

# Program Analysis for Adaptive Data Analysis

ANONYMOUS AUTHOR(S)

Data analyses are usually designed to identify some property of the population from which the data are drawn, generalizing beyond the specific data sample. For this reason, data analyses are often designed in a way that guarantees that they produce a low generalization error. That is, they are designed so that the result of a data analysis run on a sample data does not differ too much from the result one would achieve by running the analysis over the entire population.

An adaptive data analysis can be seen as a process composed by multiple queries interrogating some data, where the choice of which query to run next may rely on the results of previous queries. The generalization error of each individual query/analysis can be controlled by using an array of well-established statistical techniques. However, when queries are arbitrarily composed, the different errors can propagate through the chain of different queries and bring to high generalization error. To address this issue, data analysts are designing several techniques that not only guarantee bounds on the generalization errors of single queries, but that also guarantee bounds on the generalization error of the composed analyses. The choice of which of these techniques to use, often depends on the chain of queries that an adaptive data analysis can generate.

In this work, we consider adaptive data analyses implemented as while-like programs and we design a program analysis which can help with identifying which technique to use to control their generalization error. More specifically, we formalize the intuitive notion of *adaptivity* as a quantitative property of programs. We do this because the adaptivity level of a data analysis is a key measure to choose the right technique. Based on this definition, we design a program analysis for soundly approximating this quantity. The program analysis generates a representation of the data analysis as a weighted dependency graph, where the weight is an upper bound on the number of times each variable can be reached, and uses a path search strategy to guarantee an upper bound on the adaptivity. We implement our program analysis and show that it can help to analyze the adaptivity of several concrete data analyses with different adaptivity structures.

Additional Key Words and Phrases: Adaptive data analysis, program analysis, dependency graph

## 1 INTRODUCTION

Consider a dataset  $X$  consisting of  $n$  independent samples from some unknown population  $P$ . How can we ensure that the conclusions drawn from  $X$  *generalize* to the population  $P$ ? Despite decades of research in statistics and machine learning on methods for ensuring generalization, there is an increased recognition that many scientific findings generalize poorly (e.g. [19, 26]). While there are many reasons a conclusion might fail to generalize, one that is receiving increasing attention is *adaptivity*, which occurs when the choice of method for analyzing the dataset depends on previous interactions with the same dataset [19]. Adaptivity can arise from many common practices, such as exploratory data analysis, using the same data set for feature selection and regression, and the re-use of datasets across research projects. Unfortunately, adaptivity invalidates traditional methods for ensuring generalization and statistical validity, which assume that the method is selected independently of the data. The misinterpretation of adaptively selected results has even been blamed for a “statistical crisis” in empirical science [19].

A line of work initiated by Dwork et al. [15], Hardt and Ullman [25] posed the question: Can we design *general-purpose* methods that ensure generalization in the presence of adaptivity, together with guarantees on their accuracy? The idea that has emerged in these works is to use randomization to help ensure generalization. Specifically, these works have proposed to mediate the access of an adaptive data analysis to the data by means of queries from some pre-determined family (we will consider here a specific family of queries often called “statistical” or “linear” queries) that are sent to a *mechanism* which uses some randomized process to guarantee that the result of the query does not depend too much on the specific sampled dataset. This guarantees that the result of the queries generalizes well. This approach is described in Fig. 1. This line of work has identified many new

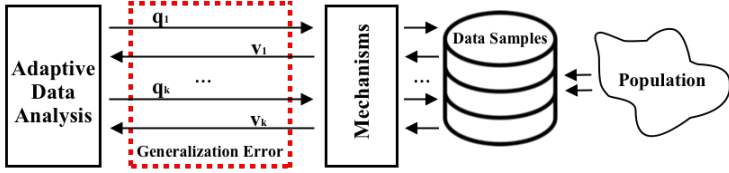


Fig. 1. Overview of our Adaptive Data Analysis model. We have a population that we are interested in studying, and a dataset containing individual samples from this population. The adaptive data analysis we are interested in running has access to the dataset through queries of some pre-determined family (e.g., statistical or linear queries) mediated by a mechanism. This mechanism uses randomization to reduce the generalization error of the queries issued to the data.

algorithmic techniques for ensuring generalization in adaptive data analysis, leading to algorithms with greater statistical power than all previous approaches. Common methods proposed by these works include, the addition of noise to the result of a query, data splitting, etc. Moreover, these works have also identified problematic strategies for adaptive analysis, showing limitations on the statistical power one can hope to achieve. Subsequent works have then further extended the methods and techniques in this approach and further extended the theoretical underpinning of this approach, e.g. [7, 13, 14, 16, 27, 31, 35, 36].

A key development in this line of work is that the best method for ensuring generalization in an adaptive data analysis depends to a large extent on the number of *rounds of adaptivity*, the depth of the chain of queries. As an informal example, the program  $x \leftarrow q_1(D); y \leftarrow q_2(D, x); z \leftarrow q_3(D, y)$  has three rounds of adaptivity, since  $q_2$  depends on  $D$  not only directly because it is one of its input but also via the result of  $q_1$ , which is also run on  $D$ , and similarly,  $q_3$  depends on  $D$  directly but also via the result of  $q_2$ , which in turn depends on the result of  $q_1$ . The works we discussed above showed that, not only does the analysis of the generalization error depend on the number of rounds, but knowing the number of rounds actually allows one to choose methods that lead to the smallest possible generalization error - we will discuss this further in Section 2.

For example, these works showed that when an adaptive data analysis uses a large number of rounds of adaptivity then a low generalization error can be achieved by a mechanism adding to the result of each query Gaussian noise scaled to the number of rounds. When instead an adaptive data analysis uses a small number of rounds of adaptivity then a low generalization error can be achieved by using more specialized methods, such as data splitting mechanism or the reusable holdout technique from Dwork et al. [15]. To better understand this idea, we show in Fig. 2 three experiments showcasing these situations. More precisely, in Fig. 2(a) we show the results of a specific analysis<sup>1</sup> with two rounds of adaptivity. This analysis can be seen as a classifier which first runs 400 non-adaptive queries on the first 400 attributes of the data, looking for correlations between the attributes and a label, and then runs one last query which depends on all these correlations. Without any mechanism the generalization error of the last query is pretty large, and the lower generalization error is achieved when the data-splitting method is used. Fig. 2(c) shows how this situation also change with the number of queries. Specifically, it shows the root mean square error of the last *adaptive* query when the numbers queries varies. This also highlight the fact that different mechanisms, for the same analysis, produce results with very different generalization error. In Fig. 2(b), we show the results of a specific analysis<sup>2</sup> with four hundreds rounds of adaptivity. At each step, this analysis runs an adaptive query based on the results of the previous ones. Without

<sup>1</sup>We will use formally a program implementing this analysis (Fig. 3) as a running example in the rest of the paper.

<sup>2</sup>We will present this analysis formally in Section 6.

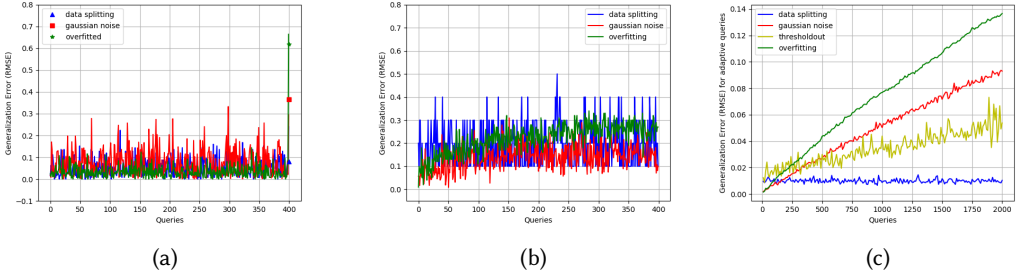


Fig. 2. The generalization errors of two adaptive data analysis examples, under different choices of mechanisms. (a) Data analysis with 2 rounds adaptivity, (b) Data analysis with 400 rounds adaptivity. (c) Same one as (a) any mechanism, the generalization error of most of the queries is pretty large, and this error can be lowered by using Gaussian noise.

This scenario motivates us to explore the design of program analysis techniques that can be used to estimate the number of *rounds of adaptivity* that a program implementing a data analysis can perform. These techniques could be used to help a data analyst in the choice of the mechanism to use, and they could ultimately be integrated into a tool for adaptive data analysis such as the *Guess and Check* framework by Rogers et al. [31].

The first problem we face is *how to formally define* a model for adaptive data analysis which is general enough to support the methods we discussed above and which would permit to formulate the notion of adaptivity these methods use. We take the approach of designing a programming framework for submitting queries to some *mechanism* giving access to the data mediated by one of the techniques we mentioned before, e.g., adding Gaussian noise, randomly selecting a subset of the data, using the reusable holdout technique, etc. In this approach, a program models an *analyst* asking a sequence of queries to the mechanism. The mechanism runs the queries on the data applying one of the methods above and returns the result to the program. The program can then use this result to decide which query to run next. Overall, we are interested in controlling the generalization of the query results returned by the mechanism, by means of the adaptivity.

The second problem we face is *how to define the adaptivity of a given program*. Intuitively, a query  $Q$  may depend on another query  $P$ , if there are two values that  $P$  can return which affect in different ways the execution of  $Q$ . For example, as shown in [14], and as we did in our example in Fig. 2(a), one can design a machine learning algorithm for constructing a classifier which first computes each feature’s correlation with the label via a sequence of queries, and then constructs the classifier based on the correlation values. If one feature’s correlation changes, the classifier depending on features is also affected. This notion of dependency builds on the execution trace as a *causal history*. In particular, we are interested in the history or provenance of a query up until this is executed, we are not then concerned about how the result is used — except for tracking whether the result of the query may further cause some other query. This is because we focus on the generalization error of queries and not their post-processing. To formalize this intuition as a quantitative program property, we use a trace semantics recording the execution history of programs on some given input — and we create a dependency graph, where the dependency between different variables (queries are also assigned to variables) is explicit and track which variable is associated with a query request. We then enrich this graph with weights describing the number of times each variable is evaluated in a program evaluation starting with an initial state. The adaptivity is then defined as the length of the walk visiting most query-related variables on this graph<sup>3</sup>. In other words, we define adaptivity as a *quantitative form of program dependency*.

<sup>3</sup>Formally, graphs will be well-defined only for terminating programs, this will guarantee that the longest walk is finite

The third problem we face is *how to estimate the adaptivity of a given program*. The adaptive data analysis model we consider and our definition of adaptivity suggest that for this task we can use a program analysis that is based on some form of dependency analysis. This analysis needs to take into consideration: 1) the fact that, in general, a query  $Q$  is not a monolithic block but rather it may depend, through the use of variables and values, on other parts of the program. Hence, it needs to consider some form of data flow analysis. 2) the fact that, in general, the decision on whether to run a query or not may depend on some other value. Hence, it needs to consider some form of control flow analysis. 3) the fact that, in general, we are not only interested in whether there is a dependency or not, but in the length of the chain of dependencies. Hence, it needs to consider some quantitative information about the program dependencies.

To address these considerations and be able to estimate a sound upper bound on the adaptivity of a program, we develop a static program analysis algorithm, named *AdaptFun*, which combines data flow and control flow analysis with reachability bound analysis [22]. This combination gives tighter bounds on the adaptivity of a program than the ones one would achieve by directly using the data and control flow analyses or the ones that one would achieve by directly using reachability bound analysis techniques alone. We evaluate *AdaptFun* on a number of examples showing that it is able to efficiently estimate precise upper bounds on the adaptivity of different programs. All the proofs and extended definitions can be found in the supplementary material.

To summarize, our work aims at the design of a static analysis for programs implementing adaptive analysis that can estimate their rounds of adaptivity. Specifically, our contributions are:

- (1) A programming framework for adaptive data analyses where programs represent analysts that can query generalization-preserving mechanisms mediating the access to some data.
- (2) A formal definition of the notion of adaptivity under the analyst-mechanism model. This definition is built on a variable-based dependency graph that is constructed using sets of program execution traces.
- (3) A static program analysis algorithm *AdaptFun* combining data flow, control flow and reachability bound analysis in order to provide tight bounds on the adaptivity of a program.
- (4) A soundness proof of the program analysis showing that the adaptivity estimated by *AdaptFun* bounds the true adaptivity of the program.
- (5) An implementation of *AdaptFun* and an experimental evaluation of the bounds this implementation provides on several examples.

## 2 OVERVIEW

### 2.1 Some results in Adaptive Data Analysis

In Adaptive Data Analysis an *analyst* is interested in studying some distribution  $P$  over some domain  $\mathcal{X}$ . Following previous works [7, 15, 25], we focus on the setting where the analyst is interested in answers to *statistical queries* (also known as *linear queries*) over the distribution. A statistical query is usually defined by some function  $\text{query} : \mathcal{X} \rightarrow [-1, 1]$  (often other codomains such as  $[0, 1]$  or  $[-R, +R]$ , for some  $R$ , are considered). The analyst wants to learn the *population mean*, which (abusing notation) is defined as  $\text{query}(P) = \mathbb{E}_{X \sim P} [\text{query}(X)]$ . We assume that the distribution  $P$  can only be accessed via a set of *samples*  $X_1, \dots, X_n$  drawn independently and identically distributed (i.i.d.) from  $P$ . These samples are held by a mechanism  $M(X_1, \dots, X_n)$  who receives the query  $\text{query}$  and computes an answer  $a \approx \text{query}(P)$ . The naïve way to approximate the population mean is to use the *empirical mean*, which (abusing notation) is defined as  $\text{query}(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n \text{query}(X_i)$ . However, the mechanism  $M$  can adopt some methods for improving the generalization error  $|a - \text{query}(P)|$ .

In this work we consider analysts that ask a sequence of  $k$  queries  $\text{query}_1, \dots, \text{query}_k$ . If the queries are all chosen in advance, independently of the answers of each other, then we say they are *non-adaptive*. If the choice of each query  $\text{query}_j$  depend on the prefix  $\text{query}_1, a_1, \dots, \text{query}_{j-1}, a_{j-1}$  then they are *fully adaptive*. An important intermediate notion is  *$r$ -round adaptive*, where the sequence can be partitioned into  $r$  batches of non-adaptive queries. Note that non-adaptive queries are 1-round and fully adaptive queries are  $k$ -round adaptive.

We now review what is known about the problem of answering  $r$ -round adaptive queries.

**THEOREM 2.1 ([7]).** (1) For any distribution  $P$ , and any  $k$  non-adaptive statistical queries,

$$\max_{j=1, \dots, k} |a_j - \text{query}_j(P)| = O\left(\sqrt{\frac{\log k}{n}}\right).$$

(2) For any distribution  $P$ , and any  $k$   $r$ -round adaptive statistical queries, with  $r \geq 2$ , the empirical mean (rounded to an appropriate number of bits of precision)<sup>4</sup> satisfies:

$$\max_{j=1, \dots, k} |a_j - \text{query}_j(P)| = O\left(\sqrt{\frac{k}{n}}\right)$$

In fact, these bounds are tight (up to constant factors) which means that even allowing one extra round of adaptivity leads to an exponential increase in the generalization error, from  $\log k$  to  $k$ .

Dwork et al. [15] and Bassily et al. [7] showed that by using carefully calibrated Gaussian noise in order to limit the dependency of a single query on the specific data instance, one can actually achieve much stronger generalization error as a function of the number of queries, specifically.

**THEOREM 2.2 ([7, 15]).** For any distribution  $P$ , any  $k$ , any  $r \geq 2$  and any  $r$ -round adaptive statistical queries, if we answer queries with carefully calibrated Gaussian noise we have:

$$\max_{j=1, \dots, k} |a_j - \text{query}_j(P)| = O\left(\frac{\sqrt{k}}{\sqrt{n}}\right)$$

More interestingly, Dwork et al. [15] also gave a refined bounds that can be achieved with different mechanisms depending on the number of rounds of adaptivity.

**THEOREM 2.3 ([15]).** For any  $r$  and  $k$ , there exists a mechanism such that for any distribution  $P$ , and any  $r \geq 2$  any  $r$ -round adaptive statistical queries, it satisfies

$$\max_{j=1, \dots, k} |a_j - \text{query}_j(P)| = O\left(\frac{r\sqrt{\log k}}{\sqrt{n}}\right)$$

Notice that Theorem 2.3 has different quantification in that the optimal choice of mechanism depends on the number of queries and number of rounds of adaptivity. This suggests that if one knows a good *a priori upper bound on the number of rounds of adaptivity*, one can choose the appropriate mechanism and get a much better guarantee in terms of generalization error. As an example, as we can see in Fig. 2, if we know that an algorithm is two rounds adaptive, we can choose data splitting as the mechanism, while if we know that an algorithm has many rounds of adaptivity we can choose Gaussian noise. It is worth to stress that by knowing the number of rounds of adaptivity one can also compute a concrete upper bound on the generalization error of a data analysis. This information allows one to have a quantitative, a priori, estimation of the effectiveness of a data analysis. This motivates us to design a static program analysis aimed at giving good *a priori* upper bounds on the number of rounds of adaptivity of a program.

## 2.2 AdaptFun formally through an example.

We illustrate the key technical components of our framework through a simple adaptive data analysis with two rounds of adaptivity. In this analysis, an analyst asks  $k+1$  queries to a mechanism

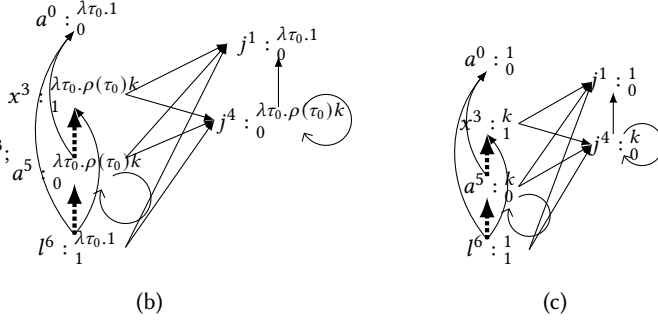
<sup>4</sup>With infinite precision even two queries may give unbounded error, when the first query's result encodes the whole data.

246  
247  
248  
249  
250  
251  
252  
253  
254

```

twoRounds(k)  $\triangleq$ 
[a  $\leftarrow$  0]0; [j  $\leftarrow$  k]1;
while [j > 0]2 do (
[x  $\leftarrow$  query( $\chi[j] \cdot \chi[k]$ )]3;
[j  $\leftarrow$  j - 1]4;
[a  $\leftarrow$  x + a]5);
[l  $\leftarrow$  query( $\chi[k] * a$ )]6

```



255

256 Fig. 3. (a) The program `twoRounds(k)`, an example with two rounds of adaptivity (b) The corresponding  
257 semantics-based dependency graph (c) The estimated dependency graph from AdaptFun.

258

259 in two phases. In the first phase, the analyst asks  $k$  queries and stores the answers that are provided  
260 by the mechanism. In the second phase, the analyst constructs a new query based on the results of  
261 the previous  $k$  queries and sends this query to the mechanism. The mechanism is abstract here  
262 and our goal is to use static analysis to provide an upper bound on adaptivity to help choose  
263 the mechanism. This data analysis assumes that the data domain  $\mathcal{X}$  contains at least  $k$  numeric  
264 attributes (every query in the first phase focuses on one), which we index just by natural numbers.  
265 The implementation of this data analysis in the language of AdaptFun is presented in Fig. 3(a).

266

267 The AdaptFun language extends a standard while language<sup>5</sup> with a query request constructor  
268 denoted `query( $\psi$ )`. Queries have the form `query( $\psi$ )`, where  $\psi$  is a special expression (see syntax in  
269 Section 3) representing a function  $\mathcal{X} \rightarrow U$  on rows. We use  $U$  to denote the codomain of queries  
270 and it could be  $[-1, 1]$ ,  $[0, 1]$  or  $[-R, +R]$ , for some  $R$  we consider. This function characterizes the  
271 linear query we are interested in running. Indeed, as we discussed in the previous section, linear  
272 queries compute the empirical mean of a function on rows – we use  $\chi$  to abstract a possible row in  
273 the database. As an example, `x  $\leftarrow$  query( $\chi[j] \cdot \chi[k]$ )` computes an approximation, according to  
274 the used mechanism, of the empirical mean of the product of the  $j^{\text{th}}$  attribute and  $k^{\text{th}}$  attribute,  
275 identified by  $\chi[j] \cdot \chi[k]$ . Notice that we don't materialize the mechanism but we assume that it is  
276 implicitly run when we execute the query. In Fig. 3(a), the queries inside the while loop correspond  
277 to the first phase of the data analysis and compute an approximation of the product of the empirical  
278 mean of the first  $k$  attributes. The query outside the loop corresponds to the second phase and  
279 computes an approximation of the empirical mean where each record is weighted by the sum of  
280 the empirical mean of the first  $k$  attributes.

280

281 This example is intuitively 2-rounds adaptive since we have two clearly distinguished phases,  
282 and the queries that we ask in the first phase do not depend on each other (the query  $\chi[j] \cdot \chi[k]$  at  
283 line 3 only relies on the counter  $j$  and input  $k$ ), while the last query (at line 6) depends on the results  
284 of all the previous queries. However, capturing this concept formally is surprisingly difficult. The  
285 difficulty comes from the fact that a query can depend on the result of another query in multiple  
286 ways, by means of data dependency or control flow dependency.

286 **2.2.1 Adaptivity definition.** The central property we are after in this work is the *adaptivity of a*  
287 *program*. We define formally this notion in three steps, which we will describe in details in Section 4.  
288 First, we define a notion of dependency, or better *may-dependency*, between variables. To do this we  
289 take inspiration from previous works on dependency analysis and information flow control and we  
290 say that a variable *may depend* on another one if changing the execution of the latter can affect the  
291 execution of the former. We can see in Fig. 3(a) that the value of the variable  $l$ , which corresponds  
292

293 <sup>5</sup>Programs components are labeled, so that we can uniquely identify every component.

294



295 to the result of the execution of the query in the second phase (in the command with label 6), is  
 296 affected by the value of the variable  $x$ , which corresponds to the result of the execution of the query  
 297 at line 3 in the first phase, via the variable  $a$ . To formally define this notion of dependency, as in  
 298 information flow control, we use the execution history of programs recorded by a trace semantics  
 299 (see Definition 3).

300 Second, we build an annotated weighted directed graph representing the possible dependencies  
 301 between labeled variables. We call this graph *semantics-based dependency graph* to stress that this  
 302 graph summarize the dependencies we could see if we knew the overall behavior of the program. The  
 303 vertices of the graph are the assigned program variables with the label of their assignments, edges are  
 304 pairs of labeled variables which satisfy the dependency relations, weights are functions associated  
 305 with vertexes and describing the number of times the assignment corresponding to the vertex is  
 306 executed when the program is run in a given starting state<sup>6</sup>, and the annotations, which we call  
 307 *query annotations*, are bits associated with vertexes and describing if the corresponding assignment  
 308 comes from a query (1) or not (0). The *semantics-based dependency graph* of the twoRounds( $k$ )  
 309 program we gave in Fig. 3(a) is described in Fig. 3(b) (we use dashed arrows for two edges that  
 310 will be highlighted in the next step, for the moment these can be considered similar to the other  
 311 edges—i.e. solid arrows). We have all the variables that are assigned in the program with their  
 312 labels, and edges representing dependency relations between them. For example, we have two  
 313 edges  $(l^6, a^5)$  and  $(a^5, x^3)$  describing the dependency between the variables assigned by queries.  
 314 The vertices  $l^6$  and  $x^3$  are the only ones with query annotation 1 (the subscript), since they are  
 315 the only two variables that are in assignments involving queries. Notice that the graph contains  
 316 cycles—in this example it contains two self-loops. These cycles capture the fact that the variables  $a^5$   
 317 and  $j^4$  are updated at every iteration of the loop using their previous value. Cycles are essential to  
 318 capture mutual dependencies like the ones that are generated in loops. Adaptivity is a quantitative  
 319 notion, so capturing this form of dependencies is not enough. This is why we also use weights. The  
 320 weight of a vertex is a function that given an initial state returns a natural number representing the  
 321 number of times the assignment corresponding to a vertex is visited during the program execution  
 322 starting in this initial state. For example, the vertex  $l^6$  has weight  $\lambda\tau.1$  since for every initial state  $\tau$   
 323 the corresponding assignment will be executed one time, the vertex  $a^5$  on the other hand has weight  
 324  $\lambda\tau.\rho(\tau)k$  since the corresponding assignment will be executed a number of times that correspond  
 325 to the value of  $k$  in the initial state  $\tau$ , and  $\rho$  is the operator reading value of  $k$  from  $\tau$ .

326 Third, we can finally define adaptivity using the semantics-based dependency graph. We actually  
 327 define this notion with respect to an initial state  $\tau$ , since different states can give very different  
 328 adaptivity. We consider the longest walk that visits each vertex  $v$  of the semantics-based dependency  
 329 graph no more than the value that the weight  $w_v$  assign to  $\tau$ , and visits as many query nodes as  
 330 possible. The number of query nodes visited is the adaptivity of the program with respect to  $\tau$ .  
 331 Looking again at Fig. 3(b), and assuming that  $\tau(k) \geq 1$ , we can see that the the walk along the  
 332 dashed arrows,  $l^6 \rightarrow a^5 \rightarrow x^3$  has two vertices with query annotation 1, and we cannot find another  
 333 walk having more than 2 vertices with query annotation 1. So the adaptivity of the program in  
 334 Fig. 3(a) with respect to  $\tau$  is 2. If we consider an initial state  $\tau$  such that  $\tau(k) = 0$  we have that the  
 335 adaptivity with respect to  $\tau$  is instead 1.

336 **2.2.2 Static analysis.** To compute statically a sound and accurate upper bound on the *adaptivity* of  
 337 a program  $c$ , we design a program analysis framework named AdaptFun which we will describe  
 338 formally in 5. The structure of AdaptFun (Fig. 4) reflects in part the definition of adaptivity we  
 339 discussed in the previous section. Specifically, AdaptFun is composed by two algorithms (the ones  
 340

341 <sup>6</sup>In our trace semantics the state is recorded in the trace, so an initial state is actually represented by an initial trace. We will  
 342 use this terminology in later sections.  
 343

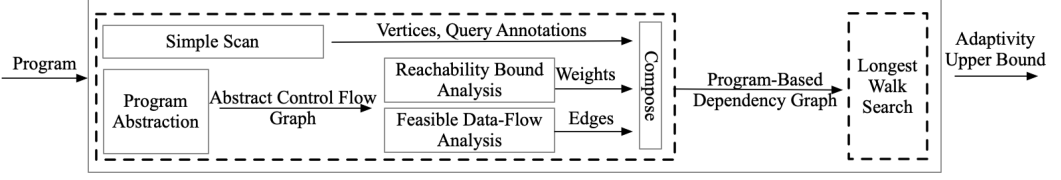


Fig. 4. The overview of AdaptFun

in dashed boxes in the figure), one for building a dependency graph, which we call *estimated dependency graph*, and the other to estimate the adaptivity from this graph. The first algorithm, which we will describe formally in Section 5, generates the *estimated dependency graph* using several program analysis techniques. Specifically, AdaptFun extracts the vertices and the query annotations by looking at the assigned variables of the program, it estimates the edges by using control flow and data flow analysis, and it estimates the weights by using symbolic reachability-bound analysis—weights in this graph are symbolic expressions over input variables. The second algorithm estimates the longest walk which respect the weights and which visit as many query nodes as possible. The two algorithm together gives us an upper bound on the program’s *adaptivity*.

We show in Fig. 3(c) the estimated dependency graph that our static analysis algorithm returns for the program `twoRounds(k)` in Fig. 3(a). Vertices and query annotations are the same as the ones in Fig. 3(b) and they are simply inferred by scanning the program. As we said before, the edges are estimated using control flow and data flow analysis. For the `twoRounds(k)` example, every edge in Fig. 3(b) is precisely inferred by our combined analysis, this is why Fig. 3(c) contains exactly the same edges. The weight of every vertex is computed using a reachability-bound estimation algorithm which output a symbolic expression over the input variables, in the example only  $k$ , representing an upper bound on the number of times each assignment is executed. For example, consider the vertex  $x^3$ , its weight is  $k$  and this provides an upper bound on the values returned by the weight function  $\lambda\tau.\rho(\tau)k$  associated with vertex  $x^3$  in Fig. 3(b) for any initial state.

The algorithm searching for the longest walk first finds a path  $l^6 : \frac{1}{1} \rightarrow a^5 : \frac{k}{1} \rightarrow x^3 : \frac{k}{1}$ , and then constructs a walk based on this path. Every vertex on this walk is visited once, and the number of vertices with query annotation 1 in this walk is 2, which is the upper bound we expect. It is worth to note here that  $x^3$  and  $a^5$  can only be visited once because there isn’t an edge to go back to them, even though they both have the weight  $k$ . In this sense, instead of simply computing the weighted length of this path  $(2k + 1)$  as adaptivity AdaptBD computes the upper bound 2. Note that 2 is not always tight, for example when  $k = 0$ .

### 3 LABELED QUERY WHILE LANGUAGE

The language of AdaptFun is a standard while language with labels to identify different components and with primitives for queries, and equipped with a trace-based operational semantics which is the main technical tool we will use to define the program’s adaptivity.

Arithmetic Expression	$a$	$::=$	$n \mid x \mid a \oplus_a a \mid \log a \mid \text{sign } a \mid \max(a, a) \mid \min(a, a)$
Boolean Expression	$b$	$::=$	$\text{true} \mid \text{false} \mid \neg b \mid b \oplus_b b \mid a \sim a$
Expression	$e$	$::=$	$v \mid a \mid b \mid [e, \dots, e]$
Value	$v$	$::=$	$n \mid \text{true} \mid \text{false} \mid [] \mid [v, \dots, v]$
Query Expression	$\psi$	$::=$	$\alpha \mid a \mid \psi \oplus_a \psi \mid \chi[a]$
Query Value	$\alpha$	$::=$	$n \mid \chi[n] \mid \alpha \oplus_a \alpha \mid n \oplus_a \chi[n] \mid \chi[n] \oplus_a n$
Label	$l$	$\in$	$\mathbb{N} \cup \{\text{in}, \text{ex}\}$
Labeled Command	$c$	$::=$	$[x \leftarrow e]^l \mid [x \leftarrow \text{query}(\psi)]^l \mid \text{while } [b]^l \text{ do } c \mid c; c \mid \text{if } ([b]^l, c, c) \mid [\text{skip}]^l$



Expressions include standard arithmetic (with value  $n \in \mathbb{N}^\infty$ ) and boolean expression, ( $a$  and  $b$ ) and extended query expressions  $\psi$ . A query expression  $\psi$  can be either a simple arithmetic expression  $a$ , an expression of the form  $\chi[a]$  where  $\chi$  represents a row of the database and  $a$  represents an index used to identify a specific attribute of the row  $\chi$ , a combination of two query expressions,  $\psi \oplus_a \psi$ , or a normal form  $\alpha$ . For example, the query expression  $\chi[3] + 5$  denotes the computation that obtains the value in the 3rd column of  $\chi$  in one row and then add 5 to it.

Command are the typical ones from while languages with an additional command  $x \leftarrow \text{query}(\psi)$  for query requests which can be used to interrogate the database and compute the linear query corresponding to  $\psi$ . Each command is annotated with a label  $l$ , we will use natural numbers as labels and we will use them to record the location of each command, so that we can uniquely identify them. We also have a set  $\mathcal{LV}$  of labeled variables, these are simply variables with a label. We denote by  $\text{LV}(c)$  the set of labeled variables which are assigned in an assignment command in the program  $c$ . We denote by  $\text{QV}(c)$  the set of labeled variables that are assigned to the result of a query in the program  $c$ .

### 3.1 Trace-based Operational Semantics

We use a trace based operational semantics tracking the history of programs execution. The operational semantics is parameterized by a database that can be access only through queries. Since this database is fixed, we omit it from the semantics but it is important to keep in mind that this database exists and it is what allow us to evaluate queries. A *trace*  $\tau$  is a list of *events* generated when executing specific commands. We denote by  $\mathcal{T}$  the set of traces and we will use list notation for traces, where  $[]$  is the empty trace, the operator  $::$  combines an event and a trace in a new event, and the operator  $++$  concatenates two traces.

We have two kinds of events: *assignment events* and *testing events*. Each event consists of a quadruple, and we use  $\mathcal{E}^{\text{asn}}$  and  $\mathcal{E}^{\text{test}}$  to denote the set of all assignment events and testing events, respectively.

$$\begin{array}{ll} \text{Event } \epsilon ::= & (x, l, v, \bullet) \mid (x, l, v, \alpha) \quad \text{Assignment Event} \\ & \mid (b, l, v, \bullet) \quad \text{Testing Event} \end{array}$$

An assignment event tracks the execution of an assignment or a query request and consists of the assigned variable, the label of the command that generates it, the value assigned to the variable, and the normal form of the query expression,  $\alpha$  if this command is a query request, otherwise a default value  $\bullet$ . A testing event tracks the execution of if and while commands and consists of the guard of the command, the label of the command, the result of evaluating the guard, while the last element is  $\bullet$ . We use the operator  $\rho(\tau)x$  to fetch the latest value assigned to  $x$  in the trace  $\tau$ . We use the operator  $\text{cnt}$  to count the occurrence of a labeled variable in the trace. We denote by  $\text{TL}(\tau) \subseteq \mathcal{L}$  the set of the labels occurring in  $\tau$ . Finally, we use  $\mathcal{T}_0(c) \subseteq \mathcal{T}$  to denote the set of *initial traces*, the ones which assign a value to the input variables.

The trace based operational semantics is described in terms of a small step evaluation relation  $\langle c, \tau \rangle \rightarrow \langle c', \tau' \rangle$  describing how a configuration program-trace evaluates to another configuration program-state. The rules for the operational semantics are described in Fig. 5. The rules for assignment and query generate assignment events, while the rules for while and if generate testing events. The rules for the standard while language constructs correspond to the usual rules extended to deal with traces. We have relations  $\langle \tau, e \rangle \Downarrow_e v$  and  $\langle \tau, b \rangle \Downarrow_b v$  to evaluate expressions and boolean expressions, respectively. Their definitions are in the supplementary material. The only rule that is non-standard is the **query** rule. When evaluating a query, the query expression  $\psi$  is first simplified to its normal form  $\alpha$  using an evaluation relation  $\langle \tau, \psi \rangle \Downarrow_q \alpha$ . Then normal form  $\alpha$  characterize the linear query that is run against the database. The query result  $v$  is the expected value of the

$$\begin{array}{c}
442 \\
443 \\
444 \\
445 \\
446 \\
447 \\
448 \\
449 \\
450 \\
451 \\
452 \\
453 \\
454 \\
455 \\
456 \\
457 \\
458 \\
459 \\
460 \\
461 \\
462 \\
463 \\
464 \\
465 \\
466 \\
467 \\
468 \\
469 \\
470 \\
471 \\
472 \\
473 \\
474 \\
475 \\
476 \\
477 \\
478 \\
479 \\
480 \\
481 \\
482 \\
483 \\
484 \\
485 \\
486 \\
487 \\
488 \\
489 \\
490
\end{array}$$

Command  $\times$  Trace  $\rightarrow$  Command  $\times$  Trace

$\langle c, \tau \rangle \rightarrow \langle c', \tau' \rangle$

$$\begin{array}{c}
\frac{\langle \tau, e \rangle \Downarrow_e v \quad \epsilon = (x, l, v, \bullet)}{\langle [x \leftarrow e]^l, \tau \rangle \rightarrow \langle [\text{skip}]^l, \tau::\epsilon \rangle} \text{assn} \qquad \frac{\tau, \psi \Downarrow_q \alpha \quad \text{query}(\alpha) = v \quad \epsilon = (x, l, v, \alpha)}{\langle [x \leftarrow \text{query}(\psi)]^l, \tau \rangle \rightarrow \langle [\text{skip}]^l, \tau::\epsilon \rangle} \text{query} \\
\frac{\tau, b \Downarrow_b \text{true} \quad \epsilon = (b, l, \text{true}, \bullet)}{\langle \text{while } [b]^l \text{ do } c, \tau \rangle \rightarrow \langle c; \text{while } [b]^l \text{ do } c, \tau::\epsilon \rangle} \text{while-t} \qquad \frac{\tau, b \Downarrow_b \text{false} \quad \epsilon = (b, l, \text{false}, \bullet)}{\langle \text{while } [b]^l, \text{ do } c, \tau \rangle \rightarrow \langle [\text{skip}]^l, \tau::\epsilon \rangle} \text{while-f} \\
\frac{\langle c_1, \tau \rangle \rightarrow \langle c'_1, \tau' \rangle}{\langle c_1; c_2, \tau \rangle \rightarrow \langle c'_1; c_2, \tau' \rangle} \text{seq1} \qquad \frac{\langle c_2, \tau \rangle \rightarrow \langle c'_2, \tau' \rangle}{\langle [\text{skip}]^l; c_2, \tau \rangle \rightarrow \langle c'_2, \tau' \rangle} \text{seq2} \qquad \frac{\tau, b \Downarrow_b \text{true} \quad \epsilon = (b, l, \text{true}, \bullet)}{\langle \text{if } ([b]^l, c_1, c_2), \tau \rangle \rightarrow \langle c_1, \tau::\epsilon \rangle} \text{if-t if-f}
\end{array}$$

Fig. 5. Trace-based Operational Semantics for Language.

function  $\lambda \chi. \alpha$  applied to each row of the dataset. We summarize this process with the notation  $\text{query}(\alpha) = v$  which we use in the rule **query**. Once the answer of the query is computed, the rules record all the needed information in the trace. As usual, we will use  $\rightarrow^*$  for the reflexive and transitive closure of  $\rightarrow$ .

The query expression evaluation relation  $\langle \tau, \psi \rangle \Downarrow_q \alpha$  is defined by the following rules which reduce a query expression to its normal form.

$$\begin{array}{c}
\frac{\langle \tau, a \rangle \Downarrow_a n}{\langle \tau, a \rangle \Downarrow_q n} \qquad \frac{\langle \tau, \psi_1 \rangle \Downarrow_q \alpha_1 \quad \langle \tau, \psi_2 \rangle \Downarrow_q \alpha_2}{\langle \tau, \psi_1 \oplus_a \psi_2 \rangle \Downarrow_q \alpha_1 \oplus_a \alpha_2} \qquad \frac{\langle \tau, a \rangle \Downarrow_a n}{\langle \tau, \chi[a] \rangle \Downarrow_q \chi[n]} \qquad \frac{}{\langle \tau, \alpha \rangle \Downarrow_q \alpha}
\end{array}$$

## 4 DEFINITION OF ADAPTIVITY

In this section, we formally present the definition of adaptivity for a given program. As we discussed in Section 2.2.1, we first define a dependency relation between program variables, we then define a semantics-based dependency graph, and finally look at longest walks in this graph.

### 4.1 May-dependency between variables

We are interested in defining a notion of dependencies between program variables since assigned variables are a good proxy to study dependencies between queries—we can recover query requests from variables associated with queries. We consider dependencies that can be generated by either data or control flow. For example, in the program

$$c_1 = [x \leftarrow \text{query}(\chi[2])]^1; [y \leftarrow \text{query}(\chi[3] + x)]^2$$

the query  $\text{query}(\chi[3] + x)$  depends on the query  $\text{query}(\chi[2])$  through a *value dependency* via  $x^1$ . Conversely, in the program

$$c_2 = [x \leftarrow \text{query}(\chi[1])]^1; \text{if } ([x > 2])^2, [y \leftarrow \text{query}(\chi[2])]^3, [\text{skip}]^4$$

the query  $\text{query}(\chi[2])$  depends on the query  $\text{query}(\chi[1])$  via the *control dependency* of the guard of the if command involving the labeled variable  $x^1$ .

To define dependency between program variables we will consider two events that are generated from the same command, hence they have the same variable name or boolean expression and label, but have either different value or different query expression, captured by the following definition.

**DEFINITION 1.** Two events  $\epsilon_1, \epsilon_2 \in \mathcal{E}$  differ in their value, or query value, denoted as  $\text{Diff}(\epsilon_1, \epsilon_2)$ , if and only if:

$$\pi_1(\epsilon_1) = \pi_1(\epsilon_2) \wedge \pi_2(\epsilon_1) = \pi_2(\epsilon_2) \tag{1a}$$

$$\wedge ((\pi_3(\epsilon_1) \neq \pi_3(\epsilon_2) \wedge \pi_4(\epsilon_1) = \pi_4(\epsilon_2) = \bullet) \vee (\pi_4(\epsilon_1) \neq \bullet \wedge \pi_4(\epsilon_2) \neq \bullet \wedge \pi_4(\epsilon_1) \neq_q \pi_4(\epsilon_2))) \tag{1b}$$

where  $\psi_1 =_q \psi_2$  denotes the semantics equivalence between query values<sup>7</sup>, and  $\pi_i$  projects the  $i$ -th element from the quadruple of an event.

We can now define when an event *may depend* on another one<sup>8</sup>.

**DEFINITION 2 (EVENT MAY-DEPENDENCY).** An event  $\epsilon_2 \in \mathcal{E}^{\text{asn}}$  may-depend on an event  $\epsilon_1 \in \mathcal{E}^{\text{asn}}$  in a program  $c$  denoted  $\text{DEP}_e(\epsilon_1, \epsilon_2, c)$ , if and only if

$$\exists \tau, \tau_0, \tau_1, \tau' \in \mathcal{T}, \epsilon'_1 \in \mathcal{E}^{\text{asn}}, c_1, c_2 \in C. \text{Diff}(\epsilon_1, \epsilon'_1) \wedge \quad (2a)$$

$$(\exists \epsilon'_2 \in \mathcal{E}. \left( \begin{array}{l} \langle c, \tau_0 \rangle \rightarrow^* \langle c_1, \tau_{1++}[\epsilon_1] \rangle \rightarrow^* \langle c_2, \tau_{1++}[\epsilon_1]_{++\tau_{++}[\epsilon_2]} \rangle \\ \wedge \langle c_1, \tau_{1++}[\epsilon'_1] \rangle \rightarrow^* \langle c_2, \tau_{1++}[\epsilon'_1]_{++\tau'_{++}[\epsilon'_2]} \rangle \\ \wedge \text{Diff}(\epsilon_2, \epsilon'_2) \wedge \text{cnt}(\tau, \pi_2(\epsilon_2)) = \text{cnt}(\tau', \pi_2(\epsilon'_2)) \end{array} \right)) \quad (2b)$$

$$\vee \left( \begin{array}{l} \exists \tau_3, \tau'_3 \in \mathcal{T}, \epsilon_b \in \mathcal{E}^{\text{test}}. \\ \langle c, \tau_0 \rangle \rightarrow^* \langle c_1, \tau_{1++}[\epsilon_1] \rangle \rightarrow^* \langle c_2, \tau_{1++}[\epsilon_1]_{++\tau_{++}[\epsilon_b]} \rangle \\ \wedge \langle c_1, \tau_{1++}[\epsilon'_1] \rangle \rightarrow^* \langle c_2, \tau_{1++}[\epsilon'_1]_{++\tau'_{++}[\neg\epsilon_b]} \rangle \\ \wedge \text{TL}(\tau_3) \cap \text{TL}(\tau'_3) = \emptyset \wedge \text{cnt}(\tau', \pi_2(\epsilon_b)) = \text{cnt}(\tau, \pi_2(\epsilon_b)) \wedge \epsilon_2 \in \tau_3 \wedge \epsilon_2 \notin \tau'_3 \end{array} \right), \quad (2c)$$

There are several components in this definition. The part with label (2a) requires that  $\epsilon_1$  and  $\epsilon'_1$  differ in their value ( $\text{Diff}(\epsilon_1, \epsilon'_1)$ ). The next two parts (2b) and (2c) capture the value dependency and control dependency, respectively. As in the literature on non-interference, and following [10], we formulate these dependencies as relational properties, i.e. in terms of two different traces of execution. We force these two traces to differ by using the event  $\epsilon_1$  in one and  $\epsilon'_1$  in the other. For the value dependency we check whether the change also create a change in the value of  $\epsilon_2$  or not. We additionally check that the two events we consider appear the same number of times in the two traces - this to make sure that if the events are generated by assignments in a loop, we consider the same iterations. For the control dependency we check whether the change in  $\epsilon_1$  affect the appearance in the computation of  $\epsilon_2$  or not. For this we require the presence of a test event whose value is affected by the change in  $\epsilon_1$  in order to guarantee that the computation goes through a control flow guard. Similarly to the previous condition, we additionally check that the two test events we consider appear the same number of times in the two traces.

We can now extend the dependency relation to variables by considering all the assignment events generated during the program's execution.

**DEFINITION 3 (VARIABLE MAY-DEPENDENCY).** A variable  $x_2^{l_2} \in \mathbb{L}\mathbb{V}(c)$  may-dependent on the variable  $x_1^{l_1} \in \mathbb{L}\mathbb{V}(c)$  in a program  $c$ ,  $\text{DEP}_{\text{var}}(x_1^{l_1}, x_2^{l_2}, c)$ , iff

$$\exists \epsilon_1, \epsilon_2 \in \mathcal{E}^{\text{asn}}, \tau \in \mathcal{T}. \pi_1(\epsilon_1)^{\pi_2(\epsilon_1)} = x_1^{l_1} \wedge \pi_1(\epsilon_2)^{\pi_2(\epsilon_2)} = x_2^{l_2} \wedge \text{DEP}_e(\epsilon_1, \epsilon_2, \tau, c)$$

Notice that in the definition above we can also have that the two variables are the same, this allow us to capture self-dependencies.

## 4.2 Semantics-based Dependency Graph

We can now define the *semantics-based dependency graph* of a program  $c$ . We want this graph to combines quantitative reachability information with dependency information.

<sup>7</sup>The formal definition is in the supplementary material

<sup>8</sup>We consider here dependencies between assignment events. This simplifies the definition and is enough for the stating the following definitions. The full definition is in the supplementary material.

540 **DEFINITION 4 (SEMANTICS-BASED DEPENDENCY GRAPH).** Given a program  $c$ , its semantics-based  
 541 dependency graph  $G_{\text{trace}}(c) = (V_{\text{trace}}(c), E_{\text{trace}}(c), W_{\text{trace}}(c), Q_{\text{trace}}(c))$  is defined as follows,

$$\begin{aligned}
 542 \text{ Vertices} \quad V_{\text{trace}}(c) &:= \{x^l \mid x^l \in \mathbb{L}\mathbb{V}(c)\} \\
 543 \text{ Directed Edges} \quad E_{\text{trace}}(c) &:= \{(x^l, y^j) \mid x^i, y^j \in \mathbb{L}\mathbb{V}(c) \wedge \text{DEP}_{\text{var}}(x^i, y^j, c)\} \\
 544 \text{ Weights} \quad W_{\text{trace}}(c) &:= \{(x^l, w) \mid w : \mathcal{T}_0(c) \rightarrow \mathbb{N} \wedge x^l \in \mathbb{L}\mathbb{V}(c) \\
 545 &\quad \wedge \forall \tau_0 \in \mathcal{T}_0(c), \tau' \in \mathcal{T}. \langle c, \tau_0 \rangle \rightarrow^* \langle \text{skip}, \tau_0 + \tau' \rangle \wedge w(\tau_0) = \text{cnt}(\tau', l)\} \\
 546 \text{ Query Annotations} \quad Q_{\text{trace}}(c) &:= \{(x^l, n) \mid x^l \in \mathbb{L}\mathbb{V}(c) \wedge (n = 1 \Leftrightarrow x^l \in \mathbb{Q}\mathbb{V}(c)) \wedge (n = 0 \Leftrightarrow x^l \notin \mathbb{Q}\mathbb{V}(c))\}
 \end{aligned}$$

548 A semantics-based dependency graph  $G_{\text{trace}}(c) = (V_{\text{trace}}(c), E_{\text{trace}}(c), W_{\text{trace}}(c), Q_{\text{trace}}(c))$  is well-formed if  
 549 and only if  $\{x^l \mid (x^l, w) \in W_{\text{trace}}(c)\} = V_{\text{trace}}(c)$ .

550 As we discussed before, vertices and query annotations are just read out from the program  $c$ .  
 551 We have an edge in  $E_{\text{trace}}(c)$  if we have a may dependency between two labeled variables in  $c$ . A  
 552 weight function  $w \in W_{\text{trace}}(c)$  is a function that for every starting trace  $\tau_0 \in \mathcal{T}_0(c)$  gives the number  
 553 of times the assignment of the corresponding vertex  $x^l$  is visited. Notice that weight functions are  
 554 total and with range  $\mathbb{N}$ . This means that if a program  $c$  has some non-terminating behavior, the set  
 555  $W_{\text{trace}}(c)$  will be empty. To rule out this situation, we consider as well-formed only graphs which  
 556 have a weight for every vertex. In the rest of the paper we will implicitly consider only well-formed  
 557 semantics-based dependency graphs.

### 559 4.3 Trace-based Adaptivity

560 We can now define the adaptivity of a program formally. This notion is formulated in terms of an  
 561 initial trace, specifying the value of the input variables, as the walk on the graph  $G_{\text{trace}}(c)$ , which  
 562 has the largest number of query requests.

564 **DEFINITION 5 (WALK ON  $G_{\text{trace}}(c)$ ).** Given the semantics-based dependency graph  $G_{\text{trace}}(c) =$   
 565  $(V_{\text{trace}}, E_{\text{trace}}, W_{\text{trace}}, Q_{\text{trace}})$  of a program  $c$ , a walk  $k : \mathcal{T}_0(c) \rightarrow \mathbb{N}$  on  $G_{\text{trace}}(c)$  is a function that  
 566 given as input an initial trace  $\tau_0$  returns a sequence of edges  $(e_1 \dots e_{n-1})$  for which there is a sequence  
 567 of vertices  $(v_1, \dots, v_n)$  such that:

- 568 •  $e_i = (v_i, v_{i+1}) \in E_{\text{trace}}$  for every  $1 \leq i < n$ .
- 569 • every  $v_i \in V_{\text{trace}}$  and  $(v_i, w_i) \in W_{\text{trace}}$ ,  $v_i$  appears in  $(v_1, \dots, v_n)$  at most  $w_i(\tau_0)$  times.

570  $(v_1, \dots, v_n)$  is the vertex sequence of  $k(\tau_0)$  and the length of  $k(\tau_0)$  is the number of vertices in its vertex  
 571 sequence, i.e.,  $|k(\tau_0)| = n$ . We denote by  $\mathcal{WK}(G_{\text{trace}}(c))$  the set of all the walks  $k$  in  $G_{\text{trace}}(c)$ .

572 Because for the adaptivity we are interested in the dependency between queries, we calculate a  
 573 special “length” of a walk, the *query length*, by counting only the vertices corresponding to queries.

575 **DEFINITION 6 (QUERY LENGTH).** Given the semantics-based dependency graph  $G_{\text{trace}}(c)$  of a  
 576 program  $c$ , and a walk  $k \in \mathcal{WK}(G_{\text{trace}}(c))$ , the query length of  $k$  is a function  $\text{len}^q(k) : \mathcal{T}_0(c) \rightarrow \mathbb{N}$   
 577 that given an initial trace  $\tau_0$  returns the number of vertices which correspond to query variables in the  
 578 vertices sequence,  $(v_1, \dots, v_n)$  as follows,

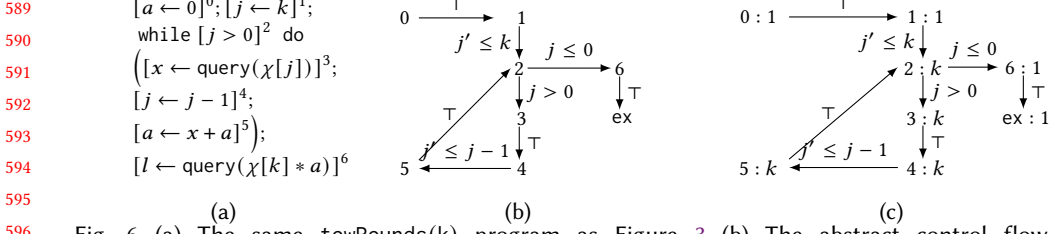
$$579 \text{len}^q(k)(\tau_0) = |\{v \mid v \in (v_1, \dots, v_n) \wedge Q(v) = 1\}|.$$

581 **DEFINITION 7 (ADAPTIVITY OF A PROGRAM).** Given a program  $c$ , its adaptivity  $A(c)$  is function  
 582  $A(c) : \mathcal{T}_0(c) \rightarrow \mathbb{N}$  such that for an initial trace  $\tau_0 \in \mathcal{T}_0(c)$ ,

$$583 A(c)(\tau_0) = \max \{\text{len}^q(k)(\tau_0) \mid k \in \mathcal{WK}(G_{\text{trace}}(c))\}$$

## 585 5 THE ADAPTIVITY ANALYSIS ALGORITHM - ADAPTFUN

586 In this section, we present our program analysis AdaptFun for computing an upper bound on the  
 587 *adaptivity* of a given program  $c$ . The high level idea behind AdaptFun is to first build an *estimated*



596 Fig. 6. (a) The same `twoRounds(k)` program as Figure 3 (b) The abstract control flow graph, `absG(twoRounds(k))` (c) `absG(twoRounds(k))` with the reachability bound.

598  
 599  
 600 *dependency graph*  $G_{\text{est}}(c)$  of a program  $c$  (Section 5.1) which overapproximates the semantics-based  
 601 *dependency graph* in two dimensions: it overapproximates the dependencies between assigned  
 602 variables (Section 5.1.2), and, it overapproximates the weights (Section 5.1.3). Then, `AdaptFun` uses  
 603 a custom algorithm to estimate the longest walk on this graph, providing in this way an upper  
 604 bound on the adaptivity of the program.

605 Given a program  $c$ , the set of vertices  $V_{\text{est}}(c)$  and query annotations  $Q_{\text{est}}(c)$  of the *estimated*  
 606 *dependency graph* can be computed by simply scanning the program  $c$ . These set can be computed  
 607 precisely and correspond to the same sets in the semantics-based *dependency graph*. This means  
 608 that  $G_{\text{est}}(c)$  has the same underlying vertex structure as the semantics-based graph  $G_{\text{trace}}(c)$ . The  
 609 differences will be in the sets of edges and weights.

## 611 5.1 Weight and Edge Estimation

612 The set of edges  $E_{\text{est}}(c)$  and the set of weights  $W_{\text{est}}(c)$  of the estimated *dependency graph* are  
 613 estimated through an analysis combining control flow, data flow, and loop bound analysis. These  
 614 analyses are naturally described over an *Abstract Transition Graph* of the input program, which we  
 615 describe next.

616  
 617 **5.1.1 Abstract Transition Graph.** We say that we have a *transition* from a program point  $l$  to a  
 618 program point  $l'$  if and only if the command with label  $l'$  can execute right after the execution of  
 619 the command with label  $l$ . The *Abstract Transition Graph*  $\text{absG}(c)$  of a program  $c$  is a graph with  
 620 the set of labels of program points in  $c$  (including a label `ex` for the exit point) as the set of vertices  
 621  $\text{absV}(c)$ , and with the set of transitions in  $c$  as the set of edges  $\text{absE}(c)$ . Each edge of the graph is  
 622 annotated with either the symbol  $\top$ , a boolean expression or a *difference constraint* [33].

623 A *difference constraint* is an inequality of the form  $x' \leq y + v$  or  $x' \leq v$  where  $x, y$  are variables  
 624 and  $v \in \mathcal{SC}$  is a symbolic constant: either a natural number, the symbol  $\infty$ , an input variable or a  
 625 symbol  $Q_m$  representing a query request. We denote by  $\mathcal{DC}$  the set of *difference constraints*.

626 A *difference constraint* on an edge,  $l \xrightarrow{x' \leq y+v} l'$  or  $l \xrightarrow{x' \leq v} l'$ , denotes that after executing the  
 627 command at location  $l$  the value of the variable  $x$  is at most the value of the expression  $y + v$  resp.  $v$   
 628 before the execution of the command  $l'$ . A boolean value  $b$  on an edge,  $l \xrightarrow{b} l'$ , denotes that after  
 629 evaluating the guard of an if or a while command with label  $l, b$  holds and the next command to be  
 630 executed is the one with label  $l'$ . A  $\top$  symbol on an edge,  $l \xrightarrow{\top} l'$  denotes that the command with  
 631 label  $l$  is a skip, and the commands that do not interfere with any loop counter variable.

632 We compute *difference constraints* and the other annotation via a simple program abstraction  
 633 method adopted from [33], described in details in the supplementary material.

634 **Example.** We show in Fig. 6(b) the abstract control flow graph, `absG(twoRounds(k))` of the  
 635 `twoRounds(k)` program we gave in Fig. 3(a) and which we also report in Fig. 6(a).

5.1.2 *Edge Estimation.* The set of edges  $E_{\text{est}}(c)$  is estimated through a combined data and control flow analysis with three components.

**Reaching definition analysis:** The first component is a reaching definition analysis computing for each label  $l$  in the graph  $\text{absG}(c)$  the set of labeled variables that may reach  $l$  as follows.

(1). For each label  $l$ , the analysis generates two initial sets of labeled variables,  $in$  and  $out$ , containing all the labeled variables  $x^l$  that are newly generated but not yet reassigned before and after executing the command  $l$ .

(2). The analysis iterates over  $\text{absG}(c)$ , and updates  $in(l)$  and  $out(l)$  until they are stable. The final  $in(l)$  is the set of reaching definitions  $\text{RD}(l, c)$  for  $l$ .

**Feasible data-flow analysis:** The second component is a **feasible data-flow analysis** computing for every pair  $x^i, y^j \in \text{LV}(c)$  whether there is a flow from  $x^i$  to  $y^j$ . This analysis is based on a relation  $\text{flowsTo}(x^i, y^j, c)$  built over the sets  $\text{RD}(l, c)$  for every location  $l$ . This relation is defined as:

DEFINITION 8 (FEASIBLE DATA-FLOW). *Given a program  $c$  and two labeled variables  $x^i, y^j$  in this program,  $\text{flowsTo}(x^i, y^j, c)$  is*

$$\begin{aligned} \text{flowsTo}(x^i, y^j, [y \leftarrow e]^l) &\triangleq \{(x^i, y^j) \mid x \in \text{FV}(e) \wedge x^i \in \text{RD}(l, [y \leftarrow e]^l)\} \\ \text{flowsTo}(x^i, y^j, [y \leftarrow \text{query}(\psi)]^l) &\triangleq \{(x^i, y^j) \mid x \in \text{FV}(\psi) \wedge x^i \in \text{RD}(l, [y \leftarrow \text{query}(\psi)]^l)\} \\ \text{flowsTo}(x^i, y^j, [\text{skip}]^l) &= \emptyset \\ \text{flowsTo}(x^i, y^j, \text{if } ([b]^l, c_1, c_2)) &\triangleq \text{flowsTo}(x^i, y^j, c_1) \vee \text{flowsTo}(x^i, y^j, c_2) \\ &\vee \{(x^i, y^j) \mid x \in \text{FV}(b) \wedge x^i \in \text{RD}(l, \text{if } ([b]^l, c_1, c_2)) \wedge y^j \in \text{LV}(c_1)\} \\ &\vee \{(x^i, y^j) \mid x \in \text{FV}(b) \wedge x^i \in \text{RD}(l, \text{if } ([b]^l, c_1, c_2)) \wedge y^j \in \text{LV}(c_2)\} \\ \text{flowsTo}(x^i, y^j, \text{while } [b]^l \text{ do } c_w) &\triangleq \text{flowsTo}(x^i, y^j, c_w) \vee \\ &\{(x^i, y^j) \mid x \in \text{FV}(b) \wedge x^i \in \text{RD}(l, \text{while } [b]^l \text{ do } c_w) \wedge y^j \in \text{LV}(c_w)\} \\ \text{flowsTo}(x^i, y^j, c_1; c_2) &\triangleq \text{flowsTo}(x^i, y^j, c_1) \vee \text{flowsTo}(x^i, y^j, c_2) \end{aligned}$$

This relation gives us an overapproximation of the *variable may-dependency* relation for direct dependencies (dependencies that do not go through other variables).

**Edge Construction:** The third component constructs an edge by computing a transitive closure (through other variables) of the  $\text{flowsTo}$  relation. There is a directed edge from  $x^i$  to  $y^j$  if and only if there is chain of variables in the  $\text{flowsTo}$  relation between  $x^i$  and  $y^j$ . This is defined as follows:

$$E_{\text{est}}(c) \triangleq \{(y^j, x^i) \mid y^j, x^i \in V_{\text{est}}(c) \wedge \exists n, z_1^r, \dots, z_n^r \in \text{LV}(c). \\ n \geq 0 \wedge \text{flowsTo}(x^i, z_1^r, c) \wedge \dots \wedge \text{flowsTo}(z_n^r, y^j, c)\}$$

We prove that the set  $E_{\text{est}}(c)$  soundly approximates the set  $G_{\text{trace}}(c)$ .

LEMMA 5.1 (MAPPING FROM EGDES OF  $G_{\text{trace}}$  TO  $G_{\text{est}}$ ). *For every program  $c$  we have:*

$$\forall e = (v_1, v_2) \in E_{\text{trace}}(c). \exists e' \in E_{\text{est}}(c). e' = (v_1, v_2)$$

**Example.** Consider Fig. 3(c), the edge  $l^6 \rightarrow a^5$  is built by the  $\text{flowsTo}(l^6, a^5, c)$  relation because  $a$  is used directly in the query expression  $\chi[k] * a$  and we also have  $a^5 \in \text{RD}(6, \text{twoRounds}(k))$  from the reaching definition analysis. The edge  $x^3 \rightarrow j^5$  represents the control flow from  $j^5$  to  $x^3$ , which is soundly approximated by our  $\text{flowsTo}$  relation. The edge  $l^6 \rightarrow x^3$  is produced by the transitivity of  $\text{flowsTo}(l^6, a^5, c)$  and  $\text{flowsTo}(a^5, x^3, c)$ .

5.1.3 *Weight Estimation.* The set  $W_{\text{est}}(c)$  of weights for the estimated dependency graph is estimated from the Abstract Transition Graph  $\text{absG}(c)$  of  $c$  using *reachability-bound analysis* [22]. Specifically, we estimate as the weight of a node with label  $l$  a symbolic upper bounds on the execution times of the command with label  $l$  obtained by reachability-bound analysis. These symbolic upper bounds are expressions with the input variables as free variables, hence they correspond to the weight functions in the semantics-based dependency graphs.

Our reachability-bound algorithm adapts to our setting ideas from previous work [32, 33, 39]. Specifically, it provides an upper bound on the number of times every command can be executed by using three steps.



- 687 (1) This step assigns to each edge  $l \xrightarrow{dc} l' \in \text{absE}(c)$  a *local bound* as follows. We look at the  
688 strongly connected components of  $\text{absG}(c)$ . If the edge does not belong to any strongly  
689 connected components, then the local bound is 1, representing the fact that the edge is not  
690 in a loop and so it gets executed at most once. If the edge belongs to a strongly connected  
691 component and one of the variables  $x$  in  $dc$  decreases, then the local bound is  $x$ . Otherwise, if  
692 the edge belongs to a strongly connected component and there is a variable  $y$  that decreases  
693 in the difference constraint of some other edge, and if by removing this other edge, the  
694 original edge does not belong anymore to the strongly connected components of  $\text{absG}(c)$ ,  
695 then the local bound is  $y$ . Otherwise, the local bound is  $\infty$ . Notice that the output is either a  
696 symbolic constant in  $\mathcal{SC}$  or a variable that is not an input variable.
- 697 (2) This step aims at determining the *reachability-bound*  $\text{TB}(e, c)$  of every edge  $e \in E_{\text{est}}(c)$ .  
698 Every bound is a symbolic expression built out of symbols in  $\mathcal{SC}$  and the operations  $+$ ,  $*$ ,  $\max$ .  
699 For every edge, if the local bound of this edge computed at the previous step is a symbol in  
700  $\mathcal{SC}$  then this is already the reachability-bound. If instead the local bound of the edge is a  
701 variable  $y$  which is not an input variable, this step will eliminate it and replace it with a  
702 symbolic expression. In order to do this, this step will compute two quantities: first, it will  
703 recursively sum the reachability-bounds of all the edges whose difference constraint may  
704 increment the variable  $y$ , plus the corresponding increment; second, it will recursively sum  
705 the reachability-bounds of all the edges whose difference constraint may reset the variable  
706  $y$  to a (symbolic) expression that doesn't depend on it, multiplied by the maximal value of  
707 this symbolic expression. The sum of these two quantities provides the symbolic expression  
708 that is an upper bound on the number of times the edge can be reached. To compute these  
709 two quantities we use two mutually recursive procedures.

710 Using the reachability-bound  $\text{TB}(e, c)$  for every edge  $e = (l, dc, l')$  we can provide a bound on the  
711 visiting times of each vertex  $x^l \in \text{absV}(c)$ . Formally:  $w = \sum \{\text{TB}(e, c) \mid e = (l, \_, \_)\}$ . Notice that  $w$   
712 is a symbolic arithmetic expression over symbols in  $\mathcal{SC}$ . In particular, it may contain the input  
713 variables and so it may effectively be used as a function of the input - and capture loop bounds in  
714 terms of these inputs.

715  
716 **THEOREM 5.1 (SOUNDNESS OF THE REACHABILITY BOUNDS ESTIMATION).** *Let  $c$  be a program and*  
717  *$W_{\text{est}}(c)$  be its estimated weight set. Then, for each  $(x^l, w) \in W_{\text{est}}(c)$ ,  $\tau_0 \in \mathcal{T}_0(c)$ ,  $\tau \in \mathcal{T}$ ,  $v \in \mathbb{N}$  we have:*

$$718 \text{ if } \langle c, \tau_0 \rangle \rightarrow^* \langle \text{skip}, \tau_0 + \tau \rangle \wedge \langle \tau_0, w \rangle \Downarrow_e v \wedge \text{ then } \text{cnt}(\tau, l) \leq v$$

719  
720 Notice that in this theorem, the evaluation  $\langle \tau_0, w \rangle \Downarrow_e v$  is needed in order to obtain a concrete  
721 value  $v$  from the symbolic weight  $w$  by specifying a value for the input variables through  $\tau_0$ .

722 **Example.** Consider again Fig. 3(c), the estimated weight for  $a^5$  is  $k$ , and this is a sound estimation.  
723 For an arbitrary  $\tau_0 \in \mathcal{T}_0(c)$ , we know that  $\langle \tau_0, k \rangle \Downarrow_e \rho(\tau_0)k$  and by the weight  $w_k$  for the vertex  $a^5$   
724 (as in Figure 3(b)) we know  $w_k(\tau_0) = \rho(\tau_0)k$ .

## 725 5.2 Adaptivity Upper Bound Computation

726  
727 We estimate the adaptivity upper bound,  $A_{\text{est}}(c)$  for a program  $c$  as the maximum query length  
728 over all finite walks in its *estimated dependency graph*,  $G_{\text{est}}(c)$ .

729 **DEFINITION 9 (ESTIMATED ADAPTIVITY).** *Given a program  $c$  and its estimated dependency graph*  
730  *$G_{\text{est}}(c)$  the estimated adaptivity for  $c$  is*

$$731 A_{\text{est}}(c) \triangleq \max\{\text{len}^q(k) \mid k \in \mathcal{WK}(G_{\text{est}}(c))\}.$$

732  
733 Notice that different from a walk on  $G_{\text{trace}}(c)$ , a walk  $k \in \mathcal{WK}(G_{\text{est}}(c))$  on the graph  $G_{\text{est}}(c)$   
734 does not rely on an initial trace. This because, similarly to what we did for the weights in  $W_{\text{est}}(c)$

in the previous section, we use symbolic expressions over input variables. Similarly, the adaptivity bound  $A_{\text{est}}(c)$  will also be a symbolic arithmetic expression over the input variables. With this symbolic expression we can prove the upper bound sound with respect to any initial trace.

**THEOREM 5.2 (SOUNDNESS OF  $A_{\text{est}}(c)$ ).** *For every program  $c$ , its estimated adaptivity is a sound upper bound of its adaptivity.*

$$\forall \tau_0 \in \mathcal{T}_0(c), v \in \mathbb{N}^\infty . \langle A_{\text{est}}(c), \tau_0 \rangle \Downarrow_e v \implies A(c)(\tau_0) \leq v.$$

Symbolic expressions as used in the weight are great to express symbolic bounds but make the computation of a maximal walk harder. Specifically, one has to face two challenges. The first is non-termination. A naive traversing strategy leads to non-termination because the weight of each vertex in  $G_{\text{est}}(c)$  is a symbolic expression containing input variables. We could try to use a depth first search strategy using the longest weighted path to approximate the longest finite walk with the weight as its visiting time. However, these approach would face the second challenge: approximation. It would consistently and considerably over-approximate the adaptivity.

To address these two challenges we design an algorithm AdaptBD combining Depth First Search and Breadth First Search. The idea of this algorithm is to reduce the task of computing the longest walk to the task of computing local versions of the adaptivity on the maximal strongly connected components (SCC) of the graph  $G_{\text{est}}(c)$  and then compose them into the program adaptivity. The algorithm uses another algorithm AdaptBD<sub>SCC</sub> recursively, in order to find the longest walk for a strong connected component (SCC) of  $G_{\text{est}}(c)$ . The pseudocode of AdaptBD<sub>SCC</sub> is given as Algorithm 1

---

### Algorithm 1 Adaptivity Bound Algorithm on An SCC (AdaptBD<sub>SCC</sub>( $c, \text{SCC}_i$ ))

---

**Require:** The program  $c$ , A strong connected component of  $G_{\text{est}}(c)$ :  $\text{SCC}_i = (V_i, E_i, W_i, Q_i)$

```

1: init.  $r_{\text{scc}}$ :  $\mathcal{A}_{\text{in}}$  List with initial value 0.
2: init.  $\text{visited}$ :  $\{0, 1\}$  List with initial value 0;  $r$ :  $\mathcal{A}_{\text{in}}$  List, initial value  $\infty$ ;
    $\text{flowcapacity}$ :  $\mathcal{A}_{\text{in}}$  List, initial value  $\infty$ ;  $\text{querynum}$ : INT List, initial value  $Q_i(v)$ .
3: if  $|V_i| = 1$  and  $|E_i| = 0$ : return  $Q(v)$ 
4: def  $\text{dfs}(G, s, \text{visited})$ :
5:   for every vertex  $v$  connected by a directed edge from  $s$ :
6:     if  $\text{visited}[v] = \text{false}$ :
7:        $\text{flowcapacity}[v] = \min(W_i(v), \text{flowcapacity}[s])$ ;    $\text{querynum}[v] = \text{querynum}[s] + Q_i(v)$ ;
8:        $r[v] = \max(r[v], \text{flowcapacity}[v] \times \text{querynum}[v])$ ;
9:        $\text{visited}[v] = 1$ ;    $\text{dfs}(G, v, \text{visited})$ ;
10:    else:  $\#$ {There is a cycle finished}
11:       $r[v] = \max(r[v], r[s] + \min(W_i(v), \text{flowcapacity}[s]) * (\text{querynum}[s] + Q_i(v)))$ ;
12:    return  $r[c]$ 
13: for every vertex  $v$  in  $V_i$ :
14:   initialize the  $\text{visited}$ ,  $r$ ,  $\text{flowcapacity}$ ,  $\text{querynum}$  with the same value at line:2.
15:    $r_{\text{scc}} = \max(r_{\text{scc}}, \text{dfs}(\text{SCC}_i, v, \text{visited}))$ ;
16: return  $r_{\text{scc}}$ 

```

---

This algorithm takes as input the program  $c$  and a  $\text{SCC}_i$  of  $G_{\text{est}}(c)$ , and outputs an adaptivity bound for  $\text{SCC}_i$ . If  $\text{SCC}_i$  contains only one vertex,  $x^l$  without any edge, AdaptBD<sub>SCC</sub> returns the query annotation of  $x^l$  as the adaptivity. If  $\text{SCC}_i$  contains at least one edge, AdaptBD<sub>SCC</sub>:

- (1) first collects all the paths in  $\text{SCC}_i$ ;
- (2) it then calculates the adaptivity of every path by the method  $\text{dfs}$ ;
- (3) in the end, it outputs the maximal adaptivity among all paths as the adaptivity of  $\text{SCC}_i$ .

By the property of SCC, the paths collected in step 1 are all simple cycles with the same starting and ending vertex. Step 2 is the key step. It recursively computes the adaptivity upper bound on

785 the fly of paths collected through a DFS procedure `dfs` (lines: 4-13). This procedure guarantees that  
 786 the visiting times of each vertex is upper bounded by its weight, and addresses the approximation  
 787 challenge, via two special lists parameters `flowcapacity` and `querynum` (lines:7-11).

788 `flowcapacity` is a list of symbolic expressions  $\mathcal{A}_{in}$  which tracks the minimum weight when  
 789 searching a path, and updates the weight when the path reaches a certain vertex. `querynum` is a list  
 790 of integer initialized with the value of the query annotation  $Q_i(v)$  for every vertex. It tracks the  
 791 total number of vertices with query annotation 1 along the path. The operation at line: 8 and line:  
 792 11 implements the operation `flowcapacity[v] × querynum[v]`. Because `flowcapacity[v]` is the  
 793 minimum weight over this path, this guarantees that every vertex on the estimated walk is allowed  
 794 to be visited at most `flowcapacity[v]`, and this walk is a valid finite walk. Then `querynum[v]`  
 795 guarantees that `flowcapacity[v] × querynum[v]` computes an accurate query length because  
 796 `querynum[v]` is only the number of the vertices with query annotation 1, giving a tight bound  
 797 without losing the soundness.

798 THEOREM 5.3 (SOUNDNESS OF AdaptBD). *For every program  $c$ , we have*

$$799 \text{AdaptBD}(G_{\text{est}}(c)) \geq A_{\text{est}}(c).$$

## 800 6 EXAMPLES

801 We illustrate here how our analysis work on two different examples. Our first example, Algorithm  
 802 `multipleRounds` in Fig. 7(a), is a simplified form of the *monitor argument* by Rogers et al. [31]. The  
 803 input  $k$  is the number of iterations. It uses a list  $I$  to track queries. Specifically at each iteration  
 804 it updates  $I$  by using the result of a query which relies on  $I$ : `query( $\chi[I]$ )`. After  $k$  iterations, the  
 805 algorithm returns the columns of the hidden database  $D$  which are not contained in the tracking list  
 806  $I$ . The functions `updnscore( $p, a$ )`, `updcscore( $p, a$ )`, `update( $I, ns, cs$ )` simplify the computations of  
 807 updating  $ns$ ,  $cs$  and  $I$ . They depends on the result of the query but they do not perform queries them-  
 808 selves. Different from the code in the example `twoRounds( $k$ )`, the query request, `[ $a \leftarrow \text{query}(I)$ ]6,  
 809 in each loop iteration depends on the tracking list  $I$ , which in turn depends on all the queries  
 810 from previous iterations. In this sense, all these  $k$  queries are fully adaptively chosen, and so the  
 811 adaptivity is  $k$ . The estimated dependency graph  $G_{\text{est}}(\text{multipleRounds}(k))$  is presented in Fig. 7(b)  
 812 and we omitted the semantics-based dependency graph  $G_{\text{trace}}(\text{multipleRounds}(k))$  because it has  
 813 the same topology and only differ in weights. Our program analysis AdaptFun provides a tight  
 814 upper bound for this example using AdaptBD(multipleRounds( $k$ )). It first finds a path on the  
 815 graph  $G_{\text{est}}(\text{multipleRounds}(k))$   $a^6 : \frac{k}{1} \rightarrow I^9 : \frac{k}{0} \rightarrow ns^7 : \frac{k}{0}$  with three weighted vertices. Then  
 816 AdaptBD algorithm transforms this path into a walk  $a^6 : \frac{k}{1} \rightarrow I^9 : \frac{k}{0} \rightarrow ns^7 : \frac{k}{0} \rightarrow a^6 : \frac{k}{1} \dots$ , where  
 817  $a^6, I^9, ns^7$  are all visited  $k$  times respectively. So  $A_{\text{est}}(\text{multipleRounds}(k)) = k$ . We know for any  
 818 initial trace  $\tau_0, \langle \tau_0, k \rangle \Downarrow_e \rho(\tau_0)k$ , i.e.,  $A(\text{multipleRounds}(k))(\tau_0) \leq \rho(\tau_0)k$  for any  $\tau_0$ , and so what  
 819 we have produced is a tight and sound bound.`

820 We want to show an example where our definition of adaptivity (Def. 7) itself over-approximates  
 821 the intuitive adaptivity. Our second example `multiRoundsS( $k$ )` in Fig. 8(a) demonstrates this over-  
 822 approximation. It is a variant of the multiple rounds strategy with input  $k$ . In each iteration, the  
 823 query `query( $\chi[y] + p$ )` in line 7 is based on previous query results stored in  $p$  and  $y$ . Different  
 824 from Ex. 6, only the query answer from the  $(k - 2)^{\text{th}}$  iteration is used in the query request  
 825 `[ $p \leftarrow \text{query}(\chi[y] + p)$ ]7. This is because the execution will reset the value of  $p$  to 0 in all the  
 826 other iterations after this query request (line 10). In this way, all the query answers stored in  $p$  are  
 827 erased and are not used in the query request at the next iteration, except the one at the  $(k - 2)^{\text{th}}$   
 828 iteration. So multiRoundsS( $k$ )'s adaptivity rounds is only 2. However, our Def. 7 fails to realize  
 829 that there is only dependency relation between  $p^7$  and  $p^7$  in one iteration, but not in others. As the  
 830  $G_{\text{trace}}(\text{multiRoundsS}(k))$  in Fig. 8(b) shows, there is an edge from  $p^7$  to itself representing Variable  
 831  
 832  
 833`

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

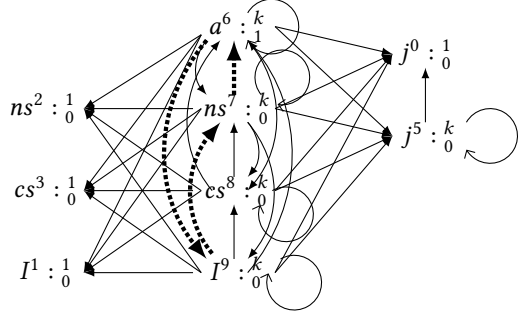
882

```

multipleRounds(k)  $\triangleq$ 
[j  $\leftarrow$  k]0; [I  $\leftarrow$  []]1;
[ns  $\leftarrow$  0]2; [cs  $\leftarrow$  0]3;
while [j > 0]4 do
([j  $\leftarrow$  j - 1]5; [a  $\leftarrow$  query(I)]6;
[ns  $\leftarrow$  updnscore(ns, a)]7;
[cs  $\leftarrow$  updcscore(cs, a)]8;
[I  $\leftarrow$  updI(I, ns, cs)]9)

```

(a)



(b)

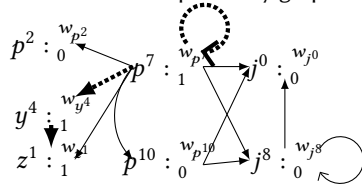
Fig. 7. (a) The simplified multiple rounds example (b) The estimated dependency graph from AdaptFun

```

multiRoundsS(k)
[j  $\leftarrow$  0]0; [z  $\leftarrow$  query(0)]1; [p  $\leftarrow$  0]2;
if ([k = 0]3, [y  $\leftarrow$  query(z)]4, [skip]5);
while [j  $\neq$  k]6 do
([p  $\leftarrow$  query( $\chi$ [y] + p)]7; [j  $\leftarrow$  j + 1]8
if ([j  $\neq$  k - 2]9, [p  $\leftarrow$  0]10, [skip]10);

```

(a)



(b)

Fig. 8. (a) The multi rounds single example (b) The semantics-based dependency graph.

*May-Dependency* of  $p^7$  on itself, and  $p^7$ 's visiting times,  $w(\tau_0)$ .  $w(\tau_0)$  counts the execution times of command  $[p \leftarrow \text{query}(\chi[y] + p)]^7$ . It equals to the loop iteration numbers, i.e.,  $k$ 's initial value. Then, as the dotted arrows, longest walk is  $p^7 \rightarrow \dots \rightarrow p^7 \rightarrow y^4 \rightarrow z^1$  computes  $2 + w_{p^7}(\tau_0)$ , instead of 2. It is worth to stress that our algorithm still compute an accurate bound w.r.t this definition, even if the definition itself is over-approximating. Indeed, the AdaptFun give us adaptivity  $2 + k$ .

## 7 IMPLEMENTATION

We implemented AdaptFun as a tool which takes a labeled program as input and outputs the upper bound on the program adaptivity and the total number of queries that the program runs. This implementation consists of a module written in OCaml for the generation of the estimated graph  $G_{\text{est}}$ , and a module written in Python for the weight estimation algorithm (Section 5.1.3) and the algorithm AdaptBD (Section 5.2). The OCaml program takes the labeled program as input and outputs a version of the graph  $G_{\text{est}}$  (without weights) and the abstract transition graph  $\text{absG}$  for the program. These two objects are then fed into the python program which computes the weights, and outputs the adaptivity bound and the query number. We evaluated this implementation on 25 examples with performances summarized in Tab. 1. The 1<sup>st</sup> column is the example name. For each example  $c$ , the 2<sup>nd</sup> column is its adaptivity rounds, AdaptFun outputs are in the 3<sup>rd</sup> and 4<sup>th</sup> columns. They are the adaptivity upper bound and  $c$ 's total query requests #. The last 4 columns are AdaptFun's performance w.r.t. the program lines. We track the running time of the OCaml code for parsing the program and generating the  $G_{\text{est}}(c)$ , and the running times of the weight analysis and the AdaptBD( $c$ ) in Python. We implemented two weight estimation methods. The first one (referred as I in Tab.1) is the one we presented formally in Section 5.1. Unfortunately, this method is accurate but slow, it doesn't performs well with big program. The second one (referred as II) is a relaxation of the first one. It is more efficient but it over-approximate complicated loops. Based on the two implementations, our AdaptFun produces two bounds on the adaptivity, corresponding to the left and right side (I | II) in the 3<sup>rd</sup>, 4<sup>th</sup> and 6<sup>th</sup> columns<sup>9</sup>. The first 5 programs are adapted

<sup>9</sup>When the method II produces the same results as I, we omit them and use the symbol -.

from classical data analysis algorithms. AdaptFun computes tight adaptivity bound for the first 3. For the forth program `multiRoundsO(k)`, AdaptFun over-approximates the adaptivity as  $1 + 2 * k$  because of its path-insensitivity. The fifth program is the one in Example 6, where AdaptFun outputs the tight bound but we give a loose definition for its actual adaptivity. The programs from Tab. 1 line:6-17 all have small size but complex structures, to test the programs under different situations including data, control dependency, the multiple paths nested loop with related counters, etc. Both implementations compute tight bounds for examples in line:6-14 and over-approximate the adaptivities for 15<sup>th</sup> and 16<sup>th</sup> due to path-insensitivity. For the 17<sup>th</sup> one, implementation I gives tight bound while II gives loose bound, so we keep both implementations. The last six programs are big but simple, to test the performance limitation. From the evaluation results, the performance bottleneck is the weight estimation algorithm. The implementation I is unable to evaluate them in a reasonable time period, denoted by \* on the left side. While the implementation II computes the *adaptivity* for them effectively on the right side.

Table 1. Experimental results of AdaptFun implementation

Program c	adaptivity	AdaptFun		lines	performance		
		AdaptBD(c) (I   II)	query# (I   II)		running time (second)		
					Ocaml	Weight	AdaptBD
twoRounds(k)	2	2 -	k + 1 -	8	0.0005	0.0017   0.0002	0.0003
multiRounds(k)	k	k  max(1, k)	k -	10	0.0012	0.0017   0.0002	0.0002
LRGD(k, r)	k	k  max(1, k)	2k -	10	0.0015	0.0072   0.0002	0.0002
mROdd(k)	1 + k	2 + max(1, 2k)  -	1 + 3k -	10	0.0015	0.0061   0.0002	0.0002
mRSingle(k)	2	1 + max(1, k)  -	1 + k 1 + k	9	0.0011	0.0075   0.0002	0.0002
iFCD()	3	3 4	3 4	5	0.0005	0.0003   0.0001	0.0001
while(k)	1 + k/2	1 + max(1, k/2)  -	1 + k/2 -	7	0.0021	0.0015   0.0001	0.0001
whileRV(k)	1 + 2k	1 + 2k 1 + max(1, 2k)	2 + 3k -	9	0.0016	0.0056   0.0002	0.0001
whileVCD(k)	1 + 2Q <sub>m</sub>	Q <sub>m</sub> + max(1, 2Q <sub>m</sub> )  -	2 + 2Q <sub>m</sub>  -	6	0.0016	0.0007   0.0002	0.0001
whileMPVCD(k)	2 + Q <sub>m</sub>	2 + Q <sub>m</sub>  -	2 + 2Q <sub>m</sub>  -	9	0.0017	0.0043   0.0002	0.0001
nestWhileVD(k)	2 + k <sup>2</sup>	3 + k <sup>2</sup>  -	1 + k + k <sup>2</sup>  -	10	0.0018	0.0126   0.0002	0.0001
nestWhileRV(k)	1 + k + k <sup>2</sup>	2 + k + k <sup>2</sup>  -	2 + k + k <sup>2</sup>  -	10	0.0017	0.0186   0.0002	0.0001
nestWhileMV(k)	1 + 2k	1 + max(1, 2k)  -	1 + k + k <sup>2</sup>  -	10	0.0016	0.0071   0.0002	0.0001
nestWhileMPRV(k)	1 + k + k <sup>2</sup>	3 + k + k <sup>2</sup>  -	2 + 2k + k <sup>2</sup>  -	10	0.019	0.0999   0.0002	0.0002
whileM(k)	1 + k	2 + max(1, 2k)  -	1 + 3k -	9	0.0017	0.0062   0.0002	0.0001
whileM2(k)	1 + k	2 + k -	1 + 3k -	9	0.0017	0.0062   0.0002	0.0001
nestWhileRC(k)	1 + 3k	1 + 3k 2 + 3k + k <sup>2</sup>	1 + 3k 1 + k + k <sup>2</sup>	11	0.019	0.2669   0.0002	0.0007
mRComplete(k, N)	k	k -	k -	27	0.0026	85.9017   0.0003	0.0004
mRCompose(k)	2k	2k -	2k -	46	0.0036	5104   0.0003	0.0013
seqCompose(k)	12	12 -	326 -	502	0.0426	1.2743   0.0003	0.0223
tRCompose(k)	2	* 2	* 1 + 5k + 2k <sup>2</sup>	42	0.0026	*   0.0003	0.0005
jumboS(k)	max(20, 8 + k <sup>2</sup> )	*  max(20, 6 + k + k <sup>2</sup> )	* 44 + k + k <sup>2</sup>	71	0.0035	*   0.0003	0.0085
jumbo(k)	max(20, 10 + k + k <sup>2</sup> )	*  max(20, 12 + k + k <sup>2</sup> )	* 286 + 26k + 10k <sup>2</sup>	502	0.0691	*   0.0009	0.018
big(k)	22 + k + k * k	* 28 + k + k <sup>2</sup>	* 121 + 11k + 4k <sup>2</sup>	214	0.0175	*   0.0004	0.002

## 8 RELATED WORK

*Dependency Definitions and Analysis.* There is a vast literature on dependency definitions and dependency analysis. We consider a semantics definition of dependencies which consider (intraprocedural) data and control dependency [8, 11, 30]. Our definition is inspired by classical works on traditional dependency analysis [12] and noninterference [20]. Formally, our definition is similar to the one by Cousot [10], which also identifies dependencies by considering differences in two execution traces. However, Cousot excludes some forms of implicit dependencies, e.g. the ones generated by empty observations, which instead we consider. Common tools to study dependencies are dependency graphs [17]. We use here a semantics-based approach to dependency graph similar, for example, to works by Austin and Sohi [5], Hammer et al. [23] and [24]. Our approach shares some similarities with the use of dependency graphs in works analyzing dependencies between events, e.g. in event programming. Memon [29] uses an event-flow graph, representing all the possible event interactions, where vertices are GUI event edges represent pairs of events that can be performed immediately one after the other. In a similar way, we use edges to track the may-dependence between variables looking at all the possible interactions. Arlt et al. [4] use a weighted edges indicating a dependency between two events, e.g. one event possibly reads data

written by the other event, with the weight showing the intensity of the dependency (the quantity of data involved). We also use weights but on vertices and with different meaning, they are functions describing the number of times the vertices can be visited given an initial state. Differently from all these previous works, we use a dependency graph with quantitative information needed to identify the length of chain of dependencies. Our weight estimation is inspired by works in complexity analysis and WCET. Specifically, it is inspired by works on reachability-bound analysis using program abstraction and invariant inference [21, 22, 34] and work on invariant inference through cost equations and ranking functions [2, 3, 9, 18].

*Generalization in Adaptive Data Analysis.* Starting from the works by Dwork et al. [15] and Hardt and Ullman [25], several works have designed methods that ensure generalization for adaptive data analyses [7, 13, 14, 16, 27, 31, 35, 36]. Several of these works drew inspiration from differential privacy, a notion of formal data privacy. By limiting the influence that an individual can have on the result of a data analysis, even in adaptive settings, differential privacy can also be used to limit the influence that a specific data sample can have on the statistical validity of a data analysis. This connection is actually in two directions, as discussed for example by Yeom et al. [38]. Considering this connection between generalization and privacy, it is not surprising that some of the works on programming language techniques for privacy-preserving data analysis are related to our work. Adaptive Fuzz [37] is a programming framework for differential privacy that is designed around the concept of adaptivity. This framework is based on a typed functional language that distinguish between several forms of adaptive and non-adaptive composition theorem with the goal of achieving better upper bounds on the privacy cost. Adaptive Fuzz uses a type system and some partial evaluation to guarantee that the programs respect differential privacy. However, it does not include any technique to bound the number of rounds of adaptivity. Lobo-Vesga et al. [28] propose a language for differential privacy where one can reason about the accuracy of programs in terms of confidence intervals on the error that the use of differential privacy can generate. These are akin to bounds on the generalization error. This language is based on a static analysis which however cannot handle adaptivity. The way we formalize the access to the data mediated by a mechanism is a reminiscence of how the interaction with an oracle is modeled in the verification of security properties. As an example, the recent works by Barbosa et al. [6] and Aguirre et al. [1] use different techniques to track the number of accesses to an oracle. However, reasoning about the number of accesses is easier than estimating the adaptivity of these calls, as we do instead here.

## 9 CONCLUSION AND FUTURE WORKS

We presented AdaptFun, a program analysis useful to provide an upper bound on the adaptivity of a data analysis, as well as on the total number of queries asked. This estimation can help data analysts to control the generalization errors of their analyses by choosing different algorithmic techniques based on the adaptivity. Besides, a key contribution of our works is the formalization of the notion of adaptivity for adaptive data analysis. We showed the applicability of our approach by implementing and experimentally evaluating our program analysis.

As future work, we plan to investigate the potential integration of AdaptFun in an adaptive data analysis framework like Guess and check by Rogers et al. [31]. As we discussed, this framework is designed to support adaptive data analyses with limited generalization error. As our experiments show, this framework could benefit from the information provided by AdaptFun to provide more precise estimate and improved confidence intervals. Another direction we will explore is to make the upper bounds provided by AdaptFun more precise by integrating our algorithm with a path-sensitive approach.



## REFERENCES

- 981  
982  
983 [1] Alejandro Aguirre, Gilles Barthe, Marco Gaboardi, Deepak Garg, Shin-ya Katsumata, and Tetsuya Sato. 2021. Higher-  
984 order probabilistic adversarial computations: Categorical semantics and program logics. *CoRR* abs/2107.01155 (2021).  
985 arXiv:2107.01155 <https://arxiv.org/abs/2107.01155>
- 986 [2] Elvira Albert, Puri Arenas, Samir Genaim, and Germán Puebla. 2008. Automatic Inference of Upper Bounds for  
987 Recurrence Relations in Cost Analysis. In *Static Analysis, 15th International Symposium, SAS 2008, Valencia, Spain,*  
988 *July 16-18, 2008. Proceedings (Lecture Notes in Computer Science, Vol. 5079)*, Maria Alpuente and Germán Vidal (Eds.).  
989 Springer, 221–237. [https://doi.org/10.1007/978-3-540-69166-2\\_15](https://doi.org/10.1007/978-3-540-69166-2_15)
- 990 [3] Christophe Alias, Alain Darté, Paul Feautrier, and Laure Gonnord. 2010. Multi-dimensional Rankings, Program  
991 Termination, and Complexity Bounds of Flowchart Programs. In *Static Analysis - 17th International Symposium, SAS*  
992 *2010, Perpignan, France, September 14-16, 2010. Proceedings (Lecture Notes in Computer Science, Vol. 6337)*, Radhia Cousot  
993 and Matthieu Martel (Eds.). Springer, 117–133. [https://doi.org/10.1007/978-3-642-15769-1\\_8](https://doi.org/10.1007/978-3-642-15769-1_8)
- 994 [4] Stephan Arlt, Andreas Podelski, Cristiano Bertolini, Martin Schäfer, Ishan Banerjee, and Atif M. Memon. 2012. Light-  
995 weight Static Analysis for GUI Testing. In *23rd IEEE International Symposium on Software Reliability Engineering, ISSRE*  
996 *2012, Dallas, TX, USA, November 27-30, 2012*. IEEE Computer Society, 301–310. <https://doi.org/10.1109/ISSRE.2012.25>
- 997 [5] Todd M. Austin and Gurindar S. Sohi. 1992. Dynamic Dependency Analysis of Ordinary Programs. In *Proceedings of*  
998 *the 19th Annual International Symposium on Computer Architecture. Gold Coast, Australia, May 1992*, Allan Gottlieb  
999 (Ed.). ACM, 342–351. <https://doi.org/10.1145/139669.140395>
- 1000 [6] Manuel Barbosa, Gilles Barthe, Benjamin Grégoire, Adrien Koutsos, and Pierre-Yves Strub. 2021. Mechanized Proofs  
1001 of Adversarial Complexity and Application to Universal Composability. *IACR Cryptol. ePrint Arch.* 2021 (2021), 156.  
1002 <https://eprint.iacr.org/2021/156>
- 1003 [7] Raef Bassily, Kobbi Nissim, Adam D. Smith, Thomas Steinke, Uri Stemmer, and Jonathan R. Ullman. 2016. Algorithmic  
1004 stability for adaptive data analysis. In *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing,*  
1005 *STOC 2016, Cambridge, MA, USA, June 18-21, 2016*, Daniel Wichs and Yishay Mansour (Eds.). ACM, 1046–1059.  
1006 <https://doi.org/10.1145/2897518.2897566>
- 1007 [8] Gianfranco Bilardi and Keshav Pingali. 1996. A framework for generalized control dependence. In *Proceedings of the*  
1008 *ACM SIGPLAN 1996 conference on Programming language design and implementation*. 291–300.
- 1009 [9] Marc Brockschmidt, Fabian Emmes, Stephan Falke, Carsten Fuhs, and Jürgen Giesl. 2016. Analyzing Runtime and Size  
1010 Complexity of Integer Programs. *ACM Trans. Program. Lang. Syst.* 38, 4 (2016), 13:1–13:50. <http://dl.acm.org/citation.cfm?id=2866575>
- 1011 [10] Patrick Cousot. 2019. Abstract Semantic Dependency. In *Static Analysis - 26th International Symposium, SAS 2019,*  
1012 *Porto, Portugal, October 8-11, 2019, Proceedings (Lecture Notes in Computer Science, Vol. 11822)*, Bor-Yuh Evan Chang  
1013 (Ed.). Springer, 389–410. [https://doi.org/10.1007/978-3-030-32304-2\\_19](https://doi.org/10.1007/978-3-030-32304-2_19)
- 1014 [11] Ron Cytron, Jeanne Ferrante, Barry K. Rosen, Mark N. Wegman, and F. Kenneth Zadeck. 1991. Efficiently Computing  
1015 Static Single Assignment Form and the Control Dependence Graph. *ACM Trans. Program. Lang. Syst.* 13, 4 (1991),  
1016 451–490. <https://doi.org/10.1145/115372.115320>
- 1017 [12] Dorothy E. Denning and Peter J. Denning. 1977. Certification of Programs for Secure Information Flow. *Commun.*  
1018 *ACM* 20, 7 (1977), 504–513. <https://doi.org/10.1145/359636.359712>
- 1019 [13] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Roth. 2015. Generalization  
1020 in Adaptive Data Analysis and Holdout Reuse. In *Advances in Neural Information Processing Systems 28: Annual*  
1021 *Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, Corinna  
1022 Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, and Roman Garnett (Eds.). 2350–2358. <https://proceedings.neurips.cc/paper/2015/hash/bad5f33780c42f2588878a9d07405083-Abstract.html>
- 1023 [14] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Roth. 2015. The reusable  
1024 holdout: Preserving validity in adaptive data analysis. *Science* 349, 6248 (2015), 636–638.
- 1025 [15] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Leon Roth. 2015. Preserving  
1026 Statistical Validity in Adaptive Data Analysis. In *Proceedings of the Forty-Seventh Annual ACM Symposium on Theory*  
1027 *of Computing, STOC 2015, Portland, OR, USA, June 14-17, 2015*, Rocco A. Servedio and Ronitt Rubinfeld (Eds.). ACM,  
1028 117–126. <https://doi.org/10.1145/2746539.2746580>
- 1029 [16] Vitaly Feldman and Thomas Steinke. 2017. Generalization for Adaptively-chosen Estimators via Stable Median.  
In *Proceedings of the 30th Conference on Learning Theory, COLT 2017, Amsterdam, The Netherlands, 7-10 July 2017*  
(*Proceedings of Machine Learning Research, Vol. 65*), Satyen Kale and Ohad Shamir (Eds.). PMLR, 728–757. <http://proceedings.mlr.press/v65/feldman17a.html>
- [17] Jeanne Ferrante, Karl J. Ottenstein, and Joe D. Warren. 1987. The Program Dependence Graph and Its Use in  
Optimization. *ACM Trans. Program. Lang. Syst.* 9, 3 (1987), 319–349. <https://doi.org/10.1145/24039.24041>
- [18] Antonio Flores-Montoya and Reiner Hähnle. 2014. Resource Analysis of Complex Programs with Cost Equations. In  
*Programming Languages and Systems - 12th Asian Symposium, APLAS 2014, Singapore, November 17-19, 2014, Proceedings*

- (*Lecture Notes in Computer Science*, Vol. 8858), Jacques Garrigue (Ed.). Springer, 275–295. [https://doi.org/10.1007/978-3-319-12736-1\\_15](https://doi.org/10.1007/978-3-319-12736-1_15)
- [19] Andrew Gelman and Eric Loken. 2014. The Statistical Crisis in Science. *Am Sci* 102, 6 (2014), 460.
- [20] Joseph A. Goguen and José Meseguer. 1982. Security Policies and Security Models. In *1982 IEEE Symposium on Security and Privacy, Oakland, CA, USA, April 26-28, 1982*. IEEE Computer Society, 11–20. <https://doi.org/10.1109/SP.1982.10014>
- [21] Sumit Gulwani, Sagar Jain, and Eric Koskinen. 2009. Control-flow refinement and progress invariants for bound analysis. In *Proceedings of the 2009 ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2009, Dublin, Ireland, June 15-21, 2009*, Michael Hind and Amer Diwan (Eds.). ACM, 375–385. <https://doi.org/10.1145/1542476.1542518>
- [22] Sumit Gulwani and Florian Zuleger. 2010. The reachability-bound problem. In *Proceedings of the 2010 ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2010, Toronto, Ontario, Canada, June 5-10, 2010*, Benjamin G. Zorn and Alexander Aiken (Eds.). ACM, 292–304. <https://doi.org/10.1145/1806596.1806630>
- [23] Christian Hammer, Martin Grimme, and Jens Krinke. 2006. Dynamic path conditions in dependence graphs. In *Proceedings of the 2006 ACM SIGPLAN Workshop on Partial Evaluation and Semantics-based Program Manipulation, 2006, Charleston, South Carolina, USA, January 9-10, 2006*, John Hatcliff and Frank Tip (Eds.). ACM, 58–67. <https://doi.org/10.1145/1111542.1111552>
- [24] Christian Hammer, Martin Grimme, and Jens Krinke. 2006. Dynamic path conditions in dependence graphs. In *Proceedings of the 2006 ACM SIGPLAN Workshop on Partial Evaluation and Semantics-based Program Manipulation, 2006, Charleston, South Carolina, USA, January 9-10, 2006*, John Hatcliff and Frank Tip (Eds.). ACM, 58–67. <https://doi.org/10.1145/1111542.1111552>
- [25] Moritz Hardt and Jonathan R. Ullman. 2014. Preventing False Discovery in Interactive Data Analysis Is Hard. In *55th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2014, Philadelphia, PA, USA, October 18-21, 2014*. IEEE Computer Society, 454–463. <https://doi.org/10.1109/FOCS.2014.55>
- [26] John PA Ioannidis. 2005. Why most published research findings are false. *PLoS medicine* 2, 8 (2005), e124.
- [27] Christopher Jung, Katrina Ligett, Seth Neel, Aaron Roth, Saeed Sharifi-Malvajerdi, and Moshe Shenfeld. 2020. A New Analysis of Differential Privacy’s Generalization Guarantees. 151 (2020), 31:1–31:17. <https://doi.org/10.4230/LIPIcs.ITCS.2020.31>
- [28] Elisabet Lobo-Vesga, Alejandro Russo, and Marco Gaboardi. 2021. A Programming Language for Data Privacy with Accuracy Estimations. *ACM Transactions on Programming Languages and Systems (TOPLAS)* 43, 2 (2021), 1–42.
- [29] Atif M. Memon. 2007. An event-flow model of GUI-based applications for testing. *Softw. Test. Verification Reliab.* 17, 3 (2007), 137–157. <https://doi.org/10.1002/stvr.364>
- [30] Lori L. Pollock and Mary Lou Soffa. 1989. An Incremental Version of Iterative Data Flow Analysis. *IEEE Trans. Software Eng.* 15, 12 (1989), 1537–1549. <https://doi.org/10.1109/32.58766>
- [31] Ryan Rogers, Aaron Roth, Adam D. Smith, Nathan Srebro, Om Thakkar, and Blake E. Woodworth. 2020. Guaranteed Validity for Empirical Approaches to Adaptive Data Analysis. In *The 23rd International Conference on Artificial Intelligence and Statistics, AISTATS 2020, 26-28 August 2020, Online [Palermo, Sicily, Italy] (Proceedings of Machine Learning Research, Vol. 108)*, Silvia Chiappa and Roberto Calandra (Eds.). PMLR, 2830–2840. <http://proceedings.mlr.press/v108/rogers20a.html>
- [32] Moritz Sinn, Florian Zuleger, and Helmut Veith. 2014. A Simple and Scalable Static Analysis for Bound Analysis and Amortized Complexity Analysis. In *Computer Aided Verification - 26th International Conference, CAV 2014, Held as Part of the Vienna Summer of Logic, VSL 2014, Vienna, Austria, July 18-22, 2014. Proceedings (Lecture Notes in Computer Science, Vol. 8559)*, Armin Biere and Roderick Bloem (Eds.). Springer, 745–761. [https://doi.org/10.1007/978-3-319-08867-9\\_50](https://doi.org/10.1007/978-3-319-08867-9_50)
- [33] Moritz Sinn, Florian Zuleger, and Helmut Veith. 2017. Complexity and resource bound analysis of imperative programs using difference constraints. *Journal of automated reasoning* 59, 1 (2017), 3–45.
- [34] Moritz Sinn, Florian Zuleger, and Helmut Veith. 2017. Complexity and Resource Bound Analysis of Imperative Programs Using Difference Constraints. *J. Autom. Reason.* 59, 1 (2017), 3–45. <https://doi.org/10.1007/s10817-016-9402-4>
- [35] Thomas Steinke and Lydia Zakyntinou. 2020. Reasoning About Generalization via Conditional Mutual Information. In *Conference on Learning Theory, COLT 2020, 9-12 July 2020, Virtual Event [Graz, Austria] (Proceedings of Machine Learning Research, Vol. 125)*, Jacob D. Abernethy and Shivani Agarwal (Eds.). PMLR, 3437–3452. <http://proceedings.mlr.press/v125/steinke20a.html>
- [36] Jonathan R. Ullman, Adam D. Smith, Kobbi Nissim, Uri Stemmer, and Thomas Steinke. 2018. The Limits of Post-Selection Generalization. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (Eds.). 6402–6411. <https://proceedings.neurips.cc/paper/2018/hash/77ee3bc58ce560b86c2b59363281e914-Abstract.html>
- [37] Daniel Winograd-Cort, Andreas Haeberlen, Aaron Roth, and Benjamin C. Pierce. 2017. A framework for adaptive differential privacy. *Proc. ACM Program. Lang.* 1, ICFP (2017), 10:1–10:29. <https://doi.org/10.1145/3110254>

Short Title

- 1079 [38] Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. 2018. Privacy Risk in Machine Learning: Analyzing  
1080 the Connection to Overfitting. In *31st IEEE Computer Security Foundations Symposium, CSF 2018, Oxford, United*  
1081 *Kingdom, July 9-12, 2018*. IEEE Computer Society, 268–282. <https://doi.org/10.1109/CSF.2018.00027>  
1082 [39] Florian Zuleger, Sumit Gulwani, Moritz Sinn, and Helmut Veith. 2011. Bound Analysis of Imperative Programs with  
1083 the Size-Change Abstraction. In *Static Analysis - 18th International Symposium, SAS 2011, Venice, Italy, September*  
1084 *14-16, 2011. Proceedings (Lecture Notes in Computer Science, Vol. 6887)*, Eran Yahav (Ed.). Springer, 280–297. [https://doi.org/10.1007/978-3-642-23702-7\\_22](https://doi.org/10.1007/978-3-642-23702-7_22)

1085 Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127