

Improving Search Algorithms by Using Intelligent Coordinates

David Wolpert, Kagan Tumer, and Esfandiar Bandari
NASA Ames Research Center, Moffett Field, CA, 94035, USA

We consider the problem of designing a set of computational agents so that as they all pursue their self-interests a global function G of the collective system is optimized. Three factors govern the quality of such design. The first relates to conventional exploration-exploitation search algorithms for finding the maxima of such a global function, e.g., simulated annealing (SA). Game-theoretic algorithms instead are related to the second of those factors, and the third is related to techniques from the field of machine learning. Here we demonstrate how to exploit all three factors by modifying the search algorithm's exploration stage so that rather than by random sampling, each coordinate of the underlying search space is controlled by an associated machine-learning-based "player" engaged in a non-cooperative game. Experiments demonstrate that this modification improves SA by up to an order of magnitude for bin-packing and for a model of an economic process run over an underlying network. These experiments also reveal novel small worlds phenomena.

PACS numbers: 89.20.Ff, 89.75.-k, 89.75.Fb, 02.60.Pn, 02.70.-c, 02.70.Tt

I. INTRODUCTION

Many distributed systems found in nature have inspired function-maximization algorithms. In some of these the coordinates of the underlying system are viewed as players engaged in a non-cooperative game, whose joint behavior (hopefully) maximizes the pre-specified global function of the entire system. Examples of such systems are auctions and clearing of markets. Typically in the computer-based algorithms inspired by such "collectives" of players, each separate coordinate of the system is controlled by an associated machine learning algorithm [3, 4, 7, 10, 16, 23], reinforcement-learning (RL) algorithms being particularly common [17, 22].

One important issue concerning such collectives is whether the payoff function g_η of each player η is sufficiently sensitive to what coordinates η controls in comparison to the other coordinates, so that η can learn how to control its coordinates to achieve high payoff. A second crucial issue is the need for all of the g_η to be "aligned" with G , so that as the players individually learn how to increase their payoffs, G also increases.

Previous work in the COllective INtelligence (COIN) framework addresses these issues. This work extends conventional game-theoretic mechanism design [8, 15] to include off-equilibrium behavior, learnability issues, g_η with non-human attributes (e.g., g_η for which incentive compatibility is irrelevant), and arbitrary G . In domains from network routing to congestion problems it outperforms traditional techniques, by up to several orders of magnitude for large systems [18, 21, 22].

Other collective systems found in nature that have inspired search algorithms do not involve players conducting a non-cooperative game. Examples include spin glasses, genomes undergoing neo-Darwinian natural selection, and eusocial insect colonies, which have been translated into simulated annealing (SA [9, 11]), genetic algorithms [1, 5], and swarm intelligence [2, 12], respectively. An important issue here is the exploration/exploitation dynamics of the overall collective.

Recent analysis reveals how G is governed by the interaction between exploration/exploitation, the alignment of the g_η and G , and the learnability of the g_η [21]. Here we use that analysis to motivate a hybrid algorithm, Intelligent Coordinates for search (IC), that addresses all three issues. It works by modifying any exploration-based search algorithm so that each coordinate being searched is made "intelligent", its exploration value being the move of a game-playing computer algorithm rather than the random sample of a probability distribution.

Like SA, IC is intended to be used as an "off the shelf" algorithm; rarely will it be the best possible algorithm for some particular domain. Rather it is designed for use in very large problems where parallelization can provide a large advantage, while there is little exploitable information concerning gradients. We present experiments comparing IC and SA on two archetypal domains: bin-packing and an economic model of people choosing formats for their home music systems.

In the bin-packing domain IC achieves a given value of G up to three orders of magnitude faster than does SA, with the improvement increasing linearly with the size of the problem. In the format choice problem G is the sum of each person's "happiness" with her format choices. Each person η 's happiness with each of her choices is set by three factors: which of her nearest neighbors on a ring network (η 's "friends") make that choice; η 's intrinsic preference for that choice; and the price of music purchased in that format, inversely proportional to the total number of players using that choice. Here again, IC improves G two orders of magnitude more quickly than does SA. We also considered an algorithm similar to the Groves mechanism of economics; IC outperformed it by over two orders of magnitude. We also modified the ring to be a small-worlds network [13, 14, 19]. This barely improved IC's performance (3%), with no effect on the other algorithms. However if G was also changed, so that each η 's happiness depends on agreeing with her friends' friends, the performance increase in changing to a small-worlds topology is significant (10%). This underscores

the multiplicity of factors behind the benefits of small-worlds networks.

II. SIMPLIFIED THEORY OF COLLECTIVES

Let $z \in \zeta$ be the joint move of all agents/players in the collective. We want the z that maximizes the provided world utility $G(z)$. In addition to G we have private utility functions $\{g_\eta\}$, one for each agent η controlling z_η . $\hat{\eta}$ refers to all agents other than η .

Intelligence “standardizes” utility functions so that the value they assign to z only reflects their ranking of z relative to some other z' . One form of it is

$$N_{\eta,U}(z) \equiv \int d\mu_{z_\eta}(z') \Theta[U(z) - U(z')], \quad (1)$$

where Θ is the Heaviside function, and where the subscript on the (normalized) measure $d\mu$ indicates it is restricted to z' such that $z'_\eta = z_\eta$.

Our uncertainty concerning the system induces a distribution $P(z)$. All attributes of the collective we can set, e.g., the private utility functions of the agents, are given by the value of the **design coordinate** s . Bayes theorem provides the **central equation**:

$$P(G | s) = \int d\vec{N}_G P(G | \vec{N}_G, s) \int d\vec{N}_g P(\vec{N}_G | \vec{N}_g, s) P(\vec{N}_g | s), \quad (2)$$

where \vec{N}_G and \vec{N}_g are the intelligence vectors for all the agents, for utilities g_η and for G , respectively. $N_{\eta,g_\eta}(z) = 1$ means that agent η 's move maximizes its utility, given the moves of the other agents. So $\vec{N}_g(z) = \vec{1}$ means z is a Nash equilibrium. Conversely, $\vec{N}_G(z') = \vec{1}$ means that the value of G cannot increase in moving from z' along any single (sic) coordinate of ζ . So if these two points are identical, then if the agents do well enough at maximizing their private utilities they must be near an (on-axis) maximizing point for G .

More formally, say for our s the third conditional probability in the integrand in the central equation (“term 3”) is peaked near $\vec{N}_g = \vec{1}$. Then s probably induces large (private utility function) intelligences (intuitively, the utilities are learnable). If in addition the second term is peaked near $\vec{N}_G = \vec{N}_g$, then \vec{N}_G will also be large (intuitively, the private utility is “aligned with G ”). This peakedness is assured if $\vec{N}_g = \vec{N}_G$ exactly $\forall z$. Such a system is said to be **factored**. Finally, if the first term in the integrand is peaked about high G when \vec{N}_G is large, then s probably induces high G , as desired.

As a trivial example, a **team game**, where $g_\eta = G \forall \eta$, is factored [7]. However team games usually have poor third terms, especially in large collectives. This is because each η has to discern how its moves affect $g_\eta = G$, given the background of the (varying) moves of the other agents whose moves comparably affect G .

Fix some $f(z_\eta)$, two moves z_η^1 and z_η^2 , a utility U , a value s , and a z_η . The associated **learnability** is

$$\Lambda_f(U; z_\eta, s, z_\eta^1, z_\eta^2) \equiv \sqrt{\frac{[E(U; z_\eta, z_\eta^1) - E(U; z_\eta, z_\eta^2)]^2}{\int dz_\eta [f(z_\eta) \text{Var}(U; z_\eta, z_\eta)]}}. \quad (3)$$

The averages and variance here are evaluated according to $P(U|n_\eta)P(n_\eta|z_\eta, z_\eta^1)$, $P(U|n_\eta)P(n_\eta|z_\eta, z_\eta)$, and $P(U|n_\eta)P(n_\eta|z_\eta, z_\eta^2)$, respectively, where n_η is η 's training set, formed by sampling U .

The denominator in Eq. 3 reflects the sensitivity of $U(z)$ to z_η , while the numerator reflects its sensitivity to z_η . So the greater the learnability of g_η , the more $g_\eta(z)$ depends only on the move of agent η , i.e., the more learnable g_η is. More formally, it can be shown that if appropriately scaled, g'_η will result in better expected intelligence for agent η than will g_η whenever $\Lambda_f(g'_\eta; z_\eta, s, z_\eta^1, z_\eta^2) > \Lambda_f(g_\eta; z_\eta, s, z_\eta^1, z_\eta^2)$ for all pairs of moves z_η^1, z_η^2 [20].

A **difference utility** is one of the form $U(z) = G(z) - D(z_\eta)$. Any difference utility is factored [20]. In addition, under usually benign approximations, the $D(z_\eta)$ that maximizes $\Lambda_f(U; z_\eta, s, z_\eta^1, z_\eta^2)$ for all pairs z_η^1, z_η^2 is $E_f(G(z) | z_\eta, s)$, where the expectation value is over z_η . The associated difference utility is called the **Aristocrat utility** (AU). If each η uses its AU as its private utility, then we have both good terms 2 and 3.

Evaluating the expectation value in AU can be difficult in practice. This motivates the **Wonderful Life Utility** (WLU), which requires no such evaluation:

$$WLU_\eta \equiv G(z) - G(z_\eta, CL_\eta), \quad (4)$$

where CL_η is the **clamping parameter**. WLU is factored, independent of the clamping parameter. Furthermore, while not matching AU , WLU typically has far better learnability than does a team game, and therefore typically results in better values of G . It is also often easier to evaluate than is G itself [18, 21].

One way to address term 1 as well as 2 and 3 is to incorporate exploration/exploitation techniques like SA.

III. EXPERIMENTS

In our version of SA, at the beginning of each time-step t a distribution $h_\eta(\zeta_\eta)$ is formed for every η by allotting probability 75% to the move η had at the end of the preceding time-step, $z_{\eta,t-1}$, and uniformly dividing probability 25% across all of its other moves. The “exploration” joint-move z_{expl} is then formed by simultaneously sampling all the h_η . If $G(z_{expl}) > G(z_{t-1})$, $z_{\eta,t}$ is set to z_{expl} . Otherwise z_t is set by sampling a Boltzmann distribution having energies $G(z_{t-1})$ and $G(z_{expl})$. Many different annealing schedules were investigated; all results below are for best schedules found.

IC is identical except that each h_η is replaced by $\frac{h_\eta(z_\eta)c_{\eta,t}(z_\eta)}{\sum_{a'} h_\eta(a'_\eta)c_{\eta,t}(a'_\eta)}$, where the distribution $c_{\eta,t}$ is set by an RL algorithm trying to optimize payoffs g_η . Here RL is done using a training set $n_{\eta,t}$ of all preceding move-payoff pairs, $\{(z_{\eta,t'}, g_\eta(z_{t'})) : t' < t\}$. For each possible move by η one forms the weighted average of the payoffs recorded in $n_{\eta,t}$ that occurred with that move, where the weights decay exponentially in $t - t'$. $c_{\eta,t}$ then is the Boltzmann distribution, parameterized by a “learning temperature” (that effectively rescales g_η) over those averages.

In all our experiments the “AU” version of IC approximated f to be uniform $\forall \eta$, and then used a mean-field approximation to pull the expectation inside G . Unless otherwise specified, the clamping elements used in WLU’s were set to $\vec{0}$.

In bin-packing N items, all of size $< c$, must be assigned into a minimal subset of N bins, without assigning a summed size $> c$ to any one bin. G of an assignment pattern is the number of occupied bins [6], and each agent controls the bin choice of one item. To improve performance all algorithms use a modified “ G ”, G_{soft} , even though their performance is measured with G :

$$G_{\text{soft}} = \begin{cases} \sum_{i=1}^N \left[\left(\frac{c}{2}\right)^2 - \left(x_i - \frac{c}{2}\right)^2 \right] & \text{if } x_i \leq c \\ \sum_{i=1}^N \left(x_i - \frac{c}{2}\right)^2 & \text{if } x_i > c \end{cases}, \quad (5)$$

where x_i is the summed size of all items in bin i . (Use of G_{soft} encourages bins to be either full or empty.)

In the IC runs learning temperature was .2, and all agents made the transition to RL-based moves after a period of 100 random z ’s used to generate the starting n_η . Exploitation temperature started at .5 for all algorithms, and was multiplied by .8 every 100 exploitation time-steps. In each SA run, the distribution h was slowly modified to generate solutions that differed in fewer items than the current solution as time progressed.

Algorithm	Ave. G	Best	Worst	% Optimum
IC WLU	3.32 ± 0.22	2	8	72 %
IC TG	7.84 ± 0.17	6	10	0 %
COIN WLU	3.52 ± 0.20	2	7	64 %
COIN TG	7.84 ± 0.15	6	9	0 %
SA	6.00 ± 0.19	4	7	0 %

TABLE I: Bin-packing G at time 200 for $N = 20, c = 12$.

In Table 1 “Best” refers to the best end-of-run G achieved by the associated algorithm, “worst” is the worst value, and “%Optimum” is the percentage of runs that were within one bin of the best value. Fig. 1 shows average performances (over 25 runs) as a function of time step. The algorithms that account for both terms 2 and 3 — IC WLU and COIN WLU — far outperform the others, with the algorithm accounting for all three terms doing best. The worst algorithms were those that accounted for only a single term (SA and COIN TG).

Linearly (i.e., optimistically) extrapolating SA’s performance from time 15000 indicates it would take over 1000 times as long as IC WLU to reach the G value IC WLU reaches at time 200. In addition the ratio of WLU’s time 1000 performance (relative to random search) to SA’s grows linearly with the size of the problem. Finally, Fig. 2 illustrates that the benefit of addressing terms 2 and 3 grows with the difficulty of the problem. In both figures SA outperforms IC - TG; this is due to there being more parameter-tuning with SA.

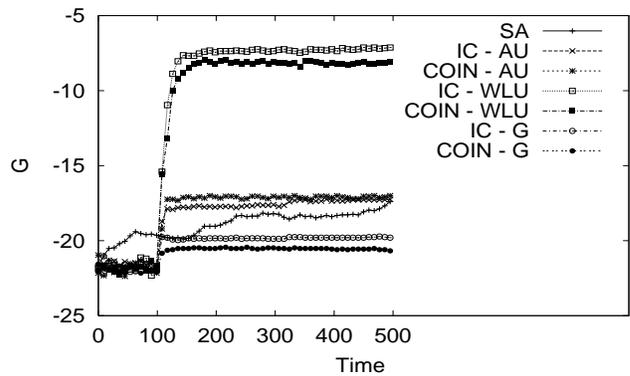


FIG. 1: Average bin-packing G for $N = 50, c = 10$. All error bars $\leq .31$ except IC - AU and COIN - AU are $\leq .57$.

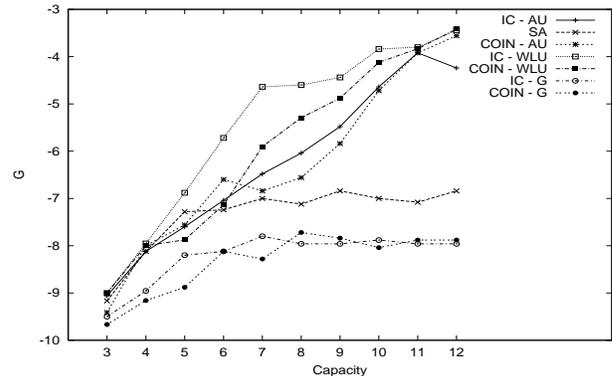


FIG. 2: G vs. c for $N = 20$ at $t = 200$. All error bars $\leq .34$.

For the format choice problem G is the sum over all N_a agents η of η ’s “happiness” with its music formats:

$$G = \sum_{\eta=1}^{N_a} \sum_{i=1}^{N_f} \sum_{\eta' \in \text{neigh}_\eta} \vartheta(i) \omega_{\eta,\eta',i} \text{pref}_{\eta,i} \quad (6)$$

where N_f is the numbers of formats; neigh_η is the set of players lying $\leq D$ hops away from player η ; $\text{pref}_{\eta,i}$ is η ’s intrinsic preference for format i (set randomly at initialization $\in [0, 1]$); $\vartheta(i)$ is the total number of players that choose format i (i.e., the inverse price for format i); and $\omega_{i,\eta,\eta'} = 1$ if the choices of players η and η' both include the format i , and 0 otherwise (each agent’s move is a selection of three of four total formats, implemented

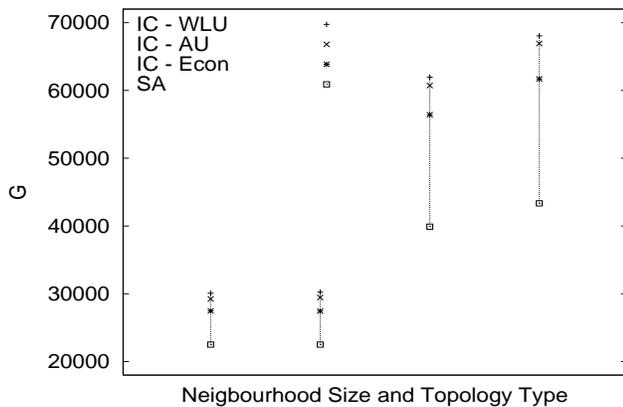


FIG. 3: $G(t = 200)$ for 100 agents. In order from left to right, $D = \{1, 1, 3, 3\}$, and topologies are $\{L, W, L, W\}$.

by choosing the one format not to be used). D values of both 1 and 3 were investigated.

In Fig. 3, “IC Econ” refers to WLU IC where clamping means the agent chooses no format whatsoever. It is essentially the game-theory Groves mechanism wherein one

sets g_η to η ’s marginal contribution to G , here rescaled and interleaved with a simulated annealing step to improve performance. “IC-WLU” instead clamps η ’s move to zero (in accord with the theory of collectives), which means that η chooses all formats. Learning temperature was now .4, and exploitation temperature was .05 (annealing provided no advantage since runs were short). Two network topologies were investigated. Both were m -node rings with an extra $.06m$ random links added, a new such set for each of the 50 runs giving a plotted average value. “Short links” (L) means that all extra links connected players two hops apart, while “small-worlds” (W) means there was no such restriction.

IC Econ’s inferior performance illustrates the shortcoming of economics-like algorithms. For $D = 1$ SA did not benefit from small worlds connections, and IC variants barely benefited (3%), despite the associated drop in average inter-node hop distance. However if D also increased, so that G directly reflected the change in the topology, then the gain with a small worlds topology grew to 10%. (See the discussion on path lengths in [14].)

The authors thank Michael New, Bill Macready, and Charlie Strauss for helpful comments.

-
- [1] T. Back, D. B. Fogel, and Z. Michalewicz, editors. *Handbook of Evolutionary Computation*. Oxford University Press, 1997.
- [2] E. Bonabeau, M. Dorigo, and G. Theraulaz. Inspiration for optimization from social insect behaviour. *Nature*, 406(6791):39–42, 2000.
- [3] G. Caldarelli, M. Marsili, and Y. C. Zhang. A prototype model of stock exchange. *Europhysics Letters*, 40:479–484, 1997.
- [4] D. Challet and N. F. Johnson. Optimal combinations of imperfect objects. *Physical Review Letters*, 89:028701, 2002.
- [5] K. Chellapilla and D.B. Fogel. Evolution, neural networks, games, and intelligence. *Proceedings of the IEEE*, pages 1471–1496, September 1999.
- [6] E. G. Coffman Jr., G. Galambos, S. Martello, and D. Vigo. Bin packing approximation algorithms: Combinatorial analysis. In *Handbook of Combinatorial Optimization*. Kluwer Academic Publishers, 1998.
- [7] R. H. Crites and A. G. Barto. Improving elevator performance using reinforcement learning. In D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo, editors, *Advances in Neural Information Processing Systems - 8*, pages 1017–1023. MIT Press, 1996.
- [8] D. Fudenberg and J. Tirole. *Game Theory*. MIT Press, Cambridge, MA, 1991.
- [9] S. Geman and D. Geman. Stochastic relaxation, Gibbs distribution and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6:721–741, 1984.
- [10] B. A. Huberman and T. Hogg. The behavior of computational ecologies. In *The Ecology of Computation*, pages 77–115. North-Holland, 1988.
- [11] S. Kirkpatrick, C. D. Jr Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science*, 220:671–680, May 1983.
- [12] M.J.B. Krieger, J.-B. Billeter, and L. Keller. Ant-like task allocation and recruitment in cooperative robots. *Nature*, 406:992–995, 2000.
- [13] R. V. Kulkarni, E. Almaas, and Stroud D. Exact results and scaling properties of small-world networks. *Physical Review E*, 61(4):4268–4271, 2000.
- [14] M. E. J. Newman, C. Moore, and D. J. Watts. Mean-field solution of the small-world network model. *Physical Review Letters*, 84(14):3201–3204, 2000.
- [15] N. Nisan and A. Ronen. Algorithmic mechanism design. *Games and Economic Behavior*, 35:166–196, 2001.
- [16] R. Savit, R. Manuca, and R. Riolo. Adaptive competition, market efficiency, phase transitions and spin-glasses. preprint cond-mat/9712006, December 1997.
- [17] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, 1998.
- [18] K. Tumer and D. H. Wolpert. Collective intelligence and Braess’ paradox. In *Proceedings of the Seventeenth National Conference on Artificial Intelligence*, pages 104–109, Austin, TX, 2000.
- [19] D. J. Watts and S. H. Strogatz. Collective dynamics of ‘small world’ networks. *Nature*, 393:440–442, 1998.
- [20] D. H. Wolpert. Theory of design of collectives. pre-print, 2003.
- [21] D. H. Wolpert and K. Tumer. Optimal payoff functions for members of collectives. *Advances in Complex Systems*, 4(2/3):265–279, 2001.
- [22] D. H. Wolpert, K. Wheeler, and K. Tumer. Collective intelligence for control of distributed dynamical systems. *Europhysics Letters*, 49(6), March 2000.
- [23] Y. C. Zhang. Modeling market mechanism with evolutionary games. *Europhysics Letters*, March/April 1998.

