# Parameter Optimization in Control Software using Statistical Fault Localization Techniques

Jyotirmoy Deshmukh
University of Southern California
jyotirmoy.deshmukh@usc.edu

Xiaoqing Jin
jinx@cs.ucr.edu

Rupak Majumdar
MPI-SWS
rupak@mpi-sws.org

Vinayak S. Prabhu
MPI-SWS
vinayak@mpi-sws.org

*Abstract*—**Embedded controllers for cyber-physical systems are often parameterized by look-up maps representing discretizations of continuous functions on metric spaces. For example, a nonlinear control action may be represented as a table of pre-computed values, and the output action of the controller for a given input computed by using interpolation. For industrial-scale control systems, several man-hours of effort are spent in tuning the values within the look-up maps. Suppose that during testing, the controller code is found to have sub-optimal performance. The parameter fault localization problem asks which parameter values in the code are potential causes of the sub-optimal behavior. We present a statistical parameter fault localization approach based on *binary similarity coefficients* and *set spectra* methods. Our approach extends previous work on (traditional) software fault localization to a quantitative setting where the parameters encode continuous functions over a metric space and the program is reactive.**

**We have implemented our approach in a simulation workflow for control systems in Simulink. Given controller code with parameters (including *look-up maps*), our framework bootstraps the simulation workflow to return a ranked list of map entries which are deemed to have most impact on the performance. On a suite of industrial case studies with seeded errors, our tool was able to precisely identify the location of the errors.**

## I. Introduction

The correct operation of modern cyber-physical systems relies on a complex software controller interacting with a physical plant (*i.e.*, a model of the physical processes that we wish to control). Controllers used in industrial-scale software present a formidable challenge for formal analysis techniques due to their scale and format. Typical controllers are defined over several modes of operation, each of which may use a customized control scheme. Control schemes based on feedforward maps that seek to statically cancel non-linearities are frequently used [BSM06], [HC13]. Even when feedback control schemes such as PID control are employed, it is common for the proportional, integral and derivative gains to vary across different modes of operation. As a result, embedded control software typically has a fixed high-level structure corresponding to the "control law" chosen by the designer. Flexibility in the algorithm is offered instead by lookup-maps representing nonlinear functions or mode-specific control rules. A typical design process then involves careful hand-tuning of these maps by engineers till the overall system meets the desired performance objectives. As a result, most industrial-scale controllers do not have closed-form symbolic representations, but have parametric representations such as

explicit *look-up maps*. A look-up map $M$ is a function from a finite subset of $\mathbb{R}^d$ to $\mathbb{R}^n$ (*e.g.*, a look-up map in two dimensions is a two dimensional table with table entries in $\mathbb{R}^n$). It defines a mathematical function $f_M^{\mathcal{I}} : C \to \mathbb{R}^n$, for $C$ a bounded subset of $\mathbb{R}^d$, by completing the map $M$ according to the specified interpolation scheme $\mathcal{I}$ (*e.g.*, linear, bilinear, bicubic interpolation etc.). Whenever controller performs a "look-up" for the value $m$, *i.e.*, wishes to compute $f_M^{\mathcal{I}}(m)$, the control software computes the value using the stored map $M$ under the interpolation scheme $\mathcal{I}$.

Suppose that we are given the code for a controller, with one or more look-up maps, and on simulating the controller in closed-loop with the plant, we find that there is a violation of some correctness or performance specification. Which map entries should we consider for investigation? Which map entries involve performance critical regions of the controller? The problem of identifying the entries most indicative of error or performance loss is called *fault localization*, and tools for effective fault localization are key to a fast design debugging and optimization process.

In this paper, we derive a (heuristics based) ranking on the estimated importance of map entries to an observed failure or to overall system performance. If map entry $m$ is ranked higher than $m'$, then performance can likely be improved by changing $m$ rather than $m'$. Our approach obtains this ranking by computing *binary similarity coefficients* and using *set spectra* based methods. A binary similarity coefficient $\phi(X,Y)$ assigns a measure of similarity between two binary categories $X$ and $Y$ for pattern classification [CC10], [JSH89]. In our setting, $X$ is the category of failed or unfavorably performing executions, and $Y$ is the category of executions that "look-up" an entry $m$ or its neighborhood. The binary similarity based ranking approach computes $\phi(X,Y)$ for each entry $m$ in the map, and orders entries based on the $\phi(X,Y)$ values. This approach has the advantage that no complicated controller analysis is performed. It integrates with existing simulation-guided analysis tools, and we only look at the similarity between failed runs and map entry access in a black box fashion. Set spectra based methods define sets of interest using set algebra methods, also in a black box fashion.

Note that the black-box assumption for industrial control systems is crucial: plant models for real-world physical systems are routinely modeled as hybrid and nonlinear differential or algebraic equations, and controllers have considerable com-

plexity. Traditional verification tools such as model checkers or reachability analysis tools are unable to digest the scale and complexity of such models.

Our approach is inspired by earlier work on software fault localization [JH05], [LNZ$^+$05], [LYF$^+$05], [JS07], [JGG08], [NLR11], [RFZO12], [YHC13] The similarity coefficient based approach is motivated by the work [JH05], [WDGL14] in software localization, and the set spectra based methods by the work [PS92], [RR03]. We extend and improve the ideas in these works in two ways to suit our context. First, we observe that compared to normal software where runs are only classified as being good or bad, executions for cyber-physical systems are ascribed numerical scores (e.g., performance measures or robustness w.r.t. logical specifications [DM10], [DFM13]). Second, in most instances, a look-up map $M$ represents a discrete approximation of the associated *continuous* function $f_M^{\mathcal{I}}$. Consequently, in most cases, the map value $M(m)$ is close to the map value $M(m')$ if $m$ and $m'$ are close. This implies that an execution can indirectly contribute to the importance level of a map entry $m$ even if that execution does not directly access $M(m)$ – this happens when indices close to $m$ are accessed. This continuity property is absent in traditional software where individual statements that are problematic can be spread out far apart from each other. We show how to cleanly extend the two Boolean software fault localization approaches to account for the additional structure in our domain by defining the rankings using certain basic quantitative functions which give a measure of the effect of a map entry on a given run.

The problem of fault localization for look-up maps can be viewed as a parameter optimization or tuning problem, where we are trying to determine values for parameters (in this case the map entries) such that the value of a certain cost function (in this case the quantitative score indicating well-behavedness of the output) is maximized. In this sense, black-box optimization approaches can be theoretically used to perform fault localization. However, even a relatively small 20x20 look-up map contains 400 entries, or from the perspective of the black-box optimizer, a large 400-dimensional space over which it needs to optimize. Such an approach, *if it scales*, would be appropriate to synthesize a new look-up map in entirety to make the design satisfy its specification, but it is not clear whether it is well-suited to determining which entries in an *existing* look-up map are responsible for making the design not satisfy its specification (identifying sub-performing regions is of key importance to engineers). Other parameter tuning tools such as the PID tuning capability within the Control System toolbox in Simulink have very specific tuning abilities and are not general enough for the problem that we address. [ADV$^+$17] is a recent work on abstraction based *verification* of look-up maps.

We test the ideas proposed using a prototype tool and present its results over case studies from the industrial processes and automotive domain. The tool is well-integrated with the simulation-guided falsification tool Breach [Don10], and is able to use a wide array of ranking heuristics to sort the entries in the specified look-up map in order of importance. The tool supports arbitrary quantitative cost functions to score executions, including the quantitative robust satisfaction function as defined for Signal Temporal Logic specifications in [DM10], [DFM13]. We demonstrate the capabilities of the tool in localizing possible sources of fault to look-up map entries, as well as in identifying promising regions of parameter spaces in order to optimize the closed-loop system performance.

## II. PROBLEM SETTING

**Software Components and Scored Executions.** In this work we focus on software components that occur in controllers of dynamical systems. For our purposes, a software component is a module that transforms inputs into outputs. Let $\mathcal{A}$ be a software component, and let the execution (or *run*) of $\mathcal{A}$ given an input $z$ be denoted as $\mathcal{A}(z)$. We consider a setting where we are given a scoring of the runs of $\mathcal{A}$ on a finite set of inputs $\mathcal{Z}$: each run $\mathcal{A}(z)$ for $z \in \mathcal{Z}$ is given a real-valued *score*, denoted $\mathsf{score}(\mathcal{A}, z)$, with negative values being "bad" (denoting *failed* or suboptimal runs), and positive values being "good" (denoting passed or *successful* runs). This extends the boolean notion where the score set can be viewed as $\{-1, 1\}$. Intuitively, the score of a run indicates our satisfaction level on how good that run is; the higher the score, the more satisfied we are with the execution. As an example, given a temporal logic (*e.g.* Signal Temporal Logic [MN13]) specification $\varphi$, the successful runs $\mathcal{A}(z)$ are those which satisfy $\varphi$, and the failed runs are those which do not satisfy $\varphi$. Quantitative scores can be given to runs based on quantitative semantics of the logic (*e.g.*, the *robustness* values, ranging over reals, corresponding to a Signal Temporal Logic formula $\varphi$ [DM10], [DFM13]).

**Look-up Maps.** Software components invoke real-valued functions in their executions. When these functions do not have a closed form representation, or have a closed form representation that is computationally inefficient, they are instead stored in the form of *look-up maps* as follows. A real-valued (finite) *map* $M$ is a function from a finite subset of $\mathbb{R}^d$ to $\mathbb{R}^n$. The domain of $M$ (which is a finite set) is denoted by $\mathsf{dom}(M)$, and we refer to elements $m \in \mathsf{dom}(M)$ as *map entries*. Let $f_M^{\mathcal{I}} : C \to \mathbb{R}^n$ be the completion of $M$ using some interpolation scheme $\mathcal{I}$, where $C$ is a bounded subset of $\mathbb{R}^d$. The stored map $M$ is a discrete representation of the function $f_M^{\mathcal{I}}$, and for a given value $m$ is used to compute $f_M^{\mathcal{I}}(m)$ using an interpolation scheme. We abuse notation and use $\mathcal{I}(M(m_1), \ldots, M(m_l))$ to denote the result of applying the interpolation scheme $\mathcal{I}$ to the values $M(m_1), \ldots, M(m_l)$. The most commonly known interpolation scheme is linear interpolation. Controller software components also use other interpolation schemes, such as nearest neighbor interpolation, bilinear and bicubic interpolation [sim07]. Multivariate interpolation schemes may offer a better approximation of function values based on the values of neighboring points [GS00].

Given $m \in C$, we define $\mathsf{depend}\,(M, \mathcal{I}, m)$ to be the values in the finite map domain $\mathsf{dom}(M)$ which are used to compute

$f_M^{\mathcal{I}}(m)$. Formally, depend $(M, \mathcal{I}, m) =$

$$\begin{cases} \{m\} & \text{if } m \in \text{dom}(M) \\ \{m_1, .., m_l\} & \text{if } f_M^{\mathcal{I}}(m) = \mathcal{I}(M(m_1), \ldots, M(m_l)). \end{cases} \quad \text{(II.1)}$$

Thus, depend $(M, \mathcal{I}, m)$ indicates which of the map domain entries $m_1, m_2 \ldots, m_l$ are used by the interpolation scheme $\mathcal{I}$ to define $f_M^{\mathcal{I}}(m)$. Consider a run $z$. Given a map entry $m$ in the finite set $\text{dom}(M)$. we say that $z$ has a *map access for $m$* or *$z$ accesses $m$* if (a) there exists a query for $f_M^{\mathcal{I}}(m)$ in $z$; or (b) there exists a query for $f_M^{\mathcal{I}}(m')$ in $z$ for some $m'$ such that $m \in$ depend $(M, \mathcal{I}, m')$ (*i.e.* the map value at $m$ is queried during an interpolation in the execution $z$).

**Parameter Localization Rankings.** In addition to look-up maps, controller software often incorporates several *parameters* – these are certain variable entities that can take values in some (usually) quantitative domain (*e.g.*, gain constants). Look-up maps can be viewed as a set of parameters (each $m \in M$ defines a parameter with value $M(m)$). The parameter localization problem is to narrow down problematic parameters, in case of sub par controller performance, amongst the total set of parameters (whether the parameters are in the form of look-up map entries, or are of other forms) in the code. Our approach to localization is to construct a *ranking* of the parameters, in decreasing order of our belief in them being problematic. Our parameter ranking approach can be applied for (i) *repair* – in this case we are given a hard requirement that must be met, and executions are scored negative iff they do not satisfy the requirement; or (ii) *robustness and optimization* – where we score the least desirable runs negatively (for example, in the case of Signal Temporal Logic where all the runs satisfy a given specification $\varphi$, we can shift the quantitative robustness values of runs by some negative constant, so that executions which originally had a low positive robustness score value get a negative score after shifting and get into the undesired class).

In the sequel, we fix $\mathcal{A}$ and $M$ and $\mathcal{I}$, and omit them (when unnecessary) for notational simplicity.

## III. Ranking based on Similarity Coefficients

Similarity coefficients [JH05], [WDGL14] have been widely used in pattern analysis problems for classification and clustering [CC10], [JSH89], [HPAW15], [PAR16]. A binary similarity coefficient $\phi(X, Y)$ assigns a measure of similarity between two binary categories $X$ and $Y$. In our setting, $X$ is the category of failed (or suboptimal) executions, and $Y$ is the category of executions that are affected by an entry $m$. In this approach, we compute a ranking $\phi(X, m)$ for each entry $m$, and sort the entries based on the $\phi(X, m)$ values. In the extant approach for traditional software debugging, the categories $X$ and $Y$ are Boolean, *i.e.*, for a given category $X$, the only relevant property is whether an instance belongs to $X$. In our case, instances are assigned a *score* measuring how well they belong to $X$. For example, real-valued scores on executions can correspond to a quantitative measure of how well an execution satisfies a given logical property, by measuring the execution's distance from the set of executions

| | | |
|---|---|---|
| $F_A^m$ | $=$ | $\sum_{z \text{ s.t. } \text{score}(z) < 0} \text{raffect}(z, m) \cdot \text{score}(z)$ |
| $P_A^m$ | $=$ | $\sum_{z \text{ s.t. } \text{score}(z) \geq 0} \text{raffect}(z, m) \cdot \text{score}(z)$ |
| $F_U^m$ | $=$ | $\sum_{z \text{ s.t. } \text{score}(z) < 0 \text{ and } \text{raffect}(z,m)=0} \text{score}(z)$ |
| $P_U^m$ | $=$ | $\sum_{z \text{ s.t. } \text{score}(z) \geq 0 \text{ and } \text{raffect}(z,m)=0} \text{score}(z)$ |
| $P$ | $=$ | $\sum_{z \text{ s.t. } \text{score}(z) \geq 0} \text{score}(z)$ |
| $F$ | $=$ | $\sum_{z \text{ s.t. } \text{score}(z) < 0} \text{score}(z)$ |

TABLE I
BUILDING BLOCKS FOR QUANT. SIMILARITY COEFFICIENTS

not satisfying the property. We first extend the standard binary similarity coefficient to this quantitative setting. For ease of presentation, we restrict ourselves to look-up maps.

### A. Basic Quantitative Approach: Preliminaries

Suppose that executions in $\mathcal{Z}$ are scored. Further, suppose that a run $z \in \mathcal{Z}$ queries the value $f_M^{\mathcal{I}}(m)$ for $m \in \text{dom}(f_M^{\mathcal{I}})$. This value is then constructed using the values for the entries in depend $(m)$. The execution of the run $z$ thus depends on the entries in depend $(m)$. For $m \in \text{dom}(f_M^{\mathcal{I}})$ and $m' \in \text{dom}(M)$ let maffect $(m, m') = 1$ if $m' \in$ depend $(m)$, and 0 otherwise. In other words, when a run $z$ attempts to compute $f_M^{\mathcal{I}}(m)$, then for each entry $m' \in \text{dom}(M)$ for which maffect$(m, m') = 1$, we say that $m'$ affects $z$.

In general, a run $z$ may query several map values $f_M^{\mathcal{I}}(m_1)$, $f_M^{\mathcal{I}}(m_2)$, ... during its execution. If any of $m_1$, $m_2$, ... are affected by a map entry $m'$, then $m'$ affects $z$. This is captured by the (Boolean for now) function

$$\text{raffect}(z, m') = \begin{cases} 1 & \text{if } z \text{ queries some entry } m \\ & \text{with maffect}(m, m') = 1 \\ 0 & \text{otherwise} \end{cases} \quad \text{(III.1)}$$

We recall that the ordering of a similarly coefficient presents a suspiciousness ranking of the indices in $M$. A similarity coefficient based rank $\phi : \text{dom}(M) \to \mathbb{R}_+$ assigns a numerical value $\phi(m)$ to each entry $m \in \text{dom}(M)$; if the rank of a particular $m \in \text{dom}(M)$ is greater than that for another entry $m'$, then $m$ is considered more suspicious. We define a few building blocks for defining quantitative binary similarity coefficients in Table I, and briefly discuss these next.

1) $F_A^m$ is the sum of (negative) scores of the failed runs which are *affected* by entry $m$. The suspiciousness of $m$ should increase with the value[1] of $|F_A^m|$.
2) $P_A^m$ is the sum of (positive) scores of the passed runs *affected* by $m$. The suspiciousness of $m$ should decrease for higher $P_A^m$.

[1]In many applications, *e.g.*, when using Signal Temporal Logic and its quantitative interpretation, a slightly negative score (which is close to 0, *e.g.* 0.04) of an execution is a noteworthy event as this means a property is violated. As defined, $F_A^m$ however effectively discards such runs as raffect $(z, m)$ is weighted by the score of the run. If this is undesired, we can "shift" the negative by a negative constant, and similarly the positive scores by a positive constant in a preprocessing step so that scores with a low absolute value do not arise.

3) $F_U^m$ is the sum of (negative) scores of the failed runs *unaffected* by entry $m$. The suspiciousness of $m$ should decrease with an increase in $|F_U^m|$, as a high $|F_U^m|$ indicates that the problem lies away from $m$.
4) $P_U^m$ is the sum of (positive) scores of the passing runs *unaffected* by entry $m$. This quantity is not important, as we are interested in executions which are problematic; not executions which are problem-free and do not access $m$.
5) $P$ is the sum of (positive) scores of the passed runs.
6) $F$ is the sum of (negative) scores of the failed runs.

### B. Basic Quantitative Similarity Coefficients

Similarity coefficients can be constructed based on the quantities $F_A^m, P_A^m, F_U^m, P_U^m, F, P$ defined previously. We refer the reader to [CC10] for 76 similarity coefficients that have been studied before in the Boolean case.

In this work we focus on three illustrative similarity coefficients – the Tarantula similarity coefficient, the Kulczynski coefficient, and the $D^*$ similarity coefficient.

**Tarantula similarity coefficient.** The Tarantula tool [JHS02], [JH05], given a test suite where each test is labeled as either passing or failing, uses similarity coefficients to rank the statements of programs in decreasing order of suspiciousness. The ranking is presented in the form of a visual map with different statements getting color and brightness levels according to the value of the similarity coefficient for that statement. The Tarantula similarity coefficient for potentially faulty map entries extended to the case where instead of pass/fail executions we have quantitative values, is given as follows:

$$R_{\text{tarantula}}(m) = \frac{F_A^m/F}{F_A^m/F \ + \ P_A^m/P} \tag{III.2}$$

**Kulczynski similarity coefficient.** This similarity coefficient extended to the quantitative setting is:

$$R_{\text{kulczynski}}(m) = \frac{|F_A^m|}{|F_U^m| \ + \ P_A^m} \tag{III.3}$$

$D^*$ **similarity coefficient.** The $D^*$ coefficient [WDGL14] is based on the Kulczynski similarity coefficient to give more importance to failed executions which are affected by an map entry $m$ compared to (i) failed executions which are not affected my $m$ (as there might be other map indices which might be to blame for those other executions), and (ii) successful executions which are affected by $m$ (as failures are more relevant than successes). This adjustment is done by raising the numerator $|F_A^m|$ in the Kulczynski similarity coefficient to a positive power $\gamma \geq 1$. The $D^*$ coefficient extended to our quantitative setting is:

$$R_{\text{dstar}}^{\gamma}(m) = \frac{|F_A^m|^{\gamma}}{|F_U^m| \ + \ P_A^m} \tag{III.4}$$

Note that we are only interested in the order relation imposed by the ranking functions, not the numerical values themselves. The work [WDGL14] found $\gamma = 2$ to be best in their experiments. They also compared this coefficient to some of the other similarity coefficients (*e.g.*, the Tarantula coefficient), and found the $D^*$ coefficient to be better than others (in the sense of buggy statements being ranked closer to the top).

### C. Utilizing Continuity and the Metric Space Structure

In many cases, the map $M$ approximates a *continuous* function $f_M^{\mathcal{I}}$. The value of $f_M^{\mathcal{I}}(m)$ hence while depending directly on the map entries in depend $(m)$, is also correlated with the map values at nearby map entries. Consider a situation where we have 51 failed executions, each scored the same, where each failed execution is affected by a different map entry (thus there are 51 potential map entries ($\text{dom}(M) = 51$). Of these, 50 map entries are clustered very closely, and the one remaining entry is far away from these 50. Intuitively then, the 50 clustered map entries should be viewed as more problematic than the other remaining entry. Additionally, design engineers are typically interested in identifying problematic regions of a look-up map. A problematic region in $f_M^{\mathcal{I}}$ might indicate that the engineers need to "refine the grid" in the map $M$ in the identified region, rather than simply changing the map values (keeping the grid resolution the same).

We account for correlation in suspiciousness of proximal map entries by introducing a dependence between map entries the decays with distance. A modular way to accomplish this goal is to build upon the approach of Subsection III-A by defining quantitative extensions of the function maffect and raffect which account for the fact that the values of neighboring map entries are related. The new similarity coefficients can then be constructed which simply use these new maffect and raffect functions. The quantitative version of maffect is defined as:

$$\text{maffect}\,(m, m') = \begin{cases} \lambda_M^{\mathcal{D}(m,m')} & \text{if } \mathcal{D}(m, m') \leq r_M \\ 0 & \text{otherwise.} \end{cases} \tag{III.5}$$

where $0 < \lambda_M < 1$ is a decay constant , and $r_M$ is a quantity denoting the radius of influence. The value $\text{maffect}\,(m, m')$ quantifies how much the map value at entry $m'$ affects the map value at entry $m$. These quantities can in the general case be dependent on $m$ – if $f_M^{\mathcal{I}}$ at $m$ is changing very fast, then $\lambda_M$ and $r_M$ at $m$ will be small.

Using the above defined maffect function we give a quantitative version of the function raffect which quantifies which map entries affect an execution as follows:

$$\text{raffect}(z, m') = \max_{\substack{m \text{ s.t. } z \text{ queries } m \\ \wedge\ \mathcal{D}(m, m') \leq r_M}} \text{maffect}\,(m, m'). \tag{III.6}$$

raffect $(z, m')$ can also be defined in other ways (*e.g.* taking a sum instead of the maximum). We compare raffect to the binary raffect function of Subsection III-A – there raffect $(z, m')$ was either 0 or 1 based on whether the map entry $m'$ was utilized or not during the execution $z$. Now, raffect $(z, m')$ gives a more nuanced estimate of the importance of the map value for entry $m'$ (accounting for the fact that neighboring map entries should have related map values).

The new ranking coefficients corresponding to $R_{\text{tarantula}}$, $R_{\text{kulczynski}}$, and $R_{\text{dstar}}^{\gamma}$ can be defined as in Equations (III.2), (III.3), and (III.4), by plugging in the new quantitative functions raffect and maffect in the defining equations.

### D. Incorporating Frequency of Access

The rankings of the previous two subsections do not incorporate the *frequency* of map accesses inside an execution. Consider a situation in which bad executions have a tendency to repeatedly access a certain portion of the map during the course of the executions. If a map region is repeatedly accessed during an execution, one natural heuristic is to give that map region more importance. For a given run $z$ and a map entry $m$, a frequency measure of accessing a particular map region around $m$ can be expressed by generalizing the raffect function as follows. Let $|z|$ denote the number of map queries in $z$, that is, if in $z$ we have the (possibly non-distinct) queries $f_M^{\mathcal{I}}(m_1), f_M^{\mathcal{I}}(m_2), \ldots f_M^{\mathcal{I}}(m_p)$ then $|z|$ is $p$. The $k$-th query in $z$ is denoted by $z[k]$. The function fraffect generalizes raffect, and is defined as

$$\text{fraffect}\,(z, m) = \frac{\sum_{k=1}^{|z|} \text{maffect}\,(z[k], m)}{|z|} \tag{III.7}$$

for $m \in \text{dom}(M)$.

Intuitively, $\text{fraffect}\,(z, m)$ can be thought of as a weighted fraction corresponding to the effect of the region around $m$ on the execution $z$. The function fraffect can be defined using the binary maffect function; or the quantitative maffect function of Equation (III.5) incorporating entry correlation. As an example, if we use the binary maffect function, and $\text{fraffect}\,(z, m)$ is $0.4$, then it means that $40\%$ of the function calls $f_M^{\mathcal{I}}(\cdot)$ during the course of $z$ were affected by $m$. Analogues to Equations (III.2), (III.3), and (III.4) can be obtained by replacing raffect by fraffect; and letting $P = P_U^m + P_A^m$, and $F = F_U^m + F_A^m$.

## IV. Set Spectra Based Methods

An alternative to the binary coefficient based ranking method is to define various *sets* of map indices for inferring possible problematic values. The work in [PS92] defines several sets of interest in software fault localization using set algebra operations. We adapt the union model, which is the most promising set-spectra based heuristic, to our quantitative setting of map lookups.

The union model [PS92], [AHLW95] for software fault localization – given a failed run – looks at statements that are executed in the buggy execution, but not in *any* successful run. The intuition behind the model is that if certain statements are executed *only* in failing runs, those statements are likely to be buggy. The term union comes from the fact that under this heuristic, the set of suspicious statements is given by

$$\bigcup_{\substack{\text{statement s executed by } z \\ \text{s.t. score}(z) < 0}} \{s\} \quad \setminus \quad \bigcup_{\substack{\text{statement s executed by } z \\ \text{s.t. score}(z) > 0}} \{s\}$$

The work in [RR03] found in their experiments that buggy statements often are executed in successful runs, and so

the union model in many cases fails to label these buggy statements as suspicious (*i.e.*, the above set is empty); however, in cases where the union model *does* label statements as suspicious, the labelled statements are almost always buggy. That is, the union model has a very low false positive rate in labelling buggy statements (and a high false negative rate). A very low false positive rate is extremely attractive in bug finding, thus, we chose to explore the performance of the union method for our look-up map setting. In case the union method labels every region as non-buggy, we can fall back on other methods, *e.g.*, the methods of the previous section.

**Utilizing Continuity of Map Values.** We modify the union model in this case as follows. We ask for the following (at a high level): is there a map access $m$ in a run $z$ with a negative score $\text{score}(z)$, such that all positively scored runs $z'$ access map indices at least $r_M$ distance away from $m$? That is, for the set of executions $\mathcal{Z}$, is there an area in $\text{dom}(f_M^{\mathcal{I}})$ that is (i) accessed only by failing runs, and (ii) is at least $r_M$ distance away from the areas accessed by successful runs?

Formally, let $\text{access}(n, m)$ be the predicate on whether a run $z$ access map entry $m$, let $M_F$ be the set of map entries accessed by faulty (negatively scored) runs, and let $M_S$ be the set of map entries accessed by successful (non-negatively scored) runs. That is,

- $M_F = \{m \,|\, \exists z \text{ s.t. } \text{score}(z) < 0 \land \text{access}(z, m) = \textsc{t}\}$;
- $M_S = \{m \,|\, \exists z \text{ s.t. } \text{score}(z) \geq 0 \land \text{access}(z, m) = \textsc{t}\}$.

For $X$ a set of map entries and $r$ a positive real number, let $\text{Ball}(X, r)$ denote the set $\{x' \in \text{dom}(M) \mid \text{ exists } x \in X \text{ such that } \mathcal{D}(x, x') \leq r\}$. The set of suspicious map indices $\text{sus}_U$ is defined to be:

$$\text{sus}_U = M_F \setminus \text{Ball}(M_S, r_M). \tag{IV.1}$$

The elements in $\text{sus}_U$ are the suspicious map entries. In most cases, a ranking on $\text{sus}_U$ will not be necessary as $\text{sus}_U$ will be a small set, and the entire set can be classified as suspicious. If desired we can rank the entries in $\text{sus}_U$ as follows. For a map entry $m \in \text{sus}_U$, the quantification of suspiciousness of $m$ depends on: (a) the negative scores of the failed runs, and (b) the distance from the map accesses in the positively scored runs. For a map entry $m \in \text{sus}_U$, we call the first quantity accessed by a negatively scored run as $s_U(m)$ and define it as follows:

$$s_U(m) = \min \{|\text{score}(z)| \mid \text{access}(z, m) = \textsc{t} \land \text{score}(z) < 0\} \tag{IV.2}$$

In other words, $s_U(m)$ tells us what is the absolute value score of the best among failing runs (*i.e.* the score which is closest to 0) which is affected by $m$. The higher the value of $s_U(m)$, the more suspicious should $m$ be. The second quantity for an entry $m \in \text{sus}_U$ is denoted $d_U(m)$ and we define it as follows:

$$d_U(m) = \min \left\{ \mathcal{D}(m, m') \;\middle|\; \begin{array}{l} \exists z' \text{ s.t. } \text{score}(z') \geq 0 \,\land \\ \text{access}(z', m') = \textsc{t} \land \mathcal{D}(m, m') > r_M \end{array} \right\}$$

Note that the above value is equivalent to minimize the distance $\{\mathcal{D}(m, m') \mid m' \in M_S\}$.

In case the set in the above equation is empty, *i.e.*, if there does not exist a run $z'$ such that $\mathsf{score}(z') \geq 0$, we let $d_U(m)$ be a high constant. Note that since $m \in \mathsf{sus}_U$, if there exists a run $z'$ such that $\mathsf{score}(z') \geq 0$, then $z'$ access $m'$ with $\mathcal{D}(m, m') > r_M$. The quantity $d_U(m)$ gives the separation distance between $m$ and $M_S$ (the set of entries affected by successful runs). The higher the value of $d_U(m)$, the more suspicious should $m$ be. The set union based ranking $R_U$ is a function of $d_U$ and $s_U$:

$$R_U(m) = s_U(m) \cdot d_U(m).$$

Note that setting $r_M = 0$ makes $R_U(m)$ be 0 iff there exists a positively scored run which accesses $m$ (like in the original union model).

The above method marks only map entries which are accessed, as opposed to the rankings of the previous section which may mark entries that are not accessed, but lie in suspicious areas of the $f_M^{\mathcal{I}}$ areas. When using this method in practice, the user thus must look at the surrounding entries in case an index $m$ is found to be buggy.

We keep only the set algebraic version of Equation (IV.1) for this model, and not define quantitative extensions based on smoothness of $f_M^{\mathcal{I}}$ (as in Subsection III-C) because this heuristic is based on set-algebra. We expect a high percentage of entries in $\mathsf{sus}_U$ to be problematic (the set $\mathsf{sus}_U$ is expected to be small), thus further quantitative refinements are not of much use.

## V. Evaluating Effectiveness of Proposed Approach

For fault localization in traditional software, the EXAM score [WGL+16] is a commonly used measure to quantify the effectiveness of the different techniques. In our context, it can be defined as

$$\text{EXAMscore} = \frac{\begin{array}{c}\text{number of look-up map entries examined}\\ \text{(in decreasing order of their scores)}\\ \text{till an entry deemed problematic is encountered}\end{array}}{\text{total number of look-up map entries}} \times 100$$

A lower score is better as it indicates that only a small number of entries are ranked above the "true" buggy ones. However, in engineering practice, a percentile relative score is not always sufficient [PO11]; thus we also define an *absolute* version of the EXAM score as

$$\text{AbsEXAMscore} = \begin{array}{c}\text{number of look-up map entries examined}\\ \text{(in decreasing order of their scores)}\\ \text{till an entry deemed problematic is encountered}\end{array}$$

In the look-up map context, EXAMscore, and AbsEXAMscore values are not always appropriate metrics. While these scores are relevant when the maps contain isolated entries that may be the cause of a fault, in typical practice, a *region* of the map (which often correspond to regions of operation of the plant and the controller) contains suboptimal entries. In essence, we want to *cluster* map entries based on their rank values, and present these clusters to the user in order to identify suboptimal map regions. Ideally, there should be a small number of clusters containing the high-ranked map entries, and moreover each

such cluster containing high-ranked map entries should contain minimal *low*-ranked entries. There is a plethora of cluster analysis algorithms and tools [AR14]. For one or two-dimensional look-up tables, one of the simplest methods is by visual inspection of a *heat-map* of the ranked entries (Matlab provides a heat-map visualization function). Our tool incorporates this heat-map visualization functionality.

## VI. Case Studies

In this section, we empirically evaluate the efficacy of the ranking heuristics on several case studies. In Subsection VI-A we investigate the ranking heuristics on smaller examples, and in Subsection VI-B we run the ranking methods on industrial case studies.

In this section, we use the term look-up table (LUT) for readers familiar with the eponymous Simulink® block for implementing $N$-dimensional look-up. Understanding how each ranking heuristic performs on toy models, can help provide CPS designers with guidelines on choosing the right ranking heuristic based on their design-type. We assume reader familiarity with temporal logics such as Signal Temporal Logic (STL). We refer the readers to [MN13], [DFM13] for STL syntax, and semantics (boolean and quantitative).

### A. Basic Models

*1) Nonlinearity Cancellator:* We designed a toy model $\mathcal{A}_{\text{nc}}$ to mimic canceling a nonlinearity in the input $u(t)$ in Simulink®. $\mathcal{A}_{\text{nc}}$ has two outputs functions:

$$y_1(t) = u(t) \cdot M(u(t)) \qquad y_2(t) = \sum_{k=1}^{\lfloor \frac{t}{\Delta} \rfloor} \Delta \cdot y_1(k\Delta) \quad \text{(VI.1)}$$

In the equation for $y_2(t)$, $\Delta$ represents the fixed time step used for simulating $\mathcal{A}_{\text{nc}}$, ($\Delta = 0.1$ sec. for this experiment). In the equation for $y_1(t)$, $M(x)$ represents a LUT representing the nonlinear function $\frac{1}{x}$. In our example, we use

$$\mathsf{dom}(M) = \{e \mid e = 0.1 \cdot z \wedge z \in \mathbb{N} \wedge z \in [1, 90]\}, \quad \text{(VI.2)}$$

(thus there are 90 entries in $M$), and for each $e$ in $\mathsf{dom}(M)$, we let $M(e) = 0.01 \cdot \mathsf{round}(\frac{100}{e})$, *i.e.*, the approximation of $\frac{1}{e}$ up to two decimal places[2].

The input signal $u(t)$ is specified in terms of 11 equally spaced control inputs between times 0 and 30 secs, where $u(t)$ is in $[0.09, 9.01]$ at each control point, and is a linear interpolation of the values at control points for all other times. The domain of the input signal is purposely chosen to exceed the domain of the LUT $M$ to exercise the extrapolation of values performed by the LUT[3]. The range of $u(t)$ and the chosen $M$ together guarantee that $\forall u(t) \in [0.09, 9.01]$, $u(t) \cdot M(u(t)) < 1.4$. We are interested in checking the model against the following STL requirements:

$$\varphi_1^{\text{nc}} \triangleq \Box_{[10,30]}(|y_1 - 1| < 0.4) \quad \varphi_2^{\text{nc}} \triangleq \Box_{[0,30]}(y_2 \leq 30) \quad \text{(VI.3)}$$

---

[2]Note $\mathsf{round}(a)$ rounds $a$ to the nearest integer.
[3] We note that the semantics of LUTs in Simulink® allow for values outside the domain of the LUT to be input to the LUT. The LUT output computed is then obtained by linear extrapolation.

We then seed the LUT $M$ with a bug, by changing $M(2)$ to 0.8 (original value 0.5). We excite the model with 100 randomly chosen piecewise linear signals $u(t)$, and then evaluate the ranking on entries produced by each of the heuristics.

*Results.* The top 3 entries deemed most important by each of the heuristics are reported in Table II.

We get no information from the $R_{\text{tarantula}}$ ranking heuristic as all entries are weighted equally. For the first property ($\varphi_1$), the ranking heuristics that take the frequency of an entry into account perform generally better than the heuristic based on a Boolean notion of access. This is expected as every violation of the requirement (which corresponds to accessing the faulty region of the LUT), has an additive effect on the frequency-based ranking heuristics. The union spectrum based methods are also ineffective because the way we designed the experiment, all LUT entries are accessed by each trace. This is achieved by supplying an input signal with the first linear segment that is a ramp from 0.09 to 9.01 within the first 3 seconds. As this example had a seeded bug, we can compute the AbsEXAMscore. This score is $-, -, 2, 1, 2, -$ corresponding to the six ranking functions in Table II (here we take the worst scores corresponding to the two specification functions $\varphi_1^{\text{nc}}$ and $\varphi_2^{\text{nc}}$). Recall that the union spectrum method works well when there are certain LUT entries accessed only by the failing traces.

*2) Two-Dimensional Nonlinear System:* In this experiment, we designed $\mathcal{A}_{\text{ff}}$, a model to represent the **F**eed**F**orward control of a nonlinear dynamical system. The plant is a 2-dimensional nonlinear (unstable) system with dynamics described by ODEs on the left side of Eqn. (VI.4).

$$\dot{x}_1 = -3x_1 + 2x_1 x_2^2 + u \qquad \varphi^{\text{ff}} \triangleq \Box_{[0.8,2]}(|x_1| < 0.8) \quad \text{(VI.4)}$$
$$\dot{x}_2 = -x_2^3 - x_2$$

The control action $u$ used is one which trivially cancels out the nonlinearity in the first state's dynamical equation ($u = -2x_1 x_2^2$). We observe that once $u$ cancels out the $2x_1 x_2^2$ term, the rest of the system is trivially asymptotically stable[4] with the $V(x_1, x_2) = x_1^2 + x_2^2$ as a Lyapunov function certifying stability. This guarantees that for any $c$, the system never leaves the set $V(x_1, x_2) < c$.

We note that the computed feedforward control action is a polynomial in the plant state that can be represented using a 2-dimensional LUT in which we add an entry at intervals of 0.5 for integer values of $x_1$ and $x_2$ in the range $[-10, 10]$. This gives a LUT with $41^2 = 1681$ entries. For the purpose of simulating the system, we pick a random initial state $x_1(0) \in [-10, 0]$ and $x_2(0) \in [0, 10]$. We then introduce bugs in 30 of the entries that correspond to $x_1$ in $[-10, -8]$, and $x_2$ in $[7.5, 10]$. The bug basically multiplies each entry by $-2$. While the original system is globally asymptotically stable, observe that the bug may cause unstable behavior in the system.

---

[4]The Lie derivative of the resulting system with $V(x_1, x_2)$ as the Lyapunov function is $-(3x_1^2 + x_2^4 + x_2^2)$, and as the term inside the parentheses is a sum-of-squares polynomial, the Lie derivative of $V(x_1, x_2)$ is negative everywhere.

The STL requirement on the right side of Eqn. (VI.4) relates to the settling time of $x_1(t)$. We run 100 randomly chosen simulations to excite the model and apply our ranking heuristics to identify the entries likeliest to be the root cause of settling time violations.

*Results.* The top 3 entries deemed most important by each of the heuristics are reported in Table II. The ranking heuristics based on quantitative similarity coefficients do not work well for this example. (There is an exception: the Tarantula rankings utilizing the metric space interpretation of the LUT). This is so because when a wrong LUT entry is accessed in computing the feedback law, the controlled signal deviates from the desired reference value, and all subsequent accesses of the LUT are indexed by the deviated values. The set spectrum method based on the union model, on the other hand, focuses on entries that are solely accessed by the failing traces, and does not get misled by a frequentist reasoning approach. The Tarantula-based binary similarity coefficients also focus on entries appearing in failing traces; however, unless additional weighting is provided to the failing entries by considering nearby entries as potential sources of failure, the Tarantula rankings are not effective. This also shows the value of using the metric space interpretation of LUT entries.

### B. Industrial Case Studies

*1) Debugging a Gain-Scheduled Control System:* In this case study, we look at $\mathcal{A}_{\text{CSTR}}$, a model of a continuously stirred tank reactor (CSTR), a common chemical system used in the industry. A Simulink® model of this process is available as a demonstration example from the Mathworks [TM]. The control objective is to ensure that the concentration of the reagent in the tank is maintained at a specified set-point. The model assumes that the controller can gets sensor readings of the residual concentration of the reagent in the tank, and is able to change the temperature of the coolant in reactor's cooling jacket to control the reaction. The control task is complex as the process dynamics are nonlinear and vary substantially with concentration. Hence, the system uses a Proportional + Integral + Derivative (PID) controller that is gain-scheduled, *i.e.*, uses different P+I+D gains for different reference set-points for the concentration. There are 8 different control regimes and 3 different lookup tables for the P, I and D gains respectively.

In this experiment, we introduce a bug in the model. We reverse the polarity of the P gain for the control regime corresponding to a concentration of 3 units (entry $M(3)$ in the P gain LUT). Outwardly, the bug is egregious; reversing the polarity of the P gain makes the closed loop system in that control regime unstable. However, once the system transitions to a different concentration regime, corrective feedback takes over, and the system eventually settles (possibly with a longer settling time). The STL requirement in (VI.5) captures that: (1) the system is excited by a step change in the reference at 2 seconds (enforced by the time-interval for the outermost $\Box$ operator), and (2) the maximum percentage deviation in the concentration signal $c(t)$ from the given settling region is less

| Ranking | Top 3 entries for $\mathcal{A}_{\text{nc}}$ (1D LUT, size 90) | | Top 3 entries for $\mathcal{A}_{\text{ff}}$ (2D LUT, size 1681) | |
| --- | --- | --- | --- | --- |
| Heuristic | Req. $\varphi_1^{\text{nc}}$ | Req. $\varphi_2^{\text{nc}}$ | Req. $\varphi^{\text{ff}}$ | AbsEXAMscore |
| $R_{\text{tarantula}}$ | − | − | (-9.5,-0.5) , (-9.5,0.0) , (0.0,-0.5) | − |
| $R_{\text{tarantula}}$ (metric) | − | − | (-10.0,7.5) , (-10.0,8.0), (-10.0,8.5) | 1 |
| $R_{\text{dstar}}^2$ | 1.7,1.4,1.1 | 2.0,1.7,2.3 | (-10.0,0.5) , (-10.0,1.0), (-10.0,4.5) | > 3 |
| $R_{\text{kulczynski}}$ (freq.) | 2.0,1.7,1.9 | 2.0,1.7,2.3 | (-9.5,0.5) , (-9.5,1.0) , (-10.0,1.5) | > 3 |
| $R_{\text{dstar}}^2$ (freq.) | 2.0,1.9,2.3 | 1.9,1.8,1.6 | (-10.0,0.5) , (-10.0,1.0), (-10.0,4.5) | > 3 |
| $R_U$ | − | − | (-10.0,10.0), (-10.0,9.0), (-10.0,9.5) | 1 |

TABLE II

THE MOST SIGNIFICANT 3 ENTRIES AS DEEMED BY THE RANKING HEURISTICS FOR EXPERIMENTS ON THE MODELS IN VI-A

| Ranking | Top 3 entries for $\mathcal{A}_{\text{CSTR}}$ (1D LUT, size 8) | Top 3 entries for $\mathcal{A}_{\text{AFC}}$ (1D LUT, size 11) | |
| --- | --- | --- | --- |
| Heuristic | Req. $\varphi_{\text{settle}}^{\text{cstr}}$ | Req. $\varphi_{\text{overshoot}}^{\text{afc}}$ | Req. $\varphi_{\text{settle}}^{\text{afc}}$ |
| $R_{\text{tarantula}}$ | 3,2,4 | 3250,3000,2750 | 3250,2750,3000 |
| $R_{\text{dstar}}^2$ | 4,3,5 | 2500,2250,2000 | 2250,2000,2500 |
| $R_{\text{kulczynski}}$ (freq.) | 4,3,5 | 2500,2250,2000 | 2250,2000,2500 |
| $R_{\text{dstar}}^2$ (freq.) | 4,3,5 | 2500,2250,2000 | 2250,2000,2500 |
| $R_U$ | − | − | − |

TABLE III

COLUMN ENTRIES CONTAIN THE TOP 3 ENTRIES DEEMED MOST IMPORTANT BY EACH RANKING HEURISTIC. THE SECOND COLUMN CONTAINS ENTRIES FOR THE CONTINUOUSLY STIRRED TANK REACTOR MODEL. THE THIRD & FOURTH COLUMNS CONTAIN RESULTS FOR AIR-FUEL RATIO CONTROL SYSTEM.

than 2% of the reference $r(t)$ after the given settling time deadline of 2.5 seconds.

$$\varphi_{\text{settle}}^{\text{cstr}} \triangleq \Box_{[2,20]}\left(\text{step}(r) \Rightarrow \Box_{[2.5,9.9]}\left(\left|\frac{c-r}{r}\right| < 0.02\right)\right) \quad \text{(VI.5)}$$

We ran 100 randomly chosen simulations and ranked entries that are the likeliest causes for failure of requirement $\varphi_{\text{settle}}^{\text{cstr}}$. *Results.* The results are shown in Table III. All the similarity coefficients except $R_U$ are able to find the seeded bug. In this case, the union spectrum based heuristic is not useful as there is no clear separation between entries accessed by failing and passing traces.

*2) Identifying Sensitive Regions in an Observer:* In this case study, we consider $\mathcal{A}_{\text{AFC}}$, a model [JDK+14] for regulating the air to fuel ratio (denoted $\lambda$) in the mixture that undergoes combustion in gasoline engines. As the peak efficiency of a catalytic converter to reduce noxious emissions in the exhaust is reached when $\lambda$ is 14.7, this is an important control problem. $\mathcal{A}_{\text{AFC}}$ has a controller that uses an observer to estimate the amount of fuel that puddles on the injection port and is thus not used for combustion in the engine. The observer is based on the Aquino model for fuel puddling in the controller [JDK+14]. This observer uses two LUTs for predicting the deposit ratio and residual ratio of the fuel from engine speed. Both LUTs have the same size and are indexed by the same quantity (engine speed), so the access patterns for entries of both LUTs are identical, and the entries shown in the results correspond to the corresponding entries for the same engine speed.

Let $\mu(t) = \frac{\lambda - 14.7}{14.7}$ and let $\theta(t)$ denote the throttle input signal from the user. Then, we wish to identify the regions of the observer LUTs that have the most impact on the maximum overshoot (5%) and settling time (1 sec.) requirements on $\lambda$ as shown in Eqns. (VI.6),(VI.7).

$$\varphi_{\text{overshoot}}^{\text{afc}} \triangleq \Box_{[2.5,10]}\left(\text{step}(\theta) \Rightarrow \Box(\mu < 0.05)\right) \quad \text{(VI.6)}$$
$$\varphi_{\text{settle}}^{\text{afc}} \triangleq \Box_{[2.5,10]}\left(\text{step}(\theta) \Rightarrow \Box_{[1,\infty]}(|\mu| < 0.01)\right) \quad \text{(VI.7)}$$

*Results.* Table III indicates that the observer region corresponding to high engine speed does not predict the state of the deposit ratio and residual ratio (two quantities to quantify the amount of fuel that puddles) accurately, causing the model to behave poorly at high engine speed. In [JDK+14], the authors indicate that the observer was designed based on a linear model at an operating point of 1000 rpm. Thus, the result above confirms the designer hypothesis that the observer performance is poor at high engine speeds.

*3) Parameter Space Optimization for a Suspension Control System:* In [MT], the authors present $\mathcal{A}_{\text{QS}}$, a quarter car model of an automotive suspension system implemented in Simulink®. $\mathcal{A}_{\text{QS}}$ considers only one of the wheels and simplifies the dynamics to that of a suspension mass suspended by two spring-damper systems. The controller uses a Proportional+Integral+Derivative (PID) scheme to provide active assistance to the suspension system. The design objective for this system is that the distance between the suspension mass and the vehicle body (denoted $y$) shows acceptable transient behavior, as sustained oscillations or a slow settling time on $y$ are a cause discomfort to the vehicle occupants.

Let the P, I and D gains be respectively $K_p, K_i$ and $K_d$. Let $K = (K_p, K_i, K_d)$. We optimize the controller performance by looking for regions in $K$-space that correlate well with good controller performance. In this model, there are no explicit LUTs, but we impose a grid on $K$-space and then use our tool to find grid regions with a high rank with respect to

the design objective. In contrast to other case studies, where we use ranking heuristics (all rankings except $R_{\text{tarantula}}$ had similar behaviour, so we picked a typical ranking function) to identify the cause of undesirable behavior, in this case, we use it to search for desirable behavior. As $K$-space does not have an explicit (user-defined) structure, this case study allows us a flexible way of exploring the (possibly nonlinear) parameter space. We explore a simple scheme, where we first search the $K$-space using a coarse grid, and then impose a finer grid on the top grid element found in the first iteration. In each iteration, we pick 200 simulations for randomly chosen values of the gains within the chosen paramter regions. This example requires large numbers for the $K_p$, $K_i$, $K_d$ gains, so for brevity in presentation, we assume that each of the gains in the following discussion are $\times 10^6$.

For the first iteration, we assume that $K_p \in [100, 500]$, $K_i \in [8, 600]$, and $K_d \in [0, 3]$. We use 21 equally spaced grid points in these intervals giving rise to a total of 9261 grid elements. In the first iteration, the tool identifies the grid element $(16, 18, 2)$ as the highest ranked grid point according to a simple majority of the rankings. This corresponds to the center of the region: $K_p \in [380, 420]$, $K_i \in [481.6, 540.8]$, $K_d \in [0, 0.3]$. We then impose a finer grid on this region, with 21 equally spaced points for each gain-parameter. After running the second iteration, we identified grid element $(13, 14, 2)$ as the top ranked element in the new grid. This corresponds to the region $K_p \in [402, 406]$, $K_i \in [517.12, 523.04]$, and $K_d \in [0, 0.03]$. Essentially, this is a narrow region for $K_p, K_i, K_d$ that correlates well with good controller performance (i.e. low settling time). This demonstrates the power of the tool for parameter optimization.

*4) **Model Predictive Control of a Diesel Engine Air Path***:
Next, we consider $\mathcal{A}_{\text{DAP}}$, an early prototype closed-loop Simulink® model from [HZBK16] of a **D**iesel engine **A**ir **P**ath controller. $\mathcal{A}_{\text{DAP}}$ has a high fidelity plant model of the air path dynamics and a model predictive controller (MPC) to regulate the intake manifold pressure and the exhaust gas recirculation (EGR) flow rate. A notable feature of $\mathcal{A}_{\text{DAP}}$ is its scale: it has more than 3000 Simulink® blocks and over 20 multi-dimensional lookup tables. $\mathcal{A}_{\text{DAP}}$ has two inputs: (1) the fuel injection rate (denoted *fr* and excited by a single step of magnitude in a given range), and the engine speed (denoted *ne*, picked from a given range, but held constant during any single simulation). There are two outputs of interest: the intake manifold pressure (denoted $p$) and the EGR flow rate (denoted *egr*). The control designers for $\mathcal{A}_{\text{DAP}}$ indicated to us that they are interested in two requirements: Requirement (VI.8) characterizes the overshoot in the intake manifold pressure. Requirement (VI.9) relates to how well the MPC scheme tracks the *egr* signal against the *egr* reference signal (denoted $egr_{\text{ref}}$). Let $\mu = \frac{egr - egr_{\text{ref}}}{egr_{\text{ref}}}$.

$$\varphi_{\text{overshoot}}^{\text{diesel}} \triangleq \square_{[2,10]}(\text{step}(fr) \Rightarrow \square_{[0,10]}(p < p_{\max})) \text{ (VI.8)}$$
$$\varphi_{\text{settle}}^{\text{diesel}} \triangleq \square_{[2,10]}(\text{step}(fr) \Rightarrow \square_{[\tau,10]}(|\mu| < v)) \text{ (VI.9)}$$

As this model is proprietary, we suppress the numeric values for the settling time ($\tau$), the settling region ($v$), and the maximum allowed overshoot ($p_{\max}$). For the LUT entries, we scale the actual values to representative integer values. We pick an LUT identified as the most important for analysis by the control designer. This is a 2D $20 \times 15$ LUT.

*Results.* We show the results obtained by applying our tool (for a typical similarity coefficient ranking) on 500 random simulation runs of $\mathcal{A}_{\text{DAP}}$ in Table IV. For this particular experiment, the LUT entries deemed most important have a direct interpretation as the inputs to the LUT are model inputs *fr* and *ne*. Thus, each entry in the LUT corresponds to a range of *fr* and *ne* values. We remark that the entry $(2, 9)$ corresponds to a low fuel injection rate and higher engine speed scenario, while the entry $(12, 1)$ corresponds to a higher fuel injection rate and very low engine speed scenario. We also remark that for the settling time requirement, we find entries accessed only by the failing trajectories. Interestingly, the top three entries, although from different regions of the LUT, all correspond to boundary (edge) cases of the LUT. This indicates that the MPC-based controller has poorer performance w.r.t. the settling time requirements at boundary conditions.

In Fig. 1, we present the results in the form of a heat-map that pinpoints hot-spots in the chosen look-up table w.r.t. a given requirement. The bottom portion of the heat-map consists of entries not accessed by any simulation trace; this is because the control designer indicated interest only in fuel injection rates less than a certain amount. Thus entries corresponding to values greater than this amount were never accessed. This also shows another value of our tool: it allows visualizing *coverage* of LUT entries or the parameter space by a given set of simulation runs.

*5) **Study of Responsiveness in a Hydrogen Fuel-Cell Vehicle***:
Next, we use a prototype Model-In-the-Loop-Simulation (MILS) model of an airpath control model from a hydrogen fuel-cell (FC) vehcile powertrain. This model has more than 7000 Simulink® blocks that give a detailed description of the physics of the airpath, along with a simplified model of the power management, and a complex controller with several look-up maps to regulate the flow of air through the fuel-cell stack. A key requirement of the closed-loop system is responsiveness: how well does the system react to a driver's request for increased torque. Internally, a torque request gets translated to a request for increased air-flow through the stack. Thus the rise time on the air-flow rate signal is a good proxy for system responsiveness (STL requirement (VI.10)). In the requirement, $r$ is the rise-time, and $\lambda$ is a suitable number in $[0.5, 1]$.
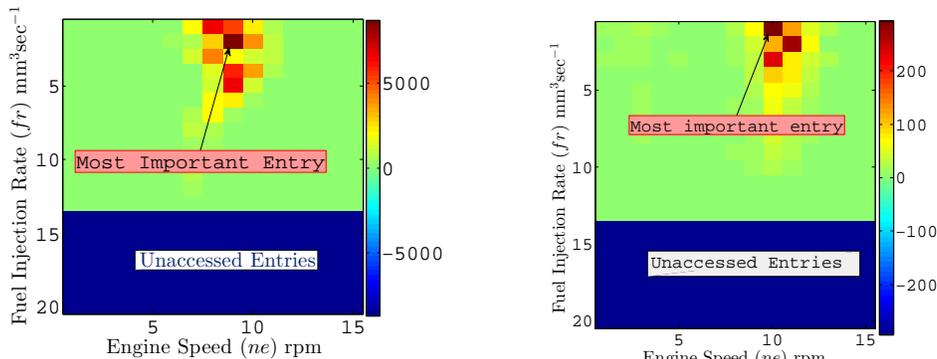
$$\varphi_{\text{rise}}^{\text{FC}} \triangleq \square_{[0,T]}\big(\text{step}(\text{AFR}_{\text{ref}}) \Rightarrow \diamondsuit_{[0,r]}(\text{AFR} > \lambda \cdot \text{AFR}_{\text{ref}})\big)(\text{VI.10})$$

In this case study, we choose three key controller look-up maps and study the correlation between accessing a certain region of the look-up map with the responsiveness of the closed-loop system. As the model is proprietary, we suppress the values on the axes.

| Ranking Heuristic | Top 3 entries for $\mathcal{A}_{\mathrm{DAP}}$ (2D LUT, size 300) | |
| --- | --- | --- |
| | Req. $\varphi_{\mathrm{overshoot}}^{\mathrm{diesel}}$ | Req. $\varphi_{\mathrm{settle}}^{\mathrm{diesel}}$ |
| $R_{\mathrm{dstar}}^2$ | (2,9),(1,9),(2,10) | (2,11),(1,11),(4,11) |
| $R_{\mathrm{kulczynski}}$ (freq.) | (2,9),(1,9),(2,10) | (2,11),(1,11),(2,10) |
| $R_{\mathrm{dstar}}^2$ (freq.) | (2,9),(1,9),(2,10) | (2,11),(1,11),(2,12) |
| $R_{\mathrm{dstar}}^2$ (freq.+metric) | (2,9),(5,9),(1,8) | (1,10),(2,11),(3,10) |
| $R_U$ | − | (12,1),(10,10),(11,10) |

TABLE IV

TOP 3 ENTRIES DEEMED THE MOST IMPORTANT BY DIFFERENT RANKING HEURISTICS FOR MODEL PREDICTIVE OF A DIESEL ENGINE AIRPATH. MODEL.



(a) Importance of LUT entries for the requirement related to overshoot on the intake manifold pressure during up-steps in fuel injection. The results shown are w.r.t. the $R_{\mathrm{dstar}}$ heuristic incorporating frequency of access and metric-space interpretation of the table.

(b) Importance of LUT entries for the requirement related to settling time on the EGR rate on down-steps in fuel injection. The results shown are w.r.t. the $R_{\mathrm{dstar}}$ heuristic incorporating Boolean access of entries and a metric-space interpretation of the table.

Fig. 1. Heatmap representation for LUT entries for the $\mathcal{A}_{\mathrm{DAP}}$ model. Entries in blue are not accessed by any trace. The color spectrum from light green (least) to dark red (most) indicates importance of the entry. The numbers on the color-bar are the rank values, that can be used to have a quantitative interpretation of the relative importance of entries. Observe that the graphical depiction may allow a broader judgement regarding problematic regions in the LUT.

*Results.* Figure 2 depicts the heat-map of the LUT rankings for a chosen LUT. All three LUTs show similar heat-maps, and we only show one due to lack of space. The actual signals used to index the look-up maps correspond to the pressure ratio across a compressor component and the air-flow through the compressor. Our analysis shows that the model is not responsive at low air-flow rates and pressure ratio values. The designers confirmed that this analysis was of interest to them, and indicated plans to improve the model performance for these conditions.
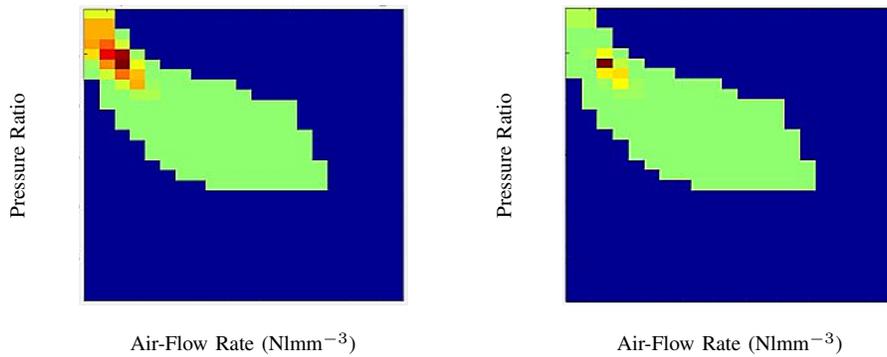
## VII. CONCLUSION

In this paper we present a set of easy-to-compute statistical correlation based rankings in order to localize parameters in control software which may be causing undesired model behavior in controlled cyber-physical systems. We empirically test the ranking heuristics provided by these methods on a number of case studies that are relevant in an industrial cyber physical systems context using our tool integrated into a simulation-guided falsification workflow for Simulink® models. It is a perfectly reasonable expectation of the reader that we suggest a single ranking scheme that should be generally used or present some guidelines to pick a ranking scheme for analysis. Unfortunately, our experiments show that each ranking scheme has its own merits, and any guidelines would

have to rely on deep knowledge of the model structure and dynamics. It is arguable that the union spectrum ranking scheme, when it gives a result should not be ignored by the designer. For the other ranking schemes, our suggestion is to use them to get an overall picture of the "problem regions" in the LUT, using graphical visualization tools such as heat-maps.

## REFERENCES

[ADV+17] Nikos Aréchiga, Sumanth Dathathri, Shashank Vernekar, Nagesh Kathare, Sicun Gao, and Shinichi Shiraishi. Osiris: A tool for abstraction and verification of control software with lookup tables. In *SCAV'17*, pages 11–18. ACM, 2017.

[AHLW95] Hiralal Agrawal, Joseph R. Horgan, Saul London, and W. Eric Wong. Fault localization using execution slices and dataflow tests. In *ISSRE 1995*, pages 143–151. IEEE, 1995.

[AR14] Charu C. Aggarwal and Chandan K. Reddy, editors. *Data Clustering: Algorithms and Applications*. CRC Press, 2014.

[BSM06] Christian Bohn, Pascal Stober, and Olaf Magnor. An optimization-based approach for the calibration of lookup tables in electronic engine control. In *Computer Aided Control System Design, 2006 IEEE*, pages 2315–2320, 2006.

[CC10] Seung Seok Choi and Sung Hyuk Cha. A survey of binary similarity and distance measures. *Journal of Systemics, Cybernetics and Informatics*, pages 43–48, 2010.

[DFM13] Alexandre Donzé, Thomas Ferrère, and Oded Maler. Efficient robust monitoring for STL. In *CAV*, LNCS 8044, pages 264–279. Springer, 2013.

(a) Importance of LUT entries for the requirement related to responsiveness of $\mathcal{A}_{\text{FC}}$. The results shown are w.r.t. the $R_{\text{dstar}}$ heuristic not incorporating the frequency of access and metric-space interpretation of the table.

(b) Importance of LUT entries for the requirement related to responsiveness of $\mathcal{A}_{\text{FC}}$. The results shown are w.r.t. the $R_{\text{dstar}}$ heuristic incorporating the frequency of access and metric-space interpretation of the table.

Fig. 2. Heatmap representation for LUT entries for the $\mathcal{A}_{\text{FC}}$ model.

[DM10]     Alexandre Donzé and Oded Maler. Robust satisfaction of temporal logic over real-valued signals. In *FORMATS 2010*, LNCS 6246, pages 92–106. Springer, 2010.

[Don10]    Alexandre Donzé. Breach, a toolbox for verification and parameter synthesis of hybrid systems. In *CAV*, pages 167–170, 2010.

[GS00]     Mariano Gasca and Thomas Sauer. Polynomial interpolation in several variables. *Advances in Computational Mathematics*, 12(4):377–410, 2000.

[HC13]     Bingzhao Gao Hong Chen. *Nonlinear Estimation and Control of Automotive Drivetrains*. Springer, 2013.

[HPAW15]   Birgit Hofer, Alexandre Perez, Rui Abreu, and Franz Wotawa. On the empirical evaluation of similarity coefficients for spreadsheets fault localization. *Autom. Softw. Eng.*, 22(1):47–74, 2015.

[HZBK16]   M. Huang, K. Zaseck, K. Butts, and I. Kolmanovsky. Rate-based model predictive controller for diesel engine air path: Design and experimental evaluation. *IEEE Trans. on Control Systems Technology*, PP(99):1–14, 2016.

[JDK+14]   Xiaoqing Jin, Jyotirmoy Deshmukh, James Kapinski, Koichi Ueda, and Ken Butts. Benchmarks for model transformations and conformance checking. In *ARCH*, 2014.

[JGG08]    Dennis Jeffrey, Neelam Gupta, and Rajiv Gupta. Fault localization using value replacement. In *ISSTA '08*, pages 167–178. ACM, 2008.

[JH05]     James A. Jones and Mary Jean Harrold. Empirical evaluation of the tarantula automatic fault-localization technique. In *ASE '05*, pages 273–282. ACM, 2005.

[JHS02]    James A. Jones, Mary Jean Harrold, and John T. Stasko. Visualization of test information to assist fault localization. In *ICSE 2002*, pages 467–477. ACM, 2002.

[JS07]     Lingxiao Jiang and Zhendong Su. Context-aware statistical debugging: From bug predictors to faulty control flow paths. In *ASE '07*, pages 184–193. ACM, 2007.

[JSH89]    Donald A. Jackson, Keith M. Somers, and Harold H. Harvey. Similarity coefficients: Measures of co-occurrence and association or simply measures of occurrence? *The American Naturalist*, 133(3):436–453, 1989.

[LNZ+05]   Ben Liblit, Mayur Naik, Alice X. Zheng, Alexander Aiken, and Michael I. Jordan. Scalable statistical bug isolation. In *PLDI 2005*, pages 15–26. ACM, 2005.

[LYF+05]   Chao Liu, Xifeng Yan, Long Fei, Jiawei Han, and Samuel P. Midkiff. Sober: Statistical model-based bug localization. *SIGSOFT Softw. Eng. Notes*, 30(5):286–295, September 2005.

[MN13]     Oded Maler and Dejan Nickovic. Monitoring properties of analog and mixed-signal circuits. *STTT*, 15(3):247–268, 2013.

[MT]       Bill Messner and Dawn Tilbury. Control tutorials for matlab and simulink.

[NLR11]    Lee Naish, Hua Jie Lee, and Kotagiri Ramamohanarao. A model for spectra-based software diagnosis. *ACM Trans. Softw. Eng. Methodol.*, 20(3):11:1–11:32, 2011.

[PAR16]    Lúcio S. Passos, Rui Abreu, and Rosaldo J. F. Rossetti. Empirical evaluation of similarity coefficients for multiagent fault localization. *IEEE Trans. Systems, Man, and Cybernetics: Systems*, To Appear, 2016.

[PO11]     Chris Parnin and Alessandro Orso. Are automated debugging techniques actually helping programmers? In *ISSTA '11*, pages 199–209. ACM, 2011.

[PS92]     Hsin Pan and Eugene H. Spafford. Heuristics for automatic localization of software faults. Technical report, Purdue University, 1992.

[RFZO12]   Jeremias Rößler, Gordon Fraser, Andreas Zeller, and Alessandro Orso. Isolating failure causes through test case generation. In *ISSTA'12*, pages 309–319. ACM, 2012.

[RR03]     Manos Renieris and Steven P. Reiss. Fault localization with nearest neighbor queries. In *(ASE 2003)*, pages 30–39. IEEE, 2003.

[sim07]    *Using Simulink*. The MathWorks, 2007.

[TM]       The Mathworks. Gain-scheduled control of a chemical reactor.

[WDGL14]   W. Eric Wong, Vidroha Debroy, Ruizhi Gao, and Yihao Li. The dstar method for effective software fault localization. *IEEE Trans. Reliability*, 63(1):290–308, 2014.

[WGL+16]   W. Eric Wong, Ruizhi Gao, Yihao Li, Rui Abreu, and Franz Wotawa. A survey on software fault localization. *IEEE Trans. Software Eng.*, 42(8):707–740, 2016.

[YHC13]    Shin Yoo, Mark Harman, and David Clark. Fault localization prioritization: Comparing information-theoretic and coverage-based approaches. *ACM Trans. Softw. Eng. Methodol.*, 22(3):19:1–19:29, July 2013.

**Basic Quantitative Similarity Coefficients** Let $\mathcal{Z}$ be the set of executions. We assume that the code is instrumented such that each execution $z$ stores a time-ordered list of map indices that are accessed during the execution. We denote by $|z|$ the number of function queries in $z$, that is, if in $z$ we get queries for $f_M^{\mathcal{I}}(m_1), f_M^{\mathcal{I}}(m_2), \dots f_M^{\mathcal{I}}(m_p)$ (the queries need not be distinct), then $|z|$ is $p$. We assume that the interpolation scheme uses a constant number of value to construct an interpolation, that is we assume $|\mathsf{depend}(m)|$ to be a constant. Thus, the time-ordered list of map indices that are accessed during the execution is of size $O(|z|)$. We let $\mathsf{size}(\mathcal{Z})$ denote $\sum_{z \in \mathcal{Z}} |z|$. Note that the total size of the time-ordered list of map indices that are accessed during the executions in $\mathcal{Z}$ is at most $\mathsf{size}(\mathcal{Z})$. The binary similarity coefficients for $m \in \mathsf{dom}(M)$ can be computed in time linear in the size of this time-ordered list. Finally, the ranking is done by sorting the indices based on the similarity coefficients. Thus, the total time required is $O\left(\mathsf{size}(\mathcal{Z}) + |M| \log(|M|)\right)$.

**Similarity Coefficients utilizing Continuity and the Metric Space structure** Let $\mathcal{Z}$ be the set of executions. In this case, we assume that the code is instrumented such that each execution $z$ stores, a time-ordered list of (i) map indices that are accessed during the execution, and (ii) function arguments to $f_M^{\mathcal{I}}$ that are queried for, *i.e.*, if the execution queries $f_M^{\mathcal{I}}(m_1), f_M^{\mathcal{I}}(m_2), \dots$, then we store $m_1, \mathsf{depend}(m_1), m_2, \mathsf{depend}(m_2), \dots$. From this list for $z$, we construct another list containing $\mathsf{raffect}(m)$ values for every $m \in \mathsf{dom}(M)$, where $\mathsf{raffect}$ is as in Equation (III.6). This can be done by first building another list which, for each $m_k$ in the original list (arising from a $f_M^{\mathcal{I}}(m_k)$ query), creates a sublist of all $\mathsf{maffect}(m_k, m)$ values for $m \in \mathsf{dom}(M)$ such that $\mathsf{maffect}(m_k, m) > 0$. In case $r_M$ is $\infty$, this takes $O(|M|)$ time for each $m_k$. In case $r_M$ is finite, this takes $O(|M|_{r_M})$ time, where $|M|_{r_M}$ denotes the maximum number of indices of $M$ in an $r_M$ sized ball (in the indices space). From this secondary list, the list of $\mathsf{raffect}$ values can be constructed in linear time, and the coefficients constructed. Thus, the total running time is $O\left(\mathsf{size}(\mathcal{Z}) \cdot |M|_{r_M} + |M| \log(|M|)\right)$. Thus, having a small radius $r_M$ allows us to avoid a (possibly quadratic) blowup in the running time as compared to the basic similarity coefficient approach.

**Similarity Coefficients incorporating frequency of accessing a map** For both definitions of $\mathsf{maffect}$, we can proceed as follows. Sort the arguments for queries in each execution $z$ (*i.e.*, sort the $m_k$ values where $f_M^{\mathcal{I}}(m_k)$ is called in $z$), based on some ordering of the indices; and then add up $\mathsf{maffect}()$ values for the same $m_k$ values. This step takes $O\left(|z| \cdot \log(|z|)\right)$ time, and gives a list of the $\mathsf{fraffect}$ values for $z$. The similarity coefficients then be calculated in linear time, and after that we need to sort $M$ elements to get a ranking. Thus, the total time required is $O\left(\sum_{z \in \mathcal{Z}} (|z| \cdot \log(|z|)) + |M| \log(|M|)\right)$.

**Union Spectrum method** The sets $M_F \subseteq M$ and $M_S \subseteq M$ can be computed in time $O(|M| + \mathsf{size}(\mathcal{Z}))$ as a sorted list. The set $\mathsf{Ball}(M_S, r_M)$ can be obtained as a sorted list in time $O\left(|M_S| \cdot |M|_{r_M}\right)$ where $|M|_{r_M}$ denotes the maximum number of indices of $M$ in an $r_M$ sized ball (in the indices space). The set difference $M_F \setminus \mathsf{Ball}(M_S, r_M)$ can be computed time $O\left(|M_F| + |\mathsf{Ball}(M_S, r_M)|\right)$. Putting everything together, we get that $\mathsf{sus}_U$ can be computed in $O\left(\mathsf{size}(\mathcal{Z}) + |M| \cdot |M|_{r_M}\right)$ time. The values $s_U(m)$ can be inferred for all $m \in \mathsf{dom}(M)$ by maintaining some additional bookkeeping in the above algorithm without increasing the running time complexity. The values $d_U(m)$ can also be computed from the sets $M_F$ and $M_S$ augmented with some additional bookkeeping in linear time. Finally, we need to sort according the values $R_U(m)$. The total running time works out to be $O\left(\mathsf{size}(\mathcal{Z}) + |M| \cdot |M|_{r_M} + |M| \cdot \log(|M|)\right)$.