# Wickedly fast Explanation Generation for Multi-Objective Optimization

Authors suppressed for blind review

Institution suppressed for blind review

**Abstract.** WICKED is a near linear-time algorithm for summarizing trade-offs in multi-objective problems. From that summary, humans can read recommendations for their systems. This paper evaluates those recommendations using data generated from (a) the POM3 model of agile selection of tasks; (b) four COCOMO-suite predictors for software development effort, months, defects and risk.

WICKED runs orders of magnitude faster than standard optimizers (NSGA-II and SPEA2). For example, for one of our larger models, WIcKET and NSGA-II terminated in 3 and 150 seconds. Further, WICKED's recommendations were just as effective at improving objective scores as the actions of standard optimizers. Hence, we recommend WICKED when some succinct summary has to be rapidly generated (e.g. in some interactive design meeting). WICKED could also be useful as post-processor to other optimizers (to generate succinct explanations of their conclusions) or as a optimizer to other optimizers (by constraining those other optimizers to only search the regions recommended by WICKED).

**Keywords:** Software engineering, explanation, optimization, multi-objective.

## 1  Introduction

> *"If you cannot- in the long run- tell everyone what you have been doing, your doing has been worthless."*
> – Erwin Schrödinger

Explaining the results of multi-objective optimization to a user can be problematic. A typical run of a multi-objective optimizer can process thousands to millions of examples. It is an overwhelming task for humans to certify the correctness of conclusions generated from so many results. Verrappa and Leiter warn that

> "..for industrial problems, these algorithms generate (many) solutions, which makes the tasks of understanding them and selecting one among them difficult and time consuming" [**?**].

Even if explanations are constrained to (say) just a few hundred examples taken from the Pareto frontier, this can still confuse the user. Valerdi notes that it can take days for panels of human experts to rigorously review even a few dozen examples [60]. For example, once had a client who disputed the results of our analysis. They demanded to audit the reasoning but when we delivered the of candidate solutions on the Pareto frontier, they were overwhelmed by the amount of information. Flustered, the client discounted the analysis and rejected our conclusions. From this experience, we learned that to better support decision making in SBSE, we must better explain SBSE results.

Other researchers have recognized the importance of explanation. It is known to be a key factor in selecting algorithms. For example, in the field of machine learning, "each time one of our favorite approaches has been applied in industry, each time the comprehensibility of the results, though ill-defined, has been a decisive factor of choice over an approach by pure statistical means, or by neural networks." [7]. Analogous terms to explainability in that community are "comprehensibility", "interpretability" [1] or "understandability" [4].

In spite of the importance attributed to the subject, explanation has not been extensively investigated in the context of SBSE. One of the few papers that does is that of Veerappa and Lieter [62] who clustered examples from the Pareto frontier (examples generated from a goal graph representation of requirements for London ambulance services). In this approach, "instead of having to inspect a large number of individual solutions, (users) can look at a much smaller number of groups of related solutions, and focus their attention on the important characteristics of the group rather than the particularities of their individual solutions" [62].

XXXX after creating. still errors in comparisons

While an innovative and insightful study, there are three open issues with that method: (a) the complexity of clustering; (b) erroneous conclusions could be generated from the users inspection of the clusters; (c) introduced by users incorrectly evaluation the value of generated recommendations. Veerappa and Lieter did not evaluate the effects of the recommendations that could be generated by users browsing their clusters. Also, their method could suffer from scalability issues since it a post-processor to a clustering algorithm (clustering is a slow process requiring say, $O(N^2)$ comparisons for the greedy agglomerate clustering algorithm used in that paper [32]).

Accordingly, in this paper, when:

1. Cluster using a near-linear time algorithm;
2. Better define the process by which recommendations are generated from clusters;
3. The generated recommendations are tested by generating more examples from the model *after* the recommendations are imposed as extra constraints on the model inputs.

XXXX One result from this work that may surprise the reader is that our "explanation" system is anctually an inference procedure as well. While our approach could be used as an add-on to other MOEA systems (to augment their reasoning with post-hoc explanations), it actually is a viable MOEA in its own right. Research results from cognitive science explain why this is so: according to Leake [**?**], explanation is actually an inference procedure in its own right where different audiences "explain" some phenonomeom in different ways according to their own goals and background knowledge

## 2   A Motivating Example

### 2.1   Our Problem

To motivate this work, we offer the real-world goal that sparked this work.

The Software Engineering Institute (SEI: http://www.sei.cmu.edu/) at Carnegie Mellon University is a centralized repository for qualitative and quantitative data collected from software development for the United States government and Department of Defense. Developers across the country looked to SEI for explanations of what factors effect their project (these explanations are used to manage their current projects and well as propose methods on how to better handle their future projects in a better man-

ner). Also, large scale policy decisions about information technology are made by the U.S. government, partially in consultation with researchers at the SEI.

It is standard for the SEI to issues "fact sheets" explaining their lessons learned as well as explaining their advise on best practice. These explanations are short reports (rarely more than both sides of one piece of paper) which are intended to give busy managers quick guidance for their projects. In the case of quantitative data, these explanations may contain some 2D plot showing how one objective (e.g. defects) changes in response to changes in one input variable (e.g. lines of code).

The veracity of these simple plots of SE data is questionable. Even for single goal reasoning such as defect reduction, there can only be poorly characterized via one input variable. Menzies et al. compared learners building models with $N = 1$ of $N > 1$ input variables: the models that used more than one input performed better [42].

Further, there are many recent SE research publications that propose multiple competing goals for SE models; e.g.

– Build software *faster* using *less* effort with *fewer* bugs [20];
– Return defect predictors that find *most* defects in the *smallest* parts of the code [6].

As we move from single goal to multiple-goal reasoning, the value of simplistic plots to explain effects in SE projects becomes even more questionable. The SBSE experience is that reasoning and trading off between multiple goals is much more complex that browsing effects related to a single isolated goal.

Worse yet, given all the context variables that can be used to describe different software projects, it is hard to believe that any *one* report can explain *all* the effects seen in all different kinds of software projects. Numerous recent reports in empirical SE offer the same *locality effect*; i.e. models generated using *all* available data perform differently, and often worse, than those learned from specific subsets [8, 40, 45, 55, 63].

## 2.2  Our Solution

New results suggest a resolution to the above problems. Firstly, work on *low dimensional approximations* of SE data shows we do not need to reason about all the context variables (since many are redundant or noisy). When applied to SBSE, low dimensional approximations can guide mutation strategies to find better solutions one to two orders of magnitude faster than standard evolutionary optimizers [34, 35].

Secondly, work on MOEA/D (multi-objective evolutionary algorithms with decomposition [65]) has shown the benefits of dividing problems into multiple cells, then optimizing each cell separately. Note that this approach is analogous (and actually predates) the locality work mentioned above.

The rest of this paper reports an experiment where we generate explanations of the forces that impact a software project by combining *decomposition* with *low dimensional approximations*. Our system is called WICKED and has six parts:

W: **W**HERE is Menzies' near-linear time clustering algorithms [40] that uses lower dimensional approximations to decompose data from models of SE processes into many small clusters,

I: **I**NFOGAIN [22] finds what numeric ranges best predict for the different clusters.

C: To find rules that distinguish the clusters, we apply Brieman's **C**ART, algorithm [11] to the data simplified by INFOGAIN.

K: **K**ILL prunes spurious leaf branches generated by CART, thus shrinking the tree. KILL outputs branches (conjunctions of *attribute=value* pairs) that lead to clusters $C_0, C_1, C_2, ...$

E: **E**NVY looks at all clusters $C_i$ trying to find a nearby cluster $C_j$ with better objective scores.

D: **D**ELTA is a contrast-set learner that reports the difference branches ending at $C_i$ and $C_j$.

It is reasonable to ask why these six particular parts, and not six other. All the above are motivated by the need for generation explanations, and keeping those explanations succinct. Kelly's Personnel construct theory [30] conjectures that humans do not explain of the world by studying all factors: rather, they find differences between things. ENVY and DELTA are our method of implementing Kelly's insight.

As to the other parts of WICKED, WHERE's dimensionality reduction lets us ignore spurious dimensions while we group the data. Once WHERE terminates, then the only ranges that are interesting are those that distinguish between the clusters (and these are found by INFOGAIN). INFOGAIN is a pre-processor to decision tree learning and KILL is a post-processor. The combination of these pre-post-processors means that the trees found by CART are very small. This, in turn, means that the contrast found by ENVY+DELTA are very succinct.

The output of WICKED is *recommendation* for how we would improve all the examples in $C_i$. Expressed in terms of evolutionary algorithms, *recommendation = mutation*; i.e. they are the suggestion on how to change examples in order to optimize multiple objectives. Return now to writing fact sheets from the SEI, we note that:

- INFOGAIN+CART+KILL generates very small trees (10 to 20 lines);
- ENVY and DELTA are very simple algorithms that can be applied by humans while manually browsing CT0's trees;
- The experiments shown below demonstrate that DELTA's recommendation are just as effective as the mutations proposed by standard optimizers (NSGA-II [17] and SPEA2 citezit02) for adjusting a population in order to improve multiple objectives.
- All of WICKED's sub-routines are near-linear time algorithms that terminate 100 times faster than standard optimizers.

Hence, we propose a modification to SEI's fact sheet: they should old the trees generated by WHERE+INFOGAIN+CART+KILL. From these, business users can apply ENVY+DELTA in order to explore for themselves the effects in SEI's data.

```
if cplx ≤ 1.1:                        Effort  Months   Defect  Risk
.. if resl ≤ 3.2:
.. .. if pvol ≤ 1.1:
.. .. .. if site ≤ 1.0:
.. .. .. ..  then: ['__25']  #             ?       ?        ?      ?
.. .. .. ..  else: ['__20']  #        d    ?       ?        ?      ?
.. .. if rely ≤ 1.0:
.. .. .. if pcap ≤ 0.9:
.. .. .. .. if ltex ≤ 1.0:
.. .. .. .. .. if site ≤ 1.0:
.. .. .. .. .. ..  then: ['__3']  #   a    ?       ?        ?      ?
.. .. .. .. .. ..  else: ['__1']  #   c    ?       ?        ?      ?
.. .. .. .. .. else: ['__6']  #       b    ?       ?        ?      ?
.. .. .. .. if ltex ≤ 1.0:
.. .. .. .. .. if pmat ≤ 3.9:
.. .. .. .. .. ..  then: ['__15']  #       ?       ?        ?      ?
.. .. .. .. .. else: ['__15']  #      e    ?       ?        ?      ?
```

## 3   Models

We have tested WICKED on numerous MOEA tasks including the standard laboratory problems (DTLZ, Schaffer, Fonseca, etc) and found it recommendations generated instances with objective scores competitive with those generated by NSGA-II or SPEA2.

Results from those standard lab problems are rarely convincing or interesting to software project managers. Hence, we show results from two business-level process models. POM3 [37, 54] implements the Boehm and Turner model [**?**, **?**, 54] of agile programming where teams select tasks as they appear in the scrum backlog. POM3 studies the implications of different ways to adjust task lists in the face of shifting priorities. XOMO [41, 47, 48] is four software process models from the University of Southern California. XOMO reports four-objective scores (which we will try to minimize): project *risk*; development *effort* and *defects*; and total *months* of development.

### 3.1   POM3

Turner and Boehm say that the agile management challenge is to strike a balance between the three objectives of *completion rates*, *idle rates*, and *overall cost* of a project. In the agile world, projects terminate after achieving a *completion rate* of $(X < 100)\%$ of its required tasks. Team members become *idle* if forced to wait for a yet-to-be-finished task from other teams. To lower *idle rate* and increase *completion rate*, management can hire staff- but this can increase *overall cost*.

When optimizing POM3, we seek changes to the controllables of Figure 1 that maximize *completion rate* while minimizing *cost* and *idle*ness. To make this task more realistic, we run POM for three different kinds of software projects, denoted POM3a, POM3b, POM3c shown in Figure 2. We make no claim that these three projects cover X% of all software projects– rather, our point here is that POM3 can handle different kinds of models.

### 3.2   XOMO

The XOMO model enables an exploration of competing factors within software projects. Ideally, management decisions can minimize all of *months*, *effort*, *defects* and *risk*. However, there are many trade-offs to be considered. For example. increasing software reliability *reduces* the number of added defects while *increasing* the software development effort. For another example, better documentation can improve team

| scale factors (exponentially decrease effort) | prec: have we done this before? flex: development flexibility resl: any risk resolution activities? team: team cohesion pmat: process maturity |
|---|---|
| upper (linearly decrease effort) | acap: analyst capability pcap: programmer capability pcon: programmer continuity aexp: analyst experience pexp: programmer experience ltex: language and tool experience tool: tool use site: multiple site development sced: length of schedule |
| lower (linearly increase effort) | rely: required reliability data: secondary memory storage requirements cplx: program complexity ruse: software reuse docu: documentation requirements time: runtime pressure stor: main memory requirements pvol: platform volatility |

Fig. 3: XOMO model decisions.

| Short name | Decision | Description | Controllable |
|---|---|---|---|
| Cult | Culture | Number (%) of requirements that change. | yes |
| Crit | Criticality | Requirements cost effect for safety critical systems. | yes |
| Crit.Mod | Criticality Modifier | Number of (%) teams affected by criticality. | yes |
| Init. Kn | Initial Known | Number of (%) initially known requirements. | no |
| Inter-D | Inter-Dependency | Number of (%) requirements that have interdependencies. Note that dependencies are requirements within the *same* tree (of requirements), but interdependencies are requirements that live in *different* trees. | no |
| Dyna | Dynamism | Rate of how often new requirements are made. | yes |
| Size | Size | Number of base requirements in the project. | no |
| Plan | Plan | Prioritization Strategy (of requirements): 0= Cost Ascending; 1= Cost Descending; 2= Value Ascending; 3= Value Descending; 4 = $\frac{Cost}{Value}$ Ascending. | yes |
| T.Size | Team Size | Number of personnel in each team | yes |

Fig. 1: List of Decisions used in POM3 (optimizers tune controllables, on right).

| | POM3a<br>A broad space<br>of projects. | POM3b<br>Highly critical<br>small projects | POM3c<br>Highly dynamic<br>large projects |
|---:|---|---|---|
| Culture | $0.10 \leq x \leq 0.90$ | $0.10 \leq x \leq 0.90$ | $0.50 \leq x \leq 0.90$ |
| Criticality | $0.82 \leq x \leq 1.26$ | $0.82 \leq x \leq 1.26$ | $0.82 \leq x \leq 1.26$ |
| Criticality Modifier | $0.02 \leq x \leq 0.10$ | $0.80 \leq x \leq 0.95$ | $0.02 \leq x \leq 0.08$ |
| Initial Known | $0.40 \leq x \leq 0.70$ | $0.40 \leq x \leq 0.70$ | $0.20 \leq x \leq 0.50$ |
| Inter-Dependency | $0.0 \leq x \leq 1.0$ | $0.0 \leq x \leq 1.0$ | $0.0 \leq x \leq 50.0$ |
| Dynamism | $1.0 \leq x \leq 50.0$ | $1.0 \leq x \leq 50.0$ | $40.0 \leq x \leq 50.0$ |
| Size | $x \in [3,10,30,100,300]$ | $x \in [3, 10, 30]$ | $x \in [30, 100, 300]$ |
| Team Size | $1.0 \leq x \leq 44.0$ | $1.0 \leq x \leq 44.0$ | $20.0 \leq x \leq 44.0$ |
| Plan | $0 \leq x \leq 4$ | $0 \leq x \leq 4$ | $0 \leq x \leq 4$ |

Fig. 2: Three classes of projects studied using POM3.

| project | | low | high | | project | | low | high | | project | | low | high |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| FL: | rely | 3 | 5 | | GR: | rely | 1 | 42 | | O2: | prec | 3 | 5 |
| JPL flight | data | 2 | 3 | | JPL ground | data | 2 | 3 | | Orbital Space | pmat | 4 | 5 |
| software | cplx | 3 | 6 | | software | cplx | 1 | 4 | | Place guidance | docu | 3 | 4 |
| | time | 3 | 4 | | | time | 3 | 4 | | navigation and | ltex | 2 | 5 |
| | stor | 3 | 4 | | | stor | 3 | 4 | | control (v2) | sced | 2 | 4 |
| | acap | 3 | 5 | | | acap | 3 | 5 | | | flex | 3 | 3 |
| | apex | 2 | 5 | | | apex | 2 | 5 | | | resl | 4 | 4 |
| | pcap | 3 | 5 | | | pcap | 3 | 5 | | | time | 3 | 3 |
| | plex | 1 | 4 | | | plex | 1 | 4 | | | stor | 3 | 3 |
| | ltex | 1 | 4 | | | ltex | 1 | 4 | | | data | 4 | 4 |
| | pmat | 2 | 3 | | | pmat | 2 | 3 | | | pvol | 3 | 3 |
| | tool | 2 | 2 | | | tool | 2 | 2 | | | reuse | 4 | 4 |
| | sced | 3 | 3 | | | sced | 3 | 3 | | | ... | | |
| | KSLOC | 7 | 418 | | | KSLOC | 11 | 392 | | | KSLOC | 75 | 125 |

Fig. 4: Three case studies used in XOMO.

communication and *decrease* the number of introduced defects. However, such increased documentation *increases* the development effort.

To explore those trade offs, XOMO uses the inputs of Figure 3 to drive four models. The *effort* model predicts for "development months" where one month is 152 work hours by one developer (and includes development and management hours):

$$effort = a \prod_i EM_i * KLOC^{b+0.01 \sum_j SF_j} \tag{1}$$

Here, *EM,SF* denote the effort multipliers and scale factors and $a, b$ are the *local calibration* parameters which in COCOMO-II have default values of 2.94 and 0.91.

The XOMO *defect* model [10] assumes that certain variable settings *add* defects while others may *subtract* (and the final defect count is the number of additions, less the number of subtractions).

Also the *months* model predicts for total development time and can be used to determine staffing levels for a software project. For example, if *effort*=200 and *months*=10, then this project needs $\frac{200}{10} = 20$ developers.

Lastly, the *risk* model comments on how sets of management decisions decrease the odds of successfully completing a project. For example suppose a manager demands *more* reliability (*rely*) while *decreasing* analyst capability (*acap*). Such a project is "risky" since it means the manager is demanding more reliability from less skilled analysts. The XOMO *risk* model contains dozens of rules that trigger on each such "risky" combinations of decisions [39].

In the following, we run three different case studies of XOMO: one for each of the NASA projects described in Figure 4.

## 4   Code

This section describes WHERE+CART+ENVY+DELTA.

### 4.1   WHERE

WHERE inputs a set of $N$ examples, each of which is a set of decisions $D$ mapped to a set of objectives $O$, so $N_i = (D, O)$ (and usually $D > 1$ and $O > 1$ and $O < D$). WHERE clusters the examples on the decisions and reports the average objective scores for each objective in each cluster.

WHERE uses a dimensionality reduction heuristic proposed by Faloutsos and Lin [21]. The method inputs $N$ examples $N_1, N_2, ...$ Next, WHERE picks any point $N_i$ at random. Thirdly, WHERE finds the point $West \in N$ that is furthest[1] from $N_i$. Finally, WHERE finds the point $East \in N$ that is furthest from $West$ (and $c = dist(West, East)$).

To recursively cluster the data, WHERE iterates over $N_i \in N$ to find $a = dist(N_i, West)$, $b = dist(N_i, East)$, $x = (a^2 + c^2 - b^2)/(2c)$. This $x$ value is the projection of $N_i$ on the line running $East$ to $West$. WHERE divides the examples on the median $x$ value, then recurses on each half. Recursion on $N$ initial examples stops when a sub-region contains less that $M$ examples (e.g. $M = \sqrt{N}$).

Note that this four-step process requires only $2N$ distance comparisons per level of recursion and one call to a sorting routine to find the median value. The total time for WHERE is some linear multiple of the sorting time used to find the median at each level. Assuming sorting takes time $O(NlogN)$, then we can say that WHERE runs in near linear time (and not the $O(N^2)$ required for other clustering algorithms such as K-Means [?]).

### 4.2   CART

### 4.3   ENVY

### 4.4   DELTA

## 5   Methods

This study ranks methods using the Scott-Knott procedure recommended by Mittas & Angelis in their 2013 IEEE TSE paper [49]. This method sorts a list of $l$ treatments with $ls$ measurements by their median score. It then splits $l$ into sub-lists $m, n$ in order to maximize the expected value of differences in the observed performances before and after divisions. E.g. for lists $l, m, n$ of size $ls, ms, ns$ where $l = m \cup n$:

$$E(\Delta) = \frac{ms}{ls} abs(m.\mu - l.\mu)^2 + \frac{ns}{ls} abs(n.\mu - l.\mu)^2$$

Scott-Knott then applies some statistical hypothesis test $H$ to check if $m, n$ are significantly different. If so, Scott-Knott then recurses on each division. For example, consider the following data collected under different treatments *rx*:

```
rx1 = [0.34, 0.49, 0.51, 0.6]
rx2 = [0.6,  0.7,  0.8,  0.9]
rx3 = [0.15, 0.25, 0.4,  0.35]
rx4= [0.6,  0.7,  0.8,  0.9]
rx5= [0.1,  0.2,  0.3,  0.4]
```

---

[1] For this work, we use the standard Euclidean measure recommended for instance-based reasoning by Aha et al. [3]; i.e. $\sqrt{\sum_i (x_i - y_i)^2}$ where $x_i, y_i$ are values normalized 0..1 for the range min..max.

After sorting and division, Scott-Knott declares:

- Ranked #1 is rx5 with median= 0.25
- Ranked #1 is rx3 with median= 0.3
- Ranked #2 is rx1 with median= 0.5
- Ranked #3 is rx2 with median= 0.75
- Ranked #3 is rx4 with median= 0.75

Note that Scott-Knott found little difference between rx5 and rx3. Hence, they have the same rank, even though their medians differ.

Scott-Knott is preferred to, say, hypothesis testing over all-pairs of methods[2]. To avoid an all-pairs comparison, Scott-Knott only calls on hypothesis tests *after* it has found splits that maximize the perfromance differences.

For this study, our hypothesis test $H$ was a conjunction of the A12 effect size test of and non-parametric bootstrap sampling; i.e. our Scott-Knott divided the data if *both* bootstrapping and an effect size test agreed that the division was statistically significant (99% confidence) and not a "small" effect ($A12 \geq 0.6$).

For a justification of the use of non-parametric bootstrapping, see Efron & Tibshirani [19, p220-223]. For a justification of the use of effect size tests see Shepperd&MacDonell [59] and Kampenes [29]. These researchers warn that even if an hypothesis test declares two populations to be "significantly" different, then that result is misleading if the "effect size" is very small[3]. Hence, to assess the performance differences we first must rule out small effects. Vargha and Delaney's non-parametric A12 effect size test explores two lists $M$ and $N$ of size $m$ and $n$. The counter $A12$ is incremented $\forall x \in M, y \in N$ as follows:

- If $x > y$ then add $1/(mn)$;
- If $x = y$ then add $0.5/(mn)$.

A12 reports the probability that numbers in one sample are bigger than in another. The A12 thresholds for "small,medium,large" effect are $\{0.56, 0.64, 0.71\}$ respectively where "small" is a euphemism for trivial or negligible effect. This test was recently endorsed by Arcuri and Briand at ICSE'11 [5].

## 6   Results

The following results come from a standard Python 2.7 interpreter (not PyPy) running on a 2.6 GHz Mac Os/X with 4 GB of ram. For NSGA-II, we used the out-of-the-box version from DEAP, https://github.com/DEAP/deap.

In the following, we compared WICKED's results with that of NSGA-II [**?**]. NSGA-II is a genetic algorithm (GA) with a highly optimized *select* operator:

- Each generation builds generation $G + 1$ by *selecting* better individuals, *combining* some of their parts, then *mutating* the results (a little).
- NSGA-II is a GA whose *select* operator uses a non-dominating sort procedure to divide the solutions into *bands* where $band_i$ dominates all of the solutions in $band_{j>i}$ (and NSGA-II favors the least-crowded solutions in the better bands).

One reason to favor NSGA-II as our comparison optimizer is *repeatability*. Many multi-objective optimizers as MOEA/D [64] and PSO [**?**] are really *frameworks* within which

---

[2] e.g. Six treatments can be compared $(6^2 - 6)/2 = 15$ ways. A 95% confidence test run 15 times total confidence $0.95^{15} = 46\%$.

[3] For example, Kocaguenli et al. [33] report on the misleading results of such hypothesis tests in software defect prediction (due to small size of the effect being explored).

an engineer has free reign to make numerous decisions (evidence: review papers list dozens of variants on PSO and MOEA/D  [**?**, 16]). Hence, in terms of *repeatability*, it can better to use precisely defined algorithms like NSGA-II.

Other reasons to use NSGA-II are that (a) it is very widely used and (b) there is no clear consensus that some other algorithm is better. When selecting a comparison algorithm, we reached out to our SBSE colleagues to find which algorithms are accepted as "best". However, no consensus was found. On the other hand, it can be shown that NSGA-II is widely used. In 2013, Sayyad and Ammar [57] surveyed 36 SBSE papers where $\frac{21}{36}$ used NSGA-II (of the others, 4 used some home-brew genetic algorithm and the remainder each used some MOEA not used by any other paper).

### 6.1   Runtimes

Standard MOEAs require at least $N^2$ comparisons between $N$ candidates, for each generation $G$ of the evolution. In theory, WICKED is much faster than that:

- The current implementation of WICKED uses $G = 1$ since its recommendations are the results of one analysis of the data (in future work, we plan to explore an iterative evolutionary version of the algorithm).
- All the sub-routines of WICKED take near linear-time (the slowest is CART that must sort all examples at each level of its trees).

On experimentation, WICKED's runtimes are consistent with its theoretical properties. Figure 5 show the effect on runtimes of increasing the initial population size for WICKED and NSGA-II. The solid lines denote WICKED's performance: note that they are always less that those seen with NSGA-II. The effect that WICKED runs faster than standard optimizers, is most pronounced in the more complex models. In the POM3 results, some of the POM3 variants take 100s of seconds to terminate. The same problems are handled by WICKER in under one minute.
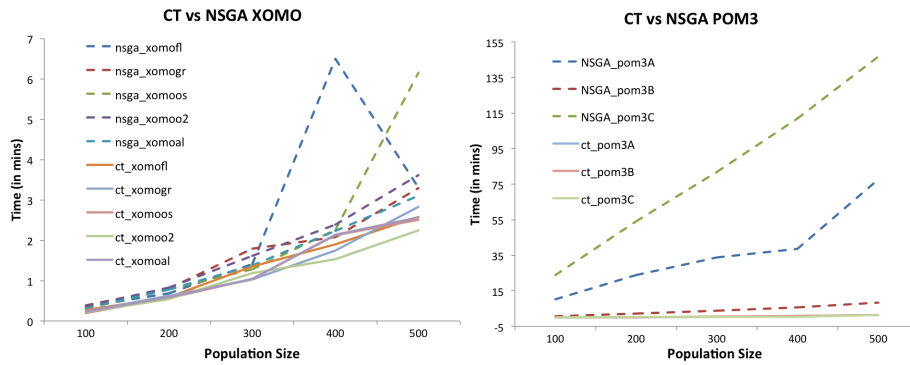


Fig. 5: Runtimes (minutes) for CT and NSGA II on POM Model (means over 20 repeats)

Note that these runtime results come from an optimized version of WICKED. Experiments are on-going with the Python profiler (to remove runtime bottlenecks). While those initial results are promising, we have nothing definitive to report at this time.

## 6.2 Optimization Improvements

To explore optimization improvements, we:
1. Collected *baseline* distributions seen in the objectives of the initial population (of 25 randomly generated individuals).
2. Using the baseline as generation $G = 1$, run NSGA-II until no improvement in any objective for three generations;
3. Using that baseline, run WICKED once. For each cluster:
   - Access the cluster items and the recommendation from the cluster;
   - Re-run the model that generated the data using constraints generated from that cluster;
4. Collected *treated* distributions from the output of steps two and three.

   Figure **??** shows the *baseline* and *treated* distributions for:
   - For all the objectives of:
     - The three variants of POM3 shown in Figure 2;
     - The three variants of XOMO shown in Figure 4;

That figure presents displays results from 20 repeated runs as horizontal quartile charts. In that figure, black dots denote median values and horizontal lines denote the 25 to 75th percentile range. To simplify readability, for each objective, all results and normalized 0..100 for the min to max values seen for that objective. Our three treatments are shown in the "Rx" column: "0" denotes the baseline; "W" denotes WICKED, and "N" denotes NSGA-II.

In all these results, *lower* values are *better* (exception: the *completion* goal in POM3 which we seek to *maximize*).

## 7   Explaining "Explanation"

As a starting point in this exploration of explanation, it is important to distinguish between the (1) problem of explaining the output of an multi-objective optimizer (discussed in this paper); from the more complex problem of (2) explaining how that output was generated. To put that another way: we seek to explain eggs, but not the chicken.

Next, a definition of "explanation" is required such that:
1. An explanation system can be designed;
2. It is possible to distinguish a "good" for a "bad" explanation.

In the SE literature, the general consensus in software engineering is that "good" explanations are succinct explanations [7, 18, 23]. On this score, MOEAs fare poorly since their output can be very verbose (hundreds or more examples from the Pareto frontier).

Also, cognitive science theory argues that there is more to "explaining" something that just showing it succinctly. According to Kelly's personal construct theory (PCT) humans explain things via "constructs" that distinguish sets of examples [30]. So, for Kelly, human explanations are not about "things" in isolation but rather the *differences between groups of things*. In data mining, finding differences between things is called *contrast set learning* [51].

Other cognitive science research studied the activities humans do during explanation generation. Leake [36] lists a dozen different tasks that humans perform when "explaining" some phenomena. Leake does not claim that the following list is complete; just that it demonstrates a wide range of goal-based purposes for explanation, including:
1. Connect event to expected/believed conditions.
2. Connect event to previously unexpected conditions.

| Fl | Rx | median | |
|---|---|---|---|
| effort | 0 | 19 | |
| | W | 4 | |
| | N | 11 | |
| months | 0 | 27 | |
| | W | 19 | |
| | N | 36 | |
| defects | 0 | 22 | |
| | W | 2 | |
| | N | 5 | |
| risks | 0 | 87 | |
| | W | 2 | |
| | N | 6 | |
| **GR** | **Rx** | **median** | |
| effort | 0 | 7 | |
| | W | 3 | |
| | N | 13 | |
| months | 0 | 21 | |
| | W | 23 | |
| | N | 39 | |
| defects | 0 | 16 | |
| | W | 1 | |
| | N | 4 | |
| risks | 0 | 74 | |
| | W | 6 | |
| | N | 5 | |
| **02** | **Rx** | **median** | |
| effort | 0 | 5 | |
| | W | 7 | |
| | N | 13 | |
| months | 0 | 25 | |
| | W | 20 | |
| | N | 43 | |
| defects | 0 | 37 | |
| | W | 1 | |
| | N | 5 | |
| risks | 0 | 79 | |
| | W | 13 | |
| | N | 6 | |

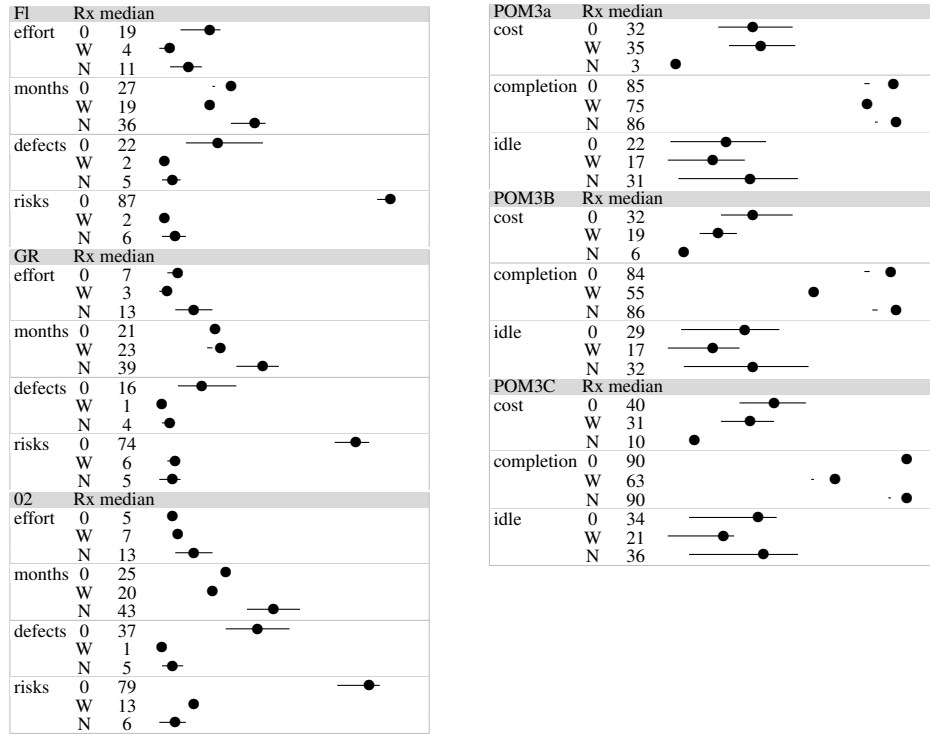| POM3a | Rx | median | |
|---|---|---|---|
| cost | 0 | 32 | |
| | W | 35 | |
| | N | 3 | |
| completion | 0 | 85 | |
| | W | 75 | |
| | N | 86 | |
| idle | 0 | 22 | |
| | W | 17 | |
| | N | 31 | |
| **POM3B** | **Rx** | **median** | |
| cost | 0 | 32 | |
| | W | 19 | |
| | N | 6 | |
| completion | 0 | 84 | |
| | W | 55 | |
| | N | 86 | |
| idle | 0 | 29 | |
| | W | 17 | |
| | N | 32 | |
| **POM3C** | **Rx** | **median** | |
| cost | 0 | 40 | |
| | W | 31 | |
| | N | 10 | |
| completion | 0 | 90 | |
| | W | 63 | |
| | N | 90 | |
| idle | 0 | 34 | |
| | W | 21 | |
| | N | 36 | |

Fig. 6: XOMO results (left); POM3 results (right); all results from 20 runs with different random seeds. Big black dots show median values. Horizontal lines show 25th to 75th percentile. All results are normalized 0..100, min..max. Except for POM3's *completion* objective, *smaller* values are *better*. In the "Rx" column, "0,W,N" denotes results from baseline, WICKED, and NSGA-II (respectively).

3. Find predictors for anomalous situation.
4. Find repair points for causes of an undesirable state.
5. Clarify current situation to predict effects or choose response.
6. Find controllable (blockable or achievable) causes.
7. Find actors contributions to outcome.
8. Find motivations for anomalous actions or decisions.
9. Find a within-theory derivation.

That is, to Leake, explanation is akin to planning where the "explainer" is showing some audience how to find or connect together information. A system that supports such explanations makes it easier to "connect the dots". In practice that means an explanation system must:

– Input a large set of axioms: e.g. examples, pieces of background knowledge;
– Output a *reduced* set of axioms: e.g. rules, model fragments, or as done by Veerappa and Lieter [62], a small number of representative examples takes from centroids of clusters on the Pareto frontier;
– Such that, in the reduces space, it is simple and quick to generate goal-based explanations including the nine kinds listed above.

### 7.1   Decision Trees as "Explanation Tools"

Decision tree learning is a widely-used framework for data mining: given a single goal (called the "class"), find some attribute value that splits the data such that the distribution of classes in each split has been simplified (where the simplest distribution is one containing examples from only one class). Decision tree learners then grow sub-trees by recursing on the data in each split. Popular decision-tree learners include:

- CART [**?**, 11] which minimizes the variance of continuous classes in each split; or
- C4.5 [56] which minimizes the information content of the discrete classes in each split.

One reason to prefer decision trees is that they can very fast to execute. Each level of recursion processes progressively less data. Also, the computation at each level of the recursion may be just a few linear passes through the data, followed by an sort of the attributes– so nothing more than $O(Nlog(N))$ at each level [**?**].

Another reason to prefer decision trees is that, as discussed below, they can operationalize much of Leake's and Kelly's cognitive models on explanation. Decision tree learners do have the disadvantage in that, as used in standard pratice, they only focus on one goal. The aim of this paper is to present a novel extension to standard decision tree learning that extends them to multi-objective optimization.

Given the above discussion, it is easy to see why that is so since decision tree learners can operationalize the above definitions of "explanation".

One reason for the popularity of decision tree learners Previously, work on contrast learning for single goal SE problems found that very succinct contrast sets could be generated as a post-processor to decision tree learning [46]:

- Building a decision tree to separate the different outcomes;
- Identifying leaves containing desired outcome $X$ and undesired outcome $Y$;
- Querying that tree to find branches $B_x$ and $B_y$ that lead to $X, Y$.
- Computing $B_x - B_y$ which selects/rejects for desired/undesired outcomes.

In one spectacularly successful demonstration of this technique [43], it was found decision trees with 6,000 nodes had much superfluous information. Specifically, when some branch point high in the tree most separated the classes, then all contrast set learning had to do is report those branch decisions that selected for branches leading to the better classes. Using that approach, a contrast set learning could report contrasts with only one to four variables in each (and when applied to test data, those contrast sets were at pruning away all the undesired outcomes). Other studies with other data sets [44] confirmed the **the law of tiny constrasts**: *the minimal constrast set between things is usually much smaller than a complete description of those things.*

For simple goal classification, one way to operationalize Leake's framework is using decision trees. Given leaves of that tree $\{X, Y, Z, etc\}$, then exists some branch $\{B_x, B_y, B_z, etc\}$ that connects the root to the leaves as a conjunction of attribute/value pairs. Given some opinion about the value of the contents of each leaf $\{U_x, U_Y, U_Z, etc\}$, then the set difference $B_x - B_y$ is the contrast set of the differences that can drive examples on $X$ over to $Y$.

"from here to there".

an explanation does not generate some single unique output. Rather, it inputs a set of axioms or examples and outputs a reduced set of axioms or examples within which it faster and simpler to generate explanations

Current MOEA algorithms are "instance-based methods" that return specific examples that perform "best" with respect to the multiple goals. The number of examples generated in this way can be overwhelming.

If a user wants to learn general principles from those examples, some secondary *explanation* process is required to group and generalize those examples. For example, Veerappa and Lieter [62] clustering examples from the Pareto frontier so users (at a minimum) need only browse the centroids of each clusters).

GAs flat vectors, not the trees explored by by ()say) Gouse et al.

Goals is performance just as good but explain better

One caveat before beginning: if the audience for the results of optimization are not human beings, then perhaps an explanation systems is not required. For example, Petke, Harman, Langdon, & Weimer [53] use evolutionary methods to rewrite code such that the new code executes faster. The audience for the rewritten code is a compiler. Such compilers do not argue or and ask questions about the code they are given to process. Hence, that rewrite system does not necessary need an explanation system. That said, a succinct and useful description of the difference between passing and failing runs of the rewrite system could be useful when (e.g.) a human is trying to debug that code rewrite system.

Yet another model of "explanation" not explored here is the "surprise modeling" approach recommended by Freitas [24], Voinea&Tulea [2] and others including Horvitz [27]. In that approach, (a) some background knowledge (e.g. summaries of prior actions by users) is used to determine "normal" behavior; (b) users are only presented results that deviated from normal expectations. In analogous research, Koegh [31] argues that *time series discords* (infrequent sequential events in a times series) are a useful way to summarize reports from complex temporal streams. The premise of surprise modeling and reporting discords is that "rare events need to be explored". In non-temporal domains, time series discords becomes *anomaly detection* [12]. For example, in the SE domian, Voinea and Telea report tools that can quickly highlight regions of unusually active debugging (and such regions should be reviewed by management) [**?**] (see also the anomaly detection work of Gruska et al. [25]).

We do not dispute the importance of exploring anomalous outliers. On the other hand, when forming policies for software projects, we need treatments that are well supported by the data. Hence, our contrast sets report changes in the data that, in our data, were *frequently* seen to lead to change.

Also, time series discords and anomaly detection are reports on some variables. Hence, they have a different goal to WICKED that strives to report recommendations on how to change the system so to remove some problem.

Further, all the systems described above [2, 25, 27, 31] are either for unsupervised learning (where no objectives are known) or for single objective systems (where only one goal is known). WICKED, on the other hand, is more ambitious since it was designed for multi-objective systems.

Another potential issue with WICKED is correlation-vs-causation conflation. The issue here is that contrast sets will be useless if they report spurious correlations and not true causal effects. Proving that some effect is truly causal is a non-trivial task. The standard Hall criteria for causal effects [52] is so strict that, outside of highly controlled lab conditions, it rarely accepts that any effect is causal. Hence, in software engineering, when researchers talk of causality [9, 15, 28, 66] they use Granger's "predictive causality"; i.e. causality is the ability of predicting values seen in the future from values seen

in the past. Elsewhere, Granger causality has been adapted to data mining by organizing cross-validations such that the test sets contain data collected at a later time than the training sets [38]. In this paper, we adapt Granger casualty to search-based methods by testing recommendations learned from $M$ simulations on a subsequent round of $N$ new simulations. Those recommendations satisfy Granger causality when the subsequent round of $N$ simulations are changed in a manner predicted by the recommendations gleaned from the original $M$ simulations.

"Data farming" is a technique used extensively by the U.S. Military [?]. Data farming builds a "landscape" of output that can be analyzed for trends, anomalies, and insights in multiple parameter dimensions. In a recent review of search-based and data mining methids in SE, we found numerous examples of data farming [?, ?, ?, ?, 13, 14, 26, 50, 58, 61].

In theory. Once a project manager can view their project on the landscape, they can use this visualization to determine

We come to this work after attending a recent seminar at the US Department of Defence's Software Engineering Institute (SEI), Pittsburgh, USA. That seminar reflected on how to best broadcast the lessons learned by SEI to a very broad audience.

In the $21^{st}$ century, it is now impossible to manually browse very large quantities of software project data. For example, as of October 2012, Mozilla Firefox had 800K reports on software projects. While it is now possible to automatically analyze such data with data miners, at some stage a group of business users will have to convene to *interpret the results* (e.g., to decide if it is wise to deploy the results as a defect reduction method within an organization). These business users are now demanding that data mining tools be augmented with tools to support business-level interpretation of that data. For example,

at a recent panel on software analytics at ICSE'12,

industrial practitioners lamented the state of the art in data mining

and software engineering [?]. Panelists commented that

"prediction is all well and good, but what about decision

making?". That is, these panelists are more interested in the interpretations

that follow the mining, rather than just the mining.

## 8   Related Work

Sayyad
    GALE

## 9   to do

Menzies has used combination of WHERE+ENVY has been used previously for finding context-specific rules for single-objective reasoning (reducing defects or software development effort [40]). What is new here is the addition of CART+CON as well as the application to multiple-objective reasoning.

## References

1.
2. Visual data mining and analysis of software repositories, 2007.
3. D. W. Aha, D. Kibler, and M. K. Albert. Instance-based learning algorithms. *Mach. Learn.*, 6(1):37–66, January 1991.

4. H. Allahyari and N. Lavesson. User-oriented assessment of classification model understandability. In *SCAI'11*, pages 11–19, 2011.

5. A. Arcuri and L. Briand. A practical guide for using statistical tests to assess randomized algorithms in software engineering. In *ICSE'11*, pages 1–10, 2011.

6. E. Arisholm and L. Briand. Predicting fault-prone components in a java legacy system. In *5th ACM-IEEE International Symposium on Empirical Software Engineering (ISESE), Rio de Janeiro, Brazil, September 21-22*, 2006. Available from `http://simula.no/research/engineering/publications/Arisholm.2006.4`.

7. I. Askira-Gelman. Knowledge discovery: Comprehensibility of the results. In *Hawaii International Conference on System Sciences*, 1998.

8. N. Bettenburg, M. Nagappan, and A. E. Hassan. Think locally, act globally: Improving defect and effort prediction models. In *MSR'12*, 2012.

9. P. Bhattacharya, M. Iliofotou, I. Neamtiu, and M. Faloutsos. Graph-based analysis and prediction for software evolution. In *Proceedings of the 34th International Conference on Software Engineering*, ICSE '12, pages 419–429, Piscataway, NJ, USA, 2012. IEEE Press.

10. B. Boehm, E. Horowitz, R. Madachy, D. Reifer, B. K. Clark, B. Steece, A. W. Brown, S. Chulani, and C. Abts. *Software Cost Estimation with Cocomo II*. Prentice Hall, 2000.

11. L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. *Classification and Regression Trees*. 1984.

12. V. Chandola, A. Banerjee, and V. Kumar. Anomaly detection: A survey. *ACM Comput. Surv.*, 41:15:1–15:58, July 2009.

13. E. Chiang and T. Menzies. Simulations for very early lifecycle quality evaluations. *Software Process: Improvement and Practice*, 7(3-4):141–159, 2003. Available from `http://menzies.us/pdf/03spip.pdf`.

14. L. Chung, B. Nixon, E. Yu, and J. Mylopoulos. *Non-Functional Requirements in Software Engineering*. Kluwer Academic Publishers, 2000.

15. C. Couto, M. Valente, P. Pires, A. Hora, N. Anquetil, and R. Bigonha. Bugmaps-granger: a tool for visualizing and predicting bugs using granger causality tests. *Journal of Software Engineering Research and Development*, 2, 1024.

16. S. Das and P. Suganthan. Differential evolution: A survey of the state-of-the-art. *Evolutionary Computation, IEEE Transactions on*, 15(1):4–31, Feb 2011.

17. K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan. A fast elitist multi-objective genetic algorithm: Nsga-ii. *IEEE Transactions on Evolutionary Computation*, 6:182–197, 2002.

18. K. Dejaeger, T. Verbraken, and B. Baesens. Toward comprehensible software fault prediction models using bayesian network classifiers. *Software Engineering, IEEE Transactions on*, 39(2):237–257, Feb 2013.

19. B. Efron and R. J. Tibshirani. *An introduction to the bootstrap*. Mono. Stat. Appl. Probab. Chapman and Hall, London, 1993.

20. O. El-Rawas and T. Menzies. A second look at faster, better, cheaper. *Innovations in Systems and Software Engineering*, 6(4):319–335, 2010.

21. C. Faloutsos and K.-I. Lin. Fastmap: a fast algorithm for indexing, data-mining and visualization of traditional and multimedia datasets. In *Proceedings of the 1995 ACM SIGMOD international conference on Management of data*, SIGMOD '95, pages 163–174, 1995.

22. U. M. Fayyad and I. H. Irani. Multi-interval discretization of continuous-valued attributes for classification learning. In *Proceedings of the Thirteenth International Joint Conference on Artificial Intelligence*, pages 1022–1027, 1993.

23. N. E. Fenton and M. Neil. Software metrics: Success, failures and new directions. *J. Syst. Softw.*, 47(2-3):149–157, July 1999.

24. A. A. Freitas. On objective measures of rule surprisingness. In *Proceedings of the Second European Conference on the Principles of Data Mining and Knowledge Discovery (PKDD'98*, pages 1–9. Springer-Verlag, 1998.

25. N. Gruska, A. Wasylkowski, and A. Zeller. Learning from 6,000 projects: lightweight cross-project anomaly detection. In *Proceedings of the 19th international symposium on Software testing and analysis*, ISSTA '10, pages 119–130. ACM, 2010.

26. W. Heaven and E. Letier. Simulating and optimising design decisions in quantitative goal models. In *Requirements Engineering Conference (RE), 2011 19th IEEE International*, pages 79–88, 2011.

27. E. Horvitz, J. Apacible, R. Sarin, and L. Liao. Prediction, expectation, and surprise: Methods, designs, and study of a deployed traffic forecasting service. In *UAI'05*, pages 275–283, 2005.

28. B. Huberman, D. Romero, and F. Wu. Crowdsourcing, attention and productivity. *Journal of Information Science*, 35(6):758–765, December 2009.

29. V. B. Kampenes, T. Dybå, J. E. Hannay, and D. I. K. Sjøberg. A systematic review of effect size in software engineering experiments. *Information & Software Technology*, 49(11-12):1073–1086, 2007.

30. G. Kelly. *The Psychology of Persona] Constructs. Volume 1: A Theory of Personality. Volume 2: Clinical Diagnosis and Psychotherapy*. Norton, 1955.

31. E. Keogh, J. Lin, and A. Fu. Hot sax: Efficiently finding the most unusual time series subsequence. In *Proceedings of the Fifth IEEE International Conference on Data Mining*, ICDM '05, pages 226–233, Washington, DC, USA, 2005. IEEE Computer Society.

32. E. Kocaguneli, T. Menzies, A. Bener, and J. Keung. Exploiting the essential assumptions of analogy-based effort estimation. *IEEE Transactions on Software Engineering*, 28:425–438, 2012. Available from `http://menzies.us/pdf/11teak.pdf`.

33. E. Kocaguneli, T. Zimmermann, C. Bird, N. Nagappan, and T. Menzies. Distributed development considered harmful? In *ICSE*, pages 882–890, 2013.

34. J. Krall and T. Menzies. Gale: Geometric active learning for search-based software engineering. *IEEE Transactions on Software Engineering (submitted)*, 2014.

35. J. Krall, T. Menzies, and M. Davies. Learning the task management space of an aircraft approach model. In *Modeling in Human Machine Systems: Challenges for Formal Verification, an AAAI 2014 Spring Symposium*, 2014.

36. D. Leake. Goal-based explanation evaluation. *Cognitive Science*, 15:509–545, 1991.

37. B. Lemon, A. Riesbeck, T. Menzies, J. Price, J. D'Alessandro, R. Carlsson, T. Prifiti, F. Peters, H. Lu, and D. Port. Applications of simulation and ai search: Assessing the relative merits of agile vs traditional software development. In *IEEE ASE'09*, 2009. Available from `http://menzies.us/pdf/09pom2.pdf`.

38. M. Lumpe, R. Vasa, T. Menzies, R. Rush, and R. Turhan. Learning better inspection optimization policies. *International Journal of Software Engineering and Knowledge Engineering*, 21(45):725–753, 2011.

39. R. Madachy. Heuristic risk assessment using cost factors. *IEEE Software*, 14(3):51–59, May 1997.

40. T. Menzies, A. Butcher, D. R. Cok, A. Marcus, L. Layman, F. Shull, B. Turhan, and T. Zimmermann. Local versus global lessons for defect prediction and effort estimation. *IEEE Trans. Software Eng.*, 39(6):822–834, 2013. Available from `http://menzies.us/pdf/12localb.pdf`.

41. T. Menzies, O. El-Rawas, J. Hihn, M. Feather, B. Boehm, and R. Madachy. The business case for automated software engineerng. In *ASE '07: Proceedings of the twenty-second IEEE/ACM international conference on Automated software engineering*, pages 303–312, New York, NY, USA, 2007. ACM. Available from `http://menzies.us/pdf/07casease-v0.pdf`.

42. T. Menzies, J. Greenwald, and A. Frank. Data mining static code attributes to learn defect predictors. *IEEE Transactions on Software Engineering*, January 2007. Available from `http://menzies.us/pdf/06learnPredict.pdf`.

43. T. Menzies and Y. Hu. Data mining for very busy people. November 2003. Available from `http://menzies.us/pdf/03tar2.pdf`.

44. T. Menzies and Y. Hu. Just enough learning (of association rules): The TAR2 treatment learner. In *Artificial Intelligence Review*, 2007. Available from `http://menzies.us/pdf/07tar2.pdf`.

45. T. Menzies and M. Shepperd. Special issue on repeatable results in software engineering prediction. *Empirical Software Engineering*, 17:1–17, 2012.

46. T. Menzies and E. Sinsel. Practical large scale what-if queries: Case studies with software risk assessment. In *Proceedings ASE 2000*, 2000. Available from `http://menzies.us/pdf/00ase.pdf`.

47. T. Menzies, S. Williams, O. El-Rawas, D. Baker, B. Boehm, J. Hihn, K. Lum, and R. Madachy. Accurate estimates without local data? *Software Process Improvement and Practice*, 14:213–225, July 2009. Available from `http://menzies.us/pdf/09nodata.pdf`.

48. T. Menzies, S. Williams, O. El-Rawas, B. Boehm, and J. Hihn. How to avoid drastic software process change (using stochastic stability). In *ICSE'09*, 2009. Available from `http://menzies.us/pdf/08drastic.pdf`.

49. N. Mittas and L. Angelis. Ranking and clustering software cost estimation models through a multiple comparisons algorithm. *IEEE Trans. Software Eng.*, 39(4):537–551, 2013.

50. I. Myrtveit, E. Stensrud, and M. Shepperd. Reliability and validity in comparative studies of software prediction models. *IEEE Trans. Softw. Eng.*, 31(5):380–391, May 2005.

51. P. K. Novak, N. Lavrač, and G. I. Webb. Supervised descriptive rule discovery: A unifying survey of contrast set, emerging pattern and subgroup mining. *J. Mach. Learn. Res.*, 10:377–403, June 2009.

52. L. Paul and N. Hall. *Causation : A Users Guide*. Oxford University Press, 2013.

53. J. Petke, M. Harman, W. Langdon, and W. Weimer. Using genetic improvement & code transplants to specialise a c++ program to a problem class. In *European Conference on Genetic Programming (EuroGP)*, 2014.

54. D. Port, A. Olkov, and T. Menzies. Using simulation to investigate requirements prioritization strategies. In *IEEE ASE'08*, 2008. Available from `http://menzies.us/pdf/08simrequire.pdf`.

55. D. Posnett, V. Filkov, and P. Devanbu. Ecological inference in empirical software engineering. In *Proceedings of ASE'11*, 2011.

56. R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufman, 1992. ISBN: 1558602380.

57. A. Sayyad and H. Ammar. Pareto-optimal search-based software engineering (posbse): A literature survey. In *RAISE'13, San Fransisco*, May 2013.

58. M. Shepperd and G. F. Kadoda. Comparing software prediction techniques using simulation. *IEEE Trans. Software Eng*, 27(11):1014–1022, 2001.

59. M. J. Shepperd and S. G. MacDonell. Evaluating prediction systems in software project estimation. *Information & Software Technology*, 54(8):820–827, 2012.

60. R. Valerdi. Convergence of expert opinion via the wideband delphi method: An application in cost estimation models. In *Incose International Symposium, Denver, USA*, 2011. Available from http://goo.gl/Zo9HT.

61. A. van Lamsweerde and E. Letier. Integrating obstacles in goal-driven requirements engineering. In *Proceedings of the 20th International Conference on Software Engineering*, pages 53–62. IEEE Computer Society Press, 1998. Available from `http://citeseer.nj.nec.com/vanlamsweerde98integrating.html`.

62. V. Veerappa and E. Letier. Understanding clusters of optimal solutions in multi-objective decision problems. In *Proceedings of the 2011 IEEE 19th International Requirements En-*

*gineering Conference*, RE '11, pages 89–98, Washington, DC, USA, 2011. IEEE Computer Society.

63. Y. Yang, Z. He, K. Mao, Q. Li, V. Nguyen, B. W. Boehm, and R. Valerdi. Analyzing and handling local bias for calibrating parametric cost estimation models. *Information & Software Technology*, 55(8):1496–1511, 2013.

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

65. Q. Zhang and H. Li. Moea/d: A multiobjective evolutionary algorithm based on decomposition. *Evolutionary Computation, IEEE Transactions on*, 11(6):712–731, Dec 2007.

66. P. Zheng, Y. Zhou, M. Lyu, and Y. Qi. Granger causality-aware prediction and diagnosis of software degradation. In *Services Computing (SCC), 2014 IEEE International Conference on*, pages 528–535, June 2014.