# Predicting Performances of Configurable Systems: the Threat of Input Sensitivity

Anonymous Author(s)

## ABSTRACT

Widely used software systems such as video encoders are by necessity highly configurable, with hundreds or even thousands of options to choose from. Their users often have a hard time finding suitable values for these options (*i.e.*, finding a proper *configuration* of the software system) to meet their goals for the tasks at hand, *e.g.*, compress a video down to a certain size. Several research work have thus proposed to resort on machine learning to predict the effect of options on such functional and performance properties, *i.e.*, build a performance model. However, while domain experts seem to be well aware of input sensitivity issues when *e.g.*,, crafting benchmarks, most of these approaches do not consider the variability of inputs to those software systems (*e.g.*, a video as input to an encoder like *x264* or a file system fed to a tool like *xz*). In this problem-statement paper, we conduct a large study over 8 configurable systems and their inputs that shows the influence of inputs on performance prediction. The results exhibit that inputs fed to systems can interact with software options, impacting their performance properties significantly. When tuning a configuration for an input, we can multiply performances up to a factor a 10. We thus warn researchers: a performance model trained on one input cannot be generalized over a large dataset of inputs.

## 1 INTRODUCTION

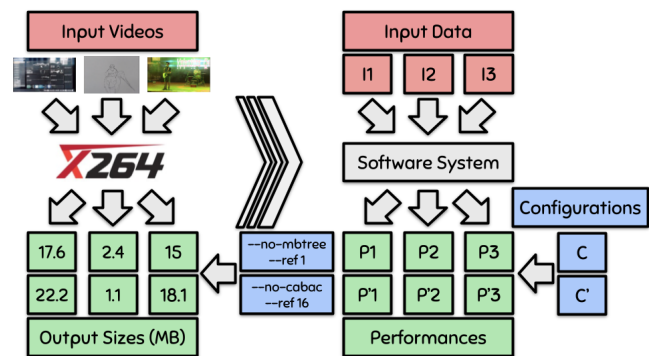Widely used software systems are by necessity highly configurable, with hundreds or even thousands of options to choose from. For example, a tool like *xz* offers various options such as –threads or –format for compressing a file system. The same applies to *Linux* kernels or video encoders such as *x264*: they all provide configuration options through compilation options, feature toggles, and command-line parameters. Software engineers often have a hard time finding suitable values for these options (*i.e.*, finding a proper configuration of the software system) to meet their goals for the tasks at hand, *e.g.*, generate code as compact as possible or compress a video down to a certain size while keeping its perceived quality. Since the number of possible configurations grows exponentially

with the number of options, even experts may end up recommending sub-optimal configurations for such complex software [38].

Research work have shown that machine learning can solve this problem and predict the performance of configurations [28, 75, 88, 98]. These works measure the performances of a configuration sample under specific settings to then build a model capable of predicting the performance of any other configuration, *i.e.*, a performance model. However, there exist cases where inputs (*e.g.*, a file fed to an archiver like *xz* or a video provided as input to an encoder like *x264*, as shown in Figure 1) can also impact the performances of a configurable system [5]. The *x264* encoder typifies this problem. For example, Kate, an engineer working for a VOD company, wants an *x264* performance model for predicting the resulting bitrate based on selected *x264* configuration options. She would build such a model using several input videos to learn about the impact of these *x264* configuration options on the encoding bitrate. Now, if Kate wants to reuse this model for new input videos, she might ask herself the following questions: Will this performance model predict an accurate value for the bitrate? Do configuration options have the same effect on the bitrate despite a different input? Do options interact in the same way no matter the inputs? These are crucial practical issues since inputs fed to configurable systems can question performance prediction models developed so far. In particular, prior research works ignore whether these prediction models generalize across different inputs and thus can be reused. If they cannot, Kate would have to train a new performance model, the worst situation being to learn as many performance models as there are inputs, making her work basically useless for a field deployment.



This paper investigates, demonstrates and quantifies what effects configurable system inputs have on performances.

**Figure 1: The diversity of inputs fed to configurable systems.**

In this work, we first conduct in-depth empirical study on the impact of inputs on configurable systems. We systematically explore the impacts of inputs on the quantitative properties of different configurations of 8 software systems. This study reveals that inputs fed to systems can interact with software options, significantly

impacting their performance properties. It also shows that a performance model fed on one input cannot be generalized over a large dataset of inputs. We then survey state-of-the-art papers on performance prediction to assess whether they address input sensitivity. We found that even when some of them acknowledge the issue, most of them do not mention or address it, which questions their practical relevance.

In summary, the contributions of this paper are as follows:

- To our best knowledge, the first large empirical study that investigates the impact of input on performances of 8 configurable systems;
- We show that **input data matter as (much as) configuration options** when predicting the performances of configurable systems processing inputs;
- An analysis of how 63 state-of-the-art research papers address the input sensibility problem in practice;
- Open science: a replication bundle that contains docker images, produced datasets of measurements and code.[1]

The remainder of this paper is organized as follows: Section 2 explains the problem of input sensitivity and the research questions addressed in this paper. Section 3 presents the experimental protocol. Section 4 details the results. Section 5 shows their significance in research. Section 6 discusses the implications of our work. Section 7 details threats to validity. Section 8 presents related work. Section 9 summarizes key insights of our paper.

**Typographic Convention.** For this paper, we adopt the following typographic convention: *emphasized* will be relative to a software system, *slanted* to its configuration options and underlined to its performance properties.

## 2  PROBLEM STATEMENT

### 2.1  Sensitivity to Inputs of Configurable Systems

Numerous works have proposed to model performance of software configurations, with several use-cases in mind for developers and users of software systems: the maintenance and understanding of configuration options and their interactions [78], the selection of an optimal configuration [24, 63, 65], the automated specialization of configurable systems [84, 85].

Configuration options of software systems can have different effects on performance (*e.g.*, runtime), but so can the input data. For example, a configurable video encoder like *x264* can process many kinds of inputs (videos) in addition to offering options on how to encode. Our hypothesis is that there is an interplay between configuration options and input data: some (combinations of) options may have different effects on performances depending on input.

Input sensitivity may have a strong impact on engineering and research work. Developers of configurable systems that process input data should be aware of this phenomenon and test their systems on a wide variety of inputs [73]. Similarly, researchers who develop learning algorithms or optimization techniques may want to benchmark them on a realistic set of inputs to draw conclusions as general as possible on configuration spaces [14]. This notably

concerns learning models that predict performances of configurations.

Researchers observed input sensitivity in multiple fields, such as SAT solvers [21, 97], compilation [16, 69], video encoding [57], data compression [46]. However, existing studies either consider a limited set of configurations (*e.g.*, only default configurations), a limited set of performance properties, or a limited set of inputs [1, 12, 22, 26, 51, 71, 81]. It limits some key insights about the input sensitivity of configurable systems. It is also a threat to validity since inputs together with options may change the underlying distributions and thus the relevance of performance models.

This work details, to the best of our knowledge, the first systematic empirical study that analyzes the influence of input data on performances for different configurable systems. Through three research questions introduced in the next section, we characterise the input sensitivity problem and study how this can affect practitioners and researchers training machine learning models to predict performances of configurable systems.

### 2.2  Research Questions

Input data have different properties that may change the software behavior and thus alter software performances of configurations. To reuse (or transfer) a performance model on multiple inputs (*i.e.*, trained on one input and tested on another input), the performance of software configurations should be similar across inputs. Hence an hypothesis is that two performance models over two different inputs are somehow related and close. In its simplest form, there is a linear relationship between these two models: their performance predictions simply increase or decrease with each other. However, more complex mappings can exist since the underlying performance distributions differ.

> **RQ$_1$ - Do software performances stay consistent across inputs?** Are the performance distributions stable from one input to another? Are the rankings of performance the same for all inputs?

But software performances are influenced by the configuration options *e.g.*, the energy consumption [13]. An option is called *influential* for a performance when its values have a strong effect on this performance [17, 40]. For example, developers might wonder whether the option they add to a configurable software has an influence on its performance. However, is an option identified as influential for some inputs still influential for other inputs? If not, it would become both tedious and time-consuming to find influential options on a per-input basis. Besides, it is unclear whether activating an option is always worth it in terms of performance; an option could improve the overall performance while reducing it for few inputs. If so, users may wonder which options to enable to improve software performances based on their input data.

> **RQ$_2$ - Do configuration option's effects change with input data?** Do the configuration options have the same effects for all inputs? Is an influential option influential for all inputs? Do the effects of configuration options vary with input data?

*RQ$_1$* and *RQ$_2$* study how inputs affect (1) performance distributions and (2) the effects of different configuration options. However,

---

[1]Available on Github:
https://anonymous.4open.science/r/df319578-8767-47b0-919d-a8e57eb67d25/

Table 1: Subject *System*s. *Domain* the area of expertise using the system. *Commit* the git commit (*i.e.*, the version) of the system. *Configs #C* the number of configurations tested per system. *Inputs I* the type of input fed to the system. *#I* the number of inputs per system. *#M* the total number of measurements, #M = #I*#C. *Performance(s) P* the performance properties measured per system. *Docker* the links to the containers to replicate the measurements. *Dataset* the links to the measurements.

| System | Domain | Commit | Configs #C | Inputs I | #I | #M | Performance(s) P | Docker | Dataset |
|--------|--------|--------|-----------|----------|-----|-----|------------------|--------|---------|
| gcc | Compilation | ccb4e07 | 80 | .c programs | 30 | 2400 | size, ctime, exec | Link | Link |
| ImageMagick | Image processing | 5ee49d6 | 100 | images | 1000 | 100 000 | size, time | Link | Link |
| lingeling | SAT solver | 7d5db72 | 100 | SAT formulae | 351 | 35 100 | #conflicts, #reductions | Link | Link |
| nodeJS | JS runtime env. | 78343bb | 50 | .js scripts | 1939 | 96 950 | #operations/s | Link | Link |
| poppler | PDF rendering | 42dde68 | 16 | .pdf files | 1480 | 23 680 | size, time | Link | Link |
| SQLite | DBMS | 53fa025 | 50 | databases | 150 | 7500 | 15 query times q1-q15 | Link | Link |
| x264 | Video encoding | e9a5903 | 201 | videos | 1397 | 280 797 | size, time, cpu, fps, bitrate | Link | Link |
| xz | Data compression | e7da44d | 30 | system files | 48 | 1440 | size, time | Link | Link |

the performance distributions could change in a negligible way, without affecting the software user's experience. Before concluding on the real impact of the input sensitivity, it is necessary to quantify how much this performance changes from one input to another.

> **RQ$_3$ - Can we ignore input sensitivity?** If we do, what is the loss in performance considering that all input data is the same and does not affect the software that processes it? Or, to put it more positively, what is the potential gain to tune a software system for its workload?

## 3 EXPERIMENTAL PROTOCOL

To answer these research questions, we have designed the following experimental protocol.

### 3.1 Data collection

We first collect performance data of configurable systems that process inputs.

   **Protocol.** Figure 2 depicts the step-by-step protocol we respect to measure performances of software systems. Each line of Table 1 should be read following Figure 2: *System* and *Domain* with Step 1; *Commit* with Step 2; *Configs #C* with Step 3; *Inputs I* and *#I* with Step 4; *#M* with Step 5; *Performance(s) P* with Step 6; *Docker* links a container for executing all the steps; *Dataset* links the results of the protocol *i.e.*, the datasets containing the performance measurements. An example with the *x264* encoder is shown in beige on Figure 2. Hereafter, we provide details for each step of the protocol.

   **Steps 1 & 2 - Software Systems.** We consider 8 software systems, open-source and well-known in various fields, that already have been studied in the literature: *gcc* [69], *ImageMagick* [82], *lingeling* [34], *nodeJS* [36], *poppler* [55], *SQLite* [84], *x264* [39] and *xz* [90]. We choose these systems because they handle different types of input data, allowing us to draw as general conclusions as possible. For each software system, we use a unique private server with the same configuration running over the same operating system.[2] We download and compile a unique version of the system, related to the git *Commit* in Table 1. All performances are measured with this version of the software.

   **Step 3 - Configuration options C.** To select the configuration options, we read the documentation of each system. We manually
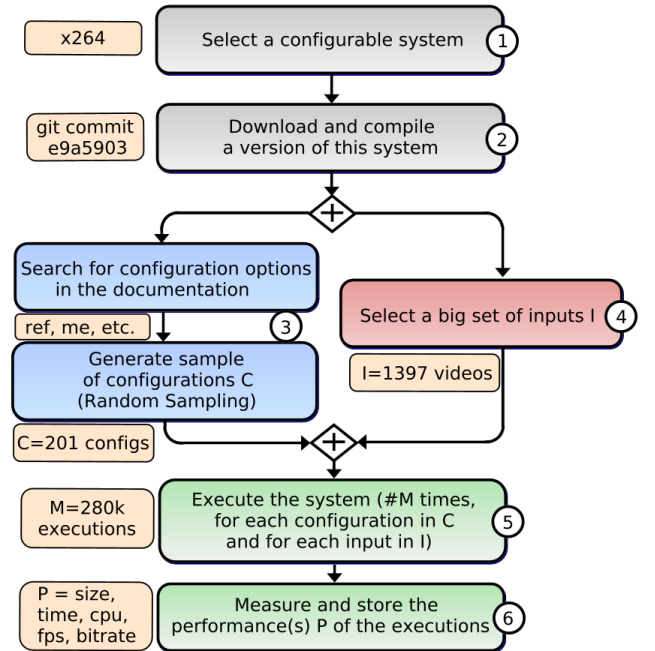


Figure 2: Measuring performances - Protocol

extracted the options affecting the performances of the system according to the documentation. We then sampled #C configurations by using random sampling [68]. We checked the uniformity of the different option values with a Kolmogorov-Smirnov test [56] applied to each configuration option.[3]

   **Step 4 - Inputs I.** For each system, we selected a different set of input data: for *gcc*, PolyBench v3.1 [72]; for *ImageMagick*, a sample of ImageNet [15] images (from 1.1 kB to 7.3 MB); for *lingeling*, the 2018 SAT competition's benchmark [34]; for *nodeJS*, its test suite; for poppler, the Trent Nelson's PDF Collection [64]; for *SQLite*, a set of generated TPC-H [70] databases (from 10 MB to 6 GB); for *x264*, the YouTube User General Content dataset [93] of videos (from 2.7 MB to 39.7 GB); for *xz*, the Canterbury corpus [54]. These are large, well-known and freely available datasets of inputs.

   **Steps 5 & 6 - Performance properties P.** For each system, we systematically executed all the configurations of C on all the

---

[2]The configurations of the running environments are available at : https://anonymous.4open.science/r/df319578-8767-47b0-919d-a8e57eb67d25/replication/Environments.md

[3]Options and tests results are available at : https://anonymous.4open.science/r/df319578-8767-47b0-919d-a8e57eb67d25/results/others/configs/sampling.md

inputs of I. For the #M resulting executions, we measured as many performance properties as possible: for *gcc*, ctime and exec the times needed to compile and execute a program and the size of the binary; for *ImageMagick*, the time to apply a *Gaussian blur* [35] to an image and the size of the resulting image; for *lingeling*, the number of reductions and conflicts found in 10 seconds of execution; for *nodeJS*, the number of operations per second (ops) executed by the script; for *poppler*, the time needed to extract the images of the pdf, and the size of the images; for *SQLite*, the time needed to answer 15 different queries q1-q15; for *x264*, the size of the compressed video, the elapsed time, the cpu usage (percentage), the bitrate (the average amount of data encoded per second) and the average number of frames encoded per second (fps); for *xz*, the size of the compressed file, and the time needed to compress it.

**Replication.** To allow researchers to easily replicate the measurement process, we provide a docker container for each system (see the links in the *Docker* column of Table 1). We also publish the resulting datasets online (see the links in the *Dataset* column) and in the companion repository with additional replication details.[4]

For the next research questions, our results are computed with Python v3.7.6 and specific versions of data science libraries.[5]

## 3.2 Performance correlations ($RQ_1$)

Based on the analysis of the data collected in Section 3.1, we can now answer the first research question: **$RQ_1$ - Do software performances stay consistent across inputs?** To check this hypothesis, we compute, analyze and compare the Spearman's rank-order correlation [45] of each couple of inputs for each system.

**Spearman correlations.** The correlations are considered as a measure of similarity between the configurations' performances over two inputs. We compute the related *p*-values: a correlation whose *p*-value is higher than the chosen threshold 0.05 is considered as null. We use the Evans rule [20] to interpret these correlations. In absolute value, we refer to correlations by the following labels; very low: 0-0.19, low: 0.2-0.39, moderate: 0.4-0.59, strong: 0.6-0.79, very strong: 0.8-1.00. A negative score tends to reverse the ranking of configurations. Very low or negative scores have practical implications: a good configuration for an input can very well exhibit bad performances for another input.

## 3.3 Effects of options ($RQ_2$)

To understand how a performance model can change based on a given input, we next study how input data interact with configuration options. **$RQ_2$ - Do configuration option's effects change with input data?** To assess the relative significance and effect of options, we use two well-known statistical methods [9, 76].

**Random Forest Importances.** The tree structure provides insights about the most essential options for prediction, because such a tree first splits w.r.t. options that provide the highest information gain. We use random forests [9], a vote between multiple decision trees: we can derive, from the forests trained on the inputs, estimates of the options importance. The computation of option importance is realized through the observation of the effect on random forest accuracy of randomly shuffling each predictor variable [58]. For a random forest, we consider that an option is influential if the median (on all inputs)o f its option importance is greater than $\frac{1}{n_{opt}}$, where $n_{opt}$ is the number of options considered in the dataset. This threshold represents the theoretic importance of options for a software having equally important options (inspired by the Kaiser rule [101]).

**Linear Regression Coefficients.** The coefficients of an Ordinary Least Square (OLS) regression [76] weight the effect of configuration options. These coefficients can be positive (resp. negative) if a bigger (resp. lower) option value results in a bigger performance. Ideally, the sign of the coefficients of a given option should remain the same for all inputs: it would suggest that the effect of an option onto performance is stable. We also provide details about coefficients related to the interactions of options (*i.e.*, feature interactions [75, 88]) in $RQ_2$ results.

## 3.4 Impact of input sensitivity ($RQ_3$)

To complete this experimental protocol, we ask whether adapting the software to its input data is worth the cost of finding the right set of parameters *i.e.*, the concrete impact of input sensitivity. **$RQ_3$ - Can we ignore input sensitivity?** To estimate how much we can lose, we first define two scenarios $S_1$ and $S_2$:

$S_1$ : *Baseline.* In this scenario, we just train a simple performance model on an input - *i.e.*, the *target* input. We choose the best configuration according to the model, configure the related software with it and execute it with the target input.

$S_2$ : *Ignoring input sensitivity.* In this scenario, let us pretend that we ignore the input sensitivity issue. We train a model on a given input *i.e.*, the *source* input, and then predict the best configuration for this source input. If we ignore the threat of input sensitivity, we can easily reuse this model for any other input, including the target input defined in $S_1$. Finally, we execute the software with the configuration predicted by our model on the *target* input.

In this part, we systematically compare $S_1$ and $S_2$ in terms of performance for all inputs, all performance properties and all software systems. For $S_1$, we repeat the scenario ten times with different sources, uniformly chosen among other inputs and consider the average performance. For both scenarios, due to the imprecision of the learning procedure, the models can recommend sub-optimal configurations. Since this imprecision can alter the results, we consider an ideal case for both scenarios and assume that the performance models always recommend the best possible configuration.

**Performance ratio.** To compare $S_1$ and $S_2$, we use a performance ratio *i.e.*, the performance obtained in $S_1$ over the performance obtained in $S_2$. If the ratio is equal to 1, there is no difference between $S_1$ and $S_2$ and the input sensitivity does not exist. A ratio of 1.4 would suggest that the performance of $S_1$ is worth 1.4 times the performance of $S_2$; therefore, it is possible to gain up to $(1.4 - 1) * 100 = 40\%$ performance by choosing $S_1$ instead of $S_2$. We also report on the standard deviation of the performance ratio distribution. A standard deviation of 0 implies that for each input, we gain or lose the same proportion of performance when picking $S_1$ over $S_2$.

---

[4]Guidelines for replication are available at: https://anonymous.4open.science/r/df319578-8767-47b0-919d-a8e57eb67d25/replication/README.md
[5]The description of the python environment is available at : https://anonymous.4open.science/r/df319578-8767-47b0-919d-a8e57eb67d25/replication/requirements.txt
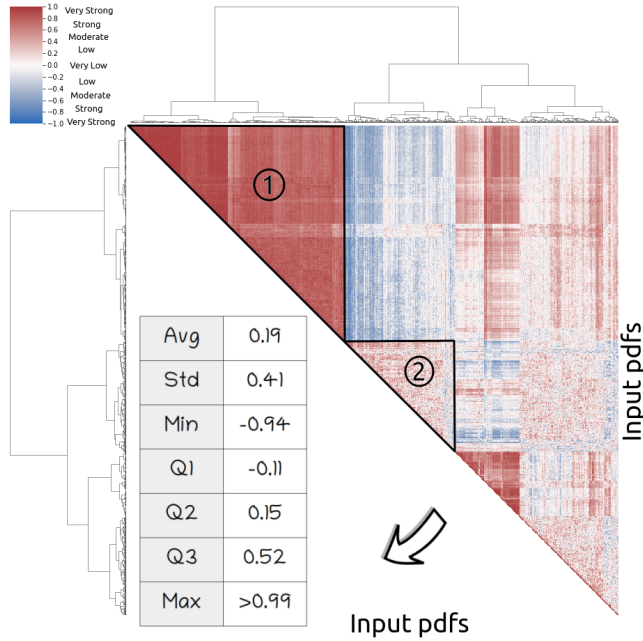
## 4 RESULTS

We now present the results obtained by following the methodology defined in Section 3.

### 4.1 Performance correlations ($RQ_1$)

We first explain the results of $RQ_1$ and their consequences on the *poppler* use case *i.e.*, an extreme case of input sensitivity, and then generalize to our other software systems.

**Extract images of input pdfs with *poppler*.** The content of input pdfs fed to *poppler* may vary *e.g.*, the pdf can be a 2-page report with a simple figure or a 300-page picture book. Depending on this content, extracting the images embedded in those files can be quick or slow. Moreover, a user can adapt different configurations for the report and not for the book (or conversely), leading to different rankings in terms of extraction time.



Each square$_{(i,j)}$ represents the Spearman correlation between the time needed to extract the images of pdfs $i$ and $j$. The color of this square respects the top-left scale: high positive correlations are red; low in white; negative in blue. Because we cannot describe each correlation individually, we added a table describing their distribution.

**Figure 3: Spearman correlogram - *poppler*, time.**

In Figure 3, we depict the Spearman rank-order correlations, in terms of extraction time, between pairs of input pdfs fed to *poppler*. We also perform hierarchical clustering [42] on *poppler* data to gather inputs having similar times distributions and visually group correlated pdfs together.[6] Results suggest a positive correlation (see dark red cells), though there are pairs of inputs with lower (see white cells) and even negative (see dark blue cells) correlations. More than a quarter of the correlations between input pdfs are positive and at least moderate - third quartile Q3 greater than 0.52. On the top-left part of the correlogram (see triangle ①), we even observe a first group of input pdfs that are highly correlated with

each others - positively, strong or very strong. In this first group, the input pdfs have similar *time* rankings; their performance react the same way to the same configurations. However, this group of pdfs is uncorrelated (very low, low) or negatively correlated (moderate, strong and very strong) with the second group of pdfs - see the triangle ②. In this case, a performance model trained on a pdf chosen in the group ① should not be reused directly on a pdf of the group ②.

**Meta-analysis.** Over the 8 systems, we observe different cases :

- There exist software systems not sensitive at all to inputs. In our experiment, *gcc*, *imagemagick* and *xz* present almost exclusively high and positive correlations between inputs *e.g.*, Q1 = 0.82 for the compressed size and *xz*. For these, un- or negatively-correlated inputs are an exception more than a rule.
- In contrast, there are software systems, namely *lingeling*, *nodeJS*, *SQLite* and *poppler*, for which performance distributions completely change and depend on input data *e.g.*, Q2 = 0.09 for *nodeJS* and ops, Q3 = 0.12 for *lingeling* and conflicts. For these, we draw similar conclusions as in the *poppler* case.
- In between, *x264* is only input-sensitive w.r.t. a performance property; it is for bitrate and size but not for cpu, fps and time *e.g.*, 0.29 as standard deviation for size and bitrate but 0.08 for the time.

> **$RQ_1$ - Do software performances stay consistent across inputs?** Performance distributions can change depending on inputs. Our systematic empirical study shows evidences about the existence of input sensitivity: (1) input sensitivity does not affect all systems; (2) input sensitivity may affect not the whole systems but some specific performance properties. So, without having scrutinized the input sensitivity of a system, one cannot develop techniques sensitive to this phenomenon.

### 4.2 Effects of options ($RQ_2$)

We first explain the results of $RQ_2$ and their concrete consequences on the bitrate of *x264* - an input-sensitive case, to then generalize to other software systems.

**Encode input videos with *x264*.** *x264* can encode different kinds of videos, such as an animation movie with many details, or a soccer game with large and monochromatic areas of grass. When encoding the soccer game, *x264* can use those fixed green areas to reduce the amount of data encoded per second (*i.e.*, the bitrate). In other words, configuration options aggregating pixels (*e.g.*, macro-block tree estimation mbtree) could both *reduce* the bitrate for the soccer game and *increase* the bitrate for the animation movie where nothing can be aggregated.

Figures 4a and 4b report on respectively the boxplots of configuration options' feature importances and effects when predicting *x264*'s bitrate for all input videos.[7] Three options are strongly influential for a majority of videos on Figure 4a: subme, mbtree and aq-mode, but their importance can differ depending on input videos:

---

[6]Detailed $RQ_1$ results for other systems available at : https://anonymous.4open.science/r/df319578-8767-47b0-919d-a8e57eb67d25/results/RQS/RQ1/RQ1.md

[7]Detailed $RQ_2$ results for other systems available at : https://anonymous.4open.science/r/df319578-8767-47b0-919d-a8e57eb67d25/results/RQS/RQ2/RQ2.md
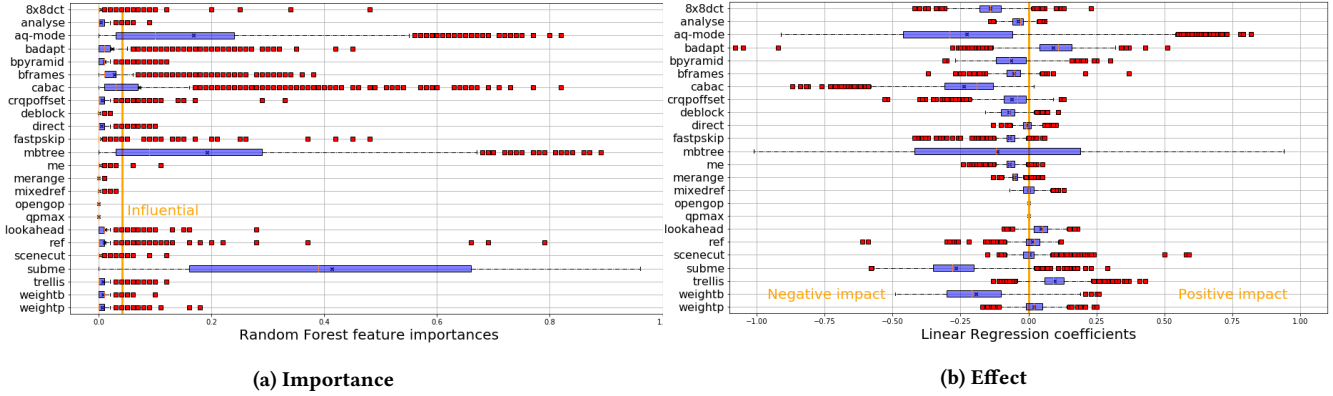
**(a) Importance**



**(b) Effect**

**Figure 4: Importance and effect of configuration options - *x264*, <u>bitrate</u>**

for instance, the importance of subme is 0.83 for video #1365 and only 0.01 for video #40. Because influential features vary with input videos for *x264*, performance models and approaches based on *feature selection* [58] may not generalize well to all input videos. Most of the options have positive and negative coefficients on Figure 4b; thus, the specific effects of options heavily depend on input videos. It is also true for influential options: mbtree can have positive and negative (influential) effects on the <u>bitrate</u> *i.e.*, activating mbtree may be worth only for few input videos. The consequence is that one cannot reliably provide end-users with a unique *x264* performance model or a *x264* default configuration whatever the input is.

Another interesting point is the link between $RQ_1$ and $RQ_2$ for *x264* and the <u>bitrate</u>; the more stable the effect of options in $RQ_2$, the more stable the distribution of performances in $RQ_1$. In fact, in a group of highly-correlated input videos (*e.g.*, like the group ① of pdfs in Figure 3, but for *x264*), the effect and importance of options are stable *i.e.*, the inputs all react the same way to the same options.[8] These different effects of influential options of *x264* may alter its encoding performances, thus explaining the different distributions pointed out in Section 4.1. Under these circumstances, training a performance model per group of inputs is probably a reasonable solution for tackling input sensitivity.

**Meta-analysis.** For *gcc*, *imagemagick* and *xz*, the importances are quite stable. As an extreme case of stability, the importances of the compressed <u>size</u> for *xz* are exactly the same, except for two inputs. For both systems, the coefficients of linear regression mostly keep the same sign across inputs *i.e.*, the effects of options do not change with inputs. For input-sensitive software systems, we always observe high variations of options' effects (*lingeling*, *poppler* or *SQLite*), sometimes coupled to high variations of options' importances (*nodeJS*). For instance, the option format for *poppler* can have an importance of 0 or 1 depending on the input. For all software systems, there exists at least one performance property whose effects are not stable for all inputs *e.g.*, one input with negative coefficient and another with a positive coefficient. For *x264*, it depends on the performance property; for cpu, fps and <u>time</u>, the effect of influential options are stable for all inputs, while for the <u>bitrate</u> and the <u>size</u>, we can draw the conclusions previously presented in this section.

> **$RQ_2$ - Do configuration option's effects change with input data?** Different inputs lead to different configuration options' significance and effects. A set of influential options with changing effects can alter the distribution of performance, thus explaining $RQ_1$ results. Therefore, a performance model should not be fixed but evolve according to the input data fed to the system.

## 4.3 Impact of input sensitivity ($RQ_3$)

This section presents the evaluation of $RQ_3$ w.r.t. the protocol of Section 3.4. In Table 2, we computed the performance ratios for the different software systems and their performance properties.[9]

For software systems whose performances are stable across inputs (*gcc*, *imagemagick* and *xz*), there are few differences between inputs. For instance, for the output <u>size</u> of *xz*, there is no variation between scenarios $S_1$ (*i.e.*, using the best configuration) and $S_2$ (*i.e.*, reusing a the best configuration of a given input for another input): all performance ratios are equals to 1 whatever the input.

For input-sensitive software systems (*lingeling*, *nodeJS*, *SQLite* and *poppler*), changing the configuration can lead to a negligible change in a few cases. For instance, for the time to answer the first query <u>q1</u> with *SQLite*, the third quartile is 1.04; in this case, *SQLite* is sensitive to inputs, but its variations of performance -less than 4 %- do not justify the complexity of tuning the software. But it can also be a huge change; for *lingeling* and solved <u>conflicts</u>, the $95^{th}$ percentile ratio is equal to 8.05 *i.e.*, a factor of 8 between $S_1$ and $S_2$. It goes up to a ratio of 10.11 for *poppler*'s extraction <u>time</u>: there exists an input pdf for which extracting its images is ten times slower when reusing a configuration, compared to the best one (*i.e.*, the fastest).

In between, *x264* is a complex case. For its low input-sensitive performances (*e.g.*, <u>cpu</u> and <u>etime</u>), it moderately impacts the performances when reusing a configuration from one input to another - average ratios at resp. 1.42 and 1.43. In this case, the rankings of performances do not change a lot with inputs, but a small ranking change does make the difference in terms of performance.

---

[8]See the detailed case of *x264* at : https://anonymous.4open.science/r/df319578-8767-47b0-919d-a8e57eb67d25/results/others/x264_groups/x264_bitrate.md

[9]Detailed $RQ_3$ results for other performance properties available at : https://anonymous.4open.science/r/df319578-8767-47b0-919d-a8e57eb67d25/results/RQS/RQ3/RQ3.md

Table 2: Performance ratio distributions across inputs, for different software systems and different performance properties. In lines, *Avg* the avegrae performance ratio. *Std* the standard deviation. $5^{th}$ the $5^{th}$ percentile. *Q1* the first quartile. *Q2* the median. *Q3* the third quartile. $95^{th}$ the $95^{th}$ percentile. Due to space constraints, we arbitrarly select few performance properties.

| System | gcc | | | lingeling | | nodeJS | poppler | | SQLite | | | x264 | | | | | xz | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Perf. P | ctime | exec | size | confl | reduc | ops | size | time | q1 | q12 | q14 | cpu | etime | fps | bitrate | size | size | time |
| Avg | 1.08 | 1.13 | 1.27 | 2.11 | 1.38 | 1.73 | 1.56 | 2.69 | 1.03 | 1.08 | 1.07 | 1.42 | 1.43 | 1.1 | 1.11 | 1.11 | 1.0 | 1.08 |
| Std | 0.07 | 0.07 | 0.36 | 2.6 | 0.79 | 1.88 | 1.27 | 3.72 | 0.02 | 0.05 | 0.05 | 1.27 | 1.45 | 0.14 | 0.13 | 0.13 | 0.0 | 0.06 |
| $5^{th}$ | 1.0 | 1.05 | 1.01 | 1.02 | 1.0 | 1.01 | 1.0 | 1.03 | 1.01 | 1.01 | 1.01 | 1.05 | 1.05 | 1.02 | 1.01 | 1.02 | 1.0 | 1.0 |
| Q1 | 1.01 | 1.11 | 1.04 | 1.05 | 1.04 | 1.08 | 1.0 | 1.14 | 1.02 | 1.03 | 1.03 | 1.12 | 1.12 | 1.04 | 1.03 | 1.05 | 1.0 | 1.02 |
| Q2 | 1.08 | 1.12 | 1.16 | 1.14 | 1.11 | 1.16 | 1.07 | 1.38 | 1.03 | 1.07 | 1.07 | 1.21 | 1.21 | 1.06 | 1.07 | 1.08 | 1.0 | 1.07 |
| Q3 | 1.11 | 1.14 | 1.32 | 1.47 | 1.25 | 1.54 | 1.51 | 2.22 | 1.04 | 1.11 | 1.09 | 1.38 | 1.37 | 1.1 | 1.15 | 1.12 | 1.0 | 1.11 |
| $95^{th}$ | 1.2 | 1.2 | 1.97 | 8.05 | 2.79 | 4.22 | 3.85 | 10.11 | 1.08 | 1.17 | 1.16 | 2.11 | 2.11 | 1.25 | 1.32 | 1.28 | 1.0 | 1.2 |

On the contrary, for the input-sensitive performances (*e.g.*, the bitrate), there are few variations of performance: we can lose $1 - \frac{1}{1.11} \simeq 9\%$ of bitrate in average. In this case, it is up to the compression experts to decide; if losing up to $1 - \frac{1}{1.32} \simeq 24\%$ of bitrate is acceptable, then we can ignore input sensitivity. Otherwise, we should consider tuning *x264* for its input video.

> **RQ$_3$ - Can we ignore input sensitivity?** There exist input-sensitive cases for which the difference of performance does not justify to consider the input sensitivity *e.g.*, 5 % change is probably negligible. However, performances can be multiplied up to a ratio of 10 if we tune other systems for their input data: we cannot ignore it.

## 5 SIGNIFICANCE OF THE THREAT

In this section, we explore the relevance and the significance of the input sensitivity problem in research. Do researchers know the input sensitivity? How do they deal with inputs in their papers? Do they ignore it? Is input sensitivity a well-known threat? What motivates us is to estimate to what extent input sensitivity is a problem in the research community, but it is also the opportunity to promote innovative and original state-of-the-art solutions.

### 5.1 Protocol

First, we aim at gathering research papers that predict performances of configurable systems.

**Gather research papers.** We focused on the publications of the last ten years. To do so, we analyzed the papers published (strictly) after 2011 from the survey of Pereira *et al.* [67] published in 2019. We completed those papers with more recent papers (2019-2021), following the same procedure as in [67]. We have only kept research work that trained performance models on software systems.

**Search for input sensitivity.** We read each selected paper and answer four different questions: Q-A. Is there a software system processing input data in the study? If not, the impact of input sensitivity in the existing research work would be relatively low. The idea of this research question is to estimate the proportion of the performance models that could be affected by input sensitivity. Q-B. Does the experimental protocol include several inputs? If not, it would suggest that the performance model only captures a partial truth, and might not generalize for other inputs fed to the software

system. Q-C. Is the problem of input sensitivity mentioned *e.g.*, in threat? This question aims to state whether researchers are aware of the input sensitivity threat, and estimate the proportion of the papers that mention it as a potential threat to validity. Q-D. Does the paper propose a solution to generalize the performance model across inputs? Finally, we check whether the paper proposes a solution managing input sensitivity *i.e.*, if the proposed approach can be applied to predict a near-optimal configuration for any input. The results were obtained by one author and validated by all other co-authors.

### 5.2 Results

Table 3 lists the 63 research papers we identified following this protocol, as well as their individual answers to Q-A→Q-D. A checked cell indicates that the answer to the corresponding question (column) for the corresponding paper (line) is *yes*. Since answering Q-B, Q-C or Q-D only makes sense if Q-A is checked, we grayed and did not consider Q-B, Q-C and Q-D if the answer of Q-A is *no*. We also provide full references and detailed justifications in the companion repository.[10] We now comment the average results for each question:

**Q-A.** Is there a software system processing input data in the study? Of the 63 papers, 59 (94 %) consider at least one configurable system processing inputs. This large proportion gives credits to input sensitivity and its potential impact on research work.

**Q-B.** Does the experimental protocol include several inputs? 63 % of the research work answering *yes* to Q-A include different inputs in their protocol, which is overall positive. But what about the other 37 %? It is understandable not to consider several inputs because of the cost of measurements. However, if we reproduce all experiments of Table 3 using other input data, will we draw the same conclusions for each paper? Based on the results of $RQ_1 \rightarrow RQ_3$, we encourage all researchers to consider at least a set of well-chosen inputs in their protocol (*e.g.*, an input per group, as shown in $RQ_1$). We give an example of such a set for *x264* in Section 6.

**Q-C.** Is the problem of input sensitivity mentioned *e.g.*, in threat? Only half (46 %) of the papers mention the threat of input sensitivity, mostly without naming it or using a domain-specific keyword (*e.g.*, workload variation [90]). For the other half, we cannot guarantee with certainty that input sensitivity concerns all papers. But we shed light on this threat to validity: ignoring input sensitivity can

---

[10]The list of papers can be consulted at https://anonymous.4open.science/r/df319578-8767-47b0-919d-a8e57eb67d25/results/RQS/RQ4/RQ4.md

Table 3: Input sensitivity in research. Paper identifier *ID* in the list. *Authors* of the paper. *Conference* in which the paper was accepted. *Year* of publication of the paper. *Title* of the paper. *Q-A.* Is there a software system processing input data in the study? *Q-B.* Does the experimental protocol include several inputs? *Q-C.* Is the problem of input sensitivity mentioned *e.g.,* in threat? *Q-D.* Does the paper propose a solution to generalize the performance model across inputs? Due to space limitation, we do not justify the answers directly in the paper, see the companion repository for full references and justifications.

| ID | Authors | Conference | Year | Title | Q-A | Q-B | Q-C | Q-D |
|---|---|---|---|---|---|---|---|---|
| 1 | Guo *et al.* [29] | ESE | 2017 | Data-efficient performance learning for configurable systems | X | | | |
| 2 | Jamshidi *et al.* [41] | SEAMS | 2017 | Transfer learning for improving model predictions [...] | X | X | X | |
| 3 | Jamshidi *et al.* [39] | ASE | 2017 | Transfer learning for performance modeling of configurable [...] | X | X | X | X |
| 4 | Oh *et al.* [65] | ESEC/FSE | 2017 | Finding near-optimal configurations in product lines by random [...] | X | | | |
| 5 | Kolesnikov *et al.* [47] | SoSyM | 2018 | Tradeoffs in modeling performance of highly configurable [...] | X | | | |
| 6 | Nair *et al.* [61] | ESEC/FSE | 2017 | Using bad learners to find good configurations | X | X | | |
| 7 | Nair *et al.* [63] | TSE | 2018 | Finding Faster Configurations using FLASH | X | X | X | |
| 8 | Murwantara *et al.* [59] | iiWAS | 2014 | Measuring Energy Consumption for Web Service Product [...] | X | X | X | X |
| 9 | Temple *et al.* [86] | SPLC | 2016 | Using Machine Learning to Infer Constraints for Product Lines | | | | |
| 10 | Temple *et al.* [84] | IEEE Software | 2017 | Learning Contextual-Variability Models | X | | X | |
| 11 | Valov *et al.* [90] | ICPE | 2017 | Transferring performance prediction models across different [...] | X | | X | X |
| 12 | Weckesser *et al.* [95] | SPLC | 2018 | Optimal reconfiguration of dynamic software product lines based [...] | | | | |
| 13 | Acher *et al.* [2] | VaMoS | 2018 | VaryLATEX: Learning Paper Variants That Meet Constraints | X | X | | X |
| 14 | Sarkar *et al.* [75] | ASE | 2015 | Cost-Efficient Sampling for Performance Prediction of [...] | X | | | |
| 15 | Temple *et al.* [83] | Report | 2018 | Towards Adversarial Configurations for Software Product Lines | | | | |
| 16 | Nair *et al.* [62] | ASE | 2018 | Faster Discovery of Faster System Configurations with [...] | X | | | |
| 17 | Siegmund *et al.* [78] | ESEC/FSE | 2015 | Performance-Influence Models for Highly Configurable Systems | X | | | |
| 18 | Valov *et al.* [88] | SPLC | 2015 | Empirical comparison of regression methods for [...] | X | | | |
| 19 | Zhang *et al.* [98] | ASE | 2015 | Performance Prediction of Configurable Software Systems [...] | X | | X | |
| 20 | Kolesnikov *et al.* [48] | ESE | 2019 | On the relation of control-flow and performance feature [...] | X | | | |
| 21 | Couto *et al.* [13] | SPLC | 2017 | Products go Green: Worst-Case Energy Consumption [...] | X | | X | |
| 22 | Van Aken *et al.* [91] | SIGMOD | 2017 | Automatic Database Management System Tuning Through [...] | X | X | X | X |
| 23 | Kaltenecker *et al.* [44] | ICSE | 2019 | Distance-based sampling of software configuration spaces | X | | | |
| 24 | Jamshidi *et al.* [40] | ESEC/FSE | 2018 | Learning to sample: exploiting similarities across environments [...] | X | X | X | X |
| 25 | Jamshidi *et al.* [38] | MASCOTS | 2016 | An Uncertainty-Aware Approach to Optimal Configuration of [...] | X | X | X | |
| 26 | Lillacka *et al.* [52] | Soft. Eng. | 2013 | Improved prediction of non-functional properties in Software [...] | X | X | X | X |
| 27 | Zuluaga *et al.* [100] | JMLR | 2016 | $\varepsilon$-pal: an active learning approach [...] | X | X | | |
| 28 | Amand *et al.* [7] | VaMoS | 2019 | Towards Learning-Aided Configuration in 3D Printing [...] | X | X | X | |
| 29 | Alipourfard *et al.* [4] | NSDI | 2017 | Cherrypick: Adaptively unearthing the best cloud configurations [...] | X | X | X | |
| 30 | Saleem *et al.* [74] | TSC | 2015 | Personalized Decision-Strategy based Web Service Selection [...] | X | X | | |
| 31 | Zhang *et al.* [99] | SPLC | 2016 | A mathematical model of performance-relevant feature interactions | X | | | |
| 32 | Ghamizi *et al.* [25] | SPLC | 2019 | Automated Search for Configurations of Deep Neural [...] | X | X | X | |
| 33 | Grebhahn *et al.* [27] | CPE | 2017 | Performance-influence models of multigrid methods [...] | | | | |
| 34 | Bao *et al.* [8] | ASE | 2018 | AutoConfig: Automatic Configuration Tuning for Distributed [...] | X | X | | |
| 35 | Guo *et al.* [28] | ASE | 2013 | Variability-aware performance prediction: A statistical learning [...] | X | | | |
| 36 | Švogor *et al.* [102] | IST | 2019 | An extensible framework for software configuration optimization [...] | X | X | | |
| 37 | El Afia *et al.* [3] | CloudTech | 2018 | Performance prediction using support vector machine for the [...] | X | X | | |
| 38 | Ding *et al.* [16] | PLDI | 2015 | Autotuning algorithmic choice for input sensitivity | X | X | X | X |
| 39 | Duarte *et al.* [19] | SEAMS | 2018 | Learning Non-Deterministic Impact Models for Adaptation | X | X | X | X |
| 40 | Thornton *et al.* [87] | KDD | 2013 | Auto-WEKA: Combined selection and hyperparameter [...] | X | X | X | |
| 41 | Siegmund *et al.* [79] | ICSE | 2012 | Predicting performance via automated feature-interaction detection | X | X | X | |
| 42 | Siegmund *et al.* [80] | SQJ | 2012 | SPL Conqueror: Toward optimization of non-functional [...] | X | X | | |
| 43 | Westermann *et al.* [96] | ASE | 2012 | Automated inference of goal-oriented performance prediction [...] | X | X | | |
| 44 | Velez *et al.* [92] | ICSE | 2021 | White-Box Analysis over Machine Learning: Modeling [...] | X | X | | |
| 45 | Pereira *et al.* [6] | ICPE | 2020 | Sampling Effect on Performance Prediction of Configurable [...] | X | X | X | |
| 46 | Shu *et al.* [77] | ESEM | 2020 | Perf-AL: Performance prediction for configurable software [...] | X | | | |
| 47 | Dorn *et al.* [18] | ASE | 2020 | Mastering Uncertainty in Performance Estimations of [...] | X | | | |
| 48 | Kaltenecker *et al.* [43] | IEEE Software | 2020 | The Interplay of Sampling and Machine Learning for Software [...] | X | | | |
| 49 | Krishna *et al.* [49] | TSE | 2020 | Whence to Learn? Transferring Knowledge in Configurable [...] | X | X | X | X |
| 50 | Weber *et al.* [94] | ICSE | 2021 | White-Box Performance-Influence Models: A Profiling [...] | X | X | | |
| 51 | Mühlbauer *et al.* [60] | ASE | 2020 | Identifying Software Performance Changes Across Variants [...] | X | X | | |
| 52 | Han *et al.* [32] | Report | 2020 | Automated Performance Tuning for Highly-Configurable [...] | X | X | | |
| 53 | Han *et al.* [33] | ICPE | 2021 | ConfProf: White-Box Performance Profiling of Configuration [...] | X | | X | |
| 54 | Valov *et al.* [89] | ICPE | 2020 | Transferring Pareto Frontiers across Heterogeneous Hardware [...] | X | | | X |
| 55 | Liu *et al.* [53] | CF | 2020 | Deffe: a data-efficient framework for performance [...] | X | X | X | X |
| 56 | Fu *et al.* [23] | NSDI | 2021 | On the Use of ML for Blackbox System Performance Prediction | X | X | X | X |
| 57 | Larsson *et al.* [50] | IFIP | 2021 | Source Selection in Transfer Learning for Improved Service [...] | X | X | X | X |
| 58 | Chen *et al.* [10] | ICSE | 2021 | Efficient Compiler Autotuning via Bayesian Optimization | X | X | X | |
| 59 | Chen *et al.* [11] | SEAMS | 2019 | All Versus One: An Empirical Comparison on Retrained [...] | X | X | | |
| 60 | Ha *et al.* [30] | ICSE | 2019 | DeepPerf: Performance Prediction for Configurable Software [...] | X | | | |
| 61 | Pei *et al.* [66] | Report | 2019 | DeepXplore: automated whitebox testing of deep learning systems | X | X | | |
| 62 | Ha *et al.* [31] | ICSME | 2019 | Performance-Influence Model for Highly Configurable Software [...] | X | | | |
| 63 | Iorio *et al.* [37] | CloudCom | 2019 | Transfer Learning for Cross-Model Regression in Performance [...] | X | X | X | X |
| | | | | Total | 59 | 37 | 27 | 15 |

prevent the generalization of performance models across inputs. This is especially true for the 37 % of papers answering *no* to Q-B *i.e.*, considering one input per system: only 14 % of these research works mention this threat in their publication.

**Q-D.** Does the paper propose a solution to generalize the performance model across inputs? We identified 15 papers that propose contributions that may help in better managing the input sensitivity problem, and that should be adapted and tested to evaluate their ability to support this problem. Out of these papers, a classification of potential solutions emerge:

*All in a dataset* [2, 23] *i.e.*, gather the measurements related to all inputs in a unique dataset. If we provide enough input data, we expect the model to predict the average performance for all inputs. This method is easy to implement, and just requires to measure few inputs. But it may lead to poor predictions for input-sensitive software systems.

*Input-aware learning* [16, 19, 91] *i.e.*, differentiate the inputs during the training of the model thanks to inputs' properties (the resolution or the duration of a video, the size of a *.pdf* file, *etc.*). For instance, [16] classifies input programs into different performance groups, and then trains a performance model per group, which can generalize to new inputs. An issue is to identify those input properties (if any) that must be both cheap to compute and discriminant enough to accurately predict performances of configurations.

*Transfer learning* [37, 39, 40, 49, 50, 53, 89, 90] uses the similarities between the source (already measured) and the target (new) to predict performances of the target based on the predictions made on the source. It should be noted that the transfer techniques were not designed to specifically address the sensitivity of inputs and were not evaluated in this context. Instead, these works have focused on changes in computing environments (*e.g.*, hardware) where it is possible to linearly transfer performance models [90]. It is unclear how transfer techniques would deal with our data.

Moreover the computational cost of learning performance models is already important for one input. Since learning from scratch each time a new input is fed to a configurable system seems impractical owing to the diversity of possible inputs, this approach proposes to adapt existing performance models, thus reducing the computational cost. However, transfer techniques still require to measure few configurations on the new input (target), which is not always realistic from a user's perspective.

To the best of our knowledge, the assessment of the underlying cost and effectiveness of these methods has never been done within the context of input sensitivity. An opportunity for researchers is to devise transfer techniques specifically designed for handling input sensitivity. We provide a large amount of data and encourage researchers to investigate this question.

> **Conclusion.** While half of the research articles mention input sensitivity, few actually address it, and most often on a single system and domain. Some state-of-the-art candidate solutions are not systematically applicable and their cost and accuracy must be assessed.

## 6 IMPACT FOR RESEARCHERS AND RESEARCH OPPORTUNITIES

This section discusses the implications of our work.

**Impacts for researchers.** We warn researchers that the effectiveness of learning strategies for a given configurable system can be biased by the inputs and the performance property used. That is, a sampling strategy, a prediction or optimisation algorithm, or a transfer technique may well be highly accurate for an input and still inaccurate for others. Most of the studies neglect either inputs or configurations, which is perfectly understandable owing to the investments required. However, the scientific community should be extremely careful with this input sensitivity issue. In view of the results of our study, new problems deserve to be tackled with associated challenges. We detail some of them hereafter.

**Sampling configurations.** With the promise to select a small and representative sample set of valid configurations, several sampling strategies have been devised in the last years [6, 43, 67] (*e.g.*, random sampling, *t*-wise sampling, distance-based sampling). As recently reported in other experimental settings [6, 43], finding the most effective combinations of sampling and learning strategies is an open problem. *Input sensitivity further exacerbates the problem.* We conjecture that some strategies for sampling configurations might be effective for specific inputs and performance properties. Pereira *et al.* [6] actually provided preliminary evidence on *x264* for 19 input videos and two performance properties. Our results show and confirm that the importance of options and their interactions is indeed sensitive to the input (see $RQ_2$), thus suggesting that some sampling strategies may not always capture them. An open issue is thus to find sampling strategies that are effective for any input.

**Tuning and performance prediction.** Numerous works aim to find optimal configurations or predict the performance of an arbitrary configuration. However, our empirical results show that the best configuration can be differently ranked (see $RQ_1$) depending on an input. The tuning or the prediction cannot be reused as such (see $RQ_3$) but should be redone or adapted whenever a system processes a new input. An open challenge is to deliver algorithms and practical tools capable of tuning a system whatever the input. Another issue is to reduce the cost of building performance models for each input (*e.g.*, through sampling or transfer learning).

**Understanding of configurable systems.** Understanding the effects of options and their interactions is hard for developers and users yet crucial for maintaining, debugging or configuring a software system. Some works (*e.g.*, [94]) have proposed to build performance models that are interpretable and capable of communicating the influence of individual options and interactions on performance (possibly back to the code). Our empirical results show that performance models and so options and their interactions are sensitive to inputs (see $RQ_2$). A first open issue is to communicate when and how options *together with input data* interact and influence performance. Another challenge is to identify a minimal set of representative inputs in such a way a configurable system can be observed and performance models learnt.

**Recommendations.** Given the state of the art and the open problems to be addressed, there is no complete solution that can be systematically employed. However, we can give two recommendations : *1. Detecting input sensitivity.* As practitioners dealing with

new inputs, we first have to to determine whether the software under study is input-sensitive w.r.t. the performance property of interest. If the input sensitivity is negligible (see $RQ_3$), we can use a single model to predict the performances of the software system. If not, measurements over multiple inputs are needed. *2. Selecting representative inputs.* To reduce the cost of measurements, the ideal would be to select a set of input data, both representative of the usage of the system and cheap to measure. We believe our work can be helpful here. On the *x264* case study, for the bitrate, we isolate 4 encoding groups of input videos (action movie - big resolution - still image - standard). Within a group, the videos share common properties, and *x264* processes them in the same way *i.e.*, same options' effects in $RQ_2$. In the companion repository, we propose to reduce the dataset of 1397 input videos [93] to a subset of 8 videos, selecting 2 cheap videos in each group of performance.[11] Automating this grouping could drastically reduce the cost of measuring configurations' performances over inputs.

## 7 THREATS TO VALIDITY

This section discusses the threats to validity related to our experimental protocol.

**Construct validity.** Due to resource constraints, we did not include all the options of the configurable systems in the experimental protocol. We may have forgotten configuration options that matter when predicting the performances of our configurable systems. However, we consider features that impact the performance properties according to the documentation, which is sufficient to show the existence of the input sensitivity threat.

**Internal Validity.** First, our results can be subject to measurement bias. We alleviated this threat by making sure only our experiment was running on the server we used to measure the performances of software systems. It has several benefits: we can guarantee we use similar hardware (both in terms of CPU and disk) for all measurements; we can control the workload of each machine (basically we force the machine to be used only by us); we can avoid networking and I/O issues by placing inputs on local folders. But it could also represent a threat: our experiments may depend on the hardware and operating system. The measurement process is launched via docker containers. If this aims at making this work reproducible, this can also alter the results of our experiment. Because of the amount of resources needed to compute all the measures, we did not repeat the process of Figure 2 several times per system. We consider that the large number of inputs under test overcomes this threat. Moreover, related work (*e.g.*, [6] for *x264*) has shown that inputs often maintain stable performances between different launches of the same configuration. Finally, the measurement process can also suffer from a lack of inputs. To limit this problem, we took relevant dataset of inputs produced and widely used in their field. For $RQ_1$-$RQ_3$, executing our code with another python environment may lead to slightly different conclusions. For $RQ_3$, we consider oracles when predicting the best configurations for both scenarios, thus neglecting the imprecision of performance models : these results might change on a real-world case. In Section 5, our

results are subject to the selection of research papers: since we use and reproduce [67], we face the same threats to validity.

**External Validity.** A threat to external validity is related to the used case studies and the discussion of the results. Because we rely on specific systems and interesting performance properties, the results may be subject to these systems and properties. To reduce this bias, we selected multiple configurable systems, used for different purposes in different domains.

## 8 RELATED WORK

In this section, we discuss other related work (see also Section 5).

**Workload performance analysis.** There are works addressing the performance analysis of software systems [12, 22, 26, 51, 71, 81] depending on different input data (also called workloads or benchmarks). However, existing studies usually consider a limited set of configurations. On the other hand, works and studies on configurable systems (see Section 5) usually neglect input data (*e.g.*, using a unique video for measuring the configurations of a video encoder). In response, we perform an in-depth, controlled study of several configurable systems to make it vary in the large, both in terms of configurations and inputs.

**Input-aware tuning.** The input sensitivity issue has been partly considered but in specific domains (SAT solvers [21, 97], compilation [16, 69], video encoding [57], data compression [46], *etc.*). It is unclear whether these ad-hoc solutions are cost-effective. As future work, we plan to systematically assess domain-specific techniques as well as generic, domain-agnostic approach (*e.g.*, transfer learning) using our dataset. Furthermore, the existence of a general solution applicable to all domains and software configurations is an open question. For example, is it always possible and effective to extract input properties for all kinds of inputs?

**Input data and other variability factors.** Most of the studies support learning models restrictive to specific static settings (*e.g.*, inputs, hardware, and version) such that a new prediction model has to be learned from scratch once the environment change [67]. Jamshidi *et al.* [39] conducted an empirical study on four configurable systems (including *SQLite* and *x264*), varying software configurations and environmental conditions, such as hardware, input, and software versions. But without isolating the individual effect of input data on software configurations, it is challenging to understand the existing interplay between the inputs and any other variability factor (*e.g.*, the hardware).

## 9 CONCLUSION

We conducted a large study over the inputs fed to 8 configurable systems that shows the significance of the input sensitivity problem on performance properties. We deliver one main message: **inputs matter as (much as) configuration options**. It appears that inputs can significantly change the performances of the configurable systems up to the point some options' values have an opposite effect depending on the input. Ignoring this lesson leads to the learning of inaccurate performance prediction models and ineffective recommendations for developers and end-users.

As future work, it is an open challenge to mitigate the threat of input sensitivity when predicting, optimising, or understanding configurable systems. We encourage researchers to confront the existing methods of the literature with our data.

---

[11]See the resulting benchmark and its construction at : https://anonymous.4open.science/r/df319578-8767-47b0-919d-a8e57eb67d25/results/others/x264_groups/x264_bitrate.md

# REFERENCES

[1] Jan De Cock, Aditya Mavlankar, Anush Moorthy, and Anne Aaron. 2016. A Large-Scale Comparison of x264, x265, and libvpx – a Sneak Peek. netflix-study.

[2] Mathieu Acher, Paul Temple, Jean-Marc Jézéquel, José A. Galindo, Jabier Martinez, and Tewfik Ziadi. 2018. VaryLATEX: Learning Paper Variants That Meet Constraints. In *Proc. of VAMOS'18*. 83–88. https://doi.org/10.1145/3168365.3168372

[3] Abdellatif El Afia and Malek Sarhani. 2017. Performance prediction using support vector machine for the configuration of optimization algorithms. In *Proc. of CloudTech'17*. 1–7. https://doi.org/10.1109/CloudTech.2017.8284699

[4] Omid Alipourfard, Hongqiang Harry Liu, Jianshu Chen, Shivaram Venkataraman, Minlan Yu, and Ming Zhang. 2017. Cherrypick: Adaptively Unearthing the Best Cloud Configurations for Big Data Analytics. In *Proc. of NSDI'17*. 469–482.

[5] A. Alourani, M.A.N. Bikas, and M. Grechanik. 2016. Input-Sensitive Profiling. In *Advances in Computers*. 31–52.

[6] Juliana Alves Pereira, Mathieu Acher, Hugo Martin, and Jean-Marc Jézéquel. 2020. Sampling Effect on Performance Prediction of Configurable Systems: A Case Study. In *Proc. of ICPE'20*. 277–288.

[7] Benoit Amand, Maxime Cordy, Patrick Heymans, Mathieu Acher, Paul Temple, and Jean-Marc Jézéquel. 2019. Towards Learning-Aided Configuration in 3D Printing: Feasibility Study and Application to Defect Prediction. In *Proc. of VAMOS'19*. https://doi.org/10.1145/3302333.3302338

[8] Liang Bao, Xin Liu, Ziheng Xu, and Baoyin Fang. 2018. AutoConfig: Automatic Configuration Tuning for Distributed Message Systems. In *Proc. of ASE'18*. 29–40. https://doi.org/10.1145/3238147.3238175

[9] Leo Breiman. 2001. Random forests. *Machine learning* 45, 1 (2001), 5–32.

[10] Junjie Chen, Ningxin Xu, Peiqi Chen, and Hongyu Zhang. 2021. Efficient Compiler Autotuning via Bayesian Optimization. In *Proc. of ICSE'21*. 1198–1209. https://doi.org/10.1109/ICSE43902.2021.00110

[11] Tao Chen. 2019. All Versus One: An Empirical Comparison on Retrained and Incremental Machine Learning for Modeling Performance of Adaptable Software. In *Proc. of SEAMS'19*. 157–168. https://doi.org/10.1109/SEAMS.2019.00029

[12] Emilio Coppa, Camil Demetrescu, Irene Finocchi, and Romolo Marotta. 2014. Estimating the Empirical Cost Function of Routines with Dynamic Workloads. In *Proc. of CGO'14*. 230:239. https://doi.org/10.1145/2581122.2544143

[13] Marco Couto, Paulo Borba, Jácome Cunha, João Paulo Fernandes, Rui Pereira, and João Saraiva. 2017. Products Go Green: Worst-Case Energy Consumption in Software Product Lines. In *Proc. of SPLC'17*. 84–93. https://doi.org/10.1145/3106195.3106214

[14] Francisco Gomes de Oliveira Neto, Richard Torkar, Robert Feldt, Lucas Gren, Carlo A. Furia, and Ziwei Huang. 2019. Evolution of statistical analysis in empirical software engineering research: Current state and steps forward. *Journal of Systems and Software* 156 (2019), 246–267. https://doi.org/10.1016/j.jss.2019.07.002

[15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*. 248–255. https://doi.org/10.1109/CVPR.2009.5206848

[16] Yufei Ding, Jason Ansel, Kalyan Veeramachaneni, Xipeng Shen, Una-May O'Reilly, and Saman Amarasinghe. 2015. Autotuning algorithmic choice for input sensitivity. In *ACM SIGPLAN Notices*, Vol. 50. ACM, 379–390.

[17] Johannes Dorn, Sven Apel, and Norbert Siegmund. 2020. Generating Attributed Variability Models for Transfer Learning. In *Proc. of VAMOS'20*.

[18] Johannes Dorn, Sven Apel, and Norbert Siegmund. 2020. Mastering Uncertainty in Performance Estimations of Configurable Software Systems. In *Proc. of ASE'20*. 684–696. https://doi.org/10.1145/3324884.3416620

[19] Francisco Duarte, Richard Gil, Paolo Romano, Antónia Lopes, and Luís Rodrigues. 2018. Learning Non-Deterministic Impact Models for Adaptation. In *Proc. of SEAMS'18*. 196–205. https://doi.org/10.1145/3194133.3194138

[20] James D Evans. 1996. *Straightforward statistics for the behavioral sciences*.

[21] Stefan Falkner, Marius Lindauer, and Frank Hutter. 2015. SpySMAC: Automated Configuration and Performance Analysis of SAT Solvers. In *Proc. of SAT'15*. 215–222.

[22] Hany FathyAtlam, Gamal Attiya, and Nawal El-Fishawy. 2013. Comparative Study on CBIR based on Color Feature. *International Journal of Computer Applications* 78, 16 (Sept. 2013), 9–15. https://doi.org/10.5120/13605-1387

[23] Silvery Fu, Saurabh Gupta, Radhika Mittal, and Sylvia Ratnasamy. 2021. On the Use of ML for Blackbox System Performance Prediction. In *Proc. of NSDI'21*. 763–784. https://www.usenix.org/conference/nsdi21/presentation/fu

[24] Wei Fu and Tim Menzies. 2017. Easy over Hard: A Case Study on Deep Learning. In *Proc. of ESEC-FSE'17*. 49–60. https://doi.org/10.1145/3106237.3106256

[25] Salah Ghamizi, Maxime Cordy, Mike Papadakis, and Yves Le Traon. 2019. Automated Search for Configurations of Convolutional Neural Network Architectures. In *Proc. of SPLC'19*. 119–130. https://doi.org/10.1145/3336294.3336306

[26] Simon F. Goldsmith, Alex S. Aiken, and Daniel S. Wilkerson. 2007. Measuring Empirical Computational Complexity. In *Proc. of ESEC-FSE'07*. 395–404.

[27] A. Grebhahn, C. Rodrigo, Norbert Siegmund, F. Gaspar, and S. Apel. 2017. Performance-influence models of multigrid methods: A case study on triangular grids. *Concurrency and Computation: Practice and Experience* 29 (2017).

[28] Jianmei Guo, Krzysztof Czarnecki, Sven Apely, Norbert Siegmundy, and Andrzej Wasowski. 2013. Variability-Aware Performance Prediction: A Statistical Learning Approach. In *Proc. of ASE'13*. 301–311. https://doi.org/10.1109/ASE.2013.6693089

[29] Jianmei Guo, Dingyu Yang, Norbert Siegmund, Sven Apel, Atrisha Sarkar, Pavel Valov, Krzysztof Czarnecki, Andrzej Wasowski, and Huiqun Yu. 2017. Data-efficient performance learning for configurable systems. *Empirical Software Engineering* 23, 3 (Nov. 2017), 1826–1867.

[30] Huong Ha and Hongyu Zhang. 2019. DeepPerf: Performance Prediction for Configurable Software with Deep Sparse Neural Network. In *Proc. of ICSE'19*. 1095–1106. https://doi.org/10.1109/ICSE.2019.00113

[31] Huong Ha and Hongyu Zhang. 2019. Performance-Influence Model for Highly Configurable Software with Fourier Learning and Lasso Regression. In *Proc. of ICSME'19*. 470–480.

[32] Xue Han and Tingting Yu. 2020. Automated Performance Tuning for Highly-Configurable Software Systems. *arXiv preprint arXiv:2010.01397* (2020).

[33] Xue Han, Tingting Yu, and Michael Pradel. 2021. ConfProf: White-Box Performance Profiling of Configuration Options. In *Proc. of ICPE'12*. 1–8.

[34] Marijn J. H. Heule, Matti Juhani Järvisalo, and Martin Suda (Eds.). 2018. *Proceedings of SAT Competition 2018: Solver and Benchmark Descriptions*. Department of Computer Science Series of Publications B, Vol. B-2018-1.

[35] Robert A Hummel, B Kimia, and Steven W Zucker. 1987. Deblurring gaussian blur. *Computer Vision, Graphics, and Image Processing* 38, 1 (1987), 66–80.

[36] Emilio Incerto, Mirco Tribastone, and Catia Trubiani. 2017. Software performance self-adaptation through efficient model predictive control. In *Proc. of ASE'17*. 485–496. https://doi.org/10.1109/ASE.2017.8115660

[37] Francesco Iorio, Ali B. Hashemi, Michael Tao, and Cristiana Amza. 2019. Transfer Learning for Cross-Model Regression in Performance Modeling for the Cloud. In *Proc. of CloudCom'19*. 9–18. https://doi.org/10.1109/CloudCom.2019.00015

[38] Pooyan Jamshidi and Giuliano Casale. 2016. An Uncertainty-Aware Approach to Optimal Configuration of Stream Processing Systems. *CoRR* (2016). http://arxiv.org/abs/1606.06543

[39] Pooyan Jamshidi, Norbert Siegmund, Miguel Velez, Christian Kästner, Akshay Patel, and Yuvraj Agarwal. 2017. Transfer Learning for Performance Modeling of Configurable Systems: An Exploratory Analysis. In *Proc. of ASE'17*. 497–508.

[40] Pooyan Jamshidi, Miguel Velez, Christian Kästner, and Norbert Siegmund. 2018. Learning to Sample: Exploiting Similarities across Environments to Learn Performance Models for Configurable Systems. In *Proc. of ESEC/FSE'18*. 71–82. https://doi.org/10.1145/3236024.3236074

[41] Pooyan Jamshidi, Miguel Velez, Christian Kästner, Norbert Siegmund, and Prasad Kawthekar. 2017. Transfer Learning for Improving Model Predictions in Highly Configurable Software. In *Proc. of SEAMS'17* (Buenos Aires). 31–41. https://doi.org/10.1109/SEAMS.2017.11

[42] Stephen C Johnson. 1967. Hierarchical clustering schemes. *Psychometrika* 32, 3 (1967), 241–254.

[43] Christian Kaltenecker, Alexander Grebhahn, Norbert Siegmund, and Sven Apel. 2020. The Interplay of Sampling and Machine Learning for Software Performance Prediction. *IEEE Softw.* 37, 4 (2020), 58–66.

[44] Christian Kaltenecker, Alexander Grebhahn, Norbert Siegmund, Jianmei Guo, and Sven Apel. 2019. Distance-Based Sampling of Software Configuration Spaces. In *Proc. of ICSE'19*. 1084–1094. https://doi.org/10.1109/ICSE.2019.00112

[45] Maurice George Kendall. 1948. *Rank correlation methods*.

[46] Mohammad Khavari Tavana, Yifan Sun, Nicolas Bohm Agostini, and David Kaeli. 2019. Exploiting Adaptive Data Compression to Improve Performance and Energy-Efficiency of Compute Workloads in Multi-GPU Systems. In *Proc. of IPDPS'19*. 664–674. https://doi.org/10.1109/IPDPS.2019.00075

[47] Sergiy Kolesnikov, Norbert Siegmund, Christian Kästner, Alexander Grebhahn, and Sven Apel. 2018. Tradeoffs in modeling performance of highly configurable software systems. *Software & Systems Modeling* 18, 3 (Feb. 2018), 2265–2283. https://doi.org/10.1007/s10270-018-0662-9

[48] Sergiy Kolesnikov, Norbert Siegmund, Christian Kästner, and Sven Apel. 2018. On the Relation of External and Internal Feature Interactions: A Case Study. arXiv:1712.07440 [cs.SE]

[49] Rahul Krishna, Vivek Nair, Pooyan Jamshidi, and Tim Menzies. 2019. Whence to Learn? Transferring Knowledge in Configurable Systems using BEETLE. *CoRR* (2019), 1–16. http://arxiv.org/abs/1911.01817

[50] Hannes Larsson, Jalil Taghia, Farnaz Moradi, and Andreas Johnsson. 2021. Source Selection in Transfer Learning for Improved Service Performance Predictions. In *Proc. of Networking'21*. 1–9. https://doi.org/10.23919/IFIPNetworking52078.2021.9472818

[51] Philipp Leitner and Jürgen Cito. 2016. Patterns in the Chaos—A Study of Performance Variation and Predictability in Public IaaS Clouds. *ACM Trans. Internet Technol.* 16, 3, Article 15 (April 2016), 23 pages. https://doi.org/10.1145/2885497

[52] Max Lillack, Johannes Müller, and Ulrich W Eisenecker. 2013. Improved prediction of non-functional properties in software product lines with domain context. *Software Engineering 2013* 1, 1 (2013), 1–14.

[53] Frank Liu, Narasinga Rao Miniskar, Dwaipayan Chakraborty, and Jeffrey S. Vetter. 2020. Deffe: A Data-Efficient Framework for Performance Characterization in Domain-Specific Computing. In *Proc. of CF'20*. 182–191. https://doi.org/10.1145/3387902.3392633

[54] Margaret Maclagan and Elizabeth Gordon. 1999. The Canterbury Corpus. *New Zealand English Journal* 13 (1999), 50–58. https://search.informit.org/doi/10.3316/informit.778572424581257

[55] D. Maiorca and B. Biggio. 2019. Digital Investigation of PDF Files: Unveiling Traces of Embedded Malware. *IEEE Security Privacy* 17, 1 (2019), 63–71.

[56] Frank J Massey Jr. 1951. The Kolmogorov-Smirnov test for goodness of fit. *Journal of the American statistical Association* 46, 253 (1951), 68–78.

[57] A. Maxiaguine, Yanhong Liu, S. Chakraborty, and Wei Tsang Ooi. 2004. Identifying "representative" workloads in designing MpSoC platforms for media processing. In *Proc. of ESTImedia'04*. 41–46. https://ieeexplore.ieee.org/document/1359702

[58] Christoph Molnar. 2020. *Interpretable Machine Learning*. Lulu. com, Munich.

[59] I Made Murwantara, Behzad Bordbar, and Leandro L. Minku. 2014. Measuring Energy Consumption for Web Service Product Configuration. In *Proc. of iiWAS'14*. 224–228. https://doi.org/10.1145/2684200.2684314

[60] Stefan Mühlbauer, Sven Apel, and Norbert Siegmund. 2020. Identifying Software Performance Changes Across Variants and Versions. In *Proc. of ASE'20*. 611–622.

[61] Vivek Nair, Tim Menzies, Norbert Siegmund, and Sven Apel. 2017. Using bad learners to find good configurations. In *Proc. of ESEC/FSE'17*. 257–267.

[62] Vivek Nair, Tim Menzies, Norbert Siegmund, and Sven Apel. 2018. Faster Discovery of Faster System Configurations with Spectral Learning. *Automated Software Engg.* 25, 2 (June 2018), 247–277. https://doi.org/10.1007/s10515-017-0225-2

[63] V. Nair, Z. Yu, T. Menzies, N. Siegmund, and S. Apel. 2020. Finding Faster Configurations Using FLASH. *IEEE Transactions on Software Engineering* 46, 7 (2020), 794–811. https://doi.org/10.1109/TSE.2018.2870895

[64] Trent Nelson. 2014. Technically-oriented PDF Collection. https://github.com/tpn/pdfs

[65] Jeho Oh, Don Batory, Margaret Myers, and Norbert Siegmund. 2017. Finding Near-Optimal Configurations in Product Lines by Random Sampling. In *Proc. of ESEC/FSE'17*. 61–71. https://doi.org/10.1145/3106237.3106273

[66] Kexin Pei, Yinzhi Cao, Junfeng Yang, and Suman Jana. 2019. DeepXplore: Automated Whitebox Testing of Deep Learning Systems. *Commun. ACM* 62, 11 (Oct. 2019), 137–145. https://doi.org/10.1145/3361566

[67] Juliana Alves Pereira, Mathieu Acher, Hugo Martin, Jean-Marc Jézéquel, Goetz Botterweck, and Anthony Ventresque. 2021. Learning software configuration spaces: A systematic literature review. *JSS* (2021), 111044. https://doi.org/10.1016/j.jss.2021.111044

[68] Quentin Plazar, Mathieu Acher, Gilles Perrouin, Xavier Devroey, and Maxime Cordy. 2019. Uniform Sampling of SAT Solutions for Configurable Systems: Are We There Yet?. In *Proc. of ICST'19*. 240–251.

[69] Dmitry Plotnikov, Dmitry Melnik, Mamikon Vardanyan, Ruben Buchatskiy, Roman Zhuykov, and Je-Hyung Lee. 2013. Automatic Tuning of Compiler Optimizations and Analysis of their Impact. *Procedia Computer Science* 18 (2013), 1312–1321. https://doi.org/10.1016/j.procs.2013.05.298

[70] Meikel Poess and Chris Floyd. 2000. New TPC Benchmarks for Decision Support and Web Commerce. *SIGMOD Rec.* 29, 4 (Dec. 2000), 64–71. https://doi.org/10.1145/369275.369291

[71] Suporn Pongnumkul, Chaiyaphum Siripanpornchana, and Suttipong Thajchayapong. 2017. Performance Analysis of Private Blockchain Platforms in Varying Workloads. In *Proc. of ICCCN'17*. 1–7. https://doi.org/10.1109/icccn.2017.8038517

[72] Louis-Noël Pouchet et al. 2012. Polybench: The polyhedral benchmark suite. *URL: http://www. cs. ucla. edu/pouchet/software/polybench* 437 (2012), 1–1.

[73] Sandra Rapps and Elaine J. Weyuker. 1985. Selecting software test data using data flow information. *IEEE transactions on software engineering* SE-11, 4 (1985), 367–375.

[74] Muhammad Suleman Saleem, Chen Ding, Xumin Liu, and Chi-Hung Chi. 2015. Personalized Decision-Strategy based Web Service Selection using a Learning-to-Rank Algorithm. *IEEE Transactions on Services Computing* 8, 5 (2015), 727–739.

[75] Atri Sarkar, Jianmei Guo, Norbert Siegmund, Sven Apel, and Krzysztof Czarnecki. 2015. Cost-Efficient Sampling for Performance Prediction of Configurable Systems. In *Proc. of ASE'30*. 342–352. https://doi.org/10.1109/ASE.2015.45

[76] George AF Seber and Alan J Lee. 2012. *Linear regression analysis*. Vol. 329.

[77] Yangyang Shu, Yulei Sui, Hongyu Zhang, and Guandong Xu. 2020. Perf-AL: Performance Prediction for Configurable Software through Adversarial Learning. In *Proc. of ESEM'20*. https://doi.org/10.1145/3382494.3410677

[78] Norbert Siegmund, Alexander Grebhahn, Sven Apel, and Christian Kästner. 2015. Performance-Influence Models for Highly Configurable Systems. In *Proc. of ESEC/FSE'15*. 284–294. https://doi.org/10.1145/2786805.2786845

[79] Norbert Siegmund, Sergiy S Kolesnikov, Christian Kästner, Sven Apel, Don Batory, Marko Rosenmüller, and Gunter Saake. 2012. Predicting performance via automated feature-interaction detection. In *Proc. of ICSE'12*. 167–177.

[80] Norbert Siegmund, Marko Rosenmüller, Martin Kuhlemann, Christian Kästner, Sven Apel, and Gunter Saake. 2012. SPL Conqueror: Toward Optimization of Non-Functional Properties in Software Product Lines. *Software Quality Journal* 20, 3–4 (Sept. 2012), 487–517. https://doi.org/10.1007/s11219-011-9152-9

[81] Urjoshi Sinha, Mikaela Cashman, and Myra B. Cohen. 2020. Using a Genetic Algorithm to Optimize Configurations in a Data-Driven Application. In *Proc. of SSBSE'20*. 137–152. https://doi.org/10.1007/978-3-030-59762-7_10

[82] Michael Still. 2006. ImageMagick.

[83] Paul Temple, Mathieu Acher, Battista Biggio, Jean-Marc Jézéquel, and Fabio Roli. 2018. Towards Adversarial Configurations for Software Product Lines. arXiv:1805.12021 [cs.LG]

[84] Paul Temple, Mathieu Acher, Jean-Marc Jezequel, and Olivier Barais. 2017. Learning Contextual-Variability Models. *IEEE Software* 34, 6 (Nov. 2017), 64–70.

[85] Paul Temple, Mathieu Acher, Jean-Marc A Jézéquel, Léo A Noel-Baron, and José A Galindo. 2017. *Learning-Based Performance Specialization of Configurable Systems*. Research Report. IRISA. https://hal.archives-ouvertes.fr/hal-01467299

[86] Paul Temple, José A. Galindo, Mathieu Acher, and Jean-Marc Jézéquel. 2016. Using machine learning to infer constraints for product lines. In *Proc. of SPLC'16*.

[87] Chris Thornton, Frank Hutter, Holger H. Hoos, and Kevin Leyton-Brown. 2013. Auto-WEKA: Combined Selection and Hyperparameter Optimization of Classification Algorithms. In *Proc. of KDD'13*. 847–855. https://doi.org/10.1145/2487575.2487629

[88] Pavel Valov, Jianmei Guo, and Krzysztof Czarnecki. 2015. Empirical Comparison of Regression Methods for Variability-Aware Performance Prediction. In *Proc. of SPLC'15*. 186–190.

[89] Pavel Valov, Jianmei Guo, and Krzysztof Czarnecki. 2020. Transferring Pareto Frontiers across Heterogeneous Hardware Environments. In *Proc. of ICPE'20*. 12–23. https://doi.org/10.1145/3358960.3379127

[90] Pavel Valov, Jean-Christophe Petkovich, Jianmei Guo, Sebastian Fischmeister, and Krzysztof Czarnecki. 2017. Transferring Performance Prediction Models Across Different Hardware Platforms. In *Proc. of ICPE'17*. 39–50.

[91] Dana Van Aken, Andrew Pavlo, Geoffrey J. Gordon, and Bohan Zhang. 2017. Automatic Database Management System Tuning Through Large-Scale Machine Learning. In *Proc. of ICMD'17*. 1009–1024.

[92] Miguel Velez, Pooyan Jamshidi, Norbert Siegmund, Sven Apel, and Christian Kästner. 2021. White-Box Analysis over Machine Learning: Modeling Performance of Configurable Systems. arXiv:2101.05362 [cs.SE]

[93] Yilin Wang, Sasi Inguva, and Balu Adsumilli. 2019. YouTube UGC Dataset for Video Compression Research. In *Proc. of MMSP'19*. 1–5. https://doi.org/10.1109/mmsp.2019.8901772

[94] Max Weber, Sven Apel, and Norbert Siegmund. 2021. White-Box Performance-Influence Models: A Profiling and Learning Approach. arXiv:2102.06395 [cs.SE]

[95] Markus Weckesser, Roland Kluge, Martin Pfannemüller, Michael Matthé, Andy Schürr, and Christian Becker. 2018. Optimal Reconfiguration of Dynamic Software Product Lines Based on Performance-Influence Models. In *Proc. of SPLC'18*.

[96] Dennis Westermann, Jens Happe, Rouven Krebs, and Roozbeh Farahbod. 2012. Automated Inference of Goal-Oriented Performance Prediction Functions. In *Proc. of ASE'12*. 190–199. https://doi.org/10.1145/2351676.2351703

[97] Lin Xu, Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. 2008. SATzilla: portfolio-based algorithm selection for SAT. *Journal of artificial intelligence research* 32 (2008), 565–606.

[98] Y. Zhang, J. Guo, E. Blais, and K. Czarnecki. 2015. Performance Prediction of Configurable Software Systems by Fourier Learning (T). In *Proc. of ASE'15*. 365–373. https://doi.org/10.1109/ASE.2015.15

[99] Yi Zhang, Jianmei Guo, Eric Blais, Krzysztof Czarnecki, and Huiqun Yu. 2016. A Mathematical Model of Performance-Relevant Feature Interactions. In *Proc. of SPLC'16*. 25–34. https://doi.org/10.1145/2934466.2934469

[100] Marcela Zuluaga, Andreas Krause, and Markus Püschel. 2016. e-PAL: An Active Learning Approach to the Multi-Objective Optimization Problem. *Journal of Machine Learning Research* 17, 104 (2016), 1–32.

[101] William R. Zwick and Wayne F. Velicer. 1982. Factors Influencing Four Rules For Determining The Number Of Components To Retain. *Multivariate Behavioral Research* 17, 2 (1982), 253–269.

[102] Ivan Švogor, Ivica Crnković, and Neven Vrček. 2019. An extensible framework for software configuration optimization on heterogeneous computing systems: Time and energy case study. *Information and Software Technology* 105 (2019).