# Transformation-based Feature Computation for Algorithm Portfolios

Barry Hurley[1], Serdar Kadioglu[2], Yuri Malitsky[1], and Barry O'Sullivan[1]

[1] Insight Centre for Data Analytics
Department of Computer Science, University College Cork, Ireland
{b.hurley|y.malitsky}@4c.ucc.ie, b.osullivan@cs.ucc.ie
[2] Oracle America Inc., Burlington, MA 01803, USA
serdark@cs.brown.edu

**Abstract.** Instance-specific algorithm configuration and algorithm portfolios have been shown to offer significant improvements over single algorithm approaches in a variety of application domains. In the SAT and CSP domains algorithm portfolios have consistently dominated the main competitions in these fields for the past five years. For a portfolio approach to be effective there are two crucial conditions that must be met. First, there needs to be a collection of complementary solvers with which to make a portfolio. Second, there must be a collection of problem features that can accurately identify structural differences between instances. This paper focuses on the latter issue: feature representation, because, unlike SAT, not every problem has well-studied features. We employ the well-known SATzilla feature set, but compute alternative sets on different SAT encodings of CSPs. We show that regardless of what encoding is used to convert the instances, adequate structural information is maintained to differentiate between problem instances, and that this can be exploited to make an effective portfolio-based CSP solver.

## 1 Introduction

Significant strides have recently been made in the application of portfolio-based algorithms in the fields of constraint satisfaction [15], quantified boolean formulae [17], and most notably in SAT [23,6]. Having a collection of solvers, these approaches compute a set of representative features about a problem instance and then use this information to decide what is the most effective solver to employ. These decisions can be made based on regression techniques [23], in which a classifier is trained to predict expected runtime of each solver and choosing the one with best predicted performance. Alternatively, a ranking algorithm can be trained to directly predict the best solver for each instance [5]. The features can also be used for clustering [7], where the best solver is chosen for each cluster of instances. In practice, regardless of the approach, portfolio algorithms have been shown to be dramatically better than using a single solver.

Algorithm portfolios also rely on a good set of features to describe the problem instance being solved. This can be seen as a major drawback since one needs

to use specific features for each problem at hand, or worse, has to come up with a set of features if none exists. If there are not enough informative features present, it is impossible to train a classifier to differentiate between classes of instances. On the other hand, if there are too many features it is possible to over fit the classifier to the training data. Furthermore, a large feature set is likely to have noisy features, which could be detrimental to the quality of the learned classifier. In the SAT domain, the features used by the solvers dominating the competitions have been thoroughly analyzed and studied over the last decade. Unfortunately, many other fields do not have such well established feature set. Even in the case of constraint satisfaction problems, where a feature set has been proposed, careful filtering can dramatically improve the quality of portfolios [9].

However, while there might not be an existing feature set, for NP-complete problems there exist polynomial-time transformations to any other NP-complete problem. In this paper we propose to take advantage of this by transforming CSP instances to SAT as a pre-processing step before computing its features. We show that such a transformation retains the necessary information needed to differentiate the classes of instances. In particular we show the effectiveness of this approach on constraint satisfaction problems. We choose the CSP domain for two reasons. First, it has a large number of solvers that can be used to make a diversified portfolio. Second, because a feature set exists for CSPs, we can compare the quality of a portfolio trained on SAT features to the domain specific CSP features.

There has been a lot of work exploring the effect of transforming CSP instances into SAT. Perhaps the most relevant work is by Ansótegui and Manyà which evaluated the performance of SAT solvers on six SAT-encodings on graph colouring, random binary CSPs, pigeon hole, and all interval series problems [1]. Solvers such as sugar [19], azucar [20], and CSP4SAT4J [11] have similarly tackled CSP problems by encoding them into SAT and then solving them with a predefined SAT solver. However, as far as we are aware, this paper represents the first time that a portfolio has been created using features gained *after* transforming a problem from one domain to another.

## 2  Encodings

There are a number of known polynomial-time transformations, or encodings, from constraint satisfaction problems to SAT [16]. In this paper we focus on three commonly used encodings: the direct, order and support encodings.

### 2.1  Direct Encoding

In the direct encoding [22] for each CSP variable $X$, with domain $\{1, \ldots, d\}$, a SAT variable is created for each domain value, i.e. $x_1, x_2, \ldots, x_d$. If $x_1$ is *true* in the resulting SAT formula, then the CSP variable $X$ is assigned the value 1 in the CSP solution. Therefore, in order to represent a solution to the CSP  exactly

one of $x_1, x_2, \ldots, x_d$ must be assigned *true*. We add an *at-least-one* clause and *at-most-one* clauses to the SAT formula for each CSP variable $X$:

$$(x_1 \vee x_2 \vee \ldots \vee x_d) \qquad \text{At Least One}$$
$$\forall v, w \in \mathrm{D}(X) : (\neg x_v \vee \neg x_w) \qquad \text{At Most One}$$

Constraints between CSP variables are represented in the direct encoding by enumerating the conflicting tuples. For a binary constraint between the pair of variables $X$ and $Y$, if the tuple $\langle X = v, Y = w \rangle$ is forbidden, then we add the conflict clause $(\neg x_v \vee \neg y_w)$.

## 2.2 Support Encoding

The support encoding [8,3] uses the same mechanism as the direct encoding to translate a CSP variable's domain into SAT. However, the support encoding differs on how the constraints between variables are encoded. Given a constraint between two variables $X$ and $Y$, for each value $v$ in the domain of $X$, let $S_{Y,X=v} \subset D(Y)$ be the subset of the values in the domain of $Y$ which are consistent with assigning $X = v$. Either $x_v$ is *false* or one of the consistent assignments from $y_1 \ldots y_d$ must be true, represented by the clause:

$$\neg x_v \vee \left( \bigvee_{i \in S_{Y,X=v}} y_i \right)$$

This must be repeated by adding clauses for each value in the domain of $Y$ and listing the values in $X$ which are consistent with each assignment.

## 2.3 Order Encoding

Unlike the direct and support encoding which model $X = v$ as a SAT variable, the order encoding creates SAT variables to represent $X \leq v$. If $X$ is less than or equal to $v$, then $X$ must also be less than or equal to $v + 1$. To enforce this across the domain we add the clauses:

$$\forall_v^{d-1} : (\neg x_{\leq v} \vee x_{\leq v+1})$$

The order encoding is naturally suited to modelling inequality constraints. To state $X \leq 3$, we would just post the unit clause $(x_{\leq 3})$. If we want to model the constraint $X = v$, we could rewrite it as $(X \leq v \wedge X \geq v)$. $X \geq v$ can then be rewritten as $\neg X \leq (v - 1)$. To state that $X = v$ under the order encoding, we would encode $(x_{\leq v} \wedge \neg x_{\leq v-1})$. A conflicting tuple between two variables, for example $\langle X = v, Y = w \rangle$ can be written in propositional logic and simplified to a CNF clause using De Morgan's Law:

$$\neg((x_{\leq v} \wedge x_{\geq v}) \wedge (y_{\leq w} \wedge y_{\geq w}))$$
$$\neg((x_{\leq v} \wedge \neg x_{\leq v-1}) \wedge (y_{\leq w} \wedge \neg y_{\leq w-1}))$$
$$\neg(x_{\leq v} \wedge \neg x_{\leq v-1}) \vee \neg(y_{\leq w} \wedge \neg y_{\leq w-1})$$
$$(\neg x_{\leq v} \vee x_{\leq v-1} \vee \neg y_{\leq w} \vee y_{\leq w-1})$$

## 3 Feature Computation

In addition to the pure direct, support and order encodings discussed in the previous section, we also consider variants of these encodings in which the clauses that encode the domains of the variables are not included. We omit the domains in order to test whether focusing only on the constraints present in a CSP is enough to differentiate the instances. We now briefly describe the features used for CSP and SAT.

**CSP Features.** We compute features for each of the original CSP instances, plus for each of the six encodings. We record 36 features directly from the CSP instance using `mistral` [4]. This includes static features such as statistics about the types of constraints used, average and maximum domain size; and dynamic statistics recorded by running `mistral` for 2 seconds: average and standard deviation of variable weights, number of nodes, number of propagations and a few others.

**SAT Features.** We use the 54 features computed using the newest feature computation tool from UBC [13]. These features include problem size features, graph-based features, balance features, proximity to horn formula features, DPLL probing features, and local search probing features.

## 4 Numerical Results

We implemented a tool to translate a CSP instance specified in XCSP format [18] into SAT (CNF). At present, it is capable of encoding inequality and binary extensional constraints using the direct, support and order encoding.

**Benchmarks.** For our evaluation, we consider CSP problem instances from the CSP solver competition.[1] Of these, we consider the instances that contain either inequality or binary extensional constraints. This presents a pool of 2,433 instances, containing Graph Colouring, Random, Quasi-random, Black Hole, Quasi-group Completion, Quasi-group With Holes, Langford, Towers of Hanoi and Pigeon Hole problems.

**Portfolio Approach.** To train our portfolios we used the ISAC methodology [7] which has been shown to work better than a regression based approaches [14]. ISAC uses the computed features to cluster the instances. Then for each cluster, the best solver in the portfolio is selected. When a new instance needs to be solved, its features are computed, it is assigned to the nearest cluster, and subsequently solved using the appropriate solver.

For our CSP solver portfolio we used: abscon [12], csp4j [11], sat4j [10], pcs [21], gecode [2], and sugar [19]. Each instance was run for 3,600 seconds. It is important to note that we include the time required for encoding the instances and computing the features as part of the computation time.

---

[1] CSP solver competition instances
http://www.cril.univ-artois.fr/~lecoutre/benchmarks.html

**Table 1.** Comparison of number of solved instances and PAR 10 score between the virtual best solver (VBS), the portfolio approach, the random cluster approach and the best single solver based on CSP and SAT features for clustering.

| PAR 10 | | | | | | |
|---|---|---|---|---|---|---|
| Approach | CSP | Direct | Direct ND | Order | Order ND | Support | Support ND |
| VBS | 1792 | 1887 | 1793 | 1806 | 1806 | 1810 | 1811 |
| Portfolio | 2066 | 3312 | 3221 | 2689 | 2077 | 2084 | **2022** |
| Random Cluster | 3806 | 3705 | 3424 | 3725 | 3797 | 3867 | 3902 |
| Best Single | 4776 | 4870 | 4777 | 4789 | 4789 | 4792 | 4792 |
| Number Solved | | | | | | |
| Approach | CSP | Direct | Direct ND | Order | Order ND | Support | Support ND |
| VBS | 2315 | 2310 | 2315 | 2315 | 2315 | 2315 | 2315 |
| Portfolio | 2297 | 2215 | 2220 | 2256 | 2297 | 2297 | **2301** |
| Random Cluster | 2180 | 2188 | 2206 | 2187 | 2182 | 2177 | 2175 |
| Best Single | 2115 | 2110 | 2115 | 2115 | 2115 | 2115 | 2115 |

We perform our experiments using stratified 10-fold cross validation. In Table 1, we present the performance for both the number of solved instances and the penalized runtime average PAR 10 which counts each time-out as taking 10 times the time-out to complete for each problem representation. The SAT encodings without the variable domains are marked with ND. We compare the portfolio performance to the best single solver as well as to the oracle Virtual Best Solver (VBS) which for every instance always selects the fastest solver. As we can see, using a portfolio approach for CSP instances is always preferable to just choosing to run a single solver. We also compare to a random clustering approach, which randomly groups the instances of the test set into the same number of clusters as the portfolio method and finds the best solver to run on each group. Note that the random clustering is trained on the same data it is evaluated on, and further that in practice one would not know to which cluster to assign a new instance. The random clustering approach is included to show that the clusters found by ISAC are indeed capturing important information about the instances. We observe this because in all cases Portfolio is better than the Random Clustering approach.

Table 1 also shows that regardless of the encoding we use, we can always close at least 50% of the performance gap between the best single solver and the virtual best one. Furthermore, we see that if we use particularly accurate encoding, which in our case is the support encoding without domain clauses, we can even achieve slightly better performance than using features that have been specifically designed for the problem domain.

## 5 Conclusion

In this paper we show that it is possible to encode an instance from one problem domain to another as a preprocessing step for feature computation. In particular, we show that even with the overhead of converting CSP instances to SAT, a CSP portfolio trained on well established SAT features can perform just as well as if it was trained on CSP specific features. These findings show that encoding techniques can retain enough information about the original instance to accurately differentiate different classes of instances. Our results serves as a proof of concept for an automated feature generation approach for NP-complete problems that do not have a well studied feature vector. We consider this as a step toward problem independent feature computation for algorithm portfolios, and we plan to analyze it further and extend its applications in the future.

## Acknowledgements

## References

1. Ansótegui, C., Manyà, F.: Mapping Problems with Finite-Domain Variables into Problems with Boolean Variables. In: The 7th International Conference on Theory and Applications of Satisfiability Testing, SAT 2004 (2004)
2. Gecode Team: Gecode: Generic Constraint Development Environment (2006), http://www.gecode.org
3. Gent, I.P.: Arc Consistency in SAT. In: Proceedings of the 15th European Conference on Artificial Intelligence, ECAI'2002. pp. 121–125 (2002)
4. Hebrard, E.: Mistral,a Constraint Satisfaction Library. In: Proceedings of the Third International CSP Solver Competition (2009)
5. Hurley, B., O'Sullivan, B.: Adaptation in a CBR-Based Solver Portfolio for the Satisfiability Problem. In: Case-Based Reasoning Research and Development - 20th International Conference, ICCBR 2012. pp. 152–166 (2012)
6. Kadioglu, S., Malitsky, Y., Sabharwal, A., Samulowitz, H., Sellmann, M.: Algorithm Selection and Scheduling. In: Proceedings of the 17th International Conference on Principles and Practice of Constraint Programming. pp. 454–469. CP'11, Springer-Verlag, Berlin, Heidelberg (2011)
7. Kadioglu, S., Malitsky, Y., Sellmann, M., Tierney, K.: ISAC - Instance-Specific Algorithm Configuration. In: Coelho, H., Studer, R., Wooldridge, M. (eds.) ECAI. Frontiers in Artificial Intelligence and Applications, vol. 215, pp. 751–756. IOS Press (2010)
8. Kasif, S.: On the Parallel Complexity of Discrete Relaxation in Constraint Satisfaction Networks. Artificial Intelligence 45(3), 275–286 (Oct 1990)

9. Kroer, C., Malitsky, Y.: Feature filtering for instance-specific algorithm configuration. In: IEEE 23rd International Conference on Tools with Artificial Intelligence, ICTAI 2011. pp. 849–855 (2011)
10. Le Berre, D., Parrain, A.: The sat4j library, release 2.2 system description. Journal on Satisfiability, Boolean Modeling and Computation 7, 59–64 (2010)
11. Le Berre, D., Lynce, I.: CSP2SAT4J: A Simple CSP to SAT Translator. In: Proceedings of the 2nd International CSP Solver Competition (2008)
12. Lecoutre, C., Tabary, S.: Abscon 112, Toward more Robustness. In: Proceedings of the Third International CSP Solver Competition (2009)
13. Lin Xu, Frank Hutter, H.H., Leyton-Brown, K.: Features for SAT (2012), http://www.cs.ubc.ca/labs/beta/Projects/SATzilla/Report_SAT_features.pdf
14. Malitsky, Y., Sabharwal, A., Samulowitz, H., Sellmann, M.: Non-model-based algorithm portfolios for sat. In: Proceedings of the 14th international conference on Theory and application of satisfiability testing. pp. 369–370. SAT'11, Springer-Verlag, Berlin, Heidelberg (2011), http://dl.acm.org/citation.cfm?id=2023474.2023517
15. O'Mahony, E., Hebrard, E., Holland, A., Nugent, C., O'Sullivan, B.: Using Case-based Reasoning in an Algorithm Portfolio for Constraint Solving. Proceeding of the 19th Irish Conference on Artificial Intelligence and Cognitive Science (2008)
16. Prestwich, S.D.: CNF Encodings. In: Handbook of Satisfiability, pp. 75–97. IOS Press (2009)
17. Pulina, L., Tacchella, A.: A multi-engine solver for quantified boolean formulas. In: Proceedings of the 13th international conference on Principles and practice of constraint programming. pp. 574–589. CP'07, Springer-Verlag, Berlin, Heidelberg (2007)
18. Roussel, O., Lecoutre, C.: XML Representation of Constraint Networks: Format XCSP 2.1. CoRR abs/0902.2362 (2009)
19. Tamura, N., Tanjo, T., Banbara, M.: System Description of a SAT-based CSP Solver Sugar. In: Proceedings of the 3rd International CSP Solver Competition. pp. 71–75 (2009)
20. Tanjo, T., Tamura, N., Banbara, M.: Azucar: A SAT-Based CSP Solver Using Compact Order Encoding — (Tool Presentation). In: Proceedings of the 15th International Conference on Theory and Applications of Satisfiability Testing (SAT 2012), LNCS 7317. pp. 456–462. Springer (2012)
21. Veksler, M., Strichman, O.: A Proof-Producing CSP Solver. In: Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2010 (2010)
22. Walsh, T.: SAT v CSP. In: Principles and Practice of Constraint Programming — CP 2000, LNCS 1894. vol. 1894, pp. 441–456. Springer-Verlag (2000)
23. Xu, L., Hutter, F., Hoos, H.H., Leyton-Brown, K.: SATzilla: Portfolio-based Algorithm Selection for SAT. Journal of Artificial Intelligence Research pp. 565–606 (June 2008)