rotation and/or reflection:

- 1) The feature's action  $a_x$  corresponds to the action  $a$ .
- 2) All elements of the feature's pattern  $p_x$  are satisfied in the game state  $s$ .

Relative positions in  $a_x$  and  $p_x$  are evaluated by resolving the walks, starting from the anchor position.

We define state-action feature vectors  $\phi(s, a)$  as binary vectors that contain a value of 1 for features  $x$  that are active for the state-action pair  $(s, a)$ , and a value of 0 for all

other features. Note that different *instantiations* (with different anchor positions, rotations, or reflections) of the same feature may be active simultaneously; in such a case, we still simply assign a feature value of 1.

#### C. Guiding the Expert using the Apprentice

We use a learned apprentice policy to guide the expert (MCTS) in its *selection* and *play-out* phases. The most common selection strategy [5] is to follow the UCB1 policy [18]. Given a current node with a state  $s$ , it selects the child node corresponding to the action  $a_{ucb1}$  given by (3).

$$a_{ucb1} = \operatorname{argmax}_a \hat{Q}(s, a) + C_{ucb1} \sqrt{\frac{\ln(\sum_{a'} N(s, a'))}{N(s, a)}}. \quad (3)$$

$\hat{Q}(s, a)$  denotes the estimated value of playing  $a$  in  $s$  based on previous MCTS iterations (i.e. the average score backpropagated through the node reached by executing  $a$  in  $s$ ),  $C_{UCB1}$  is a hyperparameter (the “exploration constant”), and  $N(s, a)$  denotes the *visit count* of  $(s, a)$  (the number of previous MCTS iterations that have selected  $a$  in the current node). The sum  $\sum_{a'} N(s, a')$  is equivalent to the total number of previous MCTS iterations that have reached the current node. This strategy only uses statistics gathered by MCTS itself, and does not use the apprentice.

In this paper, we use the apprentice to guide the selection step of MCTS using the same strategy as AlphaGo Zero [1], which selects the action  $a_{puct}$  given by (4).

$$a_{puct} = \operatorname{argmax}_a \hat{Q}(s, a) + C_{puct} p(s, a) \frac{\sqrt{\sum_{a'} N(s, a')}}{1 + N(s, a)}. \quad (4)$$

$C_{puct}$  is a hyperparameter, and  $p(s, a)$  denotes the probability of selecting  $a$  in  $s$  according to the distribution computed by the apprentice.

The most straightforward approach for using the apprentice in the play-out step is to select actions according to the probability distribution computed by the apprentice, rather than selecting actions uniformly at random. The computational overhead of computing feature vectors can be mitigated by transitioning from using the apprentice policy early in play-outs to a uniform random policy in later parts of play-outs.

#### D. Training Apprentice with Expert Iteration

We train the apprentice policy in a similar way to [1], interpreting the visit counts at the end of an MCTS search process as the target distribution. Let  $N(s, a)$  denote the number of MCTS iterations that selected action  $a$  in state  $s$ . For any node in the search tree with a state  $s$ , the quantity  $\frac{N(s, a)}{\sum_{a'} N(s, a')}$  can then be interpreted as the probability assigned to  $a$  in state  $s$  by the MCTS expert policy.

Let  $\pi(s)$  denote a vector of such probabilities for all legal actions  $a \in A(s)$ , and let  $p(s)$  denote a similar vector computed by the apprentice policy. The loss function is then given by (5):

$$\mathcal{L}(s) = -\pi(s)^\top \log p(s) + \frac{\lambda}{2} \|\theta\|^2, \quad (5)$$

which computes the cross-entropy loss between the two distributions, and an  $L_2$  regularisation penalty with a hyperparameter  $\lambda$ . A stochastic gradient descent (SGD) update to reduce this loss, based on a single example  $s$ , can be implemented according to (6):

$$\theta \leftarrow \theta - \alpha \sum_{a \in A(s)} [(p(s, a) - \pi(s, a)) \times \phi(s, a)] - \alpha \lambda \theta, \quad (6)$$

where  $p(s, a)$  and  $\pi(s, a)$  denote the entries corresponding to action  $a$  in the  $\mathbf{p}(s)$  and  $\boldsymbol{\pi}(s)$  vectors, respectively, and  $\alpha$  is a step-size hyperparameter.

#### IV. GROWING FEATURE SET DURING SELF-PLAY

In prior research using similar types of features to those used in this paper for game-playing AI applications, the complete set of features to use is typically determined before any parameters for policies or value functions using those features are learned.

For example, large sets of features were exhaustively generated in the games of Go [19], Hearts [20], and Breakthrough [21]. Such exhaustive sets can easily contain tens of thousands of features. When features are used in a pure Reinforcement Learning agent [20], or implemented for a single specific game [19], [21] which permits a highly efficient game-specific implementation (for instance by incrementally updating the set of active features [19]), it is computationally tractable to use such large feature sets. In our general game system, we find that the computational overhead of computing active features already becomes detrimental to game-playing performance in some games for feature sets containing only hundreds of features (see Section V).

Other approaches [22] for discovering sets of useful features often consist of iteratively modifying feature sets (e.g., by adding new features in some manner), and evaluating the usefulness of such modifications by running a number of evaluation games using the features. This can be slow when many evaluation games are required for an accurate evaluation. In this section, we propose an approach that starts with a small set of initial features, and gradually adds more complex features during the self-play expert iteration process, guided by the same loss function used to train the policy.

##### A. Initial Feature Set

Many ludemes used by the LUDII system to describe legal moves can easily be extended with functionality to generate patterns  $p_x$  or (parts of) features  $x$  that detect legal moves. For example, the “(to Mover (empty))” ludeme, used in Tic-Tac-Toe (Fig. 1) as well as many other games (such as Hex, Yavalath, etc.), can generate a feature  $x$  which:

- Recommends playing in the feature’s *anchor* position.
- Has a pattern  $p_x$  that requires the same anchor position (specified using a 0-length walk) to be empty.

Formally, we implement ludemes used in the specification of move rules to generate a feature set  $\mathcal{X}$  such that, for any possible game state  $s$  and legal action  $a \in A(s)$ , there exists at least one feature  $x \in \mathcal{X}$  that is active for the state-action pair

$(s, a)$ . In the worst case (for highly complex movement rules), this can simply be an “empty” feature without any restrictions on either the action or the game state (i.e. a feature that is always active).

Such features are generally not interesting features by themselves. However, they can be used to reduce the space of candidate features considered for subsequent addition to the feature set. We only allow new features to be added to the feature set if they are at least as restrictive as, and not incompatible with, at least one of the features generated by movement rule ludemes. This enables us to automatically ignore many features that would be useless due to never being active in gameplay. For example, features that require an enemy or friendly piece in the location where they recommend playing will never be considered in games that only permit playing in empty positions.

For every feature generated by move rule ludemes as described above, we generate a set of “atomic” features  $x$ , which have exactly one requirement for an element specified in their pattern  $p_x$  (in addition to any requirements that may already be there due to move rule ludemes). We exhaustively generate all such atomic features, with generated walks restricted to at most two steps. Using only atomic features (no features with more complex patterns) keeps the initial feature count down to a low number. We use the maximum number of adjacencies of any vertex in a game’s graph to determine the number of potentially meaningful rotations in a game.

##### B. Adding New Features During Expert Iteration

In the expert iteration framework, experience generated from self-play is used to update the parameters  $\theta$  of the apprentice policy, such that its output distributions  $\mathbf{p}$  more closely match the distributions  $\boldsymbol{\pi}$  of the expert policy. The error  $p(s, a) - \pi(s, a)$ , which also appears in the SGD update rule in (6), has a large absolute value for state-action pairs  $(s, a)$  for which the distributions do not yet closely match, and a low absolute value if the distributions already closely match.

We propose to use this error value as an indicator of state-action pairs  $(s, a)$  for which it is beneficial to add new features to the feature set. The intuition is that there is no need to add extra features for  $(s, a)$  pairs for which  $p(s, a)$  already closely matches  $\pi(s, a)$ , but extra features are more likely to be useful if they activate for  $(s, a)$  pairs for which  $p(s, a)$  and  $\pi(s, a)$  do not yet closely match.

Whenever we wish to add a new feature to the feature set, we sample a batch  $E = \{(s_i, A(s_i), \pi_i(s_i))\}$  of samples of experience collected from self-play. Every tuple  $(s_i, A(s_i), \pi_i(s_i))$  in this batch contains a game state  $s_i$  encountered in self-play, the list of legal actions  $A(s_i)$  in that state, and the distribution  $\pi_i(s_i)$  over the actions  $A(s_i)$  as computed by the expert at the point in time when this experience was saved. This batch is sampled without replacement from a larger experience buffer, in which we store exactly one new sample of experience (corresponding to the current game state) for every game state encountered during self-play.

Every pair of two features instances  $(x_i, x_j)$  that are active together for at least one state-action pair across the entire batch  $E$  is taken into consideration as a candidate pair that could be combined into a single new feature  $x_i x_j$ . Such a combination  $x_i x_j$  is a new feature in which the patterns of the constituents  $x_i$  and  $x_j$  are merged. Note that a combined feature  $x_i x_j$  will always be active for state-action pairs in which  $x_i$  and  $x_j$  were both active.

The candidate pair  $(x_i, x_j)$  that maximises the score given by (7) is added to the feature set as a new feature.

$$\text{score}(x_i, x_j) = |r_{err}(x_i, x_j)| \times (1 - |r_{x_i x_j}(x_i, x_j)|) \quad (7)$$

In this equation,  $r_{err}(x_i, x_j)$  denotes the Pearson correlation coefficient between errors  $p(s, a) - \pi(s, a)$ , and the event of simultaneously observing features  $x_i$  and  $x_j$  to be active for a state-action pair  $(s, a)$ . Similarly,  $r_{x_i x_j}(x_i, x_j)$  denotes the Pearson correlation coefficient between the event of simultaneously observing features  $x_i$  and  $x_j$  to be active, and the event of observing one of the constituents  $x_i$  or  $x_j$  to be active (whichever constituent leads to the strongest correlation is picked). Correlation coefficients are measured across all state-action pairs that occur in the complete batch  $E$ .

This score implements the intuition that new features are likely to be useful if they correlate strongly with observed errors between the apprentice and expert distributions, but are less likely to be useful if their activations correlate strongly with the activations of other features in the feature set. Similar intuition has also been shown to be useful for offline feature selection in supervised machine learning [23]. Ideally we would minimise correlations of candidate features between *all* features in a feature set, but this is computationally expensive. Only computing correlations between candidate features and their constituents is significantly cheaper.

## V. EXPERIMENTS

This section describes a number of experiments which evaluate the effect of features and weights learned from self-play on the game-playing performance of MCTS in a variety of board games.

### A. Setup

In the self-play process of expert iteration, experience is generated by equivalent MCTS agents playing against each other. They use (4) in the selection phase, with  $C_{puct} = \sqrt{2}$ . The first move of every play-out is sampled from the apprentice distribution  $p$ , and the corresponding node is added to the search tree. Subsequent play-out moves are selected uniformly at random, to avoid additional computational overhead of computing active features. Final moves for the “real” games are sampled from the expert distribution  $\pi$ . We use 5 seconds of “thinking time” per move. Games are terminated automatically after 100 moves, regardless of the game’s standard rules.

We use an experience buffer with a maximum capacity of 200 to store tuples of experience. Every move played in self-play results in one new tuple of experience. Old tuples are removed to make room for new tuples if necessary. We run

one SGD update to update the apprentice parameters  $\theta$  after every move, with a step-size  $\alpha = 0.05$ , and  $\lambda = 10^{-6}$  for  $L_2$  regularisation. Gradients are computed and averaged across a batch of size 20, sampled from the experience buffer.

We add one new feature after every game of self-play. A larger batch size of 30 is used in this procedure. We evaluate three simpler feature discovery strategies, in addition to the correlation-based variant described in detail in Section IV:

- 1) *Add Random*: This variant randomly selects pairs of simultaneously activated feature instances to combine. This can be viewed as an unguided baseline strategy.
- 2) *Combine Random*: Randomly combines two feature instances that activate together in the state-action pair  $(s, a)$  that maximises the absolute error  $|p(s, a) - \pi(s, a)|$ .
- 3) *Combine Max*: Combines two feature instances that activate together in the state-action pair  $(s, a)$  that maximises the absolute error  $|p(s, a) - \pi(s, a)|$ , such that one of them has the greatest absolute weight in the vector  $\theta$ , and the other is selected randomly.
- 4) *Correlation-based*: Combines feature instances such that (7) is maximised, as described in Section IV.

We evaluate the performance of a *Biased MCTS* agent (using features and weights learned from self-play) against a standard *upper confidence bounds for trees (UCT)* agent:

- *Biased MCTS*: Equal to the agent used during self-play, except that it selects moves with maximum visit counts rather than sampling moves from the  $\pi$  distribution during evaluation games.
- *UCT*: Uses (3) in the selection phase, with  $C_{ucb1} = \sqrt{2}$ . Selects moves uniformly at random in play-outs. Plays moves with maximum visit counts in evaluation games.

All MCTS agents (in self-play as well as evaluation games) use relevant parts of the search tree built up when searching for previous moves to initialise the search tree for subsequent moves. Just like self-play games, evaluation games allow for 5 seconds of thinking time per move, and are automatically declared a tie after 100 moves.

We use nine different board games with standard board sizes, all implemented in the LUDII system, for evaluation. For the game of Hex, we use a  $7 \times 7$  board in addition to the standard  $11 \times 11$  board.

### B. Results - Growing Feature Set

Fig. 4 depicts learning curves for the four different feature discovery strategies, for all ten (variants of) games. At different checkpoints (after 1, 25, 50, 100, and 200 games of self-play), we play 200 evaluation games where the *Biased MCTS* agent plays against the benchmark *UCT* agent, using the latest feature set and learned weights available at that checkpoint. Fig. 4 depicts 95% confidence intervals for the win percentage of *Biased MCTS* against *UCT* at every checkpoint. Ties count as half a win for each player.

*Biased MCTS* can quickly learn to outperform *UCT* in the majority of games; by a significant margin in Breakthrough,

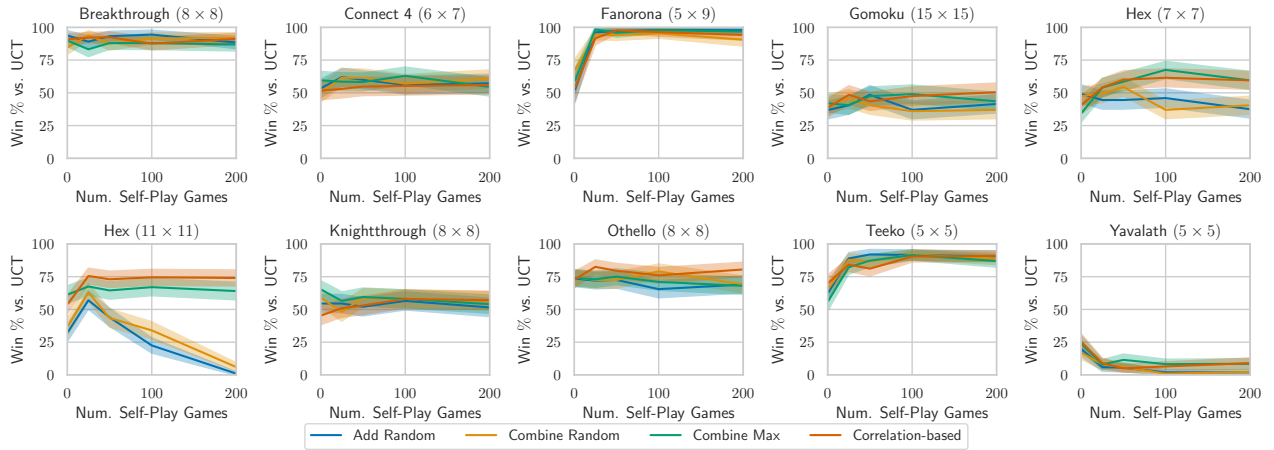


Fig. 4. Learning curves for four different feature discovery strategies, over 200 games of self-play. Shaded regions depict 95% confidence intervals for the win percentage of *Biased MCTS* vs. *UCT*. Performance evaluated by playing 200 evaluation games using feature sets and learned weights after 1, 25, 50, 100, and 200 games of self-play.

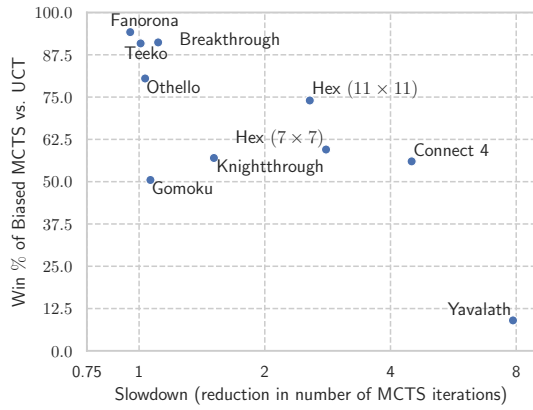


Fig. 5. Relation between win percentage of *Biased MCTS* vs. *UCT* after 200 games of self-play, and the slowdown (reduction in MCTS iteration count) due to the computational overhead of using features.

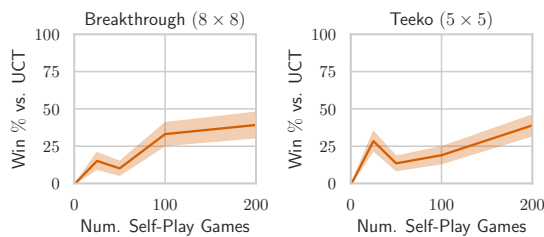


Fig. 6. Learning curves for greedy linear agent (without any search) vs. *UCT* in the games of Breakthrough and Teeko.

Fanorona, Hex ( $11 \times 11$ ), Othello, and Teeko, and by a smaller margin – but often still statistically significant for the strongest variants – in Connect 4, Hex ( $7 \times 7$ ), and Knightthrough. Note that in some of these games there already is a significant gain in playing strength from just a single game of self-play, but in others there is a clear benefit in training for a longer time and growing the feature set. In Gomoku there is no apparent change in playing strength, and in Yavalath the playing

strength is reduced by a significant margin. In most games the simpler feature discovery strategies already perform well, but we find the correlation-based feature discovery strategy to be better in some games, and never significantly worse.

Fig. 5 depicts the relation between the performance of *Biased MCTS* after 200 games of self-play against *UCT* (with win percentage on the  $y$ -axis), and the slowdown due to computing active features (reduction in number of MCTS iterations on the  $x$ -axis). The slowdown per game is computed as  $\frac{I_{uct}}{I_{biased}}$ , where  $I_{uct}$  and  $I_{biased}$  denote the average number of complete MCTS iterations performed by *UCT* and *Biased MCTS*, respectively, in the first two moves per game. We only take into account the first two moves per game, because later moves can have wildly varying iteration counts depending on the game state. These results are given for the feature set learned using the *Correlation-based* strategy. Computing features leads to the worst slowdown in Yavalath (an 8 times reduction in MCTS iteration count), which may explain the poor performance in terms of win percentage in that game. In general, the computational overhead tends to be most noticeable in games for which the game implementation itself is highly efficient in LUDII, and hardly noticeable in games where the game logic itself requires more computation.

Fig. 6 depicts learning curves for a greedy linear agent, which greedily plays actions  $a \in A(s)$  such that  $p(s, a)$  is maximised without performing any tree search at all, for the games of Breakthrough and Teeko. Surprisingly, we find that 200 games of self-play is already sufficient in these games to train a greedy agent that is competitive (reaching a win percentage of 40%) against *UCT*. Learning curves for other games (in which a simple greedy player still has a win percentage close to 0% against *UCT*, as we would expect) are omitted to save space.

### C. Results - Pruned Feature Set

To reduce the overhead of computing active features, we pruned the feature sets of all games by keeping only the 15

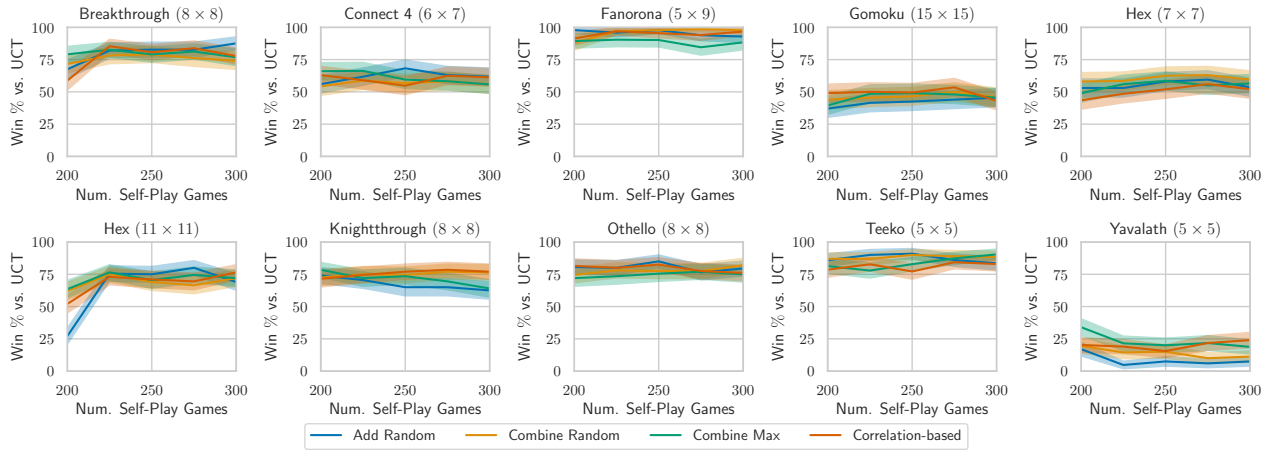


Fig. 7. Learning curves for pruned features sets, over 100 additional games of self-play. Shaded regions depict 95% confidence intervals for the win percentage of *Biased MCTS* vs. *UCT*. Performance evaluated by playing 200 evaluation games using weights learned after 0, 25, 50, 75, and 100 games of self-play (after the first sequence of 200 self-play games).

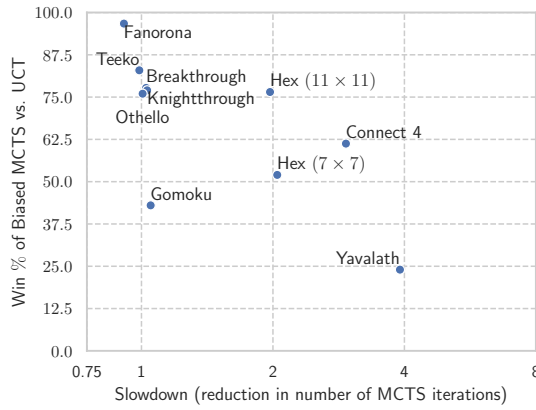


Fig. 8. Relation between win percentage of *Biased MCTS* vs. *UCT* using pruned feature sets, and the slowdown (reduction in MCTS iteration count) due to the computational overhead of using features.

features per game with the greatest absolute weights in the learned parameter vectors  $\theta$ . Starting with the weights learned from the initial 200 games of self-play, we run 100 additional games of self-play to adjust the weights (for which different values may be better now that many other features have been pruned), but freeze the feature set. Fig. 7 depicts learning curves where performance is evaluated for 0, 25, 50, 75, and 100 additional self-play games after pruning the feature set. Note that the four different feature discovery strategies may only have influence on the feature set and initial weights in this figure; they do not matter otherwise because no more features are added during these self-play games.

The relations between playing strength and slowdowns in MCTS iteration counts with pruned feature sets are depicted in Fig. 8. In comparison to the unpruned feature sets of Fig. 5, we observe the most significant changes in performance for the games of Knightthrough (where *Biased MCTS* now has a significant advantage over *UCT*), and Yavalath (where slowdowns and win percentage for *Biased MCTS* have clearly

been improved, but playing strength is still worse than *UCT*'s).

#### D. Interpreting Features in Yavalath

In the game of Yavalath, players win by making a line of four of their pieces, but lose by making a line of three of their pieces beforehand. Given these rules, it is relatively easy to construct useful features by hand. For example, Fig. 9 depicts three features that activate for moves that result in instant wins or losses. We manually constructed a feature set with only these three features, and manually assigned large weights; +3000 for the win-detecting feature, and -1000 for each of the loss-detecting features. A *Biased MCTS* player using this feature set throughout complete play-outs achieves a win percentage of 93% against *UCT*, despite a 30 $\times$  reduction in the MCTS iteration count. This indicates that the poor performance of *Biased MCTS* in Yavalath is not due to a lack of expressiveness in the feature formalisation, but rather due to poor features and/or weights being learned from self-play.

Fig. 10 depicts two of the most “important” features (with large absolute weights) found from self-play in Yavalath. These are easy to understand and seem sensible. The first feature recommends making a move to prevent the opponent from winning in their next turn. The second feature is very similar to the handcrafted win-detecting feature in Fig. 9, with the difference being that it has a seemingly unnecessary requirement for an adjacent opposing piece. This feature can still often detect immediate wins, but not all of them. Despite the poor performance in terms of win percentage in Yavalath, it appears that our proposed approach has still learned sensible – if not optimal – features in this game.

## VI. CONCLUSION

This paper describes and evaluates an approach for simultaneously learning a set of features, and corresponding weights for a linear policy function, for general games implemented in the LUDII system. The features are formalised in such a way that they are generally applicable, and easily interpretable.



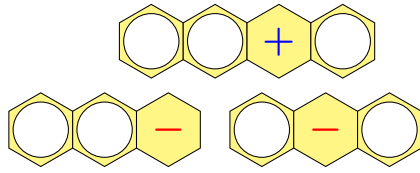


Fig. 9. Immediate win and loss features for the White player in Yavalath.

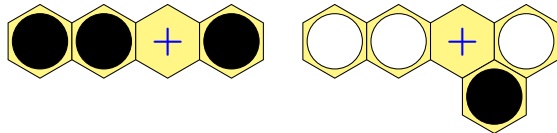


Fig. 10. Two of the learned Yavalath features with greatest absolute weights (drawn from the perspective of the White player).

The training process is also more easily applicable to general games than, for instance, Deep Neural Networks, which require game-specific knowledge to determine the numbers of input and output nodes a priori.

Using the learned features and weights to bias an MCTS agent, we demonstrate significantly improved game-playing performance over a standard UCT agent in the majority of evaluated board games. This performance is achieved with relatively few self-play games. Out of ten evaluated games, the use of features only reduced playing strength in the game of Yavalath due to computational overhead. Despite the poor performance in this game (which may also lead to poor update targets during self-play), a manual inspection of the top features learned in this game indicates that the approach still discovers sensible features.

In future research, we aim to investigate more approaches for improved feature discovery from self-play. In particular in the game of Yavalath, tests with handcrafted features indicate that it may be useful to pay extra attention to features for endgame positions. Using optimisers with momentum-based terms, rather than a simple Stochastic Gradient Descent optimiser, may enable more rapid learning of large feature weights for features that reliably detect immediate winning or losing moves. We also aim to explore transfer learning between different games, which is already facilitated by the feature representation which is shared across all games, and online fine-tuning of trained policies [24].

#### ACKNOWLEDGMENT

This research is part of the European Research Council-funded Digital Ludeme Project (ERC Consolidator Grant #771292) run by Cameron Browne at Maastricht University's Department of Data Science and Knowledge Engineering.

#### REFERENCES

- [1] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and D. Hassabis, "Mastering the game of Go without human knowledge," *Nature*, vol. 550, pp. 354–359, 2017.
- [2] T. Anthony, Z. Tian, and D. Barber, "Thinking fast and slow with deep learning and tree search," in *Adv. in Neural Inf. Process. Syst.*, 30, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., Curran Associates, Inc., 2017, pp. 5360–5370.
- [3] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, T. Lillicrap, K. Simonyan, and D. Hassabis, "A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play," *Science*, vol. 362, no. 6419, pp. 1140–1144, 2018.
- [4] D. R. Jiang, E. Ekwedike, and H. Liu, "Feedback-based tree search for reinforcement learning," in *Proc. 35th Int. Conf. Mach. Learn.*, ser. Proceedings of Machine Learning Research, J. Dy and A. Krause, Eds., vol. 80. PMLR, 2018, pp. 2284–2293.
- [5] L. Kocsis and C. Szepesvári, "Bandit based Monte-Carlo planning," in *Mach. Learn.: ECML 2006*, ser. Lecture Notes in Computer Science, J. Fürnkranz, T. Scheffer, and M. Spiliopoulou, Eds., Springer Berlin Heidelberg, 2006, vol. 4212, pp. 282–293.
- [6] R. Coulom, "Efficient selectivity and backup operators in Monte-Carlo tree search," in *Computers and Games*, ser. Lecture Notes in Computer Science, H. J. van den Herik, P. Ciancarini, and H. H. L. M. Donkers, Eds., vol. 4630. Springer Berlin Heidelberg, 2007, pp. 72–83.
- [7] C. Browne, E. Powley, D. Whitehouse, S. Lucas, P. I. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakis, and S. Colton, "A survey of Monte Carlo tree search methods," *IEEE Trans. Comput. Intell. AI Games*, vol. 4, no. 1, pp. 1–49, 2012.
- [8] C. Browne, D. J. N. J. Soemers, and E. Piette, "Strategic features for general games," in *Proc. 2nd Workshop on Knowledge Extraction from Games (KEG)*, 2019, pp. 70–75.
- [9] C. Browne, "Modern techniques for ancient games," in *Proc. 2018 IEEE Conf. Comput. Intell. Games*. IEEE, 2018, pp. 490–497.
- [10] D. Parlett, "What's a ludeme?" *Game & Puzzle Design*, vol. 2, no. 2, pp. 83–86, 2016.
- [11] M. Genesereth and N. Love, "General game playing: Overview of the AAAI competition," *AI Magazine*, vol. 26, no. 2, pp. 62–72, 2005.
- [12] M. Kirci, N. Sturtevant, and J. Schaeffer, "A GGP feature learning algorithm," *Künstliche Intelligenz*, vol. 25, no. 1, pp. 35–42, 2011.
- [13] K. Wałędzik and J. Mańdziuk, "Multigame playing by means of UCT enhanced with automatically generated evaluation functions," in *Artif. Gen. Intell.: 4th Int. Conf.*, ser. Lecture Notes in Computer Science, vol. 6830. Springer, 2011, pp. 327–332.
- [14] D. Michulke and S. Schiffel, "Distance features for general game playing agents," in *Proc. 4th Int. Conf. Agents Artif. Intell.*, 2012, pp. 127–136.
- [15] —, "Admissible distance heuristics for general games," in *Agents Artif. Intell. ICAART 2012*, ser. Communications in Computer and Inf. Science, J. Filipe and A. Fred, Eds., vol. 358. Springer, Berlin, Heidelberg, 2013.
- [16] K. Wałędzik and J. Mańdziuk, "An automatically generated evaluation function in general game playing," *IEEE Trans. Comput. Intell. AI Games*, vol. 6, no. 3, pp. 258–270, 2014.
- [17] H. Iida, M. Sakuta, and J. Rollason, "Computer shogi," *Artif. Intell.*, vol. 134, no. 1–2, pp. 121–144, 2002.
- [18] P. Auer, N. Cesa-Bianchi, and P. Fischer, "Finite-time analysis of the multiarmed bandit problem," *Mach. Learn.*, vol. 47, no. 2–3, pp. 235–256, 2002.
- [19] S. Gelly and D. Silver, "Combining online and offline knowledge in UCT," in *Proc. 24th Int. Conf. Mach. Learn.*, 2007, pp. 273–280.
- [20] N. R. Sturtevant and A. M. White, "Feature construction for reinforcement learning in Hearts," in *Computers and Games*, ser. Lecture Notes in Computer Science, H. J. van den Herik, P. Ciancarini, and H. H. L. M. Donkers, Eds., vol. 4630. Springer, 2007, pp. 122–134.
- [21] R. J. Lorentz and T. E. Zosa IV, "Machine learning in the game of Breakthrough," in *Adv. in Computer Games*, ser. Lecture Notes in Computer Science, M. H. M. Winands, H. van den Herik, and W. A. Kusters, Eds., vol. 10664. Springer, 2017, pp. 140–150.
- [22] P. Skowronski, Y. Björnsson, and M. H. M. Winands, "Automated discovery of search-extension features," in *Adv. in Computer Games*, ser. Lecture Notes in Computer Science, H. J. van den Herik and P. Spronck, Eds., vol. 6048. Springer, Berlin, Heidelberg, 2009.
- [23] M. A. Hall, "Correlation-based feature selection for machine learning," Ph.D. dissertation, Department of Computer Science, the University of Waikato, Hamilton, New Zealand, 1999.
- [24] T. Cazenave, "Playout policy adaptation with move features," *Theor. Comput. Sci.*, vol. 644, pp. 43–52, 2016.