

1 An empirical comparison of the structure of
2 configuration files of continuous integration and
3 build systems

4 Joseph Ling
jl653@kent.ac.uk



School of Computing
University of Kent
United Kingdom

Word Count: 6,100

5 January 28, 2020

7 This paper describes a simple heuristic approach to solving large-scale con-
8 straint satisfaction and scheduling problems. In this approach one starts
9 with an inconsistent assignment for a set of variables and searches through
10 the space of possible repairs. The search can be guided by a value-ordering
11 heuristic, the *min-conflicts heuristic*, that attempts to minimize the num-
12 ber of constraint violations after each step. The heuristic can be used with
13 a variety of different search strategies. We demonstrate empirically that on
14 the n -queens problem, a technique based on this approach performs orders of
15 magnitude better than traditional backtracking techniques. We also describe
16 a scheduling application where the approach has been used successfully. A
17 theoretical analysis is presented both to explain why this method works well
18 on certain types of problems and to predict when it is likely to be most
19 effective.

1 Introduction

need more introduction into CI like the travis paper probably?

For continuous integration (CI) and build systems/continuous delivery a lot of it is configured with code as configuration. The main kind of configuration format used for this is yaml (reference to what it is??) followed by xml and java based scripting formats. As CI has become increasingly more popular and practical ()......For software development teams there are roles for just looking after and managing CI and the deployment infrastructure. Understanding how code as configuration is used for CI is an emerging area of research because of this.

cite needed here to backup

The first stage in order to answer this question is to get a dataset of configuration files from different ci/cd systems. In doing so answering the questions of: - what is the spread of ci/cd systems for public github repositories? this will take into account operating system, programming language, star count, subscriber count - note: something along the lines of multiple configuration files - what naming convention do they use for the files? (in order to understand common practices)

After getting all that data then it will be need to analysed to how it is used. The key questions for this that will be focused on will be: - how are comments used in the configuration? - how are stages (named sections of the process) named for configuration files? - branch names - configuration errors when loading the config (just yaml parsing errors atm) - how scripts with the configuration files? (need to elloborate more on this one)

A key aspect is that these questions do not look too deeply into the individual implementation of each of unique piece of configuration for each type. This is because there are already some good papers looking Gallaba and McIntosh (2018) at this but in order to be able to compare the different configuration types it is important to compare similar attributes (there is also a time factor in here as well).

need to carefully plan out how to gage the size of repo, as I don't want to have too download repos

2 Previous Works

Config as code is not necessarily infrastructure as code is slightly differently. So the comparison isn't fair necessarily or accurate.

Configuration as code or infrastructure as code has been an increasing area of research over the last few years. There seems to be slightly more research in infrastructure as code Rahman, Mahdavi-Hezaveh and Williams (2019). There has been a focus on Puppet and Chef, for example in Sharma, Fragkoulis and Spinellis (2016) looks at code quality by the measure of "code smell" of Puppet code. This tackles the problem by defining by best practises and analysing the code against that. In the case of Cito et al. (2017) it uses the docker linter in order to be able to analyse the files.

15, 18, 23, 30

docker, puppet, travis analysis

3 Methodology

foo bar asdffsdf

3.1 first load

By almost any measure, the Hubble Space Telescope scheduling problem Between ten thousand and thirty thousand astronomical observations per year must be scheduled, subject to a great variety of constraints including power restrictions, observation priorities, time-dependent orbital characteristics, movement of astronomical bodies, stray light sources, etc. Because the telescope is an extremely valuable resource with a limited lifetime, efficient scheduling is a critical concern. An initial scheduling system, developed using traditional programming methods, highlighted the difficulty of the problem; it was estimated that it would take over three weeks for the system to schedule one week of observations. As described in section 5, this problem was remedied by the development of a successful constraint-based system to augment the initial system. At the heart of the constraint-based system is a neural network developed by Adorf and Johnston, the Guarded Discrete Stochastic (GDS) network,

From a computational point of view the network is interesting because

Adorf and Johnston found that it performs well on a variety of tasks, in addition to the space telescope scheduling problem. For example, the network performs significantly better on the n -queens problem than methods that were previously developed. The n -queens problem requires placing n queens on an $n \times n$ chessboard so that no two queens share a row, column or diagonal. The network has been used to solve problems of up to 1024 queens, whereas most heuristic backtracking methods encounter difficulties with problems one-tenth

In a standard Hopfield network, all connections between neurons are symmetric. In the GDS network, the main network is coupled asymmetrically to an auxiliary network of *guard neurons* which restricts the configurations that the network can assume. This modification enables the network to rapidly find a solution for many problems, even when the network is simulated on a serial machine. Unfortunately, convergence to a stable configuration is no longer guaranteed. Thus the network can fall into a local minimum involving a group of unstable states among which it will oscillate. In practice, however, if the network fails to converge after some number of neuron state transitions, it can simply be stopped and started over.

To illustrate the network architecture and updating scheme, let us consider how the network is used to solve binary constraint satisfaction problems. A problem consists of n variables, $X_1 \dots X_n$, with domains $D_1 \dots D_n$, and a set of binary constraints. Each constraint $C_\alpha(X_j, X_k)$ is a subset of $D_j \times D_k$ specifying incompatible values for a pair of variables. The goal is to find an assignment for each of the variables which satisfies the constraints. (In this paper we only consider the task of finding a single solution, rather than that of finding all solutions.) To solve a CSP using the network, each variable is represented by a separate set of neurons, one neuron for each of the variable's possible values. Each neuron is either "on" or "off", and in a solution state, every variable will have exactly one of its corresponding neurons "on", representing the value of that variable. Constraints are represented by inhibitory (i.e., negatively weighted) connections between the neurons. To

111 insure that every variable is assigned a value, there is a guard neuron for
112 each set of neurons representing a variable; if no neuron in the set is on, the
113 guard neuron will provide an excitatory input that is large enough to turn
114 one on. (Because of the way the connection weights are set up, it is unlikely
115 that the guard neuron will turn on more than one neuron.) The network is
116 updated on each cycle by randomly picking a set of neurons that represents
117 a variable, and flipping the state of the neuron in that set whose input is
118 *most inconsistent* with its current output (if any). When all neurons' states
119 are consistent with their input, a solution is achieved.

120 To solve the n -queens problem, for example, each of the $n \times n$ board posi-
121 tions is represented by a neuron whose output is either one or zero depending
122 on whether a queen is currently placed in that position or not. (Note that
123 this is a local representation rather than a distributed representation of the
124 board.) If two board positions are inconsistent, then an inhibiting connection
125 exists between the corresponding two neurons. For example, all the neurons
126 in a column will inhibit each other, representing the constraint that two
127 queens cannot be in the same column. For each row, there is a guard neuron
128 connected to each of the neurons in that row which gives the neurons in the
129 row a large excitatory input, enough so that at least one neuron in the row
130 will turn on. The guard neurons thus enforce the constraint that one queen
131 in each row must be on. As described above, the network is updated on each
132 cycle by randomly picking a row and flipping the state of the neuron in that
133 row whose input is most inconsistent with its current output. A solution is
134 realized when the output of every neuron is consistent with its input.

135 4 Why does the GDS Network Perform So 136 Well?

137 Our analysis of the GDS network was motivated by the following question:
138 “Why does the network perform so much better than traditional backtracking
139 methods on certain tasks”? In particular, we were intrigued by the results on

the n -queens problem, since this problem has received considerable attention from previous researchers. For n -queens, Adorf and Johnston found empirically that the network requires a linear number of transitions to converge. Since each transition requires linear time, the expected (empirical) time for the network to find a solution is $O(n^2)$. To check this behavior, Johnston and Adorf ran experiments with n as high as 1024, at which point memory limitations became a problem.¹

4.1 Nonsystematic Search Hypothesis

Initially, we hypothesized that the network's advantage came from the non-systematic nature of its search, as compared to the systematic organization inherent in depth-first backtracking. There are two potential problems associated with systematic depth-first search. First, the search space may be organized in such a way that poorer choices are explored first at each branch point. For instance, in the n -queens problem, depth-first search tends to find a solution more quickly when the first queen is placed in the center of the first row rather than in the corner; apparently this occurs because there are more solutions with the Queen in the center. Nevertheless, most naive algorithms tend to start in the corner simply because humans find it more natural to program that way. However, this fact by itself does not explain why nonsystematic search would work so well for n -queens. A backtracking program that randomly orders rows (and columns within rows) performs much better than the naive method, but still performs poorly relative to the GDS network.

The second potential problem with depth-first search is more significant and more subtle. As illustrated by figure 1, a depth-first search can be a disadvantage when solutions are not evenly distributed throughout the search space. In the tree at the left of the figure, the solutions are clustered

¹The network, which is programmed in Lisp, requires approximately 11 minutes to solve the 1024 queens problem on a TI Explorer II. For larger problems, memory becomes a limiting factor because the network requires approximately $O(n^2)$ space. (Although the number of connections is actually $O(n^3)$, some connections are computed dynamically rather than stored).

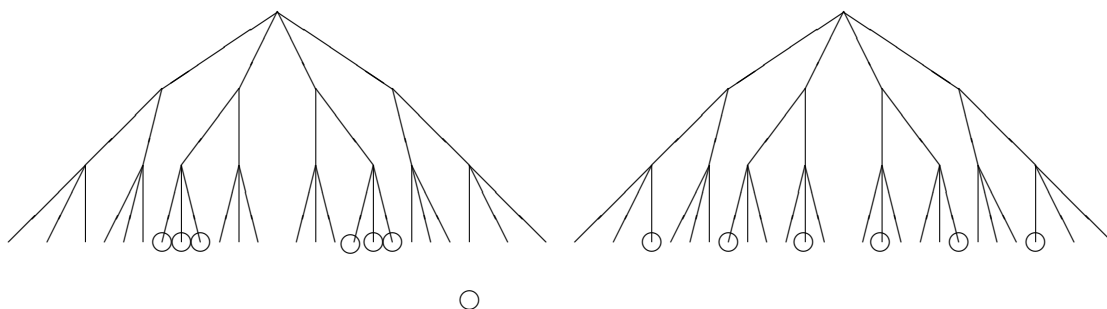


Figure 1: Solutions Clustered vs. Solutions Evenly Distributed

together. In the tree on the right, the solutions are more evenly distributed. Thus, the average distance between solutions is greater in the left tree. In a depth-first search, the average time to find the first solution increases with the average distance between solutions. Consequently depth-first search performs relatively poorly in a tree where In comparison, a search strategy which examines the leaves of the tree in random order is unaffected by solution clustering.

We investigated whether this phenomenon explained the relatively poor performance of depth-first search on n -queens by experimenting with a randomized algorithm begins by selecting a path from the root to a leaf. To select a path, the algorithm starts at the root node and chooses one of its children with equal probability. This process continues recursively until a leaf is encountered. If the leaf is a solution the algorithm terminates, if not, it starts over again at the root and selects a path. The same path may be examined more than once, since no memory is maintained between successive trials.

The Las Vegas algorithm does, in fact, perform better than simple depth-first search on n -queens However, the performance of the Las Vegas algorithm is still not nearly as good as that of the GDS network, and so we concluded that the systematicity hypothesis alone cannot explain the network's behavior.

187 4.2 Informedness Hypothesis

188 Our second hypothesis was that the network's search process uses informa-
189 tion about the current assignment that is not available to a constructive
190 backtracking program. 's use of an iterative improvement strategy guides
191 the search in a way that is not possible with a standard backtracking algo-
192 rithm. We now believe this hypothesis is correct, in that it explains why the
193 network works so well. In particular, the key to the network's performance
194 appears to be that state transitions are made so as to reduce the number of
195 outstanding inconsistencies in the network; specifically, each state transition
196 involves flipping the neuron whose output is most inconsistent with its cur-
197 rent input. From a constraint satisfaction perspective, it is as if the network
198 reassigns a value for a variable by choosing the value that violates the fewest
199 constraints. This idea is captured by the following heuristic:

200 **Min-Conflicts heuristic:**

201 *Given:* A set of variables, a set of binary constraints, and an assign-
202 ment specifying a value for each variable. Two variables *conflict* if
203 their values violate a constraint.

204 *Procedure:* Select a variable that is in conflict, and assign it a value
205 that minimizes the number of conflicts. (Break ties randomly.)

206 We have found that the network's behavior can be approximated by a
207 symbolic system that uses the min-conflicts heuristic for hill climbing. The
208 hill-climbing system starts with an initial assignment generated in a prepro-
209 cessing phase. At each choice point, the heuristic chooses a variable that is
210 currently in conflict and reassigns its value, until a solution is found. The
211 system thus searches the space of possible assignments, favoring assignments
212 with fewer total conflicts. Of course, the hill-climbing system can become
213 "stuck" in a local maximum, in the same way that the network may become
214 "stuck" in a local minimum. In the next section we present empirical evi-
215 dence to support our claim that the min-conflicts approach can account for
216 the network's effectiveness.

```

Procedure INFORMED-BACKTRACK (VARS-LEFT VARS-DONE)
  If all variables are consistent, then solution found, STOP.
  Let VAR = a variable in VARS-LEFT that is in conflict.
  Remove VAR from VARS-LEFT.
  Push VAR onto VARS-DONE.
  Let VALUES = list of possible values for VAR in ascending order according
                  to number of conflicts with variables in VARS-LEFT.
  For each VALUE in VALUES, until solution found:
    If VALUE does not conflict with any variable that is in VARS-DONE,
      then Assign VALUE to VAR.
      Call INFORMED-BACKTRACK(VARS-LEFT VARS-DONE)
    end if
  end for
end procedure

Begin program
  Let VARS-LEFT = list of all variables, each assigned an initial value.
  Let VARS-DONE = nil
  Call INFORMED-BACKTRACK(VARS-LEFT VARS-DONE)
End program

```

Figure 2: Informed Backtracking Using the Min-Conflicts Heuristic

217 There are two aspects of the min-conflicts hill-climbing method that dis-
 218 tinguish it from standard CSP algorithms. First, instead of incrementally
 219 constructing a consistent partial assignment, the min-conflicts method *re-*
 220 *pairs* a complete but inconsistent assignment by reducing inconsistencies.
 221 Thus, it uses information about the current assignment to guide its search
 222 that is not available to a standard backtracking algorithm. Second, the use
 223 of a hill-climbing strategy rather than a backtracking strategy produces a
 224 different style of search.

225 4.2.1 Repair-Based Search Strategies

226 (This is a example of a third level section.) Extracting the method from the
227 network enables us to tease apart and experiment with its different compo-
228 nents. In particular, the idea of repairing an inconsistent assignment can be
229 used with a variety of different search strategies in addition to hill climbing.
230 For example, we can backtrack through the space of possible repairs, rather
231 than using a hill-climbing strategy, as follows. Given an initial assignment
232 generated in a preprocessing phase, we can employ the min-conflicts heuristic
233 to order the choice of variables and values to consider, as described in figure
234 2. Initially, the variables are all on a list of VARS-LEFT, and as they are
235 repaired, they are pushed onto a list of VARS-DONE. The algorithm attempts
236 to find a sequence of repairs, such that no variable is repaired more than
237 once. If there is no way to repair a variable in VARS-LEFT without violat-
238 ing a previously repaired variable (a variable in VARS-DONE), the algorithm
239 backtracks.

240 Notice that this algorithm is simply a standard backtracking algorithm
241 augmented with the min-conflicts heuristic to order its choice of which vari-
242 able and value to attend to. This illustrates an important point. The back-
243 tracking repair algorithm incrementally extends a consistent partial assign-
244 ment (i.e., VARS-DONE), as does a constructive backtracking program, but
245 in addition, uses information from the initial assignment (i.e., VARS-LEFT)
246 to bias its search. Thus, it is a type of *informed backtracking*. We still char-
247 acterize it as repair-based method since its search is guided by a complete,
248 inconsistent assignment.

249 5 Experimental Results

250 [section ommitted]

251 6 A Theoretical Model

252 [section ommitted]

253 7 Discussion

254 [section ommitted]

255 8 Acknowledgement

256 The authors wish to thank Hans-Martin Adorf, Don Rosenthal, Richard
257 Franier, Peter Cheeseman and Monte Zweben for their assistance and ad-
258 vice. We also thank Ron Musick and our anonymous reviewers for their
259 comments. The Space Telescope Science Institute is operated by the Associ-
260 ation of Universities for Research in Astronomy for NASA.

261 Appendix A. Probability Distributions for N- 262 Queens

263 [section ommitted]

264 References

- 265 Cito, J., Schermann, G., Wittern, J. E., Leitner, P., Zumberi, S. and Gall,
266 H. C. (2017). An Empirical Analysis of the Docker Container Ecosystem
267 on GitHub. In *2017 IEEE/ACM 14th International Conference on Mining*
268 *Software Repositories (MSR)*, pp. 323–333, iSSN: null.
- 269 Gallaba, K. and McIntosh, S. (2018). Use and Misuse of Continuous Inte-
270 gration Features: An Empirical Study of Projects that (mis)use Travis CI.
271 *IEEE Transactions on Software Engineering*, pp. 1–1.

- 272 Rahman, A., Mahdavi-Hezaveh, R. and Williams, L. (2019). A systematic
273 mapping study of infrastructure as code research. *Information and Soft-*
274 *ware Technology*, 108, pp. 65–77.
- 275 Sharma, T., Fragkoulis, M. and Spinellis, D. (2016). Does Your Configuration
276 Code Smell? In *2016 IEEE/ACM 13th Working Conference on Mining*
277 *Software Repositories (MSR)*, pp. 189–200, iSSN: null.