

Finding Robust Solutions from Requirements Models

Gregory Gay¹, Tim Menzies¹, Omid Jalali¹, Gregory Mundy²,
Beau Gilkerson¹, Martin Feather³, and James Kiper⁴

¹ West Virginia University, Morgantown, WV, USA

² Alderson-Broaddus College, Philippi, WV

³ Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA, USA

⁴ Dept. of Computer Science and Systems Analysis, Miami University, Oxford, OH, USA

greg@greggay.com, tim@menzies.us, jalali.omid@gmail.com,
beau.gilkerson@gmail.com, mundyge@ab.edu,
martin.s.feather@jpl.nasa.gov, kiperjd@muohio.edu

Abstract. Solutions to non-linear requirements engineering problems may be “brittle”; i.e. small changes may dramatically alter solution effectiveness. Therefore, it is not enough to merely generate solutions to requirements problems—we must also assess the robustness of that solution. This paper introduces the KEYS2 algorithm that can generate *decision ordering diagrams*. Once generated, these diagrams can assess solution robustness in linear time. In experiments with real-world requirements engineering models, we show that KEYS2 can generate decision ordering diagrams in $O(N^2)$. When assessed in terms of (a) reducing inference times, (b) increasing solution quality, and (c) decreasing the variance of the generated solution, KEYS2 out-performs other search algorithms (simulated annealing, ASTAR, MaxWalkSat).

1 Introduction

Consider a “requirements model” where stakeholders write:

- Their various goals;
- Their different methods to reach those goals;
- Their view of the possible risks that compromise those goals;
- What mitigations they believe might reduce those risks.

We say that a “solution” to such a model is the the *least* cost combination of mitigations that reduce the *most* risks, thereby enabling the *most* number of requirements. In theory, software tools can be find solutions that best exploits and most satisfies the various goals of our different stakeholders. Such software might find good solutions that were missed by stakeholders in early lifecycle design discussions. Finding solutions to these requirements models is a non-linear optimization problem (*minimize* the sum of the mitigation costs while *maximizing* the number of achieved requirements).

There are many heuristic methods that can generate solutions to non-linear problems (see *Related Work*, below). However, such heuristic methods can be *brittle*; i.e. small changes may dramatically alter the effectiveness of the solution. Therefore, it important to understand the *neighborhood around the solution*. A naive approach to understanding the neighborhood might be to run a system N times then report:

- the solutions appearing more than in (say) $\frac{N}{2}$ cases;
- results with a $\pm 95\%$ confidence interval.

Note that both these approaches requires multiple runs of an analysis method. Multiple runs are undesirable since, in our experience [18], stakeholders often ask questions across a range of “scenarios”; i.e. hard-wired constraints that cannot be changed in that scenario. For example, three scenarios might be “what can be achieved assuming a maximum budget of one, two, or five million dollars?”. Scenario analysis can be time consuming. Reflecting over (say) $d = 10$ possible decisions a statistically significant number of times (e.g. $N = 20$) requires up to $20 * 2^{10} > 20,000$ repeats of the analysis.

Therefore, this paper proposes a rapid method for exploring decision neighborhoods. *Decision ordering diagrams* are a visual representation of the effects of changing a solution. We show below that:

- Using these diagrams, the region around a solution can be explored in linear time.
- A greedy Bayesian-based method called KEYS2 can generate the decision ordering diagrams in $O(N^2)$ time.
- KEYS2 yields solutions of higher quality than several other methods (simulated annealing, MaxWalkSat, ASTAR).
- Also, the variance of the solutions found by KEYS2 is less (and hence, better) than those found by the other methods.

This paper is structured as follows:

- After some background notes on the solution robustness, we describe the *DDP* requirements engineering tool used at NASA’s Jet Propulsion Laboratory (the case studies for this paper come from real-world early lifecycle DDP models);
- DDP inputs and outputs are then reviewed;
- Next, decision ordering diagrams are introduced;
- We then define and compare five different algorithms for generating solutions from DDP models: KEYS, KEYS2, Simulated Annealing, MaxWalkSat, ASTAR. Experimental results will show that KEYS2 out-performs the other methods (measured in terms of quickly generating high quality solutions that allow us to reflect over solution robustness).
- Finally, we offer some notes on related work and conclusions.

This paper extends a prior publication [37] in several ways:

- That paper did not address concerns of solution robustness.
- That paper did not explore a range of alternate algorithms.
- This paper introduces KEYS2, which is an improved version of KEYS.
- This paper offers extensive notes on related work.

2 Background

According to Harman [28], understanding the neighborhood of solutions is an open and pressing issue in search-based software engineering (SBSE). He argues that many

software engineering problems are over-constrained and no precise solution over all variables is achievable; therefore partial solutions based on heuristic search methods are preferred. *Solution robustness* is a major problem for such partial heuristic searches. The results of such partial heuristic search may be “brittle”; i.e. small changes to the search results may dramatically alter the effectiveness of the solution [29].

When offering partial solutions, it is very important to also offer insight into the space of options around the proposed solution. Such neighborhood information is very useful for managers with only partial control over their projects since it can give them confidence that even if only some of their recommendations are effected, then at least the range of outcomes is well understood. Harman [28] comments that understanding the neighborhood of our solutions is an open and pressing issue in search-based software engineering (SBSE):

“In some software engineering applications, solution robustness may be as important as solution functionality. For example, it may be better to locate an area of the search space that is rich in fit solutions, rather than identifying an even better solution that is surrounded by a set of far less fit solutions.”

“Hitherto, research on SBSE has tended to focus on the production of the fittest possible results. However, many application areas require solutions in a search space that may be subject to change. This makes robustness a natural second order property to which the research community could and should turn its attention [28].”

This paper reports a set of experiments on AI search for robust solutions. Our experiments address two important concerns. Firstly, *is demonstrating solution robustness be a time consuming task?* Secondly, is it necessary, as Harman suggests that *solution quality must be traded off against solution robustness?* That is, in our search for the conclusions that were stable within their local neighborhood, would we have to reject better solutions because they are less stable across their local neighborhood?

At least for the NASA models described in the next section, both these concerns are unfounded. KEYS2 terminates in hundredths of a second (where as our prior implementations took minutes to terminate [20]). Also, the solutions found by KEYS2 were not only of highest quality, they were also exhibited the lowest variance. Further, KEYS2 generates the decision ordering diagrams that can assess solution robustness in linear time.

3 Requirements Modeling Using “DDP”

This section introduces the DDP requirements modeling tool [13, 18]. used to interactively document the “Team-X” early lifecycle meetings at NASA’s Jet Propulsion Laboratory (JPL). These meetings are the source of the real-world requirements models used in this paper.

At “Team X” meetings, a large and diverse group of up to 30 experts from various fields (propulsion, engineering, communication, navigation, science, etc) meet for very short periods of time (hours to days) to produce a “mission concept” document. This

-
1. Requirement goals:
 - Spacecraft ground-based testing & flight problem monitoring
 - Spacecraft experiments with on-board Intelligent Systems Health Management (ISHM)
 2. Risks:
 - Obstacles to spacecraft ground-based testing & flight problem monitoring
 - Customer has no, or insufficient, money available for our use
 - Difficulty of building the models / design tools
 - ISHM Experiment is a failure (without necessarily causing flight failure)
 - Usability, User/Recipient-system interfaces undefined
 - V&V (certification path) untried and scope unknown
 - Obstacles to Spacecraft experiments with on-board ISHM
 - Bug tracking / fixes / configuration management issues, Managing revisions and upgrades (multi-center tech. development issue)
 - Concern about our technology interfering with in-flight mission
 3. Mitigations:
 - Mission-specific actions
 - Spacecraft ground-based testing & flight problem monitoring
 - Become a team member on the operations team
 - Use Bugzilla and CVS
 - Spacecraft experiments with on-board ISHM
 - Become a team member on the operations team
 - Utilize xyz's experience and guidance with certification of his technology
-

Fig. 1. Sample DDP requirements, risks, mitigations.

document may commit the project to, say, solar power rather nuclear power or a particular style of guidance software. All the subsequent work is based on the initial decisions documented in the mission concept.

DDP allows for the representation of the goals, risks, and risk-removing mitigations that belong to a specific project. During a Team X meeting, users of DDP explore combinations of mitigations that cost the least and support the most number of requirements. DDP propagate influences over matrices. For example, here is a trivial DDP model where *mitigation1* costs \$10,000 to apply and each requirement is of equal value (100). Note that the mitigation can remove 90% of the risk. Also, unless mitigated, the risk will disable 10% to 99% of requirements one and two (respectively):

$$\overbrace{\text{mitigation1}}^{\$10,000} \xrightarrow{0.9} \text{risk1} \rightarrow \begin{cases} \overbrace{\rightarrow}^{0.1} (\text{requirement1} = 100) \\ \overbrace{\rightarrow}^{0.99} (\text{requirement2} = 100) \end{cases} \quad (1)$$

The other numbers show the impact of mitigations on risks, and the impact of risks on requirements. The core of DDP is two matrices: one for *mitigations*risks* and another for *risks*requirements*.

DDP is used as follows. A dozen experts, or more, gather together for short, intensive knowledge acquisition sessions (typically, 3 to 4 half-day sessions). These sessions

must be short since it is hard to gather together these experts for more than a very short period of time. The DDP tool supports a graphical interface for the rapid entry of the assertions. Such rapid entry is essential, time is crucial and no tool should slow the debate. Therefore, DDP uses a lightweight representations for its model. Such representations are essential for early lifecycle decision making since only high-level assertions can be collected in short knowledge acquisition sessions (if the assertions get more elaborate, then experts may be unable to understand technical arguments from outside their own field of expertise). Therefore, DDP uses the following lightweight ontology:

- *Requirements* (free text) describing the objectives and constraints of the mission and its development process;
- *Weights* (numbers) of each requirements, reflecting their relative importance;
- *Risks* (free text) are events that damage requirements;
- *Mitigations*: (free text) are actions that reduce risks;
- *Costs*: (numbers) effort associated with mitigations, and repair costs for correcting Risks detected by Mitigations;
- *Mappings*: directed edges between requirements, mitigations, and risks that capture quantitative relationships among them.
- *Part-of relations* structure the collections of requirements, risks and mitigations;

Note that DDP models are the same as the “requirements models” we defined in the introduction. For examples of risks, requirements, and mitigations, see Figure 1. For an example of the network of connections between risks and requirements and mitigations, see Figure 2.

Sometimes, we are asked if the analysis of DDP requirements models is a real problem. The usual question is something like: “With these ultra-lightweight languages, aren’t all open issues just obvious?”. Such a question is usually informed by the small model fragments that appear in the ultra-lightweight modeling literature. Those sample model fragments are typically selected according to their ability to fit on a page or to succinctly illustrate some point of the authors. Real world ultra-lightweight models can be much more complex, paradoxically perhaps due to their simplicity: if a model is easy to write then it is easy to write a lot of it. Figure 2, for example, was generated in under a week by four people discussing one project. It is complex and densely-connected (a close inspection of the left and right hand sides of Figure 2 reveals the requirements and fault trees that inter-connect concepts in this model) and it is, by no means, the biggest or most complex DDP model that has ever been built.

We base our experimentation around DDP for three reasons. Firstly, one potential drawback with ultra-lightweight models is that they are excessively lightweight and contain no useful information. DDP’s models are demonstrably useful (that is, we are optimizing a real-world problem of some value). Clear project improvements have been seen from DDP sessions at JPL. Cost savings in at least two sessions have exceeded \$1 million, while savings of over \$100,000 have resulted in others [18]. Cost savings are not the only benefits of these DDP sessions. Numerous design improvements such as savings of power or mass have come out of DDP sessions. Likewise, a shifting of risks has been seen from uncertain architectural ones to more predictable and manageable ones. At some of these meetings, non-obvious significant risks have been identified and subsequently mitigated.

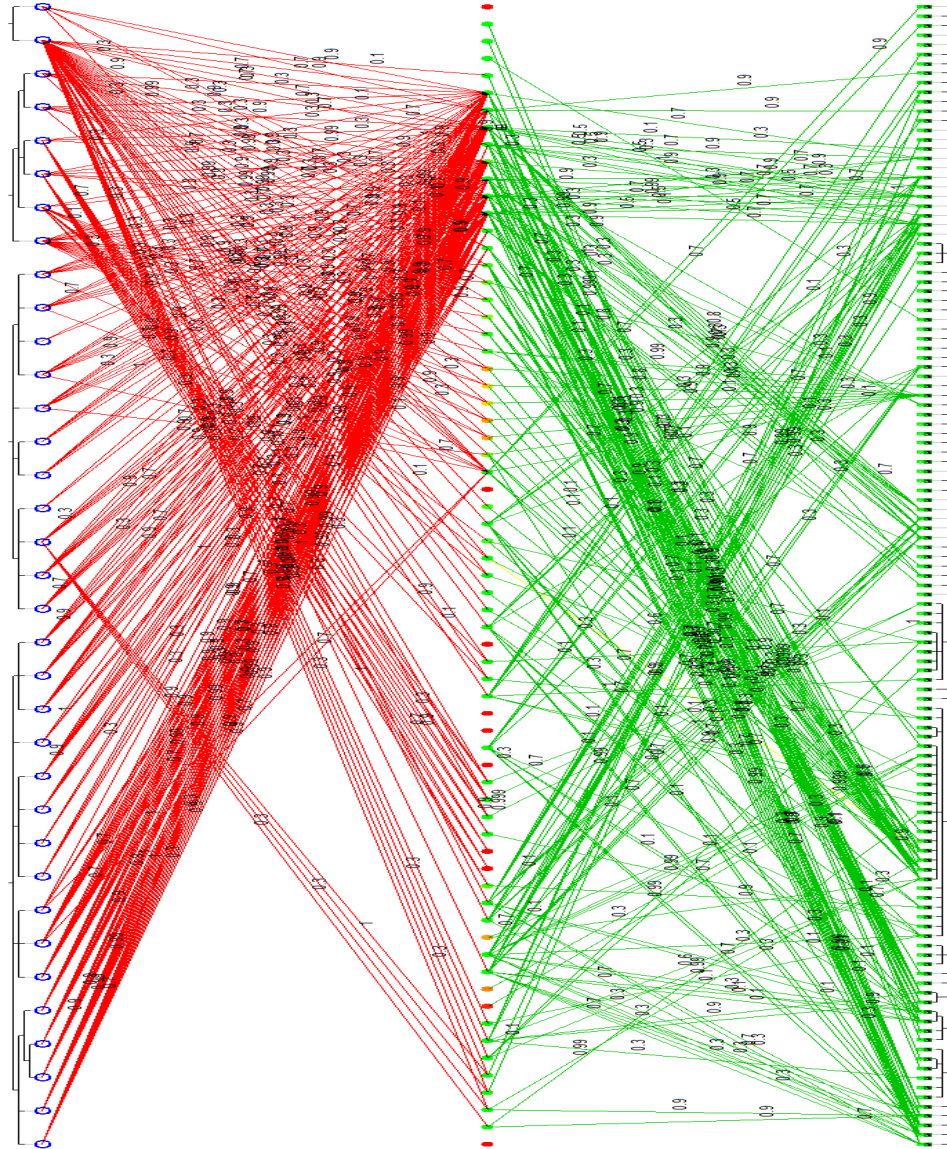


Fig. 2. An example of a model formed by the DDP tool. Red lines connect risks (middle) to requirements (left). Green lines connect mitigations (right) to the risks.

Our second reason to use DDP is that we can access numerous large real-world requirements models written in this format, both now and in the future. The DDP tool can be used to document not just final decisions but also to review the rationale that lead to those decisions. Hence, DDP remains in use at JPL: not only for its original purpose (group decision support), but also as a design rationale tool to document decisions. Recent DDP sessions included:

- An identification of the challenges of intelligent systems health management (ISHM) technology maturation (to determine the most cost-effective approach to achieving maturation) [21];
- A study on the selection and planning of deployment of prototype software [19].

Our third, and most important reason to use DDP in our research is that the tool is representative of other requirements modeling tools in widespread use. At its core, DDP is a set of influences expressed in a hierarchy, augmented with the occasional equation. Edges in the hierarchy have weights that strengthen or weaken influences that flow along those edges. At this level of abstraction, DDP is just another form of QOC [60] or a quantitative variant of Mylopoulos’ s qualitative soft goal graphs [51].

4 Model Inputs and Outputs

Before describing experimental comparisons of different methods for generating decision ordering diagrams, we first offer more details on DDP models.

4.1 Pre-Processing

To enable fast runtimes, a simple compiler exports the DDP models into a form accessible by our algorithms. This compiler stores a flattened form of the DDP requirements tree. In our compiled form, all computations are performed once and added as a constant to each reference of the requirement. For example, the compiler converts the trivial model of Equation 1 into `setupModel` and `model` functions similar to those in Figure 3. The `setupModel` function is called only once and sets several constant values. The `model` function is called whenever new cost and attainment values are needed. The topology of the mitigation network is represented as terms in equations within these functions. As our models grow more complex, so do these equations. For example, our biggest model, which contains 99 mitigations, generates 1427 lines of code. Figure 4 compares this biggest model to four other real-world DDP models.

While not a linear process, this compiler is not the bottleneck in our ability to find robust solutions to requirements models. Currently it takes about two seconds to compile a model with 50 requirements, 31 risks, and 58 mitigations. This compilation only has to happen once, after which we will run our $2^{|d|}$ what-if scenarios. Also, the current compiler is written in an unoptimized Visual Basic implementation which we can certainly significantly speed up. Finally, experts usually change a small portion of the model then run $2^{|d|}$ what-if scenarios to understand the impact of that change. An incremental compiler (that only updates changed portions) would run much faster than a full compilation of an entire DDP model.

```

#include "model.h"

#define M_COUNT 2
#define O_COUNT 3
#define R_COUNT 2

struct ddpStruct
{
    float oWeight [O_COUNT+1];
    float oAttainment [O_COUNT+1];
    float oAtRiskProp [O_COUNT+1];
    float rAPL [R_COUNT+1];
    float rLikelihood [R_COUNT+1];
    float mCost [M_COUNT+1];
    float roImpact [R_COUNT+1] [O_COUNT+1];
    float mrEffect [M_COUNT+1] [R_COUNT+1];
};

ddpStruct *ddpData;

void setupModel(void)
{
    ddpData = (ddpStruct *) malloc(sizeof(ddpStruct));
    ddpData->mCost [1]=11;
    ddpData->mCost [2]=22;
    ddpData->rAPL [1]=1;
    ddpData->rAPL [2]=1;
    ddpData->oWeight [1]=1;
    ddpData->oWeight [2]=2;
    ddpData->oWeight [3]=3;
    ddpData->roImpact [1] [1] = 0.1;
    ddpData->roImpact [1] [2] = 0.3;
    ddpData->roImpact [2] [1] = 0.2;
    ddpData->mrEffect [1] [1] = 0.9;
    ddpData->mrEffect [1] [2] = 0.3;
    ddpData->mrEffect [2] [1] = 0.4;
}

void model(float *cost, float *att, float m[])
{
    float costTotal, attTotal;
    ddpData->rLikelihood [1] = ddpData->rAPL [1] * (1 - m [1] * ddpData->mrEffect [1] [1])
        * (1 - m [2] * ddpData->mrEffect [2] [1]);
    ddpData->rLikelihood [2] = ddpData->rAPL [2] * (1 - m [1] * ddpData->mrEffect [1] [2]);
    ddpData->oAtRiskProp [1] = (ddpData->rLikelihood [1] * ddpData->roImpact [1] [1])
        + (ddpData->rLikelihood [2] * ddpData->roImpact [2] [1]);
    ddpData->oAtRiskProp [2] = (ddpData->rLikelihood [1] * ddpData->roImpact [1] [2]);
    ddpData->oAtRiskProp [3] = 0;
    ddpData->oAttainment [1] = ddpData->oWeight [1] * (1 - minValue (1, ddpData->oAtRiskProp [1]));
    ddpData->oAttainment [2] = ddpData->oWeight [2] * (1 - minValue (1, ddpData->oAtRiskProp [2]));
    ddpData->oAttainment [3] = ddpData->oWeight [3] * (1 - minValue (1, ddpData->oAtRiskProp [3]));
    attTotal = ddpData->oAttainment [1] + ddpData->oAttainment [2] + ddpData->oAttainment [3];
    costTotal = m [1] * ddpData->mCost [1] + m [2] * ddpData->mCost [2];

    *cost = costTotal;
    *att = attTotal;
}

```

Fig. 3. A trivial DDP model after knowledge compilation

Model	LOC	Objectives	Risks	Mitigations
model1.c	55	3	2	2
model2.c	272	1	30	31
model3.c	72	3	2	3
model4.c	1241	50	31	58
model5.c	1427	32	70	99

Fig. 4. Details of Five DDP Models.

4.2 Objective Function

When the `model` function is called, a pairing of the total cost of the selected mitigations and the number of reachable requirements (attainment) is returned. All of our algorithms then use that information to obtain a “score” for the current set of mitigations. The two numbers are normalized to a single score that represents the distance to a *sweet spot* of maximum requirement attainment and minimum cost:

$$score = \sqrt{cost^2 + (attainment - 1)^2} \quad (2)$$

Here, \bar{x} is a normalized value $0 \leq \frac{x - \min(x)}{\max(x) - \min(x)} \leq 1$. Hence, our scores ranges $0 \leq score \leq 1$ and *higher* scores are *better*.

4.3 Decision Ordering Diagrams

The objective function described above summarizes *one* call to a DDP model. This section describes *decision ordering diagrams*, which are a tool for summarization the results of thousands of calls to DDP models.

Consider some recommendation for changes to a project that requires decisions d of size $|d|$. In the general case, d is a subset of the space of all solutions D ($d \subseteq D$). When checking for solution robustness, or reflecting over modifications to d , a stakeholder may need to consider up to $d' \subseteq N^{|d|}$ possibilities (and $N = 2$ for binary decisions of the form “should I or should I not do this”). This can be a slow process, especially if evaluating each decision requires calling a complex and slow simulator.

Decision ordering diagrams are a linear time method for studying the robustness and neighborhood of a set of decisions. The diagrams assume that some method could offer a *linear ordering* of the decisions $x \in d$ ranked from *most-important* to *least-important*. They also assume that some method offers information on the effects of applying the top-ranked $1 \leq x \leq |d|$ decisions (e.g. the median and variance seen in the model’s objective function after applying solution $\{d_1..d_x\}$). For example, the *decision ordering diagram* of Figure 5 shows such a linear ordering (this figure presents *benefit* and *cost* results). In that figure:

- The x-axis denotes the number of decisions made.
- The y-axis shows performance statistics of an objective function seen after imposing the conjunction of decisions $1 \leq i \leq x$.

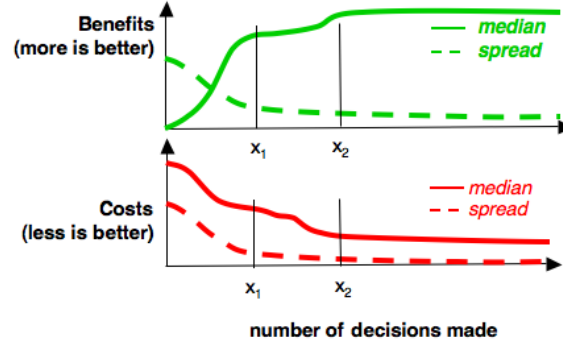


Fig. 5. A Decision Ordering Diagram. The median and spread plots show 50%-the percentile and the (75%-50)%-th percentile range (respectively) values generated from some objective function.

For performance, we run some objective function and report the median (50%-the percentile) and *spread* (the 75%-50%-the percentile range). We use median and spread to avoid any parametric assumptions.

These diagrams can comment on the robustness and neighborhood of solution $\{d_1..d_x\}$ as follows:

- By considering the variance of the performance statistics after applying $\{d_1..d_x\}$.
- By comparing the results of using the first x decisions to that of using the first $x - 1$ or $x + 1$ actions.

The neighborhood of a solution that uses decisions $\{d_1..d_x\}$ are solutions that use the decisions $\{d_1..d_{x \pm j}\}$. Since j is bounded $0 \leq |d| - 1$, this means that *reflecting over solution neighborhoods takes time linear on the number of decisions d* .

Decision ordering diagrams are a natural representation for “trade studies”. Trade studies are the activity of a multidisciplinary team to identify the most balanced technical solutions among a set of proposed viable solutions [?]. For example, minimum costs and maximum benefits are achieved at point x_2 of Figure 5. However, after applying only half the decisions (see x_1) most of the benefits could be achieved, albeit at a somewhat higher cost.

Decision ordering diagrams are *useful* under at least three conditions:

- The objective function scores are *well-behaved*; i.e. move smoothly to a plateau.
- The decisions *tame* the variance; i.e. the spread falls to value much lower than then median (otherwise, it is hard to show that decisions have any effect on the system performance).
- The are generated in a *timely* manner. Fast runtimes are required in order to keep up with fast moving discussion.

According to these definitions, Figure 5 is a *useful* decision ordering diagram if it can be generated in a *timely* manner.

It is an open issue if real worlds requirements models generate useful decision ordering diagrams. The following experiments test if, in practice, decision ordering diagrams generated from real world requirements models are *timely* to generate while being *well-behaved* and *tame*.

5 Searching for Solutions

Our experiments compare the results of numerous algorithms. We selected the comparison algorithms with much care. Numerous researchers stress the difficulties associated with comparing radically different algorithms. For example, Uribe and Stickel [63] tried to make some general statement about the value of Davis-Putnam proof (DP) procedures or binary-decision diagrams (BDD) for constraint satisfaction problems. In the end, they doubted that it was fair to compare algorithms that perform constraint satisfaction and no search (like BDDs) and methods that perform search and no constraint satisfaction (like DP). For this reason, researchers in model checking researchers like Holzmann (pers. communication) eschew companions of tools like SPIN [33], which are search-based, with tools like NuSMV [9], which are BDD-based. Hence we take care to compare algorithms similar to KEYS.

In terms of the Gu et al. survey [25], our selected algorithms (simulated annealing, ASTAR and MaxWalkSat) share four properties with KEYS and KEYS2. They are each discrete, sequential, *unconstrained* algorithms (constrained algorithms work towards a pre-determine the number of possible solutions while unconstrained methods are allowed to adjust the goal space.)

For full details on simulated annealing and ASTAR and MaxFunWalk, see below. For now, we observe that these algorithms share the property that at each step of their processing, they comment on all model inputs. KEYS2, on the other hand, explore the consequences of setting only a subset of the possible inputs.

5.1 SA: Simulated Annealing

Simulated Annealing is a classic stochastic search algorithm. It was first described in 1953 [48] and refined in 1983 [40]. The algorithm’s namesake, annealing, is a technique from metallurgy, where a material is heated, then cooled. The heat causes the atoms in the material to wander randomly through different energy states and the cooling process increases the chances of finding a state with a lower energy than the starting position.

For each round, SA “picks” a neighboring set of mitigations. To calculate this neighbor, a function traverses the mitigation settings of the current state and randomly flips those mitigations (at a 5% chance). If the neighbor has a better score, SA will move to it and set it as the current state. If it isn’t better, the algorithm will decide whether to move to it based on the mathematical function:

$$prob(w, x, y, temp(y, z)) = e^{((w-x) * \frac{y}{temp(y, z)})} \quad (3)$$

$$temp(y, z) = \frac{(z - y)}{z} \quad (4)$$

If the value of the `prob` function is greater than a randomly generated number, SA will move to that state anyways. This randomness is the very cornerstone of the Simulated

```

1. Procedure SA
2. MITIGATIONS:= set of mitigations
3. SCORE:= score of MITIGATIONS
4. while TIME < MAX_TIME && SCORE < MIN_SCORE //minScore is a constant score (threshold)
5.     find a NEIGHBOR close to MITIGATIONS
6.     NEIGHBOR_SCORE:= score of NEIGHBOR
7.     if NEIGHBOR_SCORE > SCORE
8.         MITIGATIONS:= NEIGHBOR
9.         SCORE:= NEIGHBOR_SCORE
10.    else if prob(Score, NEIGHBOR_SCORE, TIME, temp(TIME, MAX_TIME)) > RANDOM
11.        MITIGATIONS:= NEIGHBOR
12.        SCORE:= NEIGHBOR_SCORE
13.    TIME++
14. end while
15. return MITIGATIONS

```

Fig. 6. Pseudocode for SA

Annealing algorithm. Initially, the “atoms” (current solutions) will take large random jumps, sometimes to even sub-optimal new solutions. These random jumps allow simulated annealing to sample a large part of the search space, while avoiding being trapped in local minima. Eventually, the “atoms” will cool and stabilize and the search will converge to a simple hill climber.

As shown in line 4 of Figure 6, the algorithm will continue to operate until the number of tries is exhausted or a score meets the threshold requirement.

5.2 MaxFunWalk

The design of simulated annealing dates back to the 1950s. In order to benchmark our new search engine (KEYS) against a more recent state-of-the-art algorithm, we implemented the variant of MaxWalkSat described below.

MaxWalkSat is a local search method designed to address the problem of boolean satisfiability [39]. MaxWalkSat is a variant of that applies weights to each clause in a conjunctive normal form equation [58]. While WalkSat tries to satisfy the entire set of clauses, MaxWalkSat tries to maximize the sum of the weights of the satisfied clauses.

In one respect, both algorithms can be viewed as a variant of simulated annealing. Whereas simulated annealing always selects the next solution randomly, the WalkSat algorithms will *sometimes* perform random selection while, other times, conduct a local search to find the next best setting to one variable.

MaxFunWalk is a generalization of MaxWalkSat:

- MaxWalkSat is defined over CNF formulae. The success of a collection of variable settings is determined by how many clauses are “satisfiable” (defined using standard boolean truth tables).
- MaxWalkFun, on the other hand, assumes that there exist an arbitrary function that can assess a collection of variable settings. Here, we use DDP model as the assessment function.

```

1. Procedure MaxFunWalk
2. for TRIES:=1 to MAX-TRIES
3.   SELECTION:=A randomly generate assignment of mitigations
4.   for CHANGED:=1 to MAX-CHANGES
5.     if SCORE satisfies THRESHOLD return
6.     CHOSEN:= a random selection of mitigations from SELECTION
7.     with probability P
8.     flip a random setting in CHOSEN
9.     with probability (P-1)
10.    flip a setting in CHOSEN that maximizes SCORE
11.   end for
12. end for
13. return BESTSCORE

```

Fig. 7. Pseudocode for MaxFunWalk

Note that $MaxWalkFun = MaxWalkSat$ if assessment is via a logical truth table.

The MaxFunWalk procedure is shown in Figure 7. When run, the user supplies an ideal cost and attainment. This setting is normalized, scored, and set as a goal threshold. If the current setting of mitigations satisfies that threshold, the algorithm terminates.

MaxFunWalk begins by randomly setting every mitigation. From there, it will attempt to make a *single* change until the threshold is met or the allowed number of changes runs out (100 by default). A random subset of mitigations is chosen and a random number P between 0 and 1 is generated. The value of P will decide the form that the change takes:

- $P \leq \alpha$: A stochastic decision is made. A setting is changed completely at random within the set CHOSEN.
- $P > \alpha$: Local search is utilized. Each mitigation in CHOSEN is tested until one is found that improves the current score.

The best setting of α is a domain-specific. For this study, we used $\alpha = 0.3$.

If the threshold is not met by the time that the allowed number of changes is exhausted, the set of mitigations is completely reset and the algorithm starts over. This measure allows the algorithm to avoid becoming trapped in local maxima. For the DDP models, we found that the number of retries has little effect on solution quality.

If the threshold is never met, MaxFunWalk will reset and continue to make changes until the maximum number of allowed resets is exhausted. At that point, it will return the best settings found.

As an additional measure to improve the results found by MaxFunWalk, a heuristic was implemented to limit the number of mitigations that could be set at one time. If too many are set, the algorithm will turn off a few in an effort to bring the cost factor down while minimizing the effect on the attainment.

5.3 A* (ASTAR)

A* is a best-first path finding algorithm that uses distance from origin (G) and estimated cost to goal (H) to find the best path [31]. The algorithm is widely used [34, 53, 55, 62].

A* is a natural choice for DDP optimization since the scoring function described above is actually a Euclidean distance measure to the desired goal of maximum attainment and minimum costs. Hence, for H, we can make direct use of Equation 2.

The ASTAR algorithm keeps a *closed* list in order to prevent backtracking. We begin by adding the starting state to the closed list. In each “round”, a list of neighbors is populated from the series of possible states reached by making a change to a single mitigation. If that neighbor is not on the closed list, two calculations are made:

- G = Distance from the start to the current state plus the additional distance between the current state and that neighbor.
- H = Distance from that neighbor to the goal (an ideal spot, usually 0 cost and a high attainment). For DDP models, we use Equation 2 to compute H.

The *best* neighbor is the one with the lowest $F=G+H$. The algorithm “travels” to that neighbor and adds it to the closed list. Part of the optimality of the A* algorithm is that the distance to the goal is underestimated. Thus, the final goal is never actually reached by ASTAR. Our implementation terminates once it stops finding better solutions for a total of ten rounds. This number was chosen to give ample time for ASTAR to become “unstuck” if it hits a corner early on.

```

1. Procedure ASTAR
2. CURRENT_POSITION:= Starting assignment of mitigations
3. CLOSED[0]:= Add starting position to closed list
4.
5. while END:= false
6.   NEIGHBOR_LIST:=list of neighbors
7.   for each NEIGHBOR in NEIGHBOR_LIST
8.     if NEIGHBOR is not in CLOSED
9.       G:=distance from start
10.      H:=distance to goal
11.      F:=G+H
12.      if F<BEST_F
13.        BEST_NEIGHBOR:=NEIGHBOR
14.   end for
15.   CURRENT_POSITION:= BEST_NEIGHBOR
16.   CLOSED[+]:=Add new state to closed list
17.   if STUCK
18.     END:= true
19. end while
20. return CURRENT_POSITION

```

Fig. 8. Pseudocode for ASTAR

5.4 KEYS and KEYS2

The core premise of KEYS and KEYS2 is that the above algorithms perform over-elaborate searches. Suppose that the behavior of a large system is determined by a small number of *key* variables. If so, then a very rapid search for solutions can be found by (a) finding these *keys* then (b) explore the ranges of the key variables.

As documented in our *Related Work* section, this notion of *keys* has been discovered and rediscovered many times by many researchers. Historically, finding the keys has seen to be a very hard task. For example, finding the keys is analogous to finding the *minimal environments* of DeKleer’s ATMS algorithm [15]. Formally, this *logical abduction*, which is an NP-hard task [6].

Our method for finding the keys uses a Bayesian sampling method. If a model contains keys then, by definition, those variables must appear in all solutions to that model. If model outputs are scored by some oracle, then the key variables are those with ranges that occur with very different frequencies in high/low scored model outputs. Therefore, we need not search for the keys- rather, we just need to keep frequency counts on how often ranges appear in *best* or *rest* outputs.

KEYS implements this Bayesian sampling methods. It has two main components - a greedy search and the BORE ranking heuristic. The greedy search explores a space of M mitigations over the course of M “eras”. Initially, the entire set of mitigations is set randomly. During each era, one more mitigation is set to $M_i = X_j, X_j \in \{true, false\}$. In the original version of KEYS [45], the greedy search fixes one variable per era. This paper experimented with a newer version, called KEYS2, that fixes an increasing number of variables as the search progresses (see below for details).

KEYS (and KEYS2), each era e generates a set $\langle input, score \rangle$ as follows:

- 1: *MaxTries* times repeat:
 - *Selected*[1.. $(e - 1)$] are settings from previous eras.
 - *Guessed* are randomly selected values for unfixed mitigations.
 - *Input* = *selected* \cup *guessed*.
 - Call *model* to compute *score* = *ddp(input)*;
- 2: The *MaxTries* *scores* are divided into $\beta\%$ “best” and remainder become “rest”.
- 3: The *input* mitigation values are then scored using BORE (described below).
- 4: The top ranked mitigations (the default is one, but the user may fix multiple mitigations at once) are fixed and stored in *selected*[e].

```

1. Procedure KEYS
2. while FIXED_MITIGATIONS != TOTAL_MITIGATIONS
3.   for I:=1 to 100
4.     SELECTED[1...(I-1)] = best decisions up to this step
5.     GUESSED = random settings to the remaining mitigations
6.     INPUT = SELECTED + GUESSED
7.     SCORES= SCORE(INPUT)
8.   end for
9.   for J:=1 to NUM_MITIGATIONS_TO_SET
10.    TOP_MITIGATION = BORE(SCORES)
11.    SELECTED[FIXED_MITIGATIONS++] = TOP_MITIGATION
12.  end for
13. end while
14. return SELECTED

```

Fig. 9. Pseudocode for KEYS

The search moves to era $e + 1$ and repeats steps 1,2,3,4. This process stops when every mitigation has a setting. The exact settings for $MaxTries$ and β must be set via engineering judgment. After some experimentation, we used $MaxTries = 100$ and $\beta = 10$. For full details, see Figure 9.

KEYS ranks mitigations by combining a novel support-based Bayesian ranking measure. BORE [10] (short for “best or rest”) divides numeric scores seen over K runs and stores the top 10% in *best* and the remaining 90% scores in the set *rest* (the *best* set is computed by studying the delta of each score to the best score seen in any era). It then computes the probability that a value is found in *best* using Bayes theorem. The theorem uses evidence E and a prior probability $P(H)$ for hypothesis $H \in \{best, rest\}$, to calculate a posteriori probability $P(H|E) = P(E|H)P(H) / P(E)$. When applying the theorem, *likelihoods* are computed from observed frequencies. These likelihoods (called “like” below) are then normalized to create probabilities. This normalization cancels out $P(E)$ in Bayes theorem. For example, after $K = 10,000$ runs are divided into 1,000 *best* solutions and 9,000 *rest*, the value *mitigation31 = false* might appear 10 times in the *best* solutions, but only 5 times in the *rest*. Hence:

$$\begin{aligned}
E &= (mitigation31 = false) \\
P(best) &= 1000/10000 = 0.1 \\
P(rest) &= 9000/10000 = 0.9 \\
freq(E|best) &= 10/1000 = 0.01 \\
freq(E|rest) &= 5/9000 = 0.00056 \\
like(best|E) &= freq(E|best) \cdot P(best) = 0.001 \\
like(rest|E) &= freq(E|rest) \cdot P(rest) = 0.000504 \\
P(best|E) &= \frac{like(best|E)}{like(best|E) + like(rest|E)} = 0.66
\end{aligned} \tag{5}$$

Previously [10], we have found that Bayes theorem is a poor ranking heuristic since it is easily distracted by low frequency evidence. For example, note how the probability of E belonging to the best class is moderately high even though its support is very low; i.e. $P(best|E) = 0.66$ but $freq(E|best) = 0.01$.

To avoid the problem of unreliable low frequency evidence, we augment Equation 5 with a support term. Support should *increase* as the frequency of a value *increases*, i.e. $like(best|E)$ is a valid support measure. Hence, step 3 of our greedy search ranks values via

$$P(best|E) * support(best|E) = \frac{like(best|E)^2}{like(best|E) + like(rest|E)} \tag{6}$$

For each era, KEYS samples the DDP models and fixes the top $N = 1$ settings. The following observation suggested that $N = 1$ is, perhaps, an overly conservative search policy. At least for the DDP models, we have observed that the improvement in costs and attainments generally increase for each era of KEYS. This lead to the speculation that we could jump further and faster into the solution space by fixing $N \geq 1$ settings per era. Such a *jump* policy can be implemented as a small change to KEYS, which we call KEYS2.

- Standard KEYS assigns the value of one to NUM_MITIGATIONS_TO_SET (see the pseudo-code of Figure 9);

- Other variants of KEYS assigns larger values.

In era 1, KEYS2 behaves exactly the same as KEYS while in (say) era 3, KEYS2 will fix the top 3 ranked ranges. Since it sets more variables at each era, KEYS2 terminates earlier than KEYS.

Note that decision ordering diagrams could be directly generated during execution, just by collection statistics from the `SCORES` array used in line 7 of Figure 9.

6 Results

Each of the above algorithms was tested on the five models of Figure 4. Note that:

- Models one and three are trivially small and we used them to debug our code. We report our results using models two, four and five since they are large enough to stress test real-time optimization.
- Model 4 was discussed in [46] in detail. The largest, model 5 was optimized previously in [20]. At that time (2002), it took 300 seconds to generate solutions using our old, very slow, rule learning method.

We also studied how well KEYS and KEYS2 scale to larger models. Further, we instrumented KEYS and KEYS2 to generate decision ordering diagrams. The results from all these experiments are shown below.

6.1 Attainment and Costs

We ran each of our algorithms 1000 times on each model. This number was chosen because it yielded enough data points to give a clear picture of the span of results. At the same time, it is a low enough number that we can generate a set of results in a fairly short time span.

The results are pictured in Figure 10. Attainment is along the x-axis and cost (in thousands) is along the y-axis. Note that better solutions fall towards the bottom right of each plot; i.e. lower costs and higher attainment. Also better solutions exhibit less variance, i.e. are clumped tighter together.

These graphs give a clear picture of the results obtained by our various algorithms. Two methods are clearly inferior:

- SA exhibit the worst variance, lowest attainments, and highest costs.
- MaxFunWalk is better than SA (less variance, lower costs, higher attainment) but it is clearly inferior to the other methods.

As to the others:

- On larger models such as model4 and model5 KEYS and KEYS2 exhibit lower variance, lower costs, higher attainments than ASTAR.
- On smaller models such as model2, ASTAR usually produces higher attainments and lower variance than the KEYS algorithms (this advantage disappears on the larger models). However, observe the results near the (0, 0) point of model2's ASTAR results: sometimes ASTAR's heuristic search failed completely for that model

6.2 Runtime Analysis

Measured in terms of attainment and cost, there is little difference between KEYS and KEYS2. However, as shown by Figure 11, KEYS2 runs twice to three times as fast as KEYS. Interestingly, Figure 11 ranks two of the algorithms in a similar order to Figure 10:

- SA is clearly the slowest;
- MaxFunWalk is somewhat better but not as fast as the other algorithms.

As to ASTAR versus KEYS or KEYS2:

- ASTAR is faster than KEYS;
- and KEYS2 runs in time comparable to ASTAR.

Measured in terms of runtimes, there is little to recommend KEYS2 over ASTAR. However, ASTAR’s heuristic guesses were sometimes observed to be sub-optimal (recall the above discussion on the $(0, 0)$ results in model2’s ASTAR results). Such sub-optimality was never observed for KEYS2.

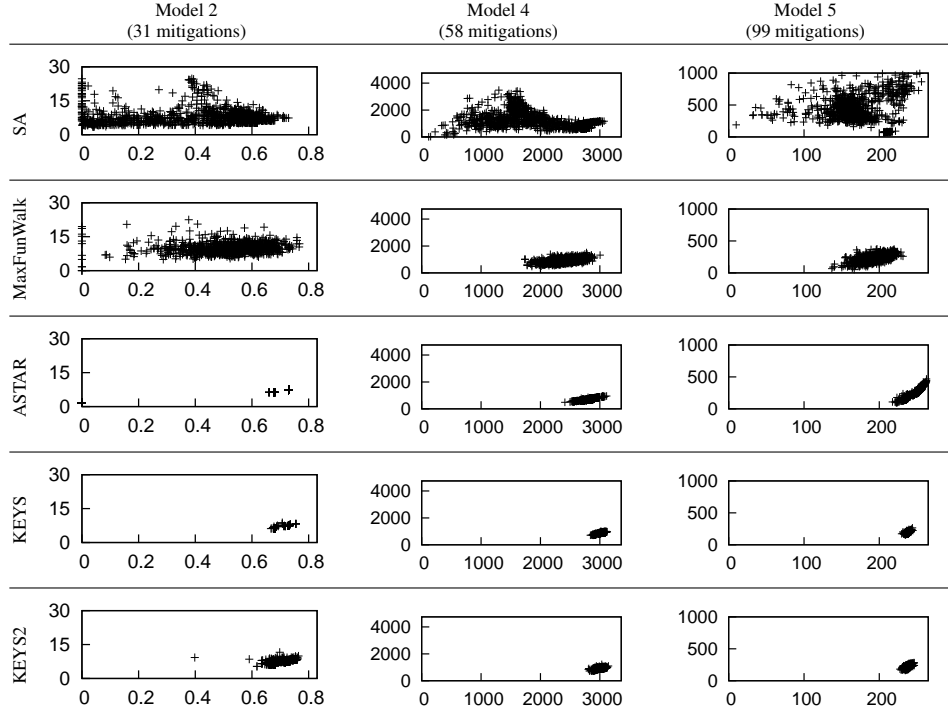


Fig. 10. 1000 results of running five algorithms on three models (15,000 runs in all). The y-axis shows cost and the x-axis shows attainment. The size of each model is measured in number of mitigations. Note that better solutions fall towards the bottom right of each plot; i.e. lower costs and higher attainment. Also better solutions exhibit less variance, i.e. are clumped tighter together.

	Model 2 (31 mitigations)	Model 4 (58 mitigations)	Model 5 (99 mitigations)
SA	0.577	1.258	0.854
MaxFunWalk	0.122	0.429	0.398
ASTAR	0.003	0.017	0.048
KEYS	0.011	0.053	0.115
KEYS2	0.006	0.018	0.038

Fig. 11. Runtimes in seconds, averaged over 100 runs, measured using the “real time” value from the Unix `times` command. The size of each model is measured in number of mitigations (and for more details on model size, see Figure 4).

model	expansion	model size	Runtime (secs)		KEYS KEYS2
			KEYS	KEYS2	
2	1	62	0.01	0.01	1.07
2	2	124	0.03	0.02	1.23
4	1	139	0.04	0.02	2.29
5	1	201	0.13	0.04	3.18
2	4	248	0.10	0.05	2.09
4	2	278	0.17	0.05	3.48
5	2	402	0.50	0.12	4.26
2	8	496	0.44	0.14	3.21
4	4	556	0.73	0.16	4.66
5	4	804	1.98	0.38	5.21
4	8	1112	2.97	0.52	5.71
5	8	1608	8.06	1.35	5.96

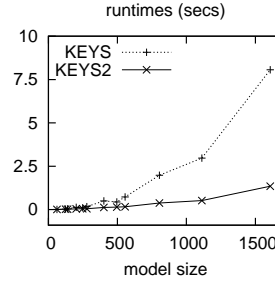


Fig. 12. Runtimes KEYS vs KEYS2 (medians over 1000 repeats) as models increase in size. The “model” number in column one corresponds to Figure 4. The “expansion factor” of column two shows how much the instance generator expanded the model. The “model sizes” of column three are the sum of mitigations, requirements, and risks seen in the expanded model.

6.3 Scale-Up Studies

Figure 12 and Figure 13 shows the effect of changing models size on the number of calls to the model and the runtimes of KEYS and KEYS2. To generate these graph, an *instance generator* was created that:

- Examined the real-world DDP models of Figure 4;
- Extracted statistics related to the different types of nodes (mitigations or risks or requirements) and the number of edges between different types of nodes;
- Used those statistics to build random models that were 1,2,4 and 8 times bigger than the original models.

Figure 14 shows the results of curve fitting to the plots of Figure 12 and Figure 13. The KEYS and KEYS2 performance curves fit a low-order polynomial (of degree two) with very high coefficients of determination ($R^2 \geq 0.98$).

Figure 14 suggests that we could scale *either* KEYS or KEYS2 to larger models. However we still recommend KEYS2. The column marked $\frac{KEYS}{KEYS2}$ in Figure 13 shows the ratio of the number of calls made by KEYS vs KEYS2. As models get larger, KEYS calls the model an order of magnitude more times than KEYS2. If applied to models with slower runtimes than DDP, then an order of magnitude more executions is highly undesirable.

model	expansion	model size	Calls to model		KEYS KEYS2
			KEYS	KEYS2	
2	1	62	3100	800	3.9
2	2	124	6200	1100	5.6
4	1	139	5800	1100	5.3
5	1	201	9900	1400	7.1
2	4	248	12400	1600	7.8
4	2	278	11600	1500	7.7
5	2	402	19800	2000	9.9
2	8	496	24800	2200	11.3
4	4	556	23200	2200	10.5
5	4	804	39600	2800	14.1
4	8	1112	46400	3000	15.5
5	8	1608	79200	4000	19.8

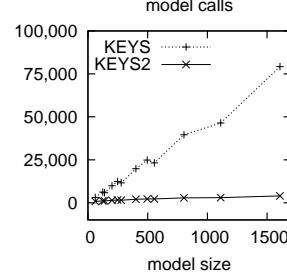


Fig. 13. Number of model calls made by KEYS vs KEYS2 (medians over 1000 results) as models increase in size. This figure uses the same column structure as Figure 12.

	KEYS		KEYS2	
	runtimes	model calls	runtimes	model calls
exponential	0.82	0.83	0.88	0.93
polynomial (of degree 2)	0.99	0.99	0.99	0.98

Fig. 14. Coefficients of determination R^2 of KEYS/KEYS2 performance figures, fitted to two different functions: exponential or polynomial of degree two. Higher values indicate a better curve fit. In all cases, the best fit is not exponential.

6.4 Decision Ordering Algorithms

The decision ordering diagrams of Figure 15 shows the effects of the decisions made by KEYS and KEYS2. For KEYS and KEYS2, at $x = 0$, all of the mitigations in the model are set at random. During each subsequent era, more mitigations are fixed (KEYS sets one at a time, KEYS2 sets one or more at a time). The lines in each of these plots show the median and spread seen in the 100 calls to the `model` function during each round.

Note that the these diagrams are *tame* and *well-behaved*:

- *Tame*: The spreads quickly shrink to values that are a small fraction of the median.
- *Well-behaved*: The median’s move smoothly to a plateau of best performance (high attainment, low costs).

On termination (at maximum value of x), KEYS and KEYS2 arrive at nearly identical median results (caveat: for `model2`, KEYS2 attains slightly more requirements at a slightly higher cost than KEYS). The spread plots for both algorithms are almost indistinguishable (exception: in `model2`, the KEYS2 spread is less than KEYS). That is, KEYS2 achieves the same results as KEYS but (as shown in Figure 11 and Figure 12) it does so in less time.

A core assumption of this work was the “keys” assumption; i.e. a small number of variables set the remaining model variables. Figure 15 offers much support for this assumption: observe how most of the improvement in costs and attainments were achieved after KEYS and KEYS2 made only a handful of decisions (ten, or less).

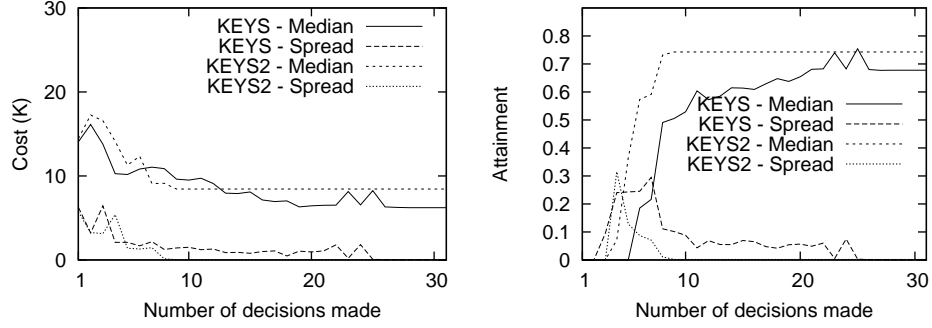


Figure 15a: Internal Decisions on Model 2.

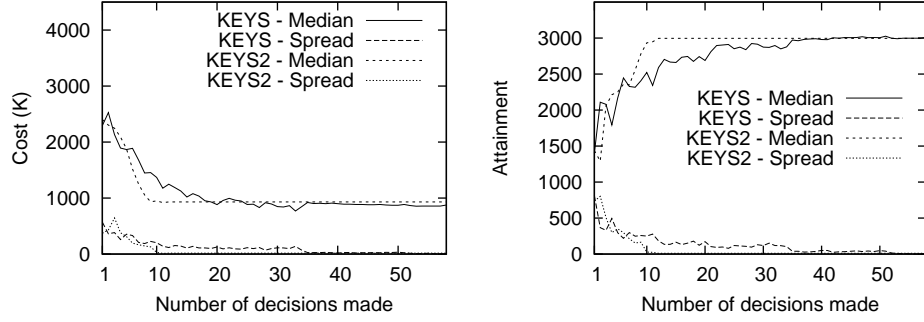


Figure 15b: Internal Decisions on Model 4.

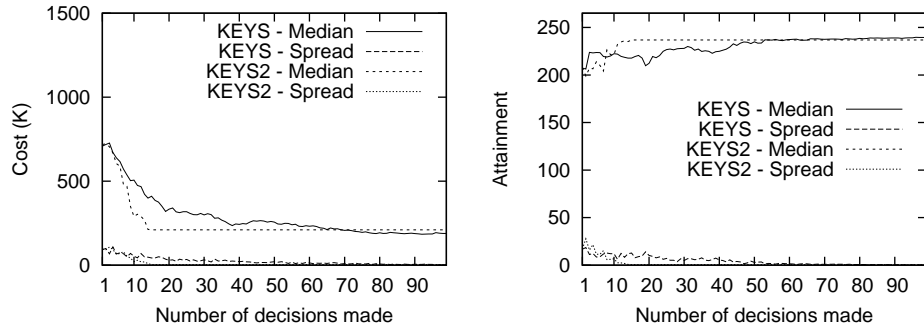


Figure 15c: Internal Decisions on Model 5.

Fig. 15. Median and spread of partial solutions learned by KEYS and KEYS2. X-axis shows the number of decisions made. “Median” shows the 50-th percentile of the measured value seen in 100 runs at each era. “Spread” shows the difference between the 75th and 50th percentile.

On another matter, it is insightful to reflect on the effectiveness of different algorithms for generating decision ordering diagrams. KEYS2 is our more direct and fastest method. As mentioned above, the information required can be collected during one execution. On the other hand, SA, ASTAR, and MaxWalkSat would require a post-processor to generate the diagrams:

- Given D possible decisions, At each era, KEYS and KEYS2 collects statistics on partial solutions where $1, 2, 3, \dots |d|$ variables are fixed (where d is the set of decisions) while the remaining $D - d$ decisions are made at random.
- ASTAR, Simulated Annealing, and MaxFunWalk work with full solutions since at each step they offer settings to all $d_i \in D$ variables. That is, in the current form, they cannot comment on partial solutions. Modified forms of these algorithms could add in extra instrumentation and extra post-processing to comment on partial solutions using methods like feature subset selection [26] or a sensitivity analysis [57].

7 Related Work

7.1 Early vs Later Life-cycle Requirements Engineering

The case studies presented in this paper come JPL’s Team X meetings. Team X conducts very early life-cycle requirement discussions.

Once a system is running, released, and being maintained/extended, another problem is *release planning*; i.e. what features to add to the next N releases. To solve this problem, an inference engine must reason about how functionality extensions to current software can best satisfy outstanding stakeholder requirements. The challenge of release planning is that the benefits of added functionality must be weighed against the cost of implementing those extensions.

Several approaches have been applied to this problem including:

- The OPTIMIZE tool of Ngo-The and Ruhe [?] that combines linear programming with genetic programming to optimize release plans for software projects
- The weighted Pareto optimal genetic algorithm approach of Zhang et al. [65]

(See also the earlier comparison of exact vs greedy algorithms by Bagnall et al. [?]).

Without further experimentation, we cannot assert that KEYS2 will work as well on later life-cycle models (such as those used in release planning) as it did above (on the earlier life-cycle Team X models). However, at this time, we can see no reason why KEYS2 would not work as a non-linear optimizer of these later life-cycle models. This could be a productive area for future work.

7.2 Other Optimizers

As documented by the search-based SE literature [11, 29, 30, 54] and Gu et al [25], there are many possible optimization methods. For example:

- Gradient descent methods assume that an objective function $F(X)$ is differentiable at any single point N . A Taylor-series approximation of $F(X)$ can be shown to decrease fastest if the negative gradient $(-\Delta F(N))$ is followed from point N .

- Sequential methods run on one CPU while parallel methods spread the work over a distributed CPU farm.
- Discrete methods assume model variables have a finite range (that may be quite small) while continuous methods assume numeric values with a very large (possibly infinite) range.
- The search-based SE literature prefers meta-heuristic methods like simulated annealing, genetic algorithms and tabu search.
- Some methods map discrete values true/false into a continuous range 1/0 and then use integer programming methods like CPLEX [50] to achieve results.
- Other methods find overlaps in goal expressions and generate a binary decision diagram (BDD) where parent nodes store the overlap of children nodes.

This list is hardly exhaustive: Gu et al. list hundreds of other methods and no single paper can experiment with them all. All the algorithms studied here are discrete and sequential. We are currently exploring parallel versions of our optimizers but, so far, the communication overhead outweighs the benefits of parallelism.

As to the general class of gradient descent methods, we do not use them since they assume the objective function being optimizing is essentially continuous. Any model with an “if” statement in it is not continuous since, at the “if”, there program’s behavior may become discontinuous. The requirements models studied here are discontinuous about every subset of every possible mitigation.

As to the more specific class of integer programming methods, we do not explore them here for two reasons. Coarfa et al. [12] found that integer programming-based approaches ran two an order of magnitude slower than discrete methods like the MaxWalkSat and KEYS2 algorithm we use below. Similar results reported by Gu et.al where discrete methods ran 100 times faster than integer programming [25].

Harman offers another reason to avoid integer programming methods. In his search-based SE manifest, Harman [29] argues that many SE problems are over-constrained and so there may exist no precise solution that covers all constraints. A complete solution over all variables is hence impossible and partial solution based on heuristic search methods are preferred. Such methods may not be complete; however, as Clarke et al remark, “...software engineers face problems which consist, not in finding *the* solution, but rather, in engineering an *acceptable* or *near-optimal solution* from a large number of alternatives.” [11]

7.3 Models of Requirements Engineering

DDP is a ultra-lightweight modeling tool. The value of ultra-lightweight ontologies in early life cycle modeling is widely recognized. For example Mylopoulos’ *soft-goal* graphs [51, 52] represent knowledge about non-functional requirements. Primitives in soft goal modeling include statements of partial influence such as *helps* and *hurts*. Another commonly used framework in the design rationale community is a “questions-options-criteria” (QOC) graph [60]. In QOC graphs:

- *Questions* suggest *options*. Deciding on one option can raise other questions;
- Options shown in a box denote *selected options*;

- Options are assessed by *criteria*;
- Criteria are gradual knowledge; i.e. they *tend/reject to support* options.

QOCs can succinctly summarize lengthy debates; e.g. the 480 sentences uttered in a debate on interface options can be displayed in a QOC graph on a single page [42]. Saaty’s Analytic Hierarchy Process (AHP) [56] is a variant of QOC.

While DDP shares many of the design aspects of softgoals & QOC & AHP, it differs in its representations and inference method. As explained above around Equation 1, where as AHP and QOC and softgoals propagate influences over hierarchies, DDP propagate influences over matrices.

7.4 Formal Models of Requirements Engineering

Zave&Jackson [?] define requirements engineering as finding the specification S for the domain assumptions K that satisfies the given requirements R ; i.e.

$$\text{find } S \text{ such that } S \vdash R \quad (7)$$

Jureta, Mylopoulos & Faulkner [?] (hereafter JMF) take issue with Equation 7, saying that it implicitly assume that K, S, R are precise and complete enough for the satisfaction relation to hold. More specifically, JMF complain that Equation 7 does not permit partial fulfillment of (some) non-functional requirements. Also, the Zave&Jackson definition does not allow any preference ordering of *specification*₁ over *specification*₂. JMF therefore offer a replacement ontology where classical inference is replaced with operators that supports the generation and ranking of subsets of domain assumptions that lead to maximal (w.r.t. size) subsets of the possible goals, and softgoal quality criteria⁵.

DDP reinterprets “ \vdash ” in Equation 7 as an inference across numeric quantities, rather than the inference over discrete logical variables suggested by Zave&Jackson. Hence, it can achieve the same goals as JMF (ranking of partial solutions with weighted goals) without requiring the JMF ontology.

7.5 Requirements Analysis Tools

There exist many powerful requirements analysis tools including continuous simulation (also called system dynamics) [1, 61], state-based simulation (including petri net and data flow approaches) [2, 27, 44], hybrid-simulation (combining discrete event simulation and systems dynamics) [17, 43, 59], logic-based and qualitative-based methods [5, chapter 20] [35], and rule-based simulations [49]. One can find these models being used in the requirements phase (i.e. the DDP tool described below), design refactoring using patterns [23], software integration [16], model-based security [38], and performance assessment [4]. Many researchers have proposed support environments to

⁵ According to JMF: “a salient characteristic of softgoals is that they cannot be satished to the ideal extent, not only because of subjectivity, but also because the ideal level of satisfaction is beyond the resources available to (and including) the system. It is therefore said that a softgoal is not satished, but satished.”

help explore the increasingly complex models that engineers are developing. Gray et al [24] have developed the Constraint-Specification Aspect Weaver (C-Saw), which uses aspect-oriented approaches [22] to help engineers in the process of model transformation. Cai and Sullivan [7] describe a formal method and tool called *Simon* that “supports interactive construction of formal models, derives and displays design structure matrices... and supports simple design impact analysis.” Other tools of note are lightweight formal methods such as ALLOY [36] and SCR [32] as well as UML tools that allow for the execution of life cycle specifications (e.g. CADENA [8]).

Many of the above tools were built to maximize the expressive power of the representation language or the constraint language used to express invariants. What distinguishes our work is that we are *willing to trade off representational or constrain expressiveness for faster runtimes*. There exists a class of ultra-lightweight model languages which, as we show above, can be processed very quickly. Any of the tools listed in the last paragraph are also candidate solutions to the problem explored in this paper, if it can be shown that their processing can generate tame and well-behaved decision ordering diagrams in a timely manner.

7.6 Other Work on “Keys”

Elsewhere [47], we have documented dozens of papers that have reported the keys effect (that a small number of variables set the rest) under different names including *narrows*, *master-variables*, *back doors*, and *feature subset selection*:

- Amarel [3] observed that search problems contain narrow sets of variables or collars that must be used in any solution. In such a search space, what matters is not so much how you get to these collars, but what decision you make when you get there. Amarel defined macros encoding paths between *narrows*, effectively permitting a search engine to jump between them.
- In a similar theoretical analysis, Menzies & Singh [47] computed the odds of a system selecting solutions to goals using complex, or simpler, sets of preconditions. In their simulations, they found that a system will naturally select for tiny sets of preconditions (a.k.a. the keys) at a very high probability.
- Numerous researchers have examined *feature subset selection*; i.e. what happens when a data miner deliberately ignores some of the variables in the training data. For example, Kohavi and John [41] showed in numerous datasets that as few as 20% of the variables are *key* - the remaining 80% of variables can be ignored without degrading a learner’s classification accuracy.
- Williams et.al. [64] discuss how to use keys (which they call “back doors”) to optimize search. Constraining these back doors also constrains the rest of the program. So, to quickly search a program, they suggest imposing some set values on the key variables. They showed that setting the keys can reduce the solution time of certain hard problems from exponential to polytime, provided that the keys can be cheaply located, an issue on which Williams et.al. are curiously silent.
- Crawford and Baker [14] compared the performance of a complete TABLEAU prover to a very simple randomized search engine called ISAMP. Both algorithms assign a value to one variable, then infer some consequence of that assignment

with forward checking. If contradictions are detected, TABLEAU backtracks while ISAMP simply starts over and re-assigns other variables randomly. Incredibly, ISAMP took *less* time than TABLEAU to find *more* solutions using just a small number of tries. Crawford and Baker hypothesized that a small set of *master variables* set the rest and that solutions are not uniformly distributed throughout the search space. TABLEAU’s depth-first search sometimes drove the algorithm into regions containing no solutions. On the other hand, ISAMP’s randomized sampling effectively searches in a smaller space.

In summary, the core assumption of our algorithms are supported in many domains.

8 Conclusion

Requirements tools as as the DDP tool (used at NASA for early lifecycle discussions) contain a shared group memory storing all the requirements, risks, and mitigations of each member of the group. Software tools can explore this shared memory to find consequences and interactions that may have been overlooked.

Studying that group memory is a non-linear optimization task: possible benefits must be traded off against the increased cost of applying various mitigations. Harman [29] caution that solutions to non-linear problems may be “brittle”; i.e. small changes to the search results may dramatically alter the effectiveness of the solution. Hence, reporting an analysis of this shared group memory, it is vitally important to comment on the robustness of the solution.

Decision ordering diagrams are a solution robustness assessment method. The diagrams rank all the possible decisions most-to-least influential. Each point x on the diagrams shows the effects on imposing the conjunction of decisions $1 \leq j \leq x$. These diagrams can comment on the robustness and neighborhood of solution $\{d_1..d_x\}$ using two operators:

1. By considering the variance of the performance statistics after applying $\{d_1..d_x\}$.
2. By comparing the results of using the first x decisions to that of using the first $x - 1$ or $x + 1$ actions.

Since the diagrams are sorted, then this analysis of robustness and neighborhood takes, at most, time linear of the number decisions. That is, theoretically, it takes linear time to *use* a decision ordering diagram (see §4.3).

Also, empirically, it take low-order polynomial time to *generate* a decision ordering diagram using KEYS2. This algorithm makes the “key” assumption (that a small group of variables set everything else) and uses Bayesian ranking mechanism to quickly find those keys. As discussed above in the *Related Work* section, this assumption holds over a wide range of models used in a wide range of domains. This “keys” assumption can be remarkably effective: our empirical results show that KEY2 can generate decision ordering diagrams faster than the other algorithms studied here. Better yet, curve fits to our empirical results show that KEYS runs on low-order polynomial time (of degree two) and so should scale to very large models.

Prior to this work, our two pre-experimental concerns were that:

- We would need to trade solution robustness against solution quality the most robust solutions may not have the highest quality.
- Demonstrating solution robustness requires multiple calls to an analysis procedure.

At least for the models studied here, neither concern were realized. KEYS2 generated the highest quality solutions (lowest cost, highest attainments) and did so in the fast runtimes.

In §4.3 it was argued that decision ordering diagrams are useful when they are *timely* to generate while being *well-behaved* and *tame*. KEYS2's results are the most *timely* (fastest to generate) of all the methods studied here. As to the other criteria, Figure 15 shows that KEYS2's decision ordering diagrams:

- Move smoothly to a plateau with only a small amount of “jitter”;
- Have very low spreads, compared to the median results.

That is, at least for the models explored here, KEYS2 generated decision ordering diagrams they are both *well-behaved* and *tame*.

In summary, we recommend KEYS2 for generating decision ordering diagrams since, apart from the (slightly slower) KEYS algorithm, we are unaware of other search-based software engineering methods that enable such a rapid reflection of solution robustness.

9 Acknowledgments

This research was conducted at West Virginia University, the Jet Propulsion Laboratory under a contract with the National Aeronautics and Space administration, Alderson-Broadbent College, and Miami University. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement by the United States Government

References

1. T. Abdel-Hamid and S. Madnick. *Software Project Dynamics: An Integrated Approach*. Prentice-Hall Software Series, 1991.
2. M. Akhavi and W. Wilson. Dynamic simulation of software process models. In *Proceedings of the 5th Software Engineering Process Group National Meeting (Held at Costa Mesa, California, April 26 - 29)*. Software engineering Institute, Carnegie Mellon University, 1993.
3. S. Amarel. Program synthesis as a theory formation task: Problem representations and solution methods. In R. S. Michalski, J. G. Carbonell, and T. M. Mitchell, editors, *Machine Learning: An Artificial Intelligence Approach: Volume II*, pages 499–569. Kaufmann, Los Altos, CA, 1986.
4. S. Balsamo, A. D. Marco, P. Inverardi, and M. Simeoni. Model-based performance prediction in software development: A survey. *IEEE Transactions on Software Engineering*, 30(5), May 2004.
5. I. Bratko. *Prolog Programming for Artificial Intelligence. (third edition)*. Addison-Wesley, 2001.

6. T. Bylander, D. Allemang, M. Tanner, and J. Josephson. The Computational Complexity of Abduction. *Artificial Intelligence*, 49:25–60, 1991.
7. Y. Cai and K. J. Sullivan. Simon: modeling and analysis of design space structures. In *ASE '05: Proceedings of the 20th IEEE/ACM international Conference on Automated software engineering*, pages 329–332, New York, NY, USA, 2005. ACM Press.
8. A. Childs, J. Greenwald, G. Jung, M. Hoosier, and J. Hatcliff. Calm and cadena: Meta-modeling for component-based product-line development. *IEEE Computer*, 39(2), February 2006. Available from <http://projects.cis.ksu.edu/docman/view.php/7/129/CALM-Cadena-IEEE-Comp%uter-Feb-2006.pdf>.
9. A. Cimatti, E. Clarke, E. Giunchiglia, F. Giunchiglia, M. Pistore, M. Roveri, R. Sebastiani, and A. Tacchella. NuSMV Version 2: An OpenSource Tool for Symbolic Model Checking. In *Proc. International Conference on Computer-Aided Verification*, 2002.
10. R. Clark. Faster treatment learning, Computer Science, Portland State University. Master's thesis, 2005.
11. J. Clarke, J. Dolado, M. Harman, R. Hierons, B. Jones, M. Lumkin, B. Mitchell, S. Mancoridis, K. Rees, M. Roper, and M. Shepperd. Reformulating software engineering as a search problem. *IEE Proceedings-Software*, 150(3):161–175, 2003.
12. C. Coarfa, D. D. Demopoulos, A. San, M. Aguirre, D. Subramanian, and M. Y. Vardi. Random 3-sat: The plot thickens. In *In Principles and Practice of Constraint Programming*, pages 143–159, 2000. Available from <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.42.3662>.
13. S. Cornford, M. Feather, and K. Hicks. DDP a tool for life-cycle risk management. In *IEEE Aerospace Conference, Big Sky, Montana*, pages 441–451, March 2001.
14. J. Crawford and A. Baker. Experimental results on the application of satisfiability algorithms to scheduling problems. In *AAAI '94*, 1994.
15. J. DeKleer. An Assumption-Based TMS. *Artificial Intelligence*, 28:163–196, 1986.
16. P. Denno, M. P. Steves, D. Libes, and E. J. Barkmeyer. Model-driven integration using existing models. *IEEE Software*, 20(5):59–63, Sept.-Oct. 2003.
17. P. Donzelli and G. Iazeolla. Hybrid simulation modelling of the software process. *Journal of Systems and Software*, 59(3), December 2001.
18. M. Feather, S. Cornford, K. Hicks, J. Kiper, and T. Menzies. Application of a broad-spectrum quantitative requirements model to early-lifecycle decision making. *IEEE Software*, 2008. Available from <http://menzies.us/pdf/08ddp.pdf>.
19. M. Feather, K. Hicks, R. Mackey, and S. Uckun. Guiding technology deployment decisions using a quantitative requirements analysis technique. In *IEEE International Conference on Requirements Engineering, Industrial Practice and Experience track Barcelona, Spain*, 2008.
20. M. Feather and T. Menzies. Converging on the optimal attainment of requirements. In *IEEE Joint Conference On Requirements Engineering ICRE'02 and RE'02, 9-13th September, University of Essen, Germany*, 2002. Available from <http://menzies.us/pdf/02re02.pdf>.
21. M. Feather, S. Uckun, and K. Hicks. Technology maturation of integrated system health management. In *Space Technology and Applications International Forum (STAIF-2008) Albuquerque, USA*, February 2008.
22. R. E. Filman. *Aspect-Oriented Software Development*. Addison-Wesley, Boston, 2004.
23. R. France, S. Ghosh, E. Song, and D. Kim. A metamodeling approach to pattern-based model refractoring. *IEEE Software*, 20(5):52–58, Sept.-Oct. 2003.
24. J. Gray, Y. Lin, and J. Zhang. Automating change evolution in model-driven engineering. *IEEE Computer*, 39(2):51–58, February 2006.

25. J. Gu, P. W. Purdom, J. Franco, and B. W. Wah. Algorithms for the satisfiability (sat) problem: A survey. In *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, pages 19–152. American Mathematical Society, 1997.
26. M. Hall and G. Holmes. Benchmarking attribute selection techniques for discrete class data mining. *IEEE Transactions On Knowledge And Data Engineering*, 15(6):1437–1447, 2003. Available from <http://www.cs.waikato.ac.nz/~mhall/HallHolmesTKDE.pdf>.
27. D. Harel. Statemate: A working environment for the development of complex reactive systems. *IEEE Transactions on Software Engineering*, 16(4):403–414, April 1990.
28. M. Harman. The current state and future of search based software engineering. In *Future of Software Engineering, ICSE'07*. 2007.
29. M. Harman and B. Jones. Search-based software engineering. *Journal of Information and Software Technology*, 43:833–839, December 2001.
30. M. Harman and J. Wegener. Getting results from search-based approaches to software engineering. In *ICSE '04: Proceedings of the 26th International Conference on Software Engineering*, pages 728–729, Washington, DC, USA, 2004. IEEE Computer Society.
31. P. Hart, N. Nilsson, and B. Raphael. A formal basis for the heuristic determination of minimum cost paths. *IEEE Transactions on Systems Science and Cybernetics*, 4:100–107, 1968.
32. C. Heitmeyer. Software cost reduction. In J. J. Marciniak, editor, *Encyclopedia of Software Engineering*, January 2002. Available from <http://chacs.nrl.navy.mil/publications/CHACS/2002/2002heitmeyer-encse.pdf>.
33. G. Holzmann. The model checker SPIN. *IEEE Transactions on Software Engineering*, 23(5):279–295, May 1997.
34. Y. Hui, E. Prakash, and N. Chaudhari. Game ai: artificial intelligence for 3d path finding. In *TENCON 2004. 2004 IEEE Region 10 Conference*, volume 2, pages 306–309, 2004.
35. Y. Iwasaki. Qualitative physics. In P. C. A. Barr and E. Feigenbaum, editors, *The Handbook of Artificial Intelligence*, volume 4, pages 323–413. Addison Wesley, 1989.
36. D. Jackson. Alloy: a lightweight object modelling notation. *ACM Trans. Softw. Eng. Methodol.*, 11(2):256–290, 2002.
37. O. Jalali, T. Menzies, and M. Feather. Optimizing requirements decisions with keys. In *Proceedings of the PROMISE 2008 Workshop (ICSE)*, 2008. Available from <http://menzies.us/pdf/08keys.pdf>.
38. J. Jerjens and J. Fox. Tools for model-based security engineering. In *ICSE '06: Proceeding of the 28th international conference on Software engineering*, pages 819–822, New York, NY, USA, 2006. ACM Press.
39. H. Kautz and B. Selman. Pushing the envelope: Planning, propositional logic and stochastic search. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence and the Eighth Innovative Applications of Artificial Intelligence Conference*, pages 1194–1201, Menlo Park, Aug. 4–8 1996. AAAI Press / MIT Press. Available from <http://www.cc.gatech.edu/~jimmyd/summaries/kautz1996.ps>.
40. S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science*, Number 4598, 13 May 1983, 220, 4598:671–680, 1983.
41. R. Kohavi and G. H. John. Wrappers for feature subset selection. *Artificial Intelligence*, 97(1-2):273–324, 1997.
42. A. MacLean, R. Young, V. Bellotti, and T. Moran. Questions, options and criteria: Elements of design space analysis. In T. Moran and J. Carroll, editors, *Design Rationale: Concepts, Techniques, and Use*, pages 53–106. Lawrence Erlbaum Associates, 1996.
43. R. Martin and R. D. M. Application of a hybrid process simulation model to a software development project. *Journal of Systems and Software*, 59(3), 2001.

44. R. Martin and D. M. Raffo. A model of the software development process using both continuous and discrete models. *International Journal of Software Process Improvement and Practice*, June/July 2000.
45. T. Menzies, O. Jalali, and M. Feather. Optimizing requirements decisions with keys. In *Proceedings PROMISE '08 (ICSE)*, 2008.
46. T. Menzies, J. Kiper, and M. Feather. Improved software engineering decision support through automatic argument reduction tools. In *SEDECS'2003: the 2nd International Workshop on Software Engineering Decision Support (part of SEKE2003)*, June 2003. Available from <http://menzies.us/pdf/03star1.pdf>.
47. T. Menzies and H. Singh. Many maybes mean (mostly) the same thing. In M. Madravian, editor, *Soft Computing in Software Engineering*. Springer-Verlag, 2003. Available from <http://menzies.us/pdf/03maybe.pdf>.
48. N. Metropolis, A. Rosenbluth, M. Rosenbluth, A. Teller, and E. Teller. Equation of state calculations by fast computing machines. *J. Chem. Phys.*, 21:1087–1092, 1953.
49. P. Mi and W. Scacchi. A knowledge-based environment for modeling and simulation software engineering processes. *IEEE Transactions on Knowledge and Data Engineering*, pages 283–294, September 1990.
50. H. Mittelman. Recent benchmarks of optimization software. In *22nd European Conference on Operational Research*, 2007.
51. J. Mylopoulos, L. Cheng, and E. Yu. From object-oriented to goal-oriented requirements analysis. *Communications of the ACM*, 42(1):31–37, January 1999.
52. J. Mylopoulos, L. Chung, and B. Nixon. Representing and using nonfunctional requirements: A process-oriented approach. *IEEE Transactions of Software Engineering*, 18(6):483–497, June 1992.
53. J. Pearl. *Heuristics: intelligent search strategies for computer problem solving*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1984.
54. L. Relia. Evolutionary computing in search-based software engineering. Master's thesis, Lappeenranta University of Technology, 2004.
55. S. J. Russell, P. Norvig, J. F. Candy, J. M. Malik, and D. D. Edwards. *Artificial intelligence: a modern approach*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 2003.
56. T. Saaty. *The Analytic Hierarchy Process: Planning, Priority Setting, Resource Allocation*. McGraw-Hill, 1980.
57. A. Saltelli, K. Chan, and E. Scott. *Sensitivity Analysis*. Wiley, 2000.
58. B. Selman, H. A. Kautz, and B. Cohen. Local search strategies for satisfiability testing. In M. Trick and D. S. Johnson, editors, *Proceedings of the Second DIMACS Challenge on Cliques, Coloring, and Satisfiability*, Providence RI, 1993.
59. S. Setamanit, W. Wakeland, and D. Raffo. Using simulation to evaluate global software development task allocation strategies. *Software Process: Improvement and Practice, (Forthcoming)*, 2007.
60. S. B. Shum and N. Hammond. Argumentation-based design rationale: What use at what cost? *International Journal of Human-Computer Studies*, 40(4):603–652, 1994.
61. H. Stermann. *Business Dynamics: Systems Thinking and Modeling for a Complex World*. Irwin McGraw-Hill, 2000.
62. B. Stout. Smart moves: Intelligent pathfinding. *Game Developer Magazine*, (7), 1997.
63. T. E. Uribe and M. E. Stickel. Ordered binary decision diagrams and the davis-putnam procedure. In *In Proc. of the 1st International Conference on Constraints in Computational Logics*, pages 34–49. Springer-Verlag, 1994.
64. R. Williams, C. Gomes, and B. Selman. Backdoors to typical case complexity. In *Proceedings of IJCAI 2003*, 2003. <http://www.cs.cornell.edu/gomes/FILES/backdoors.pdf>.

65. H. Zhang and X. Zhang. Comments on 'data mining static code attributes to learn defect predictors'. *IEEE Transactions on Software Engineering*, September 2007.

Obtaining our System

We have placed on-line all the materials required for other researchers to conduct further investigation into this problem. All the code, Makefiles, scripts, and so on used in this paper are available at <http://unbox.org/wisp/tags/ddpExperiment/install>. For security reasons, all the available JPL requirements models have been “sanitized”; i.e. all words replaced with anonymous variables.

Replies to reviewers

Thank you for your careful reviews on the first draft of this paper. On experimentation, we found that the premises of draft #1 was wrong. It turns out that KEYS2 scales very well to larger models ($O(N^2)$). Hence, the basis of the paper change from issues of runtime to other advantages of our method.

The issue of *exploring solution robustness* was not adequately discussed in draft #1. Here, we focus draft #2 on that issue. Nearly every page of this paper has been extensively rewritten and the whole paper has been shortened.

Comments from Reviewer 1

I felt that the preparatory sections built up lots of background, but the meat in the results section was skimmed over by comparison.

This draft has been extensively rewritten so that we get to the “meat” as soon as possible. New results have been added regarding how well these methods scale to larger models (see §6.3). In all, the conclusions of this paper are now based on three real-world requirements models from NASA and a dozen artificially generated models.

As to excessive amount of background material, much has been pruned from this draft but much more could also be deleted. For example, in the *Related Work* section, we could delete sections 7.3, 7.4, 7.5.

Fig 10 is too small to read properly.

That figure has been made larger.

I missed a clear and unequivocal explanation of what the fitness function and representation were for the Simulated Annealing

That has now been placed in its own section, especially marked “objective function” (see §4.2).

The work is written in a compelling and direct style (which occasionally is too direct and a little grating)....

.... I liked the novel style opening “The room is crowded and everyone is talking at once.”. However this may put off some readers. At present the paper is nice as a compelling read, but starting off like this really could distract the more traditional SE reader who have views about how scientific work should be presented...

The colorful writing style comes up again on p17; Uribe and Stickel “struggled valiantly”. I don’t think that this style is suitable for an academic journal. Maybe I am old fashioned. I think it should be toned down and the style more dispassionate. Did they really struggle “valiantly” anyway? It may well put readers off and that would be a pity.

That language has been removed, both in the introduction and around the discussion of Uribe and Stickel.

The authors make available their code which is helpful. It is not clear whether real world requirements data from JPL is also included in the on line provision. I would like some clarification on this issue in the revised version. Such provision would be a great help to follow on research, though I realize that for reasons of commercial confidentiality this may not be possible.

Appendix A now details how to obtain this material. For security reasons, all the available JPL requirements models have been “sanitized”; i.e. all words replaced with anonymous variables.

The authors statement:- " This is an exciting result since it means that current results from simulated annealing(e.g.[5,8,10,62,73]) could be greatly improved, just by switching to an alternate search engine. " really needs to be re-written.

This new draft does not include those over-inflated claims about the merits of KEYS/ KEYS2/ and simulated annealing.

Other work on SBSE for requirements that is clearly missing and needs to be included is:

- *The paper by Gunther Ruhe at FSE 08*
- *The Zhang et al (GECCO 2007) and*
- *the work by Bagnall (I&ST 2001).*

We cannot find the Ruhe paper at the FSE'08 web site⁶ or on his home pages (perhaps it was a last minute replacement for a keynote?). However, his 2009 TSE paper is clearly relevant so we used that instead. A curious omission in that paper is the runtimes of the GAs used in that paper. While this does not invalidate that paper, it does mean that we cannot make an extensive comment on that work to this paper. In any case, thanks to this review, we added a reference to his work in our section "Other Optimizers" section (see §7.2).

As to the other work on release planning, we added in notes on Ruhe, Zhang, and Bagnall's work in Related Work (see §7.1).

Some of the citations are messed up too. For example citation [5] "S M B and S. Mancoridis". S M B is Brian Mitchell I believe. Clearly some more scholarly care is required here, both to chase up relevant work and to ensure that the citations are present correctly.

Thanks for catching that- the culprit was two commas together in a .bib file. That has been fixed and the other references checked.

Oh- and that particular paper did not not make it into this draft.

Does Fig 4 add sufficient value to justify a whole page?

⁶ <https://cgi4.cc.gatech.edu/phps/conferences/fse16/program.php>

(FYI: that figure has now moved to page 6). That is, of course, for you to decide. However, we often get questions along the lines of “aren’t these models so simple you can just see the problems/solutions?”. This figure shows just how complex these models can get. We can remove that figure if you feel that we have made enough of a case in the text that the JPL are hard enough to be interesting.

In related work, some recent work by Finkelstein on multi objective requirements models for fairness is missing (this was published in RE 08)

In RE08, we found:

- Revisiting the Core Ontology and Problem in Requirements Engineering Jureta, I.J.; Mylopoulos, J.; Faulkner, S.; International Requirements Engineering, 2008. RE '08. 16th IEEE

Comments on this have been added to the “Related Work” section, see §7.4.

Comments from Reviewer 2

The novelty of the work remains unclear. The paper presents little other (of real substance) than a small performance measurement experiment, results of which are interpreted as supporting the appropriateness of the algorithm for use in automated requirements engineering.

This is a fair comment. The previous draft was poorly written and obscured the novelty of the paper.

To address this problem, this draft has been significantly reorganized and the relevant material moved to the front (and supporting material moved to a related work section at the back of the article). Consequently the novel features of this paper are not obscured:

- A completely naive search through the space of options is exponential on the number of options. We can reduce that search to linear time using *decision ordering diagrams*, described in §4.3.
- Better yet, we can generate those diagrams in low-order polynomial time using a Bayesian sampling method called KEYS2.

*The title and abstract and indeed the overall narrative of this paper are misleading. This is a paper on a small-scale runtime performance study *masquerading* as a paper on requirements analysis.*

We thank this reviewer for this comment- it inspired us to conduct the scale up study of §6.3.

We agree that the previous draft lacked insufficient analysis to merit journal publication. To address this concern, we have changed the title of the paper; changed the focus of the paper; and clarified and extended our analysis.

We would also comment that this *is* a real requirements analysis paper. In fact, this new draft’s empirical authority is *greater than* the standard

requirements engineering paper. Here, we explore three real world requirements models (from NASA) as well as a dozen artificially generated models. By contrast, a standard paper in this field uses a single small (possibly “toy”) case study.

The writing in this paper is somewhere between terrible and just plain poor: at the level of sentences, paragraphs, and overall narrative structure..... Opening scenario isn't credible. A decision as substantial as a redesign of the file system would not be made without documenting it in a way that would make everyone with an interest in that aspect of the aware of the decision....

This reviewer is quite correct in complaining about the quality of the previous draft. This draft has been extensively reorganized. The bogus opening scenario has been deleted. The core of the paper has been made clearer. The paper has been extensively proof read prior to re-submission to ASE.

Critical but unexplained and unsupported assumption: Our proposed solution to real-time requirements optimization assumes that the behavior of a large system is determined by a very small number of key variables.

This assumption has much supporting evidence from other domains (see §7.6). Also, our empirical evidence supports it (see the comment at bottom of p20).

Comments from Reviewer 3

The whole discussion of the “real time requirements optimization” problem, while a nice idea and interesting, falls flat by being bogus.

We agree. When we experimented with this scale up runtime costs using our models (and not using some bogus equation), we found the scale up factors to be far *less* than anything suggested by our formulae. So this draft has dropped that whole discussion. Instead, we focus on an issue that was poorly discussed in the previous draft: fast methods for assessing solution robustness.

First paragraph of Sec 7: Why? why is it “unfair” to compare BDD-based against search-based? Aren't they solving the same problem?

We agree this is a strange concept- and one we did not believe three years ago. But since then we have much experience at empirical software engineering conferences where authors of method1 have cried “foul!” when authors have compared model1 against model2 using a particular problem that, by definition, is not suitable for method1. The classic example is NuSMV vs SPIN: the BDDs of NuSMV make it a natural choice for models with extensive symmetry and repeated structures. Given that you know that, should you compare NuSMV and SPIN

using such models? Or are they the *last* thing you should use for a comparison?

The most telling reply we can offer to this comment is this:

- Holzmänn does not compare SPIN against NuSMV.
- Uribe and Stickel take great care when selecting comparative cases- and when the algorithms are too distant, they don't even do that.
- Holzmänn, Uribe, and Stickel have been building and assessing algorithms for a very long time. Of course, we will follow the guidance of the ASE reviewers but might it be safest to respect their experience?

Sec 7.2: what is the variable "C"?

(FYI: in the new draft, this section is appendix B.2.) "C" is "CHOSEN"
(fixed in this draft)

I'm not sure it is valid to make claims about KEYS2 performance versus MaxWalkSAT when you are actually comparing to MaxFunWalk. You have to do a better job showing how they are related. (E.g. at what level are they the same algorithm?)

Thanks for this comment- it has allowed us to clarify a confusing part of our text. MaxWalkSat is defined over CNF formulae. MaxWalkFun assumes that some function can assess a collection of variable settings. MaxWalkFun becomes MaxWalkSat when that assessment function is a logical truth table. Here, we use our DDP models as the assessment function (these notes have been added to the bottom of page 12).

Last para of Sec 9. "...some requirements models". This comes out of the blue. For which reqts models does KEYS approach fail?

This section has been removed. As to the ultimate limits of KEYS2, we don't know yet. However, given the simplicity of the algorithm's implementation and its fast runtimes, we offer the following advice: try it and see. FYI, we are using it *all* our current projects. We have found it is a good stochastic rule generator for standard UCI data and that it out-performs standard non-linear optimizers for optimizing complex FORTRAN simulators of NASA flight systems. We are also successfully applying it to effort estimation and defect prediction. All those papers are Masters projects, currently in progress (i.e. no tech reports, or publications, yet). But in summary, KEYS2 is remarkably general and remarkably useful.