

Evolutionary Computation for Reinforcement Learning

Shimon Whiteson

Abstract Algorithms for evolutionary computation, which simulate the process of natural selection to solve optimization problems, are an effective tool for discovering high-performing reinforcement-learning policies. Because they can automatically find good representations, handle continuous action spaces, and cope with partial observability, evolutionary reinforcement-learning approaches have a strong empirical track record, sometimes significantly outperforming temporal-difference methods. This chapter surveys research on the application of evolutionary computation to reinforcement learning, overviewing methods for evolving neural-network topologies and weights, hybrid methods that also use temporal-difference methods, coevolutionary methods for multi-agent settings, generative and developmental systems, and methods for on-line evolutionary reinforcement learning.

1 Introduction

Algorithms for *evolutionary computation*, sometimes known as *genetic algorithms* (Holland, 1975; Goldberg, 1989), are optimization methods that simulate the process of natural selection to find highly fit solutions to a given problem. Typically the problem assumes as input a *fitness function* $f : C \rightarrow \mathfrak{R}$ that maps C , the set of all candidate solutions, to a real-valued measure of fitness. The goal of an optimization method is to find $c^* = \arg \max_c f(c)$, the fittest solution. In some cases, the fitness function may be stochastic, in which case $f(c)$ can be thought of as a random variable and $c^* = \arg \max_c E[f(c)]$.

Evolutionary methods search for c^* by repeatedly selecting and reproducing a population of candidate solutions. The initial population is typically chosen randomly, after which each member of the population is evaluated using f and the

Shimon Whiteson
Informatics Institute, University of Amsterdam e-mail: s.a.whiteson@uva.nl

best performing ones are selected as the basis for a new population. This new population is formed via reproduction, in which the selected policies are mated (i.e., components of two different solutions are combined) and mutated (i.e., the parameter values of one solution are stochastically altered). This process repeats over many iterations, until a sufficiently fit solution has been found or the available computational resources have been exhausted.

There is an enormous number of variations on this approach, such as multi-objective methods (Deb, 2001; Coello et al, 2007) diversifying algorithms (Holland, 1975; Goldberg and Richardson, 1987; Mahfoud, 1995; Potter and De Jong, 1995; Darwen and Yao, 1996) and distribution-based methods (Larranaga and Lozano, 2002; Hansen et al, 2003; Rubinstein and Kroese, 2004). However, the basic approach is extremely general and can in principle be applied to all optimization problems for which f can be specified.

Included among these optimization problems are reinforcement-learning tasks (Moriarty et al, 1999). In this case, C corresponds to the set of possible policies, e.g., mappings from S to A , and $f(c)$ is the average cumulative reward obtained while using such a policy in a series of Monte Carlo trials in the task. In other words, in an evolutionary approach to reinforcement learning, the algorithm directly searches the space of policies for one that maximizes the expected cumulative reward.

Like many other policy-search methods, this approach reasons only about the value of entire policies, without constructing value estimates for particular state-action pairs, as temporal-difference methods do. The holistic nature of this approach is sometimes criticized. For example, Sutton and Barto write:

Evolutionary methods do not use the fact that the policy they are searching for is a function from states to actions; they do not notice which states an individual passes through during its lifetime, or which actions it selects. In some cases this information can be misleading (e.g., when states are misperceived) but more often it should enable more efficient search (Sutton and Barto, 1998, p. 9).

These facts can put evolutionary methods at a theoretical disadvantage. For example, in some circumstances, dynamic programming methods are guaranteed to find an optimal policy in time polynomial in the number of states and actions (Littman et al, 1995). By contrast, evolutionary methods, in the worst case, must iterate over an exponential number of candidate policies before finding the best one. Empirical results have also shown that evolutionary methods sometimes require more episodes than temporal-difference methods to find a good policy, especially in highly stochastic tasks in which many Monte Carlo simulations are necessary to achieve a reliable estimate of the fitness of each candidate policy (Runarsson and Lucas, 2005; Lucas and Runarsson, 2006; Lucas and Togelius, 2007; Whiteson et al, 2010b).

However, despite these limitations, evolutionary computation remains a popular tool for solving reinforcement-learning problems and boasts a wide range of empirical successes, sometimes substantially outperforming temporal-difference methods (Whitley et al, 1993; Moriarty and Miikkulainen, 1996; Stanley and Miikkulainen, 2002; Gomez et al, 2008; Whiteson et al, 2010b). There are three main reasons why.

First, evolutionary methods can cope well with partial observability. While evolutionary methods do not exploit the relationship between subsequent states that an

agent visits, this can be advantageous when the agent is unsure about its state. Since temporal-difference methods rely explicitly on the Markov property, their value estimates can diverge when it fails to hold, with potentially catastrophic consequences for the performance of the greedy policy. In contrast, evolutionary methods do not rely on the Markov property and will always select the best policies they can find for the given task. Severe partial observability may place a ceiling on the performance of such policies, but optimization within the given policy space proceeds normally (Moriarty et al, 1999). In addition, representations that use memory to reduce partial observability, such as recurrent neural networks, can be optimized in a natural way with evolutionary methods (Gomez and Miikkulainen, 1999; Stanley and Miikkulainen, 2002; Gomez and Schmidhuber, 2005a,b).

Second, evolutionary methods can make it easier to find suitable representations for the agent's solution. Since policies need only specify an action for each state, instead of the value of each state-action pair, they can be simpler to represent. In addition, it is possible to simultaneously evolve a suitable policy representation (see Sections 3 and 4.2). Furthermore, since it is not necessary to perform learning updates on a given candidate solution, it is possible to use more elaborate representations, such as those employed by *generative and developmental systems* (GDS) (see Section 6).

Third, evolutionary methods provide a simple way to solve problems with large or continuous action spaces. Many temporal-difference methods are ill-suited to such tasks because they require iterating over the action space in each state in order to identify the maximizing action. In contrast, evolutionary methods need only evolve policies that directly map states to actions. Of course, actor-critic methods (Doya, 2000; Peters and Schaal, 2008) and other techniques (Gaskett et al, 1999; Millán et al, 2002; van Hasselt and Wiering, 2007) can also be used to make temporal-difference methods suitable for continuous action spaces. Nonetheless, evolutionary methods provide a simple, effective way to address such difficulties.

Of course, none of these arguments are unique to evolutionary methods, but apply in principle to other policy-search methods too. However, evolutionary methods have proven a particularly popular way to search policy space and, consequently, there is a rich collection of algorithms and results for the reinforcement-learning setting. Furthermore, as modern methods, such as distribution-based approaches, depart further from the original genetic algorithms, their resemblance to the process of natural selection has decreased. Thus, the distinction between evolutionary methods and other policy search approaches has become fuzzier and less important.

This chapter provides an introduction to and overview of evolutionary methods for reinforcement learning. The vastness of the field makes it infeasible to address all the important developments and results. In the interest of clarity and brevity, this chapter focuses heavily on *neuroevolution* (Yao, 1999), in which evolutionary methods are used to evolve *neural networks* (Haykin, 1994), e.g., to represent policies. While evolutionary reinforcement learning is by no means limited to neural-network representations, neuroevolutionary approaches are by far the most common. Furthermore, since neural networks are a popular and well-studied representation in general, they are a suitable object of focus for this chapter.

The rest of this chapter is organized as follows. By way of introduction, Section 2 describes a simple neuroevolutionary algorithm for reinforcement learning. Section 3 considers *topology- and weight-evolving artificial neural networks* (TWEANNs), including the popular NEAT method, that automatically discover their own internal representations. Section 4 considers hybrid approaches, such as *evolutionary function approximation* and *learning classifier systems*, that integrate evolution with temporal-difference methods. Section 5 discusses *coevolution*, in which multiple competing and/or cooperating policies are evolved simultaneously. Section 6 describes *generative and developmental systems* (GDSs) such as HyperNEAT, which rely on *indirect encodings*: more complex representations in which the agent's policy must be constructed or grown from the evolved parameter values. Section 7 discusses on-line methods that strive to maximize reward during evolution instead of merely discovering a good policy quickly.

2 Neuroevolution

Neural networks (Haykin, 1994) are an extremely general-purpose way of representing complex functions. Because of their concision, they are a popular representation for reinforcement-learning policies, not only for evolutionary computation, but also for other policy-search algorithms and for temporal-difference methods. This section introduces the basics of neuroevolutionary approaches to reinforcement learning, in which evolutionary methods are used to optimize neural-network policies.

Figure 1 illustrates the basic steps of a neuroevolutionary algorithm (Yao, 1999). In each generation, each network in the population is evaluated in the task. Next, the best performing are selected, e.g., via rank-based selection, roulette wheel selection, or tournament selection (Goldberg and Deb, 1991). The selected networks are bred via crossover and mutation and reproduced (Sywerda, 1989; De Jong and Spears, 1991)) to form a new population and the process repeats.

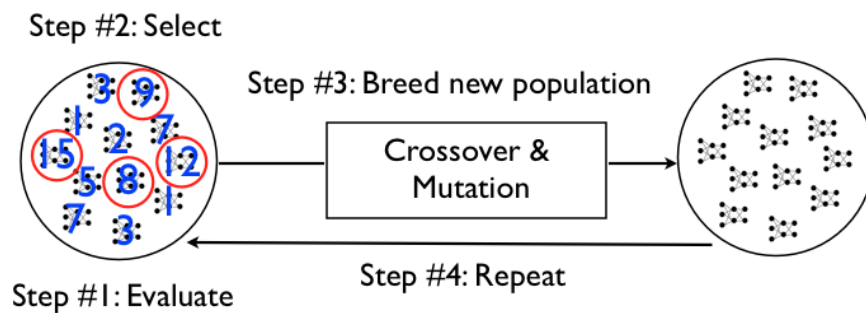


Fig. 1 The basic steps of neuroevolution.

In a reinforcement-learning setting, each input node in a network typically corresponds to a state feature, such that the value of the inputs together describe the agent’s state. There are many ways to represent actions. If the actions can also be described with features, as is common with continuous action spaces, then each output node can correspond to an action feature. In this case, the value of the outputs together describe the action to be selected for the given state. When the number of actions is small and discrete, a separate network can be used for each action, as is common in value-function approximation. In this case, each network has only one output node. The policy’s action for a given state then corresponds to the network that produces the largest output for that state.

In this chapter, we focus on a variation of this approach called *action-selection* networks. As before, we assume the set of actions is small and discrete. However, only one network is used to represent the policy. This network has one output node for each action. The policy’s action for a given state then corresponds to the node that produces the largest output for that state.

Algorithm 1 contains a high-level description of a simple neuroevolutionary method that evolves action-selection networks for an episodic reinforcement-learning problem. It begins by creating a population of random networks (line 4). In each generation, it repeatedly iterates over the current population (lines 6–7). During each step of a given episode, the agent takes whatever action corresponds to the output with the highest activation (lines 10–12). Note that s' and a' are the current state and action while s and a are the previous state and action. Neuroevolution maintains a running total of the reward accrued by the network during its evaluation (line 13). Each generation ends after e episodes, at which point each network’s average fitness is $N.fitness/N.episodes$. In stochastic domains, e typically must be much larger than $|P|$ to ensure accurate fitness estimates for each network. Neuroevolution creates a new population by repeatedly calling the BREED-NET function (line 18), which generates a new network from highly fit parents.

Note that, while the action selection described in lines 10–12 resembles greedy action selection from a value function, the network should not be interpreted as a value function.¹ Evolution does not search for the networks that best approximate the optimal value function. Instead, it searches for networks representing high performing policies. To perform well, a network need only generate more output for the optimal action than for other actions. Unlike with value functions, the scale of the outputs can be arbitrary, as well as the relative outputs of the non-selected actions.

The neural network employed by Algorithm 1 could be a simple *feed-forward* network or a more complex *recurrent* network. Recurrent neural networks can contain cycles (e.g., the output emerging from an output node can be fed into an input node). Consequently, such networks can contain internal state. In a reinforcement-learning context, this internal state can record aspects of the agent’s observation history, which can help it cope with partial observability (Wieland, 1991; Gomez and Miikkulainen, 1999; Moriarty et al, 1999; Stanley and Miikkulainen, 2002; Igel, 2003; Gomez and Schmidhuber, 2005a,b).

¹ This does not apply to the hybrid methods discussed in Section 4.

Algorithm 1 NEUROEVOLUTION(S, A, p, g, e)

```

1: //  $S$ : set of all states,  $A$ : set of all actions,  $p$ : population size
2: //  $g$ : number of generations,  $e$ : episodes per generation
3:
4:  $P \leftarrow$  INIT-POPULATION( $S, A, p$ )           // create new population  $P$  with random networks
5: for  $i \leftarrow 1$  to  $g$  do
6:   for  $j \leftarrow 1$  to  $e$  do
7:      $N, s, s' \leftarrow P[j \% p]$ , null, INIT-STATE( $S$ )           // select next network
8:     repeat
9:        $Q \leftarrow$  EVAL-NET( $N, s'$ )           // evaluate selected network on current state
10:       $a' \leftarrow \operatorname{argmax}_i Q[i]$            // select action with highest activation
11:       $s, a \leftarrow s', a'$ 
12:       $r, s' \leftarrow$  TAKE-ACTION( $a'$ )           // take action and transition to new state
13:       $N.\text{fitness} \leftarrow N.\text{fitness} + r$            // update total reward accrued by  $N$ 
14:      until TERMINAL-STATE?( $s$ )
15:       $N.\text{episodes} \leftarrow N.\text{episodes} + 1$            // update total number of episodes for  $N$ 
16:       $P' \leftarrow$  new array of size  $p$            // new array will store next generation
17:      for  $j \leftarrow 1$  to  $p$  do
18:         $P'[j] \leftarrow$  BREED-NET( $P$ )           // make a new network based on fit parents in  $P$ 
19:       $P \leftarrow P'$ 

```

While Algorithm 1 uses a traditional genetic algorithm to optimize neural-network weights, many variations are also possible. For example, *estimation of distribution algorithms* (EDAs) (Larranaga and Lozano, 2002; Hansen et al, 2003; Rubinstein and Kroese, 2004) can be used instead. EDAs, also called *probabilistic model-building genetic algorithms* (PMBGAs), do not explicitly maintain a population of candidate solutions. Instead, they maintain a distribution over solutions. In each generation, candidate solutions are sampled from this distribution and evaluated. A subset is then selected and used to update the distribution using *density estimation* techniques, unsupervised learning techniques for approximating the distribution from which a set of samples was drawn.

One of the most popular and effective EDAs is the *covariance matrix adaptation evolution strategy* (CMA-ES) (Hansen et al, 2003), a variable-metric EDA in which the distribution is a multivariate Gaussian whose covariance matrix adapts over time. When used to optimize neural networks, the resulting method, called CMA-NeuroES, has proven effective on a wide range of reinforcement-learning tasks (Igel, 2003; Heidrich-Meisner and Igel, 2008, 2009a,b,c).

3 TWEANNs

In its simplest form, Algorithm 1 evolves only neural networks with fixed representations. In such a setup, all the networks in a particular evolutionary run have the same *topology*, i.e., both the number of hidden nodes and the set of edges connecting the nodes are fixed. The networks differ only with respect to the weights of these edges, which are optimized by evolution. The use of fixed representations is

by no means unique to neuroevolution. In fact, though methods exist for automatically discovering good representations for value-functions (Mahadevan and Maggioni, 2007; Parr et al, 2007) temporal-difference methods typically also use fixed representations for function approximation.

Nonetheless, reliance on fixed representations is a significant limitation. The primary reason is that it requires the user of the algorithm to correctly specify a good representation in advance. Clearly, choosing too simple a representation will doom evolution to poor performance, since describing high quality solutions becomes impossible. However, choosing too complex a representation can be just as harmful. While such a representation can still describe good solutions, finding them may become infeasible. Since each weight in the network corresponds to a dimension of the search space, a representation with too many edges can lead to an intractable search problem.

In most tasks, the user is not able to correctly guess the right representation. Even in cases where the user possesses great domain expertise, deducing the right representation from this expertise is often not possible. Typically, finding a good representation becomes a process of trial and error. However, repeatedly running evolution until a suitable representation is found greatly increases computational costs. Furthermore, in on-line tasks (see Section 7) it also increases the real-world costs of trying out policies in the target environment.

For these reasons, many researchers have investigated ways to automate the discovery of good representations (Dasgupta and McGregor, 1992; Radcliffe, 1993; Gruau, 1994; Stanley and Miikkulainen, 2002). Evolutionary methods are well suited to this challenge because they take a direct policy-search approach to reinforcement learning. In particular, since neuroevolution already directly searches the space of network weights, it can also simultaneously search the space of network topologies. Methods that do so are sometimes called *topology- and weight-evolving artificial neural networks* (TWEANNs).

Perhaps the earliest and simplest TWEANN is the *structured genetic algorithm* (sGA) (Dasgupta and McGregor, 1992), which uses a two-part representation to describe each network. The first part represents the connectivity of the network in the form of a binary matrix. Rows and columns correspond to nodes in the network and the value of each cell indicates whether an edge exists connecting the given pair of nodes. The second part represents the weights of each edge in the network. In principle, by evolving these binary matrices along with connection weights, sGA can automatically discover suitable network topologies. However, sGA suffers from several limitations. In the following section, we discuss these limitations in order to highlight the main challenges faced by all TWEANNs.

3.1 Challenges

There are three main challenges to developing a successful TWEANN. The first is the *competing conventions* problem. In most tasks, there are multiple different

policies that have similar fitness. For example, many tasks contain symmetries that give rise to several equivalent solutions. This can lead to difficulties for evolution because of its reliance on crossover operators to breed new networks. When two networks that represent different policies are combined, the result is likely to be destructive, producing a policy that cannot successfully carry out the strategy used by either parent.

While competing conventions can arise in any evolutionary method that uses crossover, the problem is particularly severe for TWEANNs. Two parents may not only implement different policies but also have different representations. Therefore, to be effective, TWEANNs need a mechanism for combining networks with different topologies in a way that minimizes the chance of catastrophic crossover. Clearly, sGA does not meet this challenge, since the binary matrices it evolves are crossed over without regard to incompatibility in representations. In fact, the difficulties posed by the competing conventions problem were a major obstacle for early TWEANNs, to the point that some researchers simply avoided the problem by developing methods that do not perform crossover at all (Radcliffe, 1993).

The second challenge is the need to protect topological innovations long enough to optimize the associated weights. Typically, when new topological structures are introduced (e.g., the addition of a new hidden node or edge), it has a negative effect on fitness even if that structure will eventually be necessary for a good policy. The reason is that the weights associated with the new structure have not yet been optimized.

For example, consider an edge in a network evolved via sGA that is not activated, i.e., its cell in the binary matrix is set to zero. The corresponding weight for that edge will not experience any selective pressure, since it is not manifested in the network. If evolution suddenly activates that edge, the effect on fitness is likely to be detrimental, since its weight is not optimized. Therefore, if topological innovations are not explicitly protected, they will typically be eliminated from the population, causing the search for better topologies to stagnate.

Fortunately, protecting innovation is a well-studied problem in evolutionary computation. *Speciation* and *nicheing* methods (Holland, 1975; Goldberg and Richardson, 1987; Mahfoud, 1995; Potter and De Jong, 1995; Darwen and Yao, 1996) ensure diversity in the population, typically by segregating disparate individuals and/or penalizing individuals that are too similar to others. However, using such methods requires a distance metric to quantify the differences between individuals. Devising such a metric is difficult for TWEANNs, since it is not clear how to compare networks with different topologies.

The third challenge is how to evolve minimal solutions. As mentioned above, a central motivation for TWEANNs is the desire to avoid optimizing overly complex topologies. However, if evolution is initialized with a population of randomly chosen topologies, as in many TWEANNs, some of these topologies may already be too complex. Thus, at least part of the evolutionary search will be conducted in an unnecessarily high dimensional space. It is possible to explicitly reward smaller solutions by adding size penalties in the fitness function (Zhang and Muhlenbein,

1993). However, there is no principled way to determine the size of the penalties without prior knowledge about the topological complexity required for the task.

3.2 NEAT

Perhaps the most popular TWEANN is *neuroevolution of augmenting topologies* (NEAT) (Stanley and Miikkulainen, 2002). In this section, we briefly describe NEAT and illustrate how it addresses the major challenges mentioned above.

NEAT is often used to evolve action selectors, as described in Section 2. In fact, NEAT follows the framework described in Algorithm 1 and differs from traditional neuroevolution only in how INIT-POPULATION and BREED-NET are implemented.

To represent networks of varying topologies, NEAT employs a flexible genetic encoding. Each network is described by a list of *edge genes*, each of which describes an edge between two *node genes*. Each edge gene specifies the in-node, the out-node, and the weight of the edge. During mutation, new structure can be introduced to a network via special mutation operators that add new node or edge genes to the network (see Figure 2).

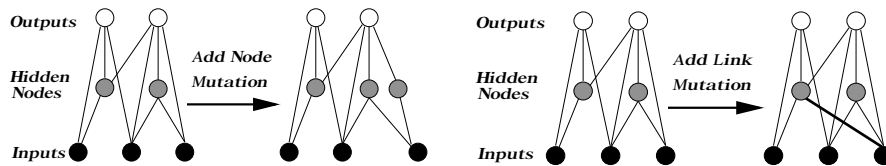


Fig. 2 Structural mutation operators in NEAT. At left, a new node is added by splitting an existing edge in two. At right, a new link (edge) is added between two existing nodes.

To avoid catastrophic crossover, NEAT relies on *innovation numbers*, which track the historical origin of each gene. Whenever a new gene appears via mutation, it receives a unique innovation number. Thus, the innovation numbers can be viewed as a chronology of all the genes produced during evolution.

During crossover, innovation numbers are used to determine which genes in the two parents correspond to each other. Genes that do not match are either *disjoint* or *excess*, depending on whether they occur within or outside the range of the other parent's innovation numbers. When crossing over, pairs of genes with the same innovation number (one from each parent) are lined up. Genes that do not match are inherited from the fitter parent. This approach makes it possible for NEAT to minimize the chance of catastrophic crossover without conducting an expensive topological analysis. Since genomes with different topologies nonetheless remain compatible throughout evolution, NEAT essentially avoids the competing conventions problem.

Innovation numbers also make possible a simple way to protect topological innovation. In particular, NEAT uses innovation numbers to speciate the population

based on topological similarity. The distance δ between two network encodings is a simple linear combination of the number of excess (E) and disjoint (D) genes, as well as the average weight differences of matching genes (\bar{W}):

$$\delta = \frac{c_1 E}{N} + \frac{c_2 D}{N} + c_3 \cdot \bar{W}$$

The coefficients c_1 , c_2 , and c_3 adjust the importance of the three factors, and the factor N , the number of genes in the larger genome, normalizes for genome size. Networks whose distance is greater than δ_t , a compatibility threshold, are placed into different species. *Explicit fitness sharing* (Goldberg, 1989), in which networks in the same species must share the fitness of their niche, is employed to protect innovative species.

To encourage the evolution of minimal solutions, NEAT begins with a uniform population of simple networks with no hidden nodes and inputs connected directly to outputs. New structure is introduced incrementally via the mutation operators that add new hidden nodes and edges to the network. Since only the structural mutations that yield performance advantages tend to survive evolution's selective pressure, minimal solutions are favored.

NEAT has amassed numerous empirical successes on difficult reinforcement-learning tasks like non-Markovian double pole balancing (Stanley and Miikkulainen, 2002), game playing (Stanley and Miikkulainen, 2004b), and robot control (Stanley and Miikkulainen, 2004a; Taylor et al, 2006; Whiteson et al, 2010b). However, Kohl and Miikkulainen (2008, 2009) have shown that NEAT can perform poorly on tasks in which the optimal action varies discontinuously across states. They demonstrate that these problems can be mitigated by providing neurons with local receptive fields and constraining topology search to cascade structures

4 Hybrids

Many researchers have investigated hybrid methods that combine evolution with supervised or unsupervised learning methods. In such systems, the individuals being evolved do not remain fixed during their fitness evaluations. Instead, they change during their 'lifetimes' by learning from the environments with which they interact.

Much of the research on hybrid methods focuses on analyzing the dynamics that result when evolution and learning interact. For example, several studies (Whitley et al, 1994; Yamasaki and Sekiguchi, 2000; Pereira and Costa, 2001; Whiteson and Stone, 2006a) have used hybrids to compare *Lamarckian* and *Darwinian* systems. In Lamarckian systems, the phenotypic effects of learning are copied back into the genome before reproduction, allowing new offspring to inherit them. In Darwinian systems, which more closely model biology, learning does not affect the genome. As other hybrid studies (Hinton and Nowlan, 1987; French and Messinger, 1994; Arita and Suzuki, 2000) have shown, Darwinian systems can indirectly transfer the results of learning into the genome by way of the *Baldwin effect* (Baldwin, 1896), in

which learning creates selective pressures favoring individuals who innately possess attributes that were previously learned.

Hybrid methods have also been employed to improve performance on supervised learning tasks (Gruau and Whitley, 1993; Boers et al, 1995; Giraud-Carrier, 2000; Schmidhuber et al, 2005, 2007). However, such methods are not directly applicable to reinforcement-learning problems because the labeled data they require is absent.

Nonetheless, many hybrid methods for reinforcement learning have been developed. To get around the problem of missing labels, researchers have employed unsupervised learning (Stanley et al, 2003), trained individuals to resemble their parents (McQuesten and Miikkulainen, 1997), trained them to predict state transitions (Nolfi et al, 1994), and trained them to teach themselves (Nolfi and Parisi, 1997). However, perhaps the most natural hybrids for the reinforcement learning setting are combination of evolution with temporal-difference methods (Ackley and Littman, 1991; Wilson, 1995; Downing, 2001; Whiteson and Stone, 2006a). In this section, we survey two such hybrids: *evolutionary function approximation* and XCS, a type of *learning classifier system*.

4.1 Evolutionary Function Approximation

Evolutionary function approximation (Whiteson and Stone, 2006a), is a way to synthesize evolutionary and temporal-difference methods into a single method that automatically selects function approximator representations that enable efficient individual learning. The main idea is that, if evolution is directed to evolve value functions instead of action selectors, then those value functions can be updated, using temporal-difference methods, during each fitness evaluation. In this way, the system can *evolve* function approximators that are better able to *learn* via temporal-difference methods. This biologically intuitive combination, which has been applied to many computational systems (Hinton and Nowlan, 1987; Ackley and Littman, 1991; Boers et al, 1995; French and Messinger, 1994; Gruau and Whitley, 1993; Nolfi et al, 1994), can yield effective reinforcement-learning algorithms. In this section, we briefly describe NEAT+Q, an evolutionary function approximation technique resulting from the combination of NEAT and Q-learning with neural-network function approximation.

To make NEAT optimize value functions instead of action selectors, all that is required is a reinterpretation of its output values. The structure of neural-network action selectors (one input for each state feature and one output for each action) is already identical to that of Q-learning function approximators. Therefore, if the weights of the networks NEAT evolves are updated during their fitness evaluations using Q-learning and backpropagation, they will effectively evolve value functions instead of action selectors. Hence, the outputs are no longer arbitrary values; they represent the long-term discounted values of the associated state-action pairs and are used, not just to select the most desirable action, but to update the estimates of other state-action pairs.

Algorithm 2 shows the inner loop of NEAT+Q, replacing lines 9–13 in Algorithm 1. Each time the agent takes an action, the network is backpropagated towards Q-learning targets (line 7) and ϵ -greedy selection occurs (lines 4–5). Figure 3 illustrates the complete algorithm: networks are selected from the population for evaluation and the Q-values they produce are used to select actions. The resulting feedback from the environment is used both to perform TD updates and to measure the network’s fitness, i.e., the total reward it accrues while learning.

Algorithm 2 NEAT+Q (inner loop)

```

1: //  $\alpha$ : learning rate,  $\gamma$ : discount factor,  $\lambda$ : eligibility decay rate,  $\epsilon$ : exploration rate
2:
3:  $Q \leftarrow \text{EVAL-NET}(N, s')$  // compute value estimates for current state
4: with-prob( $\epsilon$ )  $a' \leftarrow \text{RANDOM}(A)$  // select random exploratory action
5: else  $a' \leftarrow \text{argmax}_k Q[k]$  // or select greedy action
6: if  $s \neq \text{null}$  then
7:    $\text{BACKPROP}(N, s, a, (r + \gamma \max_k Q[k]), \alpha, \gamma, \lambda)$  // adjust weights
8:    $s, a \leftarrow s', a'$ 
9:    $r, s' \leftarrow \text{TAKE-ACTION}(a')$  // take action and transition to new state
10:   $N.\text{fitness} \leftarrow N.\text{fitness} + r$  // update total reward accrued by  $N$ 
  
```

Like other hybrid methods, NEAT+Q combines the advantages of temporal-difference methods with those of evolution. In particular, it harnesses the ability of NEAT to discover effective representations and uses it to aid neural-network value-function approximation. Unlike traditional neural-network function approximators, which put all their eggs in one basket by relying on a single manually designed network to represent the value function, NEAT+Q explores the space of such networks to increase the chance of finding a high performing representation. As a result, on certain tasks, this approach has been shown to significantly outperform both temporal-difference methods and neuroevolution on their own (Whiteson and Stone, 2006a).

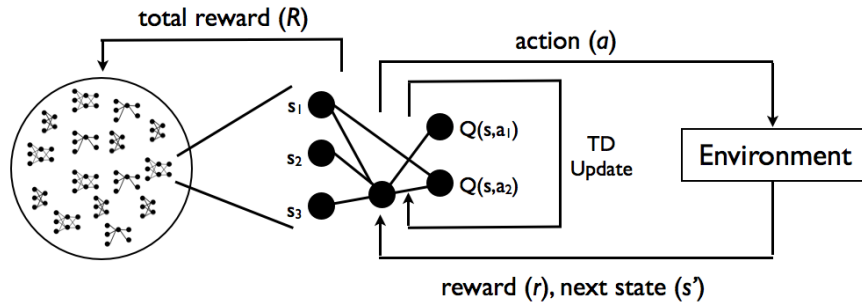


Fig. 3 The NEAT+Q algorithm.

4.2 XCS

A different type of hybrid method can be constructed using *learning classifier systems* (LCSs) (Holland, 1975; Holland and Reitman, 1977; Bull and Kovacs, 2005; Butz, 2006; Drugowitsch, 2008). An LCS is an evolutionary system that uses rules, called *classifiers*, to represent solutions. A wide range of classification, regression, and optimization problems can be tackled using such systems. These include reinforcement learning problems, in which case a set of classifiers describes a policy.

Each classifier contains a *condition* that describes the set of states to which the classifier applies. Conditions can be specified in many ways (e.g., fuzzy conditions (Bonarini, 2000) or neural network conditions (Bull and O’Hara, 2002)) but perhaps the simplest is the ternary alphabet used for binary state features: 0, 1, and #, where # indicates “don’t care”. For example, a classifier with condition 01#10 applies to the states 01010 and 01110. Classifiers also contain *actions*, which specify what action the agent should take in states to which the classifier applies, and *predictions* which estimate the corresponding action-value function.

In *Pittsburgh-style* classifier systems (De Jong et al, 1993), each individual in the evolving population consists of an entire rule set describing a complete policy. These methods can be viewed as standard evolutionary approaches to reinforcement learning with rule sets as a way to represent policies, i.e., in lieu of neural networks. In contrast, in *Michigan-style* classifier systems, each individual consists of only one rule. Thus, the entire population represents a single, evolving policy. In the remainder of this section, we give a brief overview of XCS (Wilson, 1995, 2001; Butz et al, 2008), one of the most popular Michigan-style classifiers. Since prediction updates in XCS are based on Q-learning, it can be viewed as a hybrid between temporal-difference methods and evolutionary reinforcement-learning algorithms.

In XCS, each classifier’s prediction contributes to an estimate of the Q-values of the state-action pairs to which the classifier applies. $Q(s, a)$ is the weighted average of the predictions of all the *matching* classifiers, i.e., those that apply to s and have action a . The fitness of each classifier determines its weight in the average (fitness is defined later in this section). In other words:

$$Q(s, a) = \frac{\sum_{c \in M(s, a)} c.f \cdot c.p}{\sum_{c \in M(s, a)} c.f},$$

where $M(s, a)$ is the set of all classifiers matching s and a ; $c.f$ and $c.p$ are the fitness and prediction, respectively, of classifier c .

Each time the agent is in state s , takes action a , receives reward r , and transitions to state s' , the following update rule is applied to each $c \in M(s, a)$:

$$c.p \leftarrow c.p + \beta [r + \gamma \max_{a'} Q(s', a') - c.p] \frac{c.f}{\sum_{c' \in M(s, a)} c'.f},$$

where β is a learning rate parameter. This is essentially a Q-learning update, except that the size of the update is scaled according to the relative fitness of c , since this determines its contribution to the Q-value.

Many LCS systems use *strength-based* updates, wherein each classifier's fitness is based on its prediction, i.e., classifiers that expect to obtain more reward are favored by evolution. However, this can lead to the problem of *strong over-generals* (Kovacs, 2003), which are general classifiers that have high overall value but select suboptimal actions for a subset of states. With strength-based updates, such classifiers are favored over more specific ones that select better actions for that same subset but have lower overall value.

To avoid this problem, XCS uses *accuracy-based* updates, in which a classifier's fitness is inversely proportional to an estimate of the error in its prediction. This error $c.\varepsilon$ is updated based on the absolute value of the temporal-difference error used in the Q-learning update:

$$c.\varepsilon \leftarrow c.\varepsilon + \beta(|r + \gamma \max_{a'} Q(s', a') - c.p| - c.\varepsilon)$$

Classifier accuracy is then defined in terms of this error. Specifically, when $c.\varepsilon > \varepsilon_0$, a minimum error threshold, the accuracy of c is defined as:

$$c.\kappa = \alpha(c.\varepsilon/\varepsilon_0)^{-\eta},$$

where α and η are accuracy parameters. When $c.\varepsilon \leq \varepsilon_0$, $c.\kappa = 1$.

However, fitness is computed, not with this accuracy, but instead with the *set-relative accuracy*: the accuracy divided by the sum of the accuracies of all matching classifiers. This yields the following fitness update:

$$c.f \leftarrow c.f + \beta\left(\frac{c.\kappa}{\sum_{c' \in M(s,a)} c'.\kappa} - c.f\right)$$

At every timestep, the prediction, error and fitness of each matching classifier are updated. Since XCS is a *steady-state* evolutionary method, there are no generations. Instead, the population changes incrementally through the periodic selection and reproduction of a few fit classifiers, which replace a few weak classifiers and leave the rest of the population unchanged. When selection occurs, only classifiers that match the current state and action are considered. New classifiers are created from the selected ones using crossover and mutation, and weak classifiers (chosen from the whole population, not just the matching ones) are deleted to make room.

Thanks to the Q-learning updates, the accuracy of the classifiers tends to improve over time. Thanks to steady-state evolution, the most accurate classifiers are selectively bred. General rules tend to have higher error (since they generalize over more states) and thus lower accuracy. It might seem that, as a result, XCS will evolve only highly specific rules. However, more general rules also match more often. Since only matching classifiers can reproduce, XCS balances the pressure for specific rules with pressure for general rules. Thus, it strives to learn a complete, maximally general, and accurate set of classifiers for approximating the optimal Q-function.

Though there are no convergence proofs for XCS on MDPs, it has proven empirically effective on many tasks. For example, on maze tasks, it has proven adept at automatically discovering what state features to ignore (Butz et al, 2005) and solving problems with more than a million states (Butz and Lanzi, 2009). It has also proven adept at complex sensorimotor control (Butz and Herbort, 2008; Butz et al, 2009) and autonomous robotics (Dorigo and Colombetti, 1998).

5 Coevolution

Coevolution is a concept from evolutionary biology that refers to the interactions between multiple individuals that are simultaneously evolving. In other words, coevolution occurs when the fitness function of one individual depends on other individuals that are also evolving. In nature, this can occur when different populations interact, e.g., cheetahs and the gazelles they hunt, or within a population, e.g., members of the same species competing for mates. Furthermore, coevolution can be cooperative, e.g., humans and the bacteria in our digestive systems, or competitive, e.g., predator and prey. All these forms of coevolution have been investigated and exploited in evolutionary reinforcement learning, as surveyed in this section.

5.1 Cooperative Coevolution

The most obvious application of cooperative coevolution is to cooperative multi-agent systems (Panait and Luke, 2005) which, in the context of reinforcement learning, means evolving teams of agents that coordinate their behavior to solve a sequential decision problem. In principle, such problems can be solved without coevolution by using a monolithic approach: evolving a population in which each individual specifies the policy for every agent on the team. However, such an approach quickly becomes intractable, as the size of the search space grows exponentially with respect to the number of agents.

One of the primary motivations for a coevolutionary approach is that it can help address this difficulty (Wiegand et al, 2001; Jansen and Wiegand, 2004; Panait et al, 2006). As Gomez et al. put it, “many problems may be decomposable into weakly coupled low-dimensional subspaces that can be searched semi-independently by separate species” (Gomez et al, 2008). In general, identifying these low-dimensional subspaces requires a lot of domain knowledge. However, in multi-agent problems, it is often sufficient to divide the problem up by agent, i.e., evolve one population for each agent on the team, in order to make evolution tractable. In this approach, one member of each population is selected, often randomly, to form a team that is then evaluated in the task. The total reward obtained contributes to an estimate of the fitness of each participating agent, which is typically evaluated multiple times.

While this approach often outperforms monolithic evolution and has found success in predator-prey (Yong and Miikkulainen, 2007) and robot-control (Cai and Peng, 2002) applications, it also runs into difficulties when there are large numbers of agents. The main problem is that the contribution of a single agent to the total reward accrued becomes insignificant. Thus, the fitness an agent receives depends more on which teammates it is evaluated with than on its own policy. However, it is possible to construct special fitness functions for individual agents that are much less sensitive to such effects (Agogino and Tumer, 2008). The main idea is to use *difference functions* (Wolpert and Tumer, 2002) that compare the total reward the team obtains when the agent is present to when it is absent or replaced by a fixed baseline agent. While this approach requires access to a model of the environment and increases the computational cost of fitness evaluation (so that the reward in both scenarios can be measured), it can dramatically improve the performance of cooperative coevolution.

Coevolution can also be used to simultaneously evolve multiple components of a single agent, instead of multiple agents. For example, in the task of robot soccer keepaway, domain knowledge has been used to decompose the task into different components, each representing an important skill such as running towards the ball or getting open for a pass (Whiteson et al, 2005). Neural networks for each of these components are then coevolved and together comprise a complete policy. In the keepaway task, coevolution greatly outperforms a monolithic approach.

Cooperative coevolution can also be used in a single-agent setting to facilitate neuroevolution. Rather than coevolving multiple networks, with one for each member of a team or each component of a policy, *neurons* are coevolved, which together form a single network describing the agent's policy (Potter and De Jong, 1995, 2000). Typically, networks have fixed topologies with a single hidden layer and each neuron corresponds to a hidden node, including all the weights of its incoming and outgoing edges. Just as dividing a multi-agent task up by agent often leads to simpler subproblems, so too can breaking up a neuroevolutionary task by neuron. As Moriarty and Miikkulainen say, "neuron-level evolution takes advantage of the *a priori* knowledge that individual neurons constitute basic components of neural networks" (Moriarty and Miikkulainen, 1997).

One example is *symbiotic adaptive neuroevolution* (SANE) (Moriarty and Miikkulainen, 1996, 1997) in which evolution occurs simultaneously at two levels. At the lower level, a single population of neurons is evolved. The fitness of each neuron is the average performance of the networks in which it participates during fitness evaluations. At the higher level, a population of *blueprints* is evolved, with each blueprint consisting of a vector of pointers to neurons in the lower level. The blueprints that combine neurons into the most effective networks tend to survive selective pressure. On various reinforcement-learning tasks such as robot control and pole balancing, SANE has outperformed temporal-difference methods, monolithic neuroevolution, and neuron-level evolution without blueprints.

The *enforced subpopulations* (ESP) method (Gomez and Miikkulainen, 1999) eliminates the blueprint population but segregates neurons into subpopulations, one for each hidden node. One member of each population is selected randomly to form

a network for each fitness evaluation. This encourages subpopulations to take on different specialized roles, increasing the likelihood that effective networks will be formed even without blueprints. ESP has performed particularly well on partially observable tasks, solving a non-Markovian version of the double pole-balancing problem. In addition H-ESP, a hierarchical variant that does use blueprints, has proven successful on *deep memory POMDPs*, i.e., those requiring history of hundreds or thousands of timesteps (Gomez and Schmidhuber, 2005a). ESP has even been used to evolve control systems for a finless sounding rocket (Gomez and Miikkulainen, 2003).

In *cooperative synapse neuroevolution* (CoSyNE) (Gomez et al, 2006, 2008), the idea of separate subpopulations is taken even further. Rather than a subpopulation for every neuron, which contains multiple weights, CoSyNE has a subpopulation for each edge, which has only one weight (see Figure 4). Thus the problem of finding a complete network is broken down into atomic units, solutions to which are coevolved. On several versions of the pole balancing problem, CoSyNE has been shown to outperform various temporal-difference methods and other policy-search approaches, as well as SANE, ESP, and NEAT (Gomez et al, 2006, 2008). However, in a more recent study, CMA-NeuroES (see Section 2) performed even better (Heidrich-Meisner and Igel, 2009b).

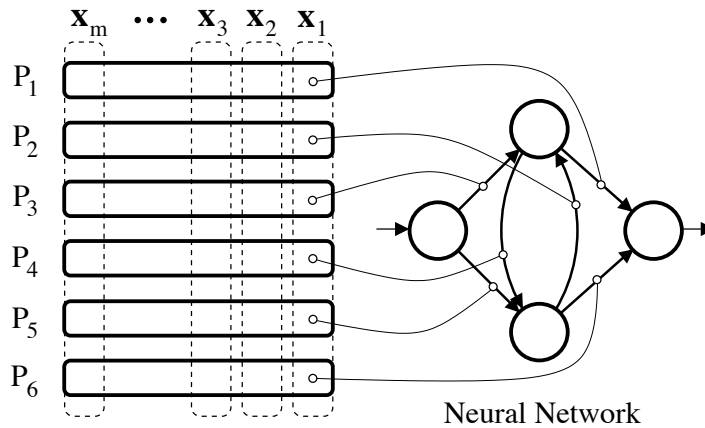


Fig. 4 The CoSyNE algorithm, using six subpopulations, each containing m weights. All the weights at a given index i form a genotype \mathbf{x}_i . Each weight is taken from a different subpopulation and describes one edge in the neural network. Figure taken with permission from (Gomez et al, 2008).

5.2 *Competitive Coevolution*

Coevolution has also proven a powerful tool for competitive settings. The most common applications are in games, in which coevolution is used to simultaneously evolve strong players and the opponents against which they are evaluated. The hope is to create an *arms race* (Dawkins and Krebs, 1979) in which the evolving agents exert continual selective pressure on each other, driving evolution towards increasingly effective policies.

Perhaps the simplest example of competitive coevolution is the work of Pollack and Blair in the game of backgammon (Pollack and Blair, 1998). Their approach relies on a simple optimization technique (essentially an evolutionary method with a population size of two) wherein a neural network plays against a mutated version of itself and the winner survives. The approach works so well that Pollack and Blair hypothesize that Tesauro’s great success with TD-Gammon (Tesauro, 1994) is due more to the nature of backgammon than the power of temporal-difference methods.²

Using larger populations, competitive coevolution has also found success in the game of checkers. The Blondie24 program uses the *minimax* algorithm (Von Neumann, 1928) to play checkers, relying on neuroevolution to discover an effective evaluator of board positions (Chellapilla and Fogel, 2001). During fitness evaluations, members of the current population play games against each other. Despite the minimal use of human expertise, Blondie24 evolved to play at a level competitive with human experts.

Competitive coevolution can also have useful synergies with TWEANNs. In fixed-topology neuroevolution, arms races may be cut short when additional improvement requires an expanded representation. Since TWEANNs can automatically expand their representations, coevolution can give rise to *continual complexification* (Stanley and Miikkulainen, 2004a).

The methods mentioned above evolve only a single population. However, as in cooperative coevolution, better performance is sometimes possible by segregating individuals into separate populations. In the *host/parasite* model (Hillis, 1990), one population evolves *hosts* and another *parasites*. Hosts are evaluated based on their robustness against parasites, e.g., how many parasites they beat in games of checkers. In contrast, parasites are evaluated based on their uniqueness, e.g., how many hosts they can beat that other parasites cannot. Such fitness functions can be implemented using *competitive fitness sharing* (Rosin and Belew, 1997).

In *Pareto coevolution*, the problem is treated as a multi-objective one, with each opponent as an objective (Ficici and Pollack, 2000, 2001). The goal is thus to find a *Pareto-optimal solution*, i.e., one that cannot be improved with respect to one objective without worsening its performance with respect to another. Using this approach, many methods have been developed that maintain *Pareto archives* of oppo-

² Tesauro, however, disputes this claim, pointing out that the performance difference between Pollack and Blair’s approach and his own is quite significant, analogous to that between an average human player and a world-class one (Tesauro, 1998).

nents against which to evaluate evolving solutions (De Jong, 2004; Monroy et al, 2006; De Jong, 2007; Popovici et al, 2010).

6 Generative and Developmental Systems

All of the evolutionary reinforcement-learning methods described above rely on *direct encodings* to represent policies. In such representations, evolution optimizes a *genotype* (e.g., a vector of numbers specifying the weights of a neural network) that can be trivially transformed into a *phenotype* (e.g., the neural network itself). While the simplicity of such an approach is appealing, it has limited scalability. Since the genotype is always as large as the phenotype, evolving the complex policies required to solve many realistic tasks requires searching a high dimensional space.

Generative and developmental systems (Gruau, 1994; Hornby and Pollack, 2002; Stanley and Miikkulainen, 2003; Stanley et al, 2009) is a subfield of evolutionary computation that focuses on evolving *indirect encodings*. In such representations, the phenotype is ‘grown’ via a developmental process specified by the genotype. In many cases, good policies possess simplifying regularities such as symmetry and repetition, which allow for genotypes that are much smaller than the phenotypes they produce. Consequently, searching genotype space is more feasible, making it possible to scale to larger reinforcement-learning tasks.

Furthermore, indirect encodings often provide a natural way to exploit a task’s *geometry*, i.e., the spatial relationship between state features. In most direct encodings, such geometry cannot be exploited because it is not captured in the representation. For example, consider a neural network in which each input describes the current state of one square on a chess board. Since these inputs are treated as an unordered set, the distance between squares is not captured in the representation. Thus, structures for exploiting the relationship between squares must be evolved, complicating the task. In contrast, an indirect encoding can describe a network where the structure for processing each square’s state is a function of that square’s position on the board, with the natural consequence that nearby squares are processed similarly.

Like other evolutionary methods, systems using indirect encodings are inspired by analogies with biological systems: e.g., human beings have trillions of cells in their bodies but develop from genomes containing only tens of thousands of genes. Thus, it is not surprising that many indirect encodings are built on models of natural development. For example, *L-systems* (Lindenmayer, 1968), which are formal grammars for describing complex recursive structures, have been used to evolve both the morphology and control system for autonomous robots, greatly outperforming direct encodings (Hornby and Pollack, 2002). Similarly, *cellular encodings* (Gruau and Whitley, 1993; Gruau, 1994) evolve graph grammars for generating modular neural networks composed of simpler subnetworks. This approach allows evolution to exploit regularities in the solution structure by instantiating multiple copies of the same subnetwork in order to build a complete network.

More recently, the HyperNEAT method (Stanley et al, 2009) has been developed to extend NEAT to use indirect encodings. This approach is based on *compositional pattern producing networks* (CPPNs). CPPNs are neural networks for describing complex patterns. For example, a two-dimensional image can be described by a CPPN whose inputs correspond to an x - y position in the image and whose output corresponds to the color that should appear in that position. The image can then be generated by querying the CPPN at each x - y position and setting that position's color based on the output. Such CPPNs can be evolved by NEAT, yielding a developmental system with the CPPN as the genotype and the image as the phenotype.

In HyperNEAT, the CPPN is used to describe a neural network instead of an image. Thus, both the genotype and phenotype are neural networks. As illustrated in Figure 5, the nodes of the phenotypic network are laid out on a *substrate*, i.e., a grid, such that each has a position. The CPPN takes as input two positions instead of one and its output specifies the weight of the edge connecting the two corresponding nodes. As before, these CPPNs can be evolved by NEAT based on the fitness of the resulting phenotypic network, e.g., its performance as a policy in a reinforcement learning task. The CPPNs can be interpreted as describing a spatial pattern in a four-dimensional hypercube, yielding the name HyperNEAT. Because the developmental approach makes it easy to specify networks that exploit symmetries and regularities in complex tasks, HyperNEAT has proven an effective tool for reinforcement learning, with successful applications in domains such as checkers (Gauci and Stanley, 2008, 2010), keepaway soccer (Verbancsics and Stanley, 2010), and multi-agent systems (D'Ambrosio et al, 2010).

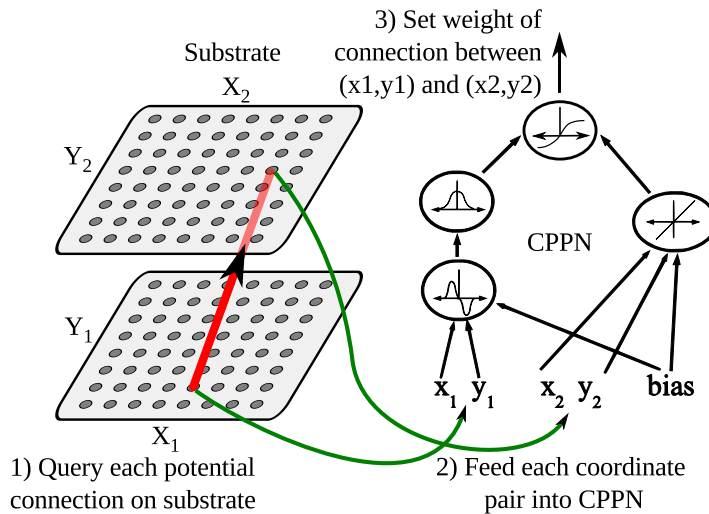


Fig. 5 The HyperNEAT algorithm. Figure taken with permission from (Gauci and Stanley, 2010).

7 On-Line Methods

While evolutionary methods have excelled in many challenging reinforcement-learning problems, their empirical success is largely restricted to *off-line* scenarios, in which the agent learns, not in the real-world, but in a ‘safe’ environment like a simulator. In other words, the problem specification usually includes a fitness function that requires only computational resources to evaluate, as in other optimization problems tackled with evolutionary computation.

In off-line scenarios, an agent’s only goal is to learn a good policy as quickly as possible. How much reward it obtains *while it is learning* is irrelevant because those rewards are only hypothetical and do not correspond to real-world costs. If the agent tries disastrous policies, only computation time is lost.

While efficient off-line learning is an important goal, it has limited practical applicability. In many cases, no simulator is available because the dynamics of the task are unknown, e.g., when a robot explores an unfamiliar environment or a chess player plays a new opponent. Other times, the dynamics of the task are too complex to accurately simulate, e.g., user behavior on a large computer network or the noise in a robot’s sensors and actuators.

Therefore, many researchers consider *on-line* learning a fundamental challenge in reinforcement learning. In an on-line learning scenario, the agent must maximize the reward it accrues while it is learning because those rewards correspond to real-world costs. For example, if a robot learning on-line tries a policy that causes it to drive off a cliff, then the negative reward the agent receives is not hypothetical; it corresponds to the real cost of fixing or replacing the robot.

Evolutionary methods have also succeeded on-line (Steels, 1994; Nordin and Banzhaf, 1997; Schroder et al, 2001), especially in evolutionary robotics (Meyer et al, 1998; Floreano and Urzelai, 2001; Floreano and Mondada, 2002; Pratihar, 2003; Kernbach et al, 2009; Zufferey et al, 2010), and some research has investigated customizing such methods to on-line settings (Floreano and Urzelai, 2001; Whiteson and Stone, 2006a,b; Priesterjahn et al, 2008; Tan et al, 2008; Cardamone et al, 2009, 2010).

Nonetheless, in most applications, researchers typically report performance using only the off-line measures common for optimization problems, e.g., the number of fitness evaluations needed to find a policy achieving a threshold performance or the performance of the best policy found after a given number of fitness evaluations. Therefore, determining how best to use evolutionary methods for reinforcement learning in on-line settings, i.e., how to maximize cumulative reward during evolution, remains an important and under-explored research area.

7.1 Model-Based Methods

One possible approach is to use evolution, not as a complete solution method, but as a component in a model-based method. In model-based algorithms, the agent’s

interactions with its environment are used to learn a model, to which planning methods are then applied. As the agent gathers more samples from the environment, the quality of the model improves, which, in turn, improves the quality of the policy produced via planning. Because planning is done off-line, the number of interactions needed to find a good policy is minimized, leading to strong on-line performance.

In such an approach, planning is typically conducted using dynamic programming methods like value iteration. However, many other methods can be used instead; if the model is continuous and/or high dimensional, evolutionary or other policy-search methods may be preferable. Unfortunately, most model-based methods are designed only to learn tabular models for small, discrete state spaces. Still, in some cases, especially when considerable domain expertise is available, more complex models can be learned.

For example, linear regression has been used learn models of helicopter dynamics, which can then be used for policy-search reinforcement learning (Ng et al, 2004). The resulting policies have successfully controlled real model helicopters. A similar approach was used to maximize on-line performance in the helicopter-hovering events in recent Reinforcement Learning Competitions (Whiteson et al, 2010a): models learned via linear regression were used as fitness functions for policies evolved off-line via neuroevolution (Koppejan and Whiteson, 2009).

Alternatively, evolutionary methods can be used for the model-learning component of a model-based solution. In particular, *anticipatory learning classifier systems* (Butz, 2002; Gerard et al, 2002, 2005; Sigaud et al, 2009), a type of LCS, can be used to evolve models of the environment that are used for planning in a framework similar to Dyna-Q (Sutton, 1990).

7.2 On-Line Evolutionary Computation

Another possible solution is *on-line evolutionary computation* (Whiteson and Stone, 2006a,b). The main idea is to borrow exploration strategies commonly used to select actions in temporal-difference methods and use them to select policies for evaluation in evolution. Doing so allows evolution to balance exploration and exploitation in a way that improves on-line performance.

Of course, evolutionary methods already strive to balance exploration and exploitation. In fact, this is one of the main motivations originally provided for genetic algorithms (Holland, 1975). However, this balance typically occurs only *across* generations, not *within* them. Once the members of each generation have been determined, they all typically receive the same evaluation time.

This approach makes sense in deterministic domains, where each member of the population can be accurately evaluated in a single episode. However, many real-world domains are stochastic, in which case fitness evaluations must be averaged over many episodes. In these domains, giving the same evaluation time to each member of the population can be grossly suboptimal because, within a generation, it is purely exploratory.

Instead, on-line evolutionary computation exploits information gained earlier in the generation to systematically give more evaluations to more promising policies and avoid re-evaluating weaker ones. This is achieved by employing temporal-difference exploration strategies to select policies for evaluation in each generation.

For example, ϵ -greedy selection can be used at the beginning of each episode to select a policy for evaluation. Instead of iterating through the population, evolution selects a policy randomly with probability ϵ . With probability $1 - \epsilon$, the algorithm selects the best policy discovered so far in the current generation. The fitness of each policy is just the average reward per episode it has received so far. Each time a policy is selected for evaluation, the total reward it receives is incorporated into that average, which can cause it to gain or lose the rank of best policy.

For the most part, ϵ -greedy selection does not alter evolution's search but simply interleaves it with exploitative episodes that increase average reward during learning. However, softmax selection can also be used to focus exploration on the most promising alternatives. At the beginning of each generation, each individual is evaluated for one episode, to initialize its fitness. Then, the remaining $e - |P|$ episodes are allocated according to a Boltzmann distribution.

Neither ϵ -greedy nor softmax consider the uncertainty of the estimates on which they base their selections, a shortcoming that can be addressed with interval estimation (Kaelbling, 1993). When used in temporal-difference methods, interval estimation computes a $(100 - \alpha)\%$ confidence interval for the value of each available action. The agent always takes the action with the highest upper bound on this interval. This strategy favors actions with high estimated value and also focuses exploration on promising but uncertain actions. The α parameter controls the balance between exploration and exploitation, with smaller values generating greater exploration.

The same strategy can be employed within evolution to select policies for evaluation. At the beginning of each generation, each individual is evaluated for one episode, to initialize its fitness. Then, the remaining $e - |P|$ episodes are allocated to the policy that currently has the highest upper bound on its confidence interval.

All three of these implementations of on-line evolutionary computation have been shown to substantially improve on-line performance, both in conjunction with NEAT (Whiteson and Stone, 2006b; Cardamone et al, 2009, 2010) and NEAT+Q (Whiteson and Stone, 2006a).

8 Conclusion

Evolutionary methods are a powerful tool for tackling challenging reinforcement learning problems. They are especially appealing for problems that include partial observability, have continuous action spaces, or where effective representations cannot be manually specified. Particularly in the area of neuroevolution, sophisticated methods exist for evolving neural-network topologies, decomposing the task based on network structure, and exploiting indirect encodings. Thanks to hybrid methods, the use of evolutionary computation does not require forgoing the power of

temporal-difference methods. Furthermore, coevolutionary approaches extend the reach of evolution to multi-agent reinforcement learning, both cooperative and competitive. While most work in evolutionary computation has focused on off-line settings, promising research exists in developing evolutionary methods for on-line reinforcement learning, which remains a critical and exciting challenge for future work.

Acknowledgments

Thanks to Ken Stanley, Risto Miikkulainen, Jürgen Schmidhuber, Martin Butz, Julian Bishop, and the anonymous reviewers for their invaluable input regarding the state of the art in evolutionary reinforcement learning.

References

- Ackley D, Littman M (1991) Interactions between learning and evolution. *Artificial Life II, SFI Studies in the Sciences of Complexity* 10:487–509
- Agogino AK, Tumer K (2008) Efficient evaluation functions for evolving coordination. *Evolutionary Computation* 16(2):257–288
- Arita T, Suzuki R (2000) Interactions between learning and evolution: The outstanding strategy generated by the Baldwin Effect. *Artificial Life* 7:196–205
- Baldwin JM (1896) A new factor in evolution. *The American Naturalist* 30:441–451
- Boers E, Borst M, Sprinkhuizen-Kuyper I (1995) Evolving Artificial Neural Networks using the “Baldwin Effect”. In: *Artificial Neural Nets and Genetic Algorithms, Proceedings of the International Conference in Ales, France*
- Bonarini A (2000) An introduction to learning fuzzy classifier systems. *Learning Classifier Systems* pp 83–104
- Bull L, Kovacs T (2005) Foundations of learning classifier systems: An introduction. *Foundations of Learning Classifier Systems* pp 1–17
- Bull L, O’Hara T (2002) Accuracy-based neuro and neuro-fuzzy classifier systems. In: *Proceedings of the Genetic and Evolutionary Computation Conference*, pp 905–911
- Butz M (2002) *Anticipatory learning classifier systems*. Kluwer Academic Publishers
- Butz M (2006) *Rule-based evolutionary online learning systems: A principled approach to LCS analysis and design*. Springer Verlag
- Butz M, Herbot O (2008) Context-dependent predictions and cognitive arm control with XCSF. In: *Proceedings of the 10th Annual Conference on Genetic and Evolutionary Computation, ACM*, pp 1357–1364
- Butz M, Lanzi P (2009) Sequential problems that test generalization in learning classifier systems. *Evolutionary Intelligence* 2(3):141–147
- Butz M, Goldberg D, Lanzi P (2005) Gradient descent methods in learning classifier systems: Improving XCS performance in multistep problems. *IEEE Transactions on Evolutionary Computation* 9(5)
- Butz M, Lanzi P, Wilson S (2008) Function approximation with XCS: Hyperellipsoidal conditions, recursive least squares, and compaction. *IEEE Transactions on Evolutionary Computation* 12(3):355–376

- Butz M, Pedersen G, Stalsh P (2009) Learning sensorimotor control structures with XCSF: Redundancy exploitation and dynamic control. In: Proceedings of the 11th Annual conference on Genetic and evolutionary computation, pp 1171–1178
- Cai Z, Peng Z (2002) Cooperative coevolutionary adaptive genetic algorithm in path planning of cooperative multi-mobile robot systems. *Journal of Intelligent and Robotic Systems* 33(1):61–71
- Cardamone L, Loiacono D, Lanzi P (2009) On-line neuroevolution applied to the open racing car simulator. In: Proceedings of the Congress on Evolutionary Computation (CEC), pp 2622–2629
- Cardamone L, Loiacono D, Lanzi PL (2010) Learning to drive in the open racing car simulator using online neuroevolution. *Computational Intelligence and AI in Games, IEEE Transactions on* 2(3):176–190
- Chellapilla K, Fogel D (2001) Evolving an expert checkers playing program without using human expertise. *IEEE Transactions on Evolutionary Computation* 5(4):422–428
- Coello C, Lamont G, Van Veldhuizen D (2007) *Evolutionary algorithms for solving multi-objective problems*. Springer-Verlag
- D’Ambrosio D, Lehman J, Risi S, Stanley KO (2010) Evolving policy geometry for scalable multi-agent learning. In: Proceedings of the Ninth International Conference on Autonomous Agents and Multiagent Systems (AAMAS-2010), pp 731–738
- Darwen P, Yao X (1996) Automatic modularization by speciation. In: Proceedings of the 1996 IEEE International Conference on Evolutionary Computation (ICEC96), pp 88–93
- Dasgupta D, McGregor D (1992) Designing application-specific neural networks using the structured genetic algorithm. In: Proceedings of the International Conference on Combinations of Genetic Algorithms and Neural Networks, pp 87–96
- Dawkins R, Krebs J (1979) Arms races between and within species. *Proceedings of the Royal Society of London Series B, Biological Sciences* 205(1161):489–511
- De Jong E (2004) The incremental Pareto-coevolution archive. In: *Genetic and Evolutionary Computation—GECCO 2004*, Springer, pp 525–536
- De Jong E (2007) A monotonic archive for Pareto-coevolution. *Evolutionary computation* 15(1):61–93
- De Jong K, Spears W (1991) An analysis of the interacting roles of population size and crossover in genetic algorithms. *Parallel problem solving from nature* pp 38–47
- De Jong K, Spears W, Gordon D (1993) Using genetic algorithms for concept learning. *Machine learning* 13(2):161–188
- Deb K (2001) *Multi-objective optimization using evolutionary algorithms*. Wiley
- Dorigo M, Colombetti M (1998) *Robot shaping: An experiment in behavior engineering*. The MIT Press
- Downing KL (2001) Reinforced genetic programming. *Genetic Programming and Evolvable Machines* 2(3):259–288
- Doya K (2000) Reinforcement learning in continuous time and space. *Neural Computation* 12(1):219–245
- Drugowitsch J (2008) *Design and analysis of learning classifier systems: A probabilistic approach*. Springer Verlag
- Ficici S, Pollack J (2000) A game-theoretic approach to the simple coevolutionary algorithm. In: *Parallel Problem Solving from Nature PPSN VI*, Springer, pp 467–476
- Ficici S, Pollack J (2001) Pareto optimality in coevolutionary learning. *Advances in Artificial Life* pp 316–325
- Floreano D, Mondada F (2002) Evolution of homing navigation in a real mobile robot. *IEEE Transactions on Systems, Man, and Cybernetics, Part B* 26(3):396–407
- Floreano D, Urzelai J (2001) Evolution of plastic control networks. *Autonomous Robots* 11(3):311–317
- French R, Messinger A (1994) Genes, phenes and the Baldwin effect: Learning and evolution in a simulated population. *Artificial Life* 4:277–282
- Gaskett C, Wettergreen D, Zelinsky A (1999) Q-learning in continuous state and action spaces. *Advanced Topics in Artificial Intelligence* pp 417–428

- Gauci J, Stanley KO (2008) A case study on the critical role of geometric regularity in machine learning. In: Proceedings of the Twenty-Third AAAI Conference on Artificial Intelligence (AAAI-2008)
- Gauci J, Stanley KO (2010) Autonomous evolution of topographic regularities in artificial neural networks. *Neural Computation* 22(7):1860–1898
- Gerard P, Stolzmann W, Sigaud O (2002) YACS: a new learning classifier system using anticipation. *Soft Computing-A Fusion of Foundations, Methodologies and Applications* 6(3):216–228
- Gerard P, Meyer J, Sigaud O (2005) Combining latent learning with dynamic programming in the modular anticipatory classifier system. *European Journal of Operational Research* 160(3):614–637
- Giraud-Carrier C (2000) Unifying learning with evolution through Baldwinian evolution and Lamarckism: A case study. In: Proceedings of the Symposium on Computational Intelligence and Learning (CoIL-2000), pp 36–41
- Goldberg D (1989) *Genetic Algorithms in Search, Optimization and Machine Learning*. Addison-Wesley
- Goldberg D, Deb K (1991) A comparative analysis of selection schemes used in genetic algorithms. *Foundations of genetic algorithms* 1:69–93
- Goldberg D, Richardson J (1987) Genetic algorithms with sharing for multimodal function optimization. In: Proceedings of the Second International Conference on Genetic Algorithms and their Application, p 49
- Gomez F, Miikkulainen R (1999) Solving non-Markovian control tasks with neuroevolution. In: Proceedings of the International Joint Conference on Artificial Intelligence, pp 1356–1361
- Gomez F, Miikkulainen R (2003) Active guidance for a finless rocket using neuroevolution. In: GECCO-03: Proceedings of the Genetic and Evolutionary Computation Conference
- Gomez F, Schmidhuber J (2005a) Co-evolving recurrent neurons learn deep memory POMDPs. In: GECCO-05: Proceedings of the Genetic and Evolutionary Computation Conference, pp 491–498
- Gomez F, Schmidhuber J (2005b) Evolving modular fast-weight networks for control. *Artificial Neural Networks: Formal Models and Their Applications-ICANN 2005* pp 383–389
- Gomez F, Schmidhuber J, Miikkulainen R (2006) Efficient non-linear control through neuroevolution. In: Proceedings of the European Conference on Machine Learning
- Gomez F, Schmidhuber J, Miikkulainen R (2008) Accelerated neural evolution through cooperatively coevolved synapses. *Journal of Machine Learning Research* 9:937–965
- Gruau F (1994) Automatic definition of modular neural networks. *Adaptive Behavior* 3(2):151
- Gruau F, Whitley D (1993) Adding learning to the cellular development of neural networks: Evolution and the Baldwin effect. *Evolutionary Computation* 1:213–233
- Hansen N, Müller S, Koumoutsakos P (2003) Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (CMA-ES). *Evolutionary Computation* 11(1):1–18
- van Hasselt H, Wiering M (2007) Reinforcement learning in continuous action spaces. In: *IEEE International Symposium on Approximate Dynamic Programming and Reinforcement Learning, 2007. ADPRL 2007*, pp 272–279
- Haykin S (1994) *Neural networks: a comprehensive foundation*. Prentice Hall
- Heidrich-Meisner V, Igel C (2008) Variable metric reinforcement learning methods applied to the noisy mountain car problem. *Recent Advances in Reinforcement Learning* pp 136–150
- Heidrich-Meisner V, Igel C (2009a) Hoeffding and Bernstein races for selecting policies in evolutionary direct policy search. In: Proceedings of the 26th Annual International Conference on Machine Learning, pp 401–408
- Heidrich-Meisner V, Igel C (2009b) Neuroevolution strategies for episodic reinforcement learning. *Journal of Algorithms* 64(4):152–168
- Heidrich-Meisner V, Igel C (2009c) Uncertainty handling CMA-ES for reinforcement learning. In: Proceedings of the 11th Annual Conference on Genetic and Evolutionary Computation, pp 1211–1218

- Hillis W (1990) Co-evolving parasites improve simulated evolution as an optimization procedure. *Physica D: Nonlinear Phenomena* 42(1-3):228–234
- Hinton GE, Nowlan SJ (1987) How learning can guide evolution. *Complex Systems* 1:495–502
- Holland J, Reitman J (1977) Cognitive systems based on adaptive algorithms. *ACM SIGART Bulletin* 63:49–49
- Holland JH (1975) *Adaptation in Natural and Artificial Systems: An Introductory Analysis with Applications to Biology, Control and Artificial Intelligence*. University of Michigan Press
- Hornby G, Pollack J (2002) Creating high-level components with a generative representation for body-brain evolution. *Artificial Life* 8(3):223–246
- Igel C (2003) Neuroevolution for reinforcement learning using evolution strategies. In: *Congress on Evolutionary Computation*, vol 4, pp 2588–2595
- Jansen T, Wiegand RP (2004) The cooperative coevolutionary (1+1) EA. *Evolutionary Computation* 12(4):405–434
- Kaelbling LP (1993) *Learning in Embedded Systems*. MIT Press
- Kernbach S, Meister E, Scholz O, Humza R, Liedke J, Ricotti L, Jemai J, Havlik J, Liu W (2009) Evolutionary robotics: The next-generation-platform for on-line and on-board artificial evolution. In: *CEC'09: IEEE Congress on Evolutionary Computation*, pp 1079–1086
- Kohl N, Miikkulainen R (2008) Evolving neural networks for fractured domains. In: *Proceedings of the Genetic and Evolutionary Computation Conference*, pp 1405–1412
- Kohl N, Miikkulainen R (2009) Evolving neural networks for strategic decision-making problems. *Neural Networks* 22:326–337, special issue on Goal-Directed Neural Systems.
- Koppejan R, Whiteson S (2009) Neuroevolutionary reinforcement learning for generalized helicopter control. In: *GECCO 2009: Proceedings of the Genetic and Evolutionary Computation Conference*, pp 145–152
- Kovacs T (2003) *Strength or accuracy: credit assignment in learning classifier systems*. Springer-Verlag
- Larranaga P, Lozano J (2002) *Estimation of distribution algorithms: A new tool for evolutionary computation*. Springer Netherlands
- Lindenmayer A (1968) Mathematical models for cellular interactions in development II. Simple and branching filaments with two-sided inputs. *Journal of Theoretical Biology* 18(3):300–315
- Littman ML, Dean TL, Kaelbling LP (1995) On the complexity of solving Markov decision processes. In: *Proceedings of the Eleventh International Conference on Uncertainty in Artificial Intelligence*, pp 394–402
- Lucas SM, Runarsson TP (2006) Temporal difference learning versus co-evolution for acquiring othello position evaluation. In: *IEEE Symposium on Computational Intelligence and Games*
- Lucas SM, Togelius J (2007) Point-to-point car racing: an initial study of evolution versus temporal difference learning. In: *IEEE Symposium on Computational Intelligence and Games*, pp 260–267
- Mahadevan S, Maggioni M (2007) Proto-value functions: A Laplacian framework for learning representation and control in Markov decision processes. *Journal of Machine Learning Research* 8:2169–2231
- Mahfoud S (1995) A comparison of parallel and sequential niching methods. In: *Conference on Genetic Algorithms*, vol 136, p 143
- McQuesten P, Miikkulainen R (1997) Culling and teaching in neuro-evolution. In: *Proceedings of the Seventh International Conference on Genetic Algorithms*, pp 760–767
- Meyer J, Husband P, Harvey I (1998) Evolutionary robotics: A survey of applications and problems. In: *Evolutionary Robotics*, Springer, pp 1–21
- Millán J, Posenato D, Dedieu E (2002) Continuous-action Q-learning. *Machine Learning* 49(2):247–265
- Monroy G, Stanley K, Miikkulainen R (2006) Coevolution of neural networks using a layered Pareto archive. In: *Proceedings of the 8th Annual Conference on Genetic and Evolutionary Computation*, p 336
- Moriarty D, Miikkulainen R (1997) Forming neural networks through efficient and adaptive co-evolution. *Evolutionary Computation* 5(4):373–399

- Moriarty DE, Miikkulainen R (1996) Efficient reinforcement learning through symbiotic evolution. *Machine Learning* 22(11):11–33
- Moriarty DE, Schultz AC, Grefenstette JJ (1999) Evolutionary algorithms for reinforcement learning. *Journal of Artificial Intelligence Research* 11:199–229
- Ng AY, Coates A, Diel M, Ganapathi V, Schulte J, Tse B, Berger E, Liang E (2004) Inverted autonomous helicopter flight via reinforcement learning. In: *Proceedings of the International Symposium on Experimental Robotics*
- Nolfi S, Parisi D (1997) Learning to adapt to changing environments in evolving neural networks. *Adaptive Behavior* 5(1):75–98
- Nolfi S, Elman JL, Parisi D (1994) Learning and evolution in neural networks. *Adaptive Behavior* 2:5–28
- Nordin P, Banzhaf W (1997) An on-line method to evolve behavior and to control a miniature robot in real time with genetic programming. *Adaptive Behavior* 5(2):107
- Panait L, Luke S (2005) Cooperative multi-agent learning: The state of the art. *Autonomous Agents and Multi-Agent Systems* 11(3):387–434
- Panait L, Luke S, Harrison JF (2006) Archive-based cooperative coevolutionary algorithms. In: *GECCO '06: Proceedings of the 8th Annual Conference on Genetic and Evolutionary Computation*, pp 345–352
- Parr R, Painter-Wakefield C, Li L, Littman M (2007) Analyzing feature generation for value-function approximation. In: *Proceedings of the 24th International Conference on Machine Learning*, p 744
- Pereira FB, Costa E (2001) Understanding the role of learning in the evolution of busy beaver: A comparison between the Baldwin Effect and a Lamarckian strategy. In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO-2001)*
- Peters J, Schaal S (2008) Natural actor-critic. *Neurocomputing* 71(7-9):1180–1190
- Pollack J, Blair A (1998) Co-evolution in the successful learning of backgammon strategy. *Machine Learning* 32(3):225–240
- Popovici E, Bucci A, Wiegand P, De Jong E (2010) Coevolutionary principles. In: *Rozenberg G, Baeck T, Kok J (eds) Handbook of Natural Computing*, Springer-Verlag, Berlin
- Potter MA, De Jong KA (1995) Evolving neural networks with collaborative species. In: *Summer Computer Simulation Conference*, pp 340–345
- Potter MA, De Jong KA (2000) Cooperative coevolution: An architecture for evolving coadapted subcomponents. *Evolutionary Computation* 8:1–29
- Pratihari D (2003) Evolutionary robotics: A review. *Sadhana* 28(6):999–1009
- Priesterjahn S, Weimer A, Eberling M (2008) Real-time imitation-based adaptation of gaming behaviour in modern computer games. In: *Proceedings of the Genetic and Evolutionary Computation Conference*, pp 1431–1432
- Radcliffe N (1993) Genetic set recombination and its application to neural network topology optimisation. *Neural Computing & Applications* 1(1):67–90
- Rosin CD, Belew RK (1997) New methods for competitive coevolution. *Evolutionary Computation* 5(1):1–29
- Rubinstein R, Kroese D (2004) *The cross-entropy method: a unified approach to combinatorial optimization, Monte-Carlo simulation, and machine learning*. Springer-Verlag
- Runarsson TP, Lucas SM (2005) Co-evolution versus self-play temporal difference learning for acquiring position evaluation in small-board go. *IEEE Transactions on Evolutionary Computation* 9:628–640
- Schmidhuber J, Wierstra D, Gomez FJ (2005) Evolino: Hybrid neuroevolution / optimal linear search for sequence learning. In: *Proceedings of the Nineteenth International Joint Conference on Artificial Intelligence*, pp 853–858
- Schmidhuber J, Wierstra D, Gagliolo M, Gomez F (2007) Training recurrent networks by evolino. *Neural computation* 19(3):757–779
- Schroder P, Green B, Grum N, Fleming P (2001) On-line evolution of robust control systems: an industrial active magnetic bearing application. *Control Engineering Practice* 9(1):37–49

- Sigaud O, Butz M, Kozlova O, Meyer C (2009) Anticipatory Learning Classifier Systems and Factored Reinforcement Learning. *Anticipatory Behavior in Adaptive Learning Systems* pp 321–333
- Stanley K, Miikkulainen R (2003) A taxonomy for artificial embryogeny. *Artificial Life* 9(2):93–130
- Stanley KO, Miikkulainen R (2002) Evolving neural networks through augmenting topologies. *Evolutionary Computation* 10(2):99–127
- Stanley KO, Miikkulainen R (2004a) Competitive coevolution through evolutionary complexification. *Journal of Artificial Intelligence Research* 21:63–100
- Stanley KO, Miikkulainen R (2004b) Evolving a roving eye for go. In: *Proceedings of the Genetic and Evolutionary Computation Conference*
- Stanley KO, Bryant BD, Miikkulainen R (2003) Evolving adaptive neural networks with and without adaptive synapses. In: *Proceedings of the 2003 Congress on Evolutionary Computation (CEC 2003)*, vol 4, pp 2557–2564
- Stanley KO, D’Ambrosio DB, Gauci J (2009) A hypercube-based indirect encoding for evolving large-scale neural networks. *Artificial Life* 15(2):185–212
- Steels L (1994) Emergent functionality in robotic agents through on-line evolution. In: *Artificial Life IV: Proceedings of the Fourth International Workshop on the Synthesis and Simulation of Living Systems*, pp 8–16
- Sutton RS (1990) Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In: *Proceedings of the Seventh International Conference on Machine Learning*, pp 216–224
- Sutton RS, Barto AG (1998) *Reinforcement Learning: An Introduction*. MIT Press
- Sywerda G (1989) Uniform crossover in genetic algorithms. In: *Proceedings of the Third International Conference on Genetic Algorithms*, pp 2–9
- Tan C, Ang J, Tan K, Tay A (2008) Online adaptive controller for simulated car racing. In: *Congress on Evolutionary Computation (CEC)*, pp 2239–2245
- Taylor ME, Whiteson S, Stone P (2006) Comparing evolutionary and temporal difference methods in a reinforcement learning domain. In: *GECCO 2006: Proceedings of the Genetic and Evolutionary Computation Conference*, pp 1321–1328
- Tesauro G (1994) TD-gammon, a self-teaching backgammon program achieves master-level play. *Neural Computation* 6:215–219
- Tesauro G (1998) Comments on co-evolution in the successful learning of backgammon strategy. *Machine Learning* 32(3):241–243
- Verbancsics P, Stanley K (2010) Evolving Static Representations for Task Transfer. *Journal of Machine Learning Research* 11:1737–1769
- Von Neumann J (1928) Zur Theorie der Gesellschaftsspiele *Math. Annalen* 100:295–320
- Whiteson S, Stone P (2006a) Evolutionary function approximation for reinforcement learning. *Journal of Machine Learning Research* 7:877–917
- Whiteson S, Stone P (2006b) On-line evolutionary computation for reinforcement learning in stochastic domains. In: *GECCO 2006: Proceedings of the Genetic and Evolutionary Computation Conference*, pp 1577–1584
- Whiteson S, Kohl N, Miikkulainen R, Stone P (2005) Evolving keepaway soccer players through task decomposition. *Machine Learning* 59(1):5–30
- Whiteson S, Tanner B, White A (2010a) The reinforcement learning competitions. *AI Magazine* 31(2):81–94
- Whiteson S, Taylor ME, Stone P (2010b) Critical factors in the empirical performance of temporal difference and evolutionary methods for reinforcement learning. *Autonomous Agents and Multi-Agent Systems* 21(1):1–27
- Whitley D, Dominic S, Das R, Anderson CW (1993) Genetic reinforcement learning for neuro-control problems. *Machine Learning* 13:259–284
- Whitley D, Gordon S, Mathias K (1994) Lamarckian evolution, the Baldwin effect and function optimization. In: *Parallel Problem Solving from Nature - PPSN III*, pp 6–15

- Wiegand R, Liles W, De Jong K (2001) An empirical analysis of collaboration methods in cooperative coevolutionary algorithms. In: Proceedings of the Genetic and Evolutionary Computation Conference (GECCO), pp 1235–1242
- Wieland A (1991) Evolving neural network controllers for unstable systems. In: International Joint Conference on Neural Networks, vol 2, pp 667–673
- Wilson S (1995) Classifier fitness based on accuracy. *Evolutionary computation* 3(2):149–175
- Wilson S (2001) Function approximation with a classifier system. In: GECCO-01: Proceedings of the Genetic and Evolutionary Computation Conference, pp 974–982
- Wolpert D, Tumer K (2002) Optimal payoff functions for members of collectives. *Modeling complexity in economic and social systems* p 355
- Yamasaki K, Sekiguchi M (2000) Clear explanation of different adaptive behaviors between Darwinian population and Lamarckian population in changing environment. In: Proceedings of the Fifth International Symposium on Artificial Life and Robotics, vol 1, pp 120–123
- Yao X (1999) Evolving artificial neural networks. *Proceedings of the IEEE* 87(9):1423–1447
- Yong CH, Miikkulainen R (2007) Coevolution of role-based cooperation in multi-agent systems. Tech. Rep. AI07-338, Department of Computer Sciences, The University of Texas at Austin
- Zhang B, Muhlenbein H (1993) Evolving optimal neural networks using genetic algorithms with Occam’s razor. *Complex Systems* 7(3):199–220
- Zufferey J, Floreano D, Van Leeuwen M, Merenda T (2010) Evolving vision-based flying robots. In: *Biologically Motivated Computer Vision*, Springer, pp 13–29