

Program Analysis for Adaptivity Analysis

In this appendix, we present the full details of the 2 languages: while language and the SSA language.

Contents

| | | |
|----------|--|-----------|
| 1 | System Overview | 1 |
| 2 | Labeled While Language | 3 |
| 2.1 | Syntax and Semantics | 3 |
| 2.2 | Trace-based Operational Semantics | 4 |
| 2.3 | Trace-based Adaptivity | 8 |
| 3 | Labeled SSA Language | 13 |
| 3.1 | The Limit of While Language | 13 |
| 3.2 | SSA form Language | 14 |
| 3.3 | Trace-based Operational Semantics for SSA Language | 15 |
| 3.4 | Event and Trace | 16 |
| 3.5 | Trace-based Adaptivity | 21 |
| 3.6 | SSA Transformation and Soundness of Transformation | 26 |
| 4 | AdaptFun | 27 |
| 4.1 | Notations | 27 |
| 4.2 | Algorithmic Analysis Rules | 29 |
| 4.3 | Adaptivity Based on Program Analysis in AdaptFun | 33 |
| 4.4 | [[Soundness of the AdaptFun]] | 34 |
| 5 | [[Examples]] | 38 |
| 6 | Non Determinism | 42 |
| 7 | Analysis of Generalization Error | 43 |

1 System Overview

In adaptive data analysis, a data analysis can depend on the results of previous analysis over the same data. This dependency may affect the *generalization properties of the data analysis*. To study this phenomenon in a formal way, we consider the *statistical query model*. In this model, a dataset D consisting of d attributes (columns) and n individuals' data (rows) can be accessed only through an

interface to which one can submit statistical queries. More precisely, suppose that the type of a row is R (as an example, a row with d binary attributes would have type $R = \{0, 1\}^d$). Then, in the statistical query model one can access the dataset only by submitting a query to the interface, in the form of a function $p : D \rightarrow [0, 1]$ where D represents dataset. The collected answer of the asked query is the average result of p on each row in the dataset D . For example, the result is the value $\frac{1}{n} \sum_{i=1}^n p(D_i)$ where D_i is the row of index i in D . While this model is rather simple, in fact it supports sufficient statistics one may be interested.

We are interested in the adaptivity of mechanisms in the model, which is straightforward supported by a high level language. In this language, queries are allowed to carry arguments to simulate the process of submitting a query to the interface in the model, for example, the expression $\text{query}(\psi)$ tells us the argument ψ is consumed to construct the query. To be precise, one submitted query who needs the average of answers of previous queries is expressed as $\text{query}(x)$, where the variable x stores the expected average results. This makes these mechanisms quite straightforward to express in the high level language. However, this convenience pays at the price that the adaptivity A of a mechanism P becomes quite tricky to estimate because the definition of dependency between two queries becomes vague in the high level language.

```

 $x \leftarrow \text{query}(0);$ 
 $\text{if } (x_1 > 0)$ 
 $y \leftarrow \text{query}(x)$ 
```

The dependency between two query submissions is the essential of the adaptivity of a mechanism. To study the dependency, we first study its dual, independence between two queries, which is defined to be: one query $\text{query}(0)$ does not depend on another query $\text{query}(x)$ when the result of $\text{query}(0)$ remains the same regardless of the modification of the result of $\text{query}(x)$. Hence, it becomes hard to distinguish whether the variance of result of $\text{query}(0)$ comes from the control flow or the argument of queries. Since we know that the result of one query from a specific D may vary under different contexts in the high level language.

[[To resolve the dilemma, we translate any program(mechanism) into its counterpart in a low level language, which mimics the high level one except its only allowing atomic queries, $\text{-- query}(0)$ -- . That is to say, given a data base D , the result of the query from D becomes deterministic. We need to show the two programs P and P^* are observably equivalent over the translation. In this way, we can define the adaptivity of a program under this model only based on the control flow. To be specific, the adaptivity A of a program P is defined based on graphs, called dependency graph, which comes from the semantics of the low level program.]] The dependency graph is constructed using a trace of queries generated along with the semantics: The queries in the trace consists of the nodes in the graph while the edge represents dependency. If there is no dependency between two node(queries), there will be no edge. Intuitively, we want to give an approximation of the adaptivity by static analysis. To this end, we propose **AdaptFun**, which estimates an upper bound on the program.

[[The adaptivity A of arbitrary high level program c is defined to be the minimal of the adaptivity A of all the possible c via various valid translations. Being valid means the programs before and after the translation are observably equivalent. Naturally, following this definition, the upper bound estimated by **AdaptFun is sound with respect to its low level adaptivity A , hence the high level one A .]]**

Finally we extend the language to support the probabilistic program and extend the adaptivity definition accordingly.

The key component of the system is a program analysis tool, which provides an upper bound on the adaptivity of the program.

2 Labeled While Language

2.1 Syntax and Semantics

| | | | |
|------------------------|---------------|-------|---|
| Arithmetic Operators | \oplus_a | $::=$ | $+ \mid - \mid \times \mid \div$ |
| Boolean Operators | \oplus_b | $::=$ | $\vee \mid \wedge$ |
| Relational Operators | \sim | $::=$ | $< \mid \leq \mid ==$ |
| Label | l | \in | \mathbb{N} |
| Arithmetic Expressions | a | $::=$ | $n \mid x \mid a \oplus_a a \mid$ |
| Boolean Expressions | b | $::=$ | $\text{true} \mid \text{false} \mid \neg b \mid b \oplus_b b \mid a \sim a$ |
| Value | v | $::=$ | $n \mid \text{true} \mid \text{false} \mid [] \mid [v, \dots, v]$ |
| Expression | e | $::=$ | $v \mid a \mid b \mid [e, \dots, e]$ |
| Query Value | α | $::=$ | $n \mid \chi[n] \oplus_a \chi[n] \mid n \oplus_a \chi[n] \mid \chi[n] \oplus_a n$ |
| Query Expression | ψ | $::=$ | $\alpha \mid a \mid \psi \oplus_a \psi$ |
| Labeled Command | c | $::=$ | $[\text{skip}]^l \mid [x \leftarrow e]^l \mid [x \leftarrow \text{query}(\psi)]^l$ $\mid \text{while } [b]^l, n \text{ do } c \mid c; c \mid \text{if } ([b]^l, c, c)$ |
| Memory | m | $::=$ | $[] \mid (x^l \rightarrow v) :: m$ |
| Annotated Variable | av | $::=$ | (x, v, l, n) |
| Variable Trace | τ | $::=$ | $[] \mid \text{av} :: \tau$ |
| Variable Counter | vcnt | $::=$ | $\mathcal{VAR} \rightarrow \mathbb{N}$ |

We use following notations to represent the set of corresponding definitions:

| | | |
|--------------------------|---|----------------------------|
| \mathcal{VAR} | : | Set of Variables |
| \mathcal{VAL} | : | Set of Values |
| \mathcal{AV} | : | Set of Annotated Variables |
| \mathcal{M} | : | Set of Memories |
| \mathcal{DB} | : | Set of Databases |
| $\mathcal{QD} = [-1, 1]$ | : | Domain of Query Results |

[[

Definition 1 (Assigned Variables (aVar)). *Given a program c , its assigned variables aVar_c is a vector containing all variables newly assigned in the program preserving the order, $\forall \mathbf{x} \in \text{aVar}, \mathbf{x} \in \mathcal{VAR}$. It is defined as follows:*

$$\text{aVar}_c \triangleq \begin{cases} [\mathbf{x}] & \mathbf{c} = [\mathbf{x} \leftarrow \mathbf{e}]^{(l, w)} \\ [\mathbf{x}] & \mathbf{c} = [\mathbf{x} \leftarrow \text{query}(\psi)]^{(l, w)} \\ \text{aVar}_{\mathbf{c}_1} + \text{aVar}_{\mathbf{c}_2} & \mathbf{c} = \mathbf{c}_1; \mathbf{c}_2 \\ \text{aVar}_{\mathbf{c}_1} + \text{aVar}_{\mathbf{c}_2} & \mathbf{c} = \text{if } ([b]^{(l, w)}, \mathbf{c}_1, \mathbf{c}_2) \\ \text{aVar}_{\mathbf{c}'} & \mathbf{c} = \text{while } [b]^{(l, w)}, n \text{ do } \mathbf{c}' \end{cases}$$

Definition 2 (Query Variables (qVar)). .

Given a program c , its query variables qVar is a vector containing all variables newly assigned by a

query in the program, $\text{qVar} \subset \mathcal{VAR}$. It is defined as follows:

$$\text{qVar}_c \triangleq \begin{cases} [] & \mathbf{c} = [\mathbf{x} \leftarrow \mathbf{e}]^{(l,w)} \\ [\mathbf{x}] & \mathbf{c} = [\mathbf{x} \leftarrow \text{query}(\psi)]^{(l,w)} \\ \text{qVar}_{c_1} ++ \text{qVar}_{c_2} & \mathbf{c} = \mathbf{c}_1; \mathbf{c}_2 \\ \text{qVar}_{c_1} ++ \text{qVar}_{c_2} & \mathbf{c} = \text{if } ([\mathbf{b}]^{(l,w)}, c_1, c_2) \\ \text{qVar}_{c'} & \mathbf{c} = \text{while } [b]^{(l,w)}, n \text{ do } c' \end{cases}$$

We are abusing the notations and operators from list here. The notation $[]$ represents an empty vector and $x :: A$ represents add an element x to the head of the vector A . The concatenation operation between 2 vectors A_1 and A_2 , i.e., $A_1 ++ A_2$ is mimic the standard list concatenation operations as follows:

$$A_1 ++ A_2 \triangleq \begin{cases} A_2 & A_1 = [] \\ x :: (A'_1 ++ A_2) & A_1 = x :: A'_1 \end{cases} \quad (1)$$

We use index within parenthesis to denote the access to the element of corresponding location, $A(i)$ denotes the element at location i in the vector A and $M(i, j)$ denotes the element at location i -th row, j -th column in the matrix M .

The variable counter vcnt_c maps every assigned variables aVar_c to a natural number $n \in \mathbb{N}$ in a certain execution. This natural number n represents the visiting times of this variable in this certain execution.

We use variable name x within parenthesis to denote the access to the associated natural number of this variable in the variable counter vcnt_c , $\text{vcnt}_c(x)$ denote the visiting times of variable x .

Definition 3 (Initial Variable Counter vcnt_c^0). Given a program c with its assigned variables aVar_c of length N , its initial variable counter vcnt_c^0 maps all the variable to 0, i.e.:

$$\text{vcnt}_c^0(x) = 0, x = \text{aVar}_c(i) \forall i = 1, \dots, N$$

]]

2.2 Trace-based Operational Semantics

We evaluate programs in the `While` language by means of a trace-based operational semantics, to capture the dependency between queries. For distinguishing elements in the trace, we add a label to commands in the `While` language as defined in the syntax. Each command is labeled with a label l , a natural number standing for the line of code where the command appears. Notice that we associate the label l to the conditional predicate b in the `if` statement, and to the guard b in the `while` statement. Some non-standard syntax is explained as follows:

A memory is standard, a map from labeled variables to values. Queries can be uniquely annotated as defined in \mathcal{AQ} , and the annotation (l, w) considers the location of the query by line number l and which iteration the query is at when it appears in a `while` statement, specified by w .

A configuration, $\langle m, c, t, w \rangle$, contains four elements: a memory m , the command c to be evaluated, a starting trace t , a starting while map w . Most of the time, the while maps remains empty until the evaluation goes into `while` statements.

[[The annotated query $\text{av} = (\alpha, v, l, n)$ is a tuple contains 4 elements.]]

Definition 4 (Order of Annotated Variables). .

Given 2 annotated queries $\text{av}_1 = (x_1, v_1, l_1, n_1), \text{av}_2 = (x_2, v_2, l_2, n_2)$:

$$\text{av}_1 <_{\text{av}} \text{av}_2 \triangleq \begin{cases} n_1 < n_2 & l_1 = l_2 \\ w_1 <_w w_2 & \text{o.w.} \end{cases} \text{[WQ:hard:-]}$$

$\text{av}_1 \geq_{\text{av}} \text{av}_2$ is defined vice versa.

A variable trace τ is a list of annotated queries accumulated along the execution of the program. A trace can be regarded as the program history, where this history consists of all the queries asked by the analyst during the execution of the program. We collect the trace with a trace-based small-step operational semantics based on transitions of the program configuration $\langle m, c, \tau, \text{vcnt} \rangle$, of form $\langle m, c, \tau, \text{vcnt} \rangle \rightarrow \langle m', \text{skip}, \tau', \text{vcnt}' \rangle$. The evaluation rules for arithmetic and boolean expressions are standard. They have the form $\langle m, a \rangle \rightarrow_a a'$, evaluating an arithmetic expression a in the memory m , and similar for the boolean expressions $\langle m, b \rangle \rightarrow_b b'$, defined as follows:

$$\langle m, a \rangle \rightarrow_a a' : \text{Memory} \times \text{AExpr} \Rightarrow \text{AExpr}$$

$$\langle m, b \rangle \rightarrow_b b' : \text{Memory} \times \text{BExpr} \Rightarrow \text{BExpr}$$

Given the evaluation for the arithmetic and boolean expression, we defined the evaluation rules for query expression ψ correspondingly as follows:

$$\langle m, \psi \rangle \rightarrow_q \psi' : \text{Memory} \times \text{QExpr} \rightarrow_q \text{QExpr}$$

$$\begin{array}{c} \frac{\langle m, n \oplus_a n \rangle \rightarrow_a n'}{\langle m, n \oplus_a n \rangle \rightarrow_q n'} \quad \frac{\langle m, \psi \rangle \rightarrow_q \psi'}{\langle m, \psi \oplus_a \alpha \rangle \rightarrow_q \psi' \oplus_a \alpha} \quad \frac{\langle m, \psi_2 \rangle \rightarrow_q \psi'_2}{\langle m, \psi_1 \oplus_a \psi_2 \rangle \rightarrow_q \psi_1 \oplus_a \psi'_2} \\[10pt] \frac{\langle m, a \rangle \rightarrow_a a'}{\langle m, \chi[a] \rangle \rightarrow_q \chi[a']} \quad \frac{\langle m, a \rangle \rightarrow_a a'}{\langle m, a \rangle \rightarrow_q a'} \end{array}$$

Given the evaluation rules for query expression, we can define its equivalence relation, as follows in Definition 17.

Definition 5 (Equivalence of Query). . Given a memory m and 2 query expressions ψ_1, ψ_2 s.t., $FV(\psi_1) \in \text{dom}(m)$ and $FV(\psi_2) \in \text{dom}(m)$:

$$\psi_1 =_q^m \psi_2 \triangleq \begin{cases} \text{true} & \exists \alpha_1, \alpha_2. (\langle m, \psi_1 \rangle \rightarrow_q \alpha_1 \wedge \langle m, \psi_2 \rangle \rightarrow_q \alpha_2) \\ & \wedge (\forall r \in \mathcal{QD}. \exists v. \text{s.t.}, \langle m, \alpha_1[r/\chi] \rangle \rightarrow_a v \wedge \langle m, \alpha_2[r/\chi] \rangle \rightarrow_a v) \\ \text{false} & \text{o.w.} \end{cases}$$

, where $FV(\psi)$ is the set of free variables in the query expression ψ . $\psi_1 \neq_q^m \psi_2$ is defined vice versa. We use $=_q$ and \neq_q as the shorthands for $=_q^\square$ and \neq_q^\square .

Then, we have the corresponding equivalence relation between 2 annotated queries defined in Definition 18: \square

Definition 6 (Equivalence of Annotated Variables). Given 2 annotated queries $\text{av}_1 = (x_1, v_1, l_1, n_1), \text{av}_2 = (x_2, v_2, l_2, n_2)$:

$$\text{av}_1 =_{\text{av}} \text{av}_2 \triangleq (l_1 = l_2 \wedge w_1 =_w w_2 \wedge \alpha_1 =_q \alpha_2)$$

$\text{av}_1 \neq_{\text{av}} \text{av}_2$ is defined vice versa.

]] [[Given an annotated variable av and a trace t , the appending operation $av :: t$ is the standard list appending operation, appends av to the head of trace t . The concatenation operation between 2 traces t_1 and t_2 , i.e., $t_1 ++ t_2$ is the standard list concatenation operation as follows:

$$t_1 ++ t_2 \triangleq \begin{cases} t_2 & t_1 = [] \\ av :: (t'_1 ++ t_2) & t_1 = av :: t'_1 \end{cases} \quad (2)$$

The subtraction operation between 2 traces t_1 and t_2 , i.e., $t_1 - t_2$ is defined as follows:

$$t_1 - t_2 \triangleq t_3 \text{ s.t., } t_2 ++ t_3 = t_1 \quad (3)$$

Given an annotated query av , av belongs to a trace t , i.e., $av \in_{av} t$ are defined as follows:

$$av \in_{av} t \triangleq \begin{cases} \text{false} & t = [] \\ \text{true} & t = av' :: t' \quad av =_{av} av' \\ av \in t' & t = av' :: t' \quad av \neq_{av} av' \end{cases} \quad (4)$$

]] [[

Definition 7 (Equivalence of Program). Given 2 programs c_1 and c_2 :

$$c_1 =_c c_2 \triangleq \begin{cases} \text{true} & c_1 = \text{skip} \wedge c_2 = \text{skip} \\ \forall m, e_1. \langle m, e_1 \rangle \rightarrow_a^* v \wedge \langle m, e_1 \rangle \rightarrow_a^* v & c_1 = x \leftarrow e_1 \wedge c_2 = x \leftarrow e_2 \\ \psi_1 =_q \psi_2 & c_1 = x \leftarrow \text{query}(\psi_1) \wedge c_2 = x \leftarrow \text{query}(\psi_2) \\ c_1^f =_c c_2^f \wedge c_1^t =_c c_2^t & c_1 = \text{if } (b, c_1^t, c_1^f) \wedge c_2 = \text{if } (b, c_2^t, c_2^f) \\ c_1' =_c c_2' & c_1 = \text{while } b \text{ do } c_1' \wedge c_2 = \text{while } b \text{ do } c_2' \\ c_1^h =_c c_2^h \wedge c_1^t =_c c_2^t & c_1 = c_1^h; c_1^t \wedge c_2 = c_2^h; c_2^t \end{cases}$$

$c_1 \neq_c c_2$ is defined vice versa.

Given 2 programs c and c' , c' is a sub-program of c , i.e., $c' \in_c c$ is defined as:

$$c' \in_c c \triangleq \exists c_1, c_2, c''. \text{ s.t., } c =_c c_1; c''; c_2 \wedge c' =_c c'' \quad (5)$$

]]

The small-step transition states that a configuration $\langle m, c, t, w \rangle$ evaluates to another configuration with the trace and while map updated along with the evaluation of the command c to the normal form of the command skip . We define rules of the trace-based operational semantics in Figure 1. The rule **query-c** evaluates the argument of a query request. When the argument is in normal form, this query will be answered. The rule **query-v** modifies the starting memory m to $m[\alpha/x]$ using the answer α of the query $\text{query}(\alpha)$ from the mechanism, with the trace expanded by appending the query $\text{query}(\alpha)$ with the current annotation (l, w) . The rule for assignment is standard and the trace remains unchanged. The sequence rule keeps tracking the modification of the trace, and the evaluation rule for if conditional goes into one branch based on the result of the conditional predicate b . The rules for while modify the while map w . In the rule **ifw-true**, the while map w is updated by $w + l$ because the execution goes into another iteration when the condition $n > 0$ is satisfied. When n reaches 0, the loop exits and the while map w eliminates the label l of this while statement by $w - l$ in the rule **ifw-false**. With the operational semantics and relations between annotated queries, we restrict the well-formed trace w.r.t. the execution of a program c in Definition 23. [WQ: we define map update as follows, if we update the map m with value v at its key k , we denote $m[x \rightarrow v]$ or $m[v/x]$.]

[WQ: Evaluation context $E ::= [\cdot] | x \leftarrow E | \text{if } E \text{ } c_1, c_2 | x \leftarrow \text{query}(E) | \text{while } E \text{ do } c | E; c | \text{skip}; E$

$$\overline{\langle m, E[e], \tau, \text{vcnt} \rangle} \rightarrow \langle m, E[e'] \tau, \text{vcnt} \rangle$$

]

[JL:

$$\boxed{Memory \times Command \times VTrace \times VCounter \rightarrow Memory \times Command \times VTrace \times VCounter}$$

$$\boxed{\langle m, c, \tau, \text{vcnt} \rangle \rightarrow \langle m', c', \tau', \text{vcnt}' \rangle}$$

$$\frac{\langle m, e \rangle \rightarrow \langle m, e' \rangle}{\langle m, [x \leftarrow e]^l, \tau, \text{vcnt} \rangle \rightarrow \langle m, [x \leftarrow e']^l, \tau, \text{vcnt} \rangle} \text{ assn-e}$$

$$\frac{n = \text{vcnt}[x] + 1 \quad \text{av} = (x, v, l, n)}{\langle m, [x \leftarrow v]^l, \tau, \text{vcnt} \rangle \rightarrow \langle m[v/x], [\text{skip}]^l, \tau + +[\text{av}], \text{vcnt}[x \rightarrow n] \rangle} \text{ assn-v}$$

$$\frac{}{\langle m, \text{while } [b]^l, n \text{ do } c, \tau, \text{vcnt} \rangle \rightarrow \langle m, c; \text{if } (b, c; \text{while } [b]^l, (n+1) \text{ do } c, \text{skip}), \tau, \text{vcnt} \rangle} \text{ while-b}$$

$$\frac{\langle m, \psi \rangle \rightarrow_q \psi'}{\langle m, [x \leftarrow \text{query}(\psi)]^l, \tau, \text{vcnt} \rangle \rightarrow \langle m, [x \leftarrow \text{query}(\psi')]^l, \tau, \text{vcnt} \rangle} \text{ query-e}$$

$$\frac{\text{query}(\alpha) = v \quad n = \text{vcnt}[x] + 1 \quad \text{av} = (x, \alpha, l, n)}{\langle m, [x \leftarrow \text{query}(\alpha)]^l, \tau, \text{vcnt} \rangle \rightarrow \langle m[v/x], \text{skip}, \tau + +[\text{av}], \text{vcnt}[x \rightarrow n] \rangle} \text{ query-v}$$

$$\frac{\langle m, c_1, \tau, \text{vcnt} \rangle \rightarrow \langle m', c'_1, \tau', \text{vcnt}' \rangle}{\langle m, c_1; c_2, \tau, \text{vcnt} \rangle \rightarrow \langle m', c'_1; c_2, \tau', \text{vcnt}' \rangle} \text{ seq1} \quad \frac{}{\langle m, [\text{skip}]^l; c_2, \tau, \text{vcnt} \rangle \rightarrow \langle m, c_2, \tau, \text{vcnt} \rangle} \text{ seq2}$$

$$\frac{\langle m, b \rangle \rightarrow_b b'}{\langle m, \text{if } ([b]^l, c_1, c_2), \tau, \text{vcnt} \rangle \rightarrow \langle m, \text{if } ([b']^l, c_1, c_2), \tau, \text{vcnt} \rangle} \text{ if-b}$$

$$\frac{}{\langle m, \text{if } ([\text{true}]^l, c_1, c_2), \tau, \text{vcnt} \rangle \rightarrow \langle m, c_1, \tau, \text{vcnt} \rangle} \text{ if-t}$$

$$\frac{}{\langle m, \text{if } ([\text{false}]^l, c_1, c_2), \tau, \text{vcnt} \rangle \rightarrow \langle m, c_2, \tau, \text{vcnt} \rangle} \text{ if-f}$$

]

Figure 1: Trace-based Operational Semantics of While Language.

2.3 Trace-based Adaptivity

We define adaptivity through a query-based dependency graph. In our model, an *analyst* asks a sequence of queries to the mechanism, and the analyst receives the answers to these queries from the mechanism. A query is adaptively chosen by the analyst when the choice of this query is affected by answers from previous queries. In this model, the adaptivity we are interested in is the length of the longest sequence of such adaptively chosen queries, among all the queries the data analyst asks to the mechanism. Also, when the analyst asks a query, the only information the analyst will have will be the answers to previous queries and the state of the program. It means that when we want to know if this query is adaptively chosen, we only need to check whether the choice of this query will be affected by changes of answers to previous queries. There are two possible situations that can affect the choice of a query, either the query argument directly uses the results of previous queries (data dependency), or the control flow of the program with respect to a query (whether to ask this query or not) depends on the results of previous queries (control flow dependency).

As a first step, we give a definition of when one query may depend on a previous query, which is supposed to consider both control dependency and data dependency. We first look at two possible candidates:

1. One query may depend on a previous query if and only if a change of the answer to the previous query may also change the result of the query.
2. One query may depend on a previous query if and only if a change of the answer to the previous query may also change the appearance of the query.

The first candidate works well by witnessing the result of one query according to the change of the answer of another query. We can easily find that the two queries have nothing to do with each other in a simple example $c = x \leftarrow \text{query}(\chi(1)); y \leftarrow \text{query}(\chi(2))$. This candidate definition works well with respect to data dependency. However, it fails to handle control dependency since it just monitors the changes to the answer of a query when the answer of previous queries returned change. The key point is that this query may also not be asked because of an analyst decision which depend on the answers of previous queries. An example of this situation is shown in program c_1 as follows.

$$c_1 = x \leftarrow \text{query}(\chi(1)); \text{if } (x > 2, y \leftarrow \text{query}(\chi(2)), \text{skip})$$

We choose the second candidate, which performs well by witnessing the appearance of one query $\text{query}(\chi(2))$ upon the change of the result of one previous query $\text{query}(\chi(1))$ in c_1 . It considers the control dependency, and at the same, does not miss the data dependency. In particular, the arguments of a query characterizes it. In this sense, if the data used in the arguments changes due to a different answer to a certain previous query, the appearance of the query may change as well. This situation is also captured by our definition. Let us look at another variant of program c , p_2 , in which the queries equipped with functions using previously assigned variables storing answer of its previous query.

$$c_2 = x \leftarrow \text{query}(\chi(2)); y \leftarrow \text{query}(x + \chi(3))$$

As a reminder, in the `While` language, the query request is composed by two components: a symbol query representing a linear query type and the argument e , which represents the function specifying what the query asks. So we do think $\text{query}(\chi(1))$ is different from $\text{query}(\chi(2))$. Informally, we think $\text{query}(x + \chi(3))$ may depend on the query $\text{query}(\chi(2))$, because equipped function of the former $x + \chi(3)$ depend on the data assigned with $\text{query}(\chi(2))$. We can see the appearance definition catches

data dependency in such a way, since $\text{query}(x + \chi(2))$ will not be the same query if the value of x is changed.

We give a formal definition of variable may dependency based on the trace-based operational semantics as follows. [JL]:

Definition 8 (Annotated Variables May Dependency). .

One annotated variable $\text{av}_2 = (x_2, v_2, l_2, n_2)$ may depend on another one $\text{av}_1 = (x_1, v_1, l_1, n_1)$ in a program c , with a starting memory m and hidden database D , denoted as $\text{DEP}_{\text{av}}(\text{av}_1, \text{av}_2, c, m, D)$ is defined below.

$$\begin{aligned} & \exists \mathbf{m}_1, \mathbf{m}_3, \tau_1, \tau_3, \mathbf{c}_2, v_1, (\alpha_1 \vee \mathbf{e}_1). \\ & \left(\begin{aligned} & \langle \mathbf{m}, \mathbf{c}, [], [] \rangle \rightarrow^* \langle \mathbf{m}_1, [\mathbf{x}_1 \leftarrow \text{query}(\alpha_1)(/\mathbf{e}_1)]^{l_1}; \mathbf{c}_1, \text{qt}_1, \tau_1, w_1 \rangle \xrightarrow{\text{ssa-query-v} (/ \text{assn-v})} \\ & \langle \mathbf{m}_1[v_1/\mathbf{x}], \mathbf{c}_2, \text{qt}'_1, \tau_1 + +[\text{av}_1], w_1 \rangle \rightarrow^* \langle \mathbf{m}_3, \text{skip}, \text{qt}_3, \tau_3, w_3 \rangle \\ & \wedge \left(\begin{aligned} & \text{av}_2 \in (\tau'_3 - (\tau_1 + +[\text{av}_1])) \\ & \Rightarrow \exists v \in \mathcal{QD}, v \neq v_1, \mathbf{m}'_3, \text{qt}'_3, \tau'_3, w'_3. \langle \mathbf{m}_1[v/\mathbf{x}], \mathbf{c}_2, \text{qt}'_1, \tau_1 + +[\text{av}_1], w_1 \rangle \\ & \rightarrow^* (\langle \mathbf{m}'_3, \text{skip}, \text{qt}'_3, \tau'_3, w'_3 \rangle \\ & \wedge \text{av}_2 \notin (\tau'_3 - (\tau_1 + +[\text{av}_1]))) \\ & \wedge \left(\begin{aligned} & \text{av}_2 \notin (\tau_3 - (\tau_1 + +[\text{av}_1])) \\ & \Rightarrow \exists v \in \mathcal{QD}, v \neq v_1, \mathbf{m}'_3, \text{qt}'_3, \tau'_3, w'_3. \langle \mathbf{m}_1[v/\mathbf{x}], \mathbf{c}_2, \text{qt}'_1, \tau_1 + +[\text{av}_1], w_1 \rangle \\ & \rightarrow^* (\langle \mathbf{m}'_3, \text{skip}, \text{qt}'_3, \tau'_3, w'_3 \rangle \\ & \wedge \text{av}_2 \in (\tau'_3 - (\tau_1 + +[\text{av}_1]))) \end{aligned} \right) \end{aligned} \right) \end{aligned} \end{aligned}$$

Definition 9 (Annotated Variables May Dependency – Version 2). .

One annotated variable $\text{av}_2 = (x_2, v_2, l_2, n_2)$ may depend on another one $\text{av}_1 = (x_1, v_1, l_1, n_1)$ in a program \mathbf{c} , with a starting memory \mathbf{m} and hidden database D , denoted as $\text{DEP}_{\text{av}}(\text{av}_1, \text{av}_2, c, m, D)$ is defined below.

$$\begin{aligned} & \exists \mathbf{m}, \mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_3, \mathbf{m}'_2, \mathbf{m}'_3, \tau_1, \tau_2, \tau'_2, t_1, t_2, t'_2, \mathbf{c}_1, \mathbf{c}_2, v'_1. \\ & \left(\begin{aligned} & \langle \mathbf{m}, \mathbf{c}, [] \rangle \rightarrow^* \langle \mathbf{m}_1, [\mathbf{x}_1 \leftarrow \text{query}(\alpha_1)(/\mathbf{e}_1)]^{l_1}; \mathbf{c}_1, \tau_1, t_1 \rangle \\ & \wedge \langle \mathbf{m}_1[v_1/\mathbf{x}_1], \mathbf{c}_1, \tau_1 + +[\text{av}_1], t_1[\mathbf{x}_1] + + \rangle \rightarrow^* \langle \mathbf{m}_2, [x_2 \leftarrow \text{query}(\alpha_2)(/\mathbf{e}_2)]^{l_2}; \mathbf{c}_2, \tau_2, t_2 \rangle \\ & \xrightarrow{\text{ssa-query-v} (/ \text{assn-v})} \langle \mathbf{m}_3, \mathbf{c}_2, \tau_2 + +[\text{av}_2], t_2[\mathbf{x}_2] + + \rangle \\ & \wedge \langle \mathbf{m}_1[v'_1/\mathbf{x}_1], \mathbf{c}_1, \tau_1, t_1 \rangle \rightarrow^* \langle \mathbf{m}'_2, \mathbf{c}_2, \tau'_2, t'_2 \rangle \\ & \wedge \text{av}_2 \notin \tau'_2 \end{aligned} \right) \end{aligned}$$

Definition 10 (Variable May Dependency). .

Given a program \mathbf{c} with its assigned variables $\text{aVar}_{\mathbf{c}}$, one variable $\mathbf{x}_2 \in \text{aVar}_{\mathbf{c}}$ may depend on another variable $\mathbf{x}_1 \in \text{aVar}_{\mathbf{c}}$ in \mathbf{c} denoted as $\text{DEP}_{\text{var}}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{c})$ is defined below.

$$\exists v_1, v_2, n_1, n_2, m, D. \text{av}_1 = (x_1, v_1, l_1, n_1) \wedge \text{av}_2 = (x_2, v_2, l_2, n_2) \wedge \text{DEP}_{\text{av}}(\text{av}_1, \text{av}_2, c, m, D)$$

Definition 11 (Execution Based Dependency Graph). .

Given a program c , a database D , a starting memory m with its assigned variables aVar_c and initial variable counter vcnt_c^0 with its corresponding execution: $\langle m, c, [], \text{vcnt}_c^0 \rangle \rightarrow^* \langle \mathbf{m}', \text{skip}, \tau, \text{vcnt} \rangle$, the dependency graph $\mathbf{G}_{\text{trace}}(\mathbf{c}, \mathbf{m}, \mathbf{D}) = (V, E, W, \text{qF})$ is defined as:

$$\begin{aligned} \text{Vertices } V & := \{x \in \mathcal{VAR} \mid x = \text{aVar}_c(i); i = 0, \dots, |\text{aVar}_c|\} \\ \text{Directed Edges } E & := \{(x, x') \in \mathcal{VAR} \times \mathcal{VAR} \mid \text{DEP}_{\text{var}}(x, x', c); x = \text{aVar}_c(i); x' = \text{aVar}_c(j); i, j = 0, \dots, |\text{aVar}_c|\} \\ \text{Weights } W & := \{(x, n) \in \mathcal{VAR} \times \mathbb{N} \mid n = \text{vcnt}(x); x = \text{aVar}_c(i); i = 0, \dots, |\text{aVar}_c|\} \\ \text{Query Flags } \text{qF} & := \left\{ (x, n) \in V \times \{0, 1\} \mid \left\{ \begin{array}{ll} n = 1 & x \in \text{qVar}_c \\ n = 0 & \text{o.w.} \end{array} \right\}; x = \text{aVar}_c(i); i = 0, \dots, |\text{aVar}_c| \right\} \end{aligned}$$

Definition 12 (Finite Walk (k)). .

Given a labeled weighted graph $G = (V, E, W, qF)$, a finite walk k in G is a sequence of edges $(e_1 \dots e_{n-1})$ for which there is a sequence of vertices (v_1, \dots, v_n) such that:

- $e_i = (v_i, v_{i+1})$ for every $1 \leq i < n$.
- every vertex $v \in V$ appears in this vertices sequence (v_1, \dots, v_n) of k at most $W(v)$ times.

(v_1, \dots, v_n) is the vertex sequence of this walk.

Length of this finite walk k is the number of vertices in its vertex sequence, i.e., $\text{len}(k) = n$.

Given a labeled weighted graph $G = (V, E, W, qF)$, we use $\mathcal{WALK}(G)$ to denote a set containing all finite walks k in G ; and $k_{v_1 \rightarrow v_2} \in \mathcal{WALK}(G)$ where $v_1, v_2 \in V$ denotes the walk from vertex v_1 to v_2 .

Definition 13 (Length of Finite Walk w.r.t. Query (len_q)). .

Given a labeled weighted graph $G = (V, E, W, qF)$ and a finite walk k in G with its vertex sequence (v_1, \dots, v_n) , the length of k w.r.t query is defined as:

$$\text{len}_q(k) = \text{len}(v \mid v \in (v_1, \dots, v_n) \wedge F(v) = 1)$$

, where $(v \mid v \in (v_1, \dots, v_n) \wedge F(v) = 1)$ is a subsequence of k 's vertex sequence.

Given a program c with a starting memory m and database D , we generate its program-based graph $\mathbf{G}_{\text{trace}}(\mathbf{c}, \mathbf{m}, \mathbf{D}) = (V, E, W, qF)$. Then the adaptivity bound based on program analysis for \mathbf{c} is the number of query vertices on a finite walk in $\mathbf{G}_{\text{prog}}(\mathbf{c})$. This finite walk satisfies:

- the number of query vertices on this walk is maximum
- the visiting times of each vertex v on this walk is bound by its weight $W(v)$.

It is formally defined in 30.

Definition 14 (Adaptivity of A Program). .

Given a program \mathbf{c} in SSA language, its adaptivity is defined for all possible starting SSA memory \mathbf{m} and database D as follows:

$$A(c) = \max\{\text{len}_q(k) \mid \mathbf{m} \in \mathcal{SM}, D \in \mathcal{DB}, k \in \mathcal{WALK}(\mathbf{G}_{\text{trace}}(\mathbf{c}, \mathbf{m}, \mathbf{D}))\}$$

]

We proved some useful properties for our language.

[[

Definition 15 (Well-formed Trace). A trace t is well formed if and only if it preserves the following two properties:

- (Uniqueness) $\forall \mathbf{av}_1, \mathbf{av}_2 \in_{\text{av}} t. (\mathbf{av}_1 \neq_{\text{av}} \mathbf{av}_2)$
- (Ordering) $\forall \mathbf{av}_1, \mathbf{av}_2 \in_{\text{av}} t. (\mathbf{av}_1 <_{\text{av}} \mathbf{av}_2) \iff \exists t_1, t_2, t_3, \mathbf{av}'_1, \mathbf{av}'_2. \text{ s.t., } (\mathbf{av}_1 =_{\text{av}} \mathbf{av}'_1) \wedge (\mathbf{av}_2 =_{\text{av}} \mathbf{av}'_2) \wedge t_1 + +[\mathbf{av}'_1] + + t_2 + +[\mathbf{av}'_2] + + t_3 = t$

]]

[[

Theorem 2.1 (Variable Trace Generated from Operational Semantics is Well-formed). .

Given a program c , with arbitrary starting memory m , trace τ and variable counter vcnt if $\langle m, c, \tau, \text{vcnt} \rangle \rightarrow^* \langle m', \text{skip}, \tau', \text{vcnt}' \rangle$, then $(\tau' - \tau)$ is a well formed trace with respect to program c , m and w , denoted as $m, c \models \tau' - \tau$.

Proof. Proof in File: "thm_os_wf_trace.tex". □

]] [[

Lemma 1 (While Map Remains Unchanged (Invariant)). Given a program c with a starting memory m , trace t and while map w , s.t., $\langle m, c, t, w \rangle \rightarrow^* \langle m', \text{skip}, t', w' \rangle$ and $\text{Labels}(c) \cap \text{Keys}(w) = \emptyset$, then

$$w = w'$$

Proof of Lemma 8. Proof in File: "lem_wunchange.tex" ■

]] [[

Lemma 2 (Trace is Written Only). Given a program c with starting trace t_1 and t_2 , for arbitrary starting memory m and while map w , if there exist evaluations

$$\langle m, c, t_1, w \rangle \rightarrow^* \langle m'_1, \text{skip}, t'_1, w'_1 \rangle$$

$$\langle m, c, t_2, w \rangle \rightarrow^* \langle m'_2, \text{skip}, t'_2, w'_2 \rangle$$

then:

$$m'_1 = m'_2 \wedge w'_1 = w'_2$$

Proof of Lemma 9. Proof in File: "lem_twriteonly.tex" ■

]] [[

Lemma 3 (Trace Uniqueness). Given a program c with a starting memory m , [WQ: a while map w ,] for any starting trace t_1 and t_2 , if there exist evaluations

$$\langle m, c, t_1, w \rangle \rightarrow^* \langle m'_1, \text{skip}, t'_1, w'_1 \rangle$$

$$\langle m, c, t_2, w \rangle \rightarrow^* \langle m'_2, \text{skip}, t'_2, w'_2 \rangle$$

then:

$$t'_1 - t_1 = t'_2 - t_2$$

Proof of Lemma 10. Proof in File: "lem_tunique.tex" ■

]] [[

Corollary 2.1.1.

$$\text{av} \in_{\text{av}} t \implies \exists t_1, t_2, \text{av}' . \text{ s.t.}, (\text{av} =_{\text{av}} \text{av}') \wedge t_1 ++ [\text{av}'] ++ t_2 = t$$

Proof. Proof in File: "coro_aqintrace.tex" ■

]]

Lemma 4 (Trace Non-Decreasing). .

[JL: For any program c with a starting memory m , trace t and while map w :

$$\langle m, c, t, w \rangle \rightarrow \langle m, c', t', w' \rangle \implies \exists t'', \text{ s.t., } t ++ t'' = t'$$

]

Proof. Proof is obvious by induction on the operational semantic rules applied in the transition . By induction on the operational semantic rules applied in the transition $\langle m, c, t, w \rangle \rightarrow \langle m, c', t', w' \rangle$, we have cases for each rule. By observation on the rules, the trace t remains unchanged in all the rules except the only one **query-v**. So, the rule **query-v** is the only interesting case to be discussed as following.

• **case:**

$$\frac{\text{query}(\alpha) = v}{\langle m, [x \leftarrow \text{query}(\alpha)]^l, t, w \rangle \rightarrow \langle m, \text{skip}, t ++ [(\alpha, l, w)], w \rangle} \text{query-v}$$

In this case, we have $c' = \text{skip}$, $t' = t ++ [(\alpha, l, w)]$, $m' = m[v/x]$ and $w' = w$.

Let $t'' = [(\alpha, l, w)]$, we have $t ++ [(\alpha, l, w)] = t'$, i.e., $t ++ t'' = t'$. This case is proved.

□

[[The following lemma describes a property of the trace-based dependency graph. For any program c with a database D and a starting memory m , the directed edges in its trace-based dependency graph can only be constructed from nodes representing smaller annotated queries to annotated queries of greater order. There doesn't exist backward edges with direction from greater annotated queries to smaller ones.]]

Lemma 5. [Edges are Forwarding Only].

Given a program c , a database D , a starting memory m and the corresponding trace-based dependency graph $G(c, D, m) = (V, E)$, for any directed edge $(av', av) \in E$, this is not the case that:

$$av' \geq_{av} av$$

Proof. Proof in File: "edge_forward.tex".

□

Lemma 6. [Trace-based Dependency Graph is Directed Acyclic].

Every trace-based dependency graph is a directed acyclic graph.

Proof. Proof is obvious based on the Lemma 12.

□

Lemma 7 (Adaptivity is Bounded). .

Given the program c with a certain database D and starting memory m , the $A(c)$ w.r.t. the D and m is bounded, i.e.,:

$$\langle m, c, [], [] \rangle \rightarrow^* \langle m', \text{skip}, t', w' \rangle \implies A_{D,m}(c) \leq |t'|$$

Proof. Proof is obvious based on the Lemma 13.

□

3 Labeled SSA Language

3.1 The Limit of While Language

we see the power of the labelled loop language to achieve the adaptivity semantically, from its being capable to express many adaptive data analysis algorithm, allowing the construction of the query-based dependency graph using traces from the execution, and so on. However, it is not powerful enough to reach the adaptivity syntactically. The main difficulty is its implicit control flow which raises extra complexity to figure out where some variables used come from. We use three simple but relevant examples to show why the loop language suffers. We use $\text{query}(0), \text{query}(1)$ to represent linear queries.

$$\begin{array}{lll}
\begin{array}{l} [x \leftarrow \text{query}(0)]^1; \\ \text{if } [(x < 0)]^2 \\ c_1 = \text{then } [x \leftarrow \text{query}(1)]^3 \\ \text{else } [\text{skip}]^4; \\ [y \leftarrow \text{query}(x + \chi(3))]^5 \end{array} &
\begin{array}{l} [x \leftarrow \text{query}(0)]^1; \\ \text{if } [(x < 0)]^2 \\ c_2 = \text{then } [x \leftarrow \text{query}(1)]^3 \\ \text{else } [x \leftarrow \text{query}(2)]^4; \\ [y \leftarrow \text{query}(x + \chi(3))]^5 \end{array} &
\begin{array}{l} [x \leftarrow \text{query}(0)]^1; \\ \text{if } [(x < 0)]^2 \\ c_3 = \text{then } [z \leftarrow \text{query}(1)]^3 \\ \text{else } [\text{skip}]^4; \\ [y \leftarrow \text{query}(x + \chi(3))]^5 \end{array}
\end{array}$$

In these three examples, the variable x at line 5 is implicit. In program c_1 , it refers to the either x at line 1, or x at line 3, which means the result of query request $\text{query}(x + \chi(3))$ assigned to the variable y may depend on $\text{query}(0)$ (bound to x at line 1) or $\text{query}(1)$ (x at line 3). When we have a look at the other two programs c_2 and c_3 , it is another talk. We think $\text{query}(x + \chi(3))$ may depend on either $\text{query}(1)$ (x at line 3) or $\text{query}(x + \chi(3))$ (x at line 4) in c_2 , while $\text{query}(x + \chi(3))$ only depends on $\text{query}(0)$ at line 1 in program c_3 . These three examples are structural similar in loop language, however, the dependency between variables are quite dissimilar. We consider variables here because query request is also bound to variables. To solve this dilemma, we move to single static assignment as follows.

$$\begin{array}{lll}
\begin{array}{l} [\mathbf{x}_1 \leftarrow \text{query}(0)]^1; \\ \text{if } [(\mathbf{x}_1 < 0)]^2 \\ c_1^{ssa} = ([], [\mathbf{x}_3, \mathbf{x}_1, \mathbf{x}_2], []) \\ \text{then } [\mathbf{x}_2 \leftarrow \text{query}(1)]^3 \\ \text{else } [\text{skip}]^4; \\ [\mathbf{y}_1 \leftarrow \text{query}(\mathbf{x}_3 + \chi(3))]^5 \end{array} &
\begin{array}{l} [\mathbf{x}_1 \leftarrow \text{query}(0)]^1; \\ \text{if } [(\mathbf{x}_1 < 0)]^2, \\ c_2^{ssa} = ([\mathbf{x}_4, \mathbf{x}_2, \mathbf{x}_3], [], []) \\ \text{then } [\mathbf{x}_2 \leftarrow \text{query}(1)]^3 \\ \text{else } [\mathbf{x}_3 \leftarrow \text{query}(2)]^4; \\ [\mathbf{y}_1 \leftarrow \text{query}(\mathbf{x}_4 + \chi(3))]^5 \end{array} &
\begin{array}{l} [\mathbf{x}_1 \leftarrow \text{query}(0)]^1; \\ \text{if } [(\mathbf{x}_1 < 0)]^2 \\ c_3^{ssa} = ([], [], []) \\ \text{then } [\mathbf{z}_1 \leftarrow \text{query}(1)]^3 \\ \text{else } [\text{skip}]^4; \\ [\mathbf{y}_1 \leftarrow \text{query}(\mathbf{x}_1 + \chi(3))]^5 \end{array}
\end{array}$$

To distinguish between the loop language and in ssa form, we denote the ssa variable \mathbf{x}_1 in bold. As we can see, the data flow becomes explicit in ssa form and the analysis on the dependency between variables in the program becomes much clear now. Considering this advantage, we aim to estimate the adaptivity through an analysis on program in ssa form.

3.2 SSA form Language

| | | |
|---------------------------|---------------|---|
| Arithmetic Operators | \oplus_a | $::= + \mid - \mid \times \mid \div$ |
| Boolean Operators | \oplus_b | $::= \vee \mid \wedge$ |
| Relational Operators | \sim | $::= < \mid \leq \mid ==$ |
| Label | l | $::= \mathbb{N}$ |
| SSA Arithmetic Expression | \mathbf{a} | $::= n \mid \mathbf{x} \mid \mathbf{a} \oplus_a \mathbf{a}$ |
| SSA Boolean Expression | \mathbf{b} | $::= \text{true} \mid \text{false} \mid \neg \mathbf{b} \mid \mathbf{b} \oplus_b \mathbf{b} \mid \mathbf{a} \sim \mathbf{a}$ |
| SSA Query Expression | ψ | $::= \alpha \mid \mathbf{a} \mid \psi \oplus_a \psi$ |
| Query Value | α | $::= n \mid \chi[n] \mid \chi[n] \oplus_a \chi[n] \mid n \oplus_a \chi[n] \mid \chi[n] \oplus_a n$ |
| Value | v | $::= n \mid \text{true} \mid \text{false} \mid [] \mid [v, \dots, v]$ |
| SSA Expression | \mathbf{e} | $::= v \mid \mathbf{a} \mid \mathbf{b} \mid [e, \dots, e]$ |
| Labeled SSA Command | \mathbf{c} | $::= [\mathbf{x} \leftarrow \mathbf{e}]^l \mid [\mathbf{x} \leftarrow \text{query}(\psi)]^l \mid \text{ifvar}(\bar{\mathbf{x}}, \bar{\mathbf{x}}') \mid$ $\text{while } [\mathbf{b}]^l, n, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] \text{ do } \mathbf{c}$ $\mid \mathbf{c}; \mathbf{c} \mid [\text{if } (\mathbf{b}, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], \mathbf{c}, \mathbf{c})]^l \mid [\text{skip}]^l$ |
| Event | ϵ | $::= (\mathbf{x}, v, l, n) \mid (\mathbf{x}, \alpha, l, n) \mid (\mathbf{b}, v, l, n)$ |
| Trace | τ | $::= [] \mid \epsilon :: \tau$ |
| Environment | θ | $::= \tau \rightarrow v$ |
| Variable Counter | vcnt | $::= \mathcal{SVAR} \rightarrow \mathbb{N}$ |

We use following notations to represent the set of corresponding terms:

| | | |
|--------------------------|---|-------------------------|
| \mathcal{SVAR} | : | Set of Variables |
| \mathcal{VAL} | : | Set of Values |
| \mathcal{E} | : | Set of Events |
| \mathcal{SM} | : | Set of SSA Memories |
| \mathcal{DB} | : | Set of Databases |
| $\mathcal{QD} = [-1, 1]$ | : | Domain of Query Results |

Consistences to the While Language Each command is labeled with a label l , a natural number standing for the line of code where the command appears. Notice that we associate the label l to the conditional predicate b in the if statement, and to the guard b in the while statement.

A memory is standard, a map from labeled variables to values.

The variable counter vcnt maps every variable to a natural number $n \in \mathbb{N}$ in a certain execution of program \mathbf{c} . This natural number n represents the visiting times of this variable in this certain execution. We use variable name x within parenthesis to denote the access to the associated natural number of this variable in the variable counter vcnt_c , $\text{vcnt}_c(x)$ denote the visiting times of variable x .

The annotated variable / event $\text{av} = (x, l, n, v)$ or $\text{av} = (x, l, n, \alpha)$ is a quaternary tuple contains 4 elements. α is a query value representing the corresponding query request $x \leftarrow \text{query}(\alpha)$ during the execution of the program.

A variable trace τ is a list of annotated queries accumulated along the execution of the program. A trace can be regarded as the program history, where this history consists of all the queries asked by the analyst during the execution of the program. **[[Given an annotated variable av and a trace t , the appending operation $\text{av} :: t$ is the standard list appending operation, appends av to the head of trace t . The concatenation operation between 2 traces t_1 and t_2 , i.e., $t_1 ++ t_2$ is the standard list concatenation operation as follows:**

$$t_1 ++ t_2 \triangleq \begin{cases} t_2 & t_1 = [] \\ \text{av} :: (t'_1 ++ t_2) & t_1 = \text{av} :: t'_1 \end{cases} \quad (6)$$

The subtraction operation between 2 traces t_1 and t_2 , i.e., $t_1 - t_2$ is defined as follows:

$$t_1 - t_2 \triangleq t_3 \text{ s.t., } t_2 ++ t_3 = t_1 \quad (7)$$

]] A configuration, $\langle \mathbf{m}, \mathbf{c}, \tau, \text{vcnt} \rangle$, contains four elements: a SSA memory \mathbf{m} , the command \mathbf{c} to be evaluated, a trace t and variable vcnt .

We collect the trace with a trace-based small-step operational semantics based on transitions of the program configuration $\langle m, c, \tau, \text{vcnt} \rangle$, of form $\langle m, c, \tau, \text{vcnt} \rangle \rightarrow \langle m', \text{skip}, \tau', \text{vcnt}' \rangle$.

Differences in SSA form Language We use \mathbf{a} to express arithmetic expressions which now contains ssa variable $\mathbf{x} \in \mathcal{SVAR}$, and the boolean expression as \mathbf{b} . The ssa expression can be either \mathbf{a} and \mathbf{b} . We also have the ssa variables annotated in a similar way as the annotated queries in the while language. The labeled commands \mathbf{c} are now in the ssa form. In the assignment command $[\mathbf{x} \leftarrow \mathbf{e}]^l$ and query request command $[\mathbf{x} \leftarrow \text{query}(\psi)]^l$, the expression \mathbf{e} and query expression ψ is now in their corresponding ssa forms.

The if command now contains the extra part $([\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2])$, which helps to track the dependency of new assigned variables in both branches $([\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2])$, then branch $[\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2]$, and else branch $[\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2]$. The $\bar{\mathbf{x}}$ is a list of ssa variables, in which every element \mathbf{x} may depends on the corresponding element \mathbf{x}_1 from $\bar{\mathbf{x}}_1$ collected in the then branch or the corresponding element \mathbf{x}_2 from $\bar{\mathbf{x}}_2$ collected in the else branch. Every tuple $(\mathbf{x}, \mathbf{x}_1, \mathbf{x}_2)$ from $[\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]$ can be understood as $\mathbf{x} = \phi(\mathbf{x}_1, \mathbf{x}_2)$ in the normal ssa form. The previous example c_2^2 can be used for reference. The second part $[\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2]$ focuses on the then branch. The list of ssa variables \mathbf{y}_1 stores the assigned ssa variables before the if command, whose non-ssa version (variables in the while language) will be modified only in the then branch. We can look at program c_1 as a reference, in which x at line 1 may be modified only in the then branch at line 3. The list $\bar{\mathbf{y}}_2$ tracks the ssa variables assigned only in the then branch. If the variables are assigned in both branches such as in the program c_2 , they goes into $[\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]$. Then we think every ssa variable in $\bar{\mathbf{y}}$ may come from the corresponding variable \mathbf{y}_1 in $\bar{\mathbf{y}}_1$ before the if command or \mathbf{y}_2 in $\bar{\mathbf{y}}_2$ in the then branch. In this sense, we can also regard every tuple $(\mathbf{y}, \mathbf{y}_1, \mathbf{y}_2)$ from $[\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2]$ as $\mathbf{y} = \phi(\mathbf{y}_1, \mathbf{y}_2)$. The rest part $[\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2]$ focus on the else branch and can be understood similarly.

The while command also has similar part $[\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]$, focusing on the while body. The new command `ifvar`($\bar{\mathbf{x}}, \bar{\mathbf{x}}'$) does not have explicit label because it is only used for evaluation internally, we will discuss more about it when in the small-step operational semantics for SSA language.

The SSA memory \mathbf{m} is a map from SSA variables \mathbf{x} to values.

The others remain the same in SSA form language as in the `While` language.

3.3 Trace-based Operational Semantics for SSA Language

The small-step transition states that a configuration $\langle m, c, \tau, \text{vcnt} \rangle$ evaluates to another configuration with the trace and while map updated along with the evaluation of the command c to the normal form of the command `skip`. We define rules of the trace-based operational semantics in Figure 1. The rule **query-e** evaluates the argument of a query request. When the argument is in normal form, this query will be answered. The rule **query-v** modifies the starting memory m to $m[\alpha/x]$ using the answer α of the query $\text{query}(\alpha)$ from the mechanism, with the trace expanded by appending the query $\text{query}(\alpha)$ with the current annotation (l, w) . The rule for assignment is standard and the trace remains unchanged. The sequence rule keeps tracking the modification of the trace, and the evaluation rule for if conditional goes into one branch based on the result of the conditional predicate b . The rules for

while modify the while map w . In the rule **ifw-true**, the while map w is updated by $w + l$ because the execution goes into another iteration when the condition $n > 0$ is satisfied. When n reaches 0, the loop exits and the while map w eliminates the label l of this while statement by $w - l$ in the rule **ifw-false**. With the operational semantics and relations between annotated queries, we restrict the well-formed trace w.r.t. the execution of a program c in Definition 23. When switching to the SSA language, we show that we are still able to achieve what we can get in Section 2. The operational semantics of the SSA language mimics its counterpart, of the form $\langle \mathbf{m}, \mathbf{c}, t, w \rangle \rightarrow \langle \mathbf{m}', \text{skip}, t', w' \rangle$. The SSA memory \mathbf{m} is a map from SSA variable \mathbf{x} to values. It still uses a trace to track the query requests during the execution, starting from an SSA configuration with an SSA memory \mathbf{m} and a program in its SSA form \mathbf{c} , which allows a similar construction of the query-based dependency graph in the SSA language as in the `While` language. We show the evaluation rules in Figure 2. The command `ifvar(\bar{x}, \bar{x}')` stores the variable map during the run time, which is a map from ssa variable \mathbf{x} to variable \mathbf{x}_i . This map is designed for `if` command, when the variable may comes from two branches and this command records which branch the variable comes from. The key idea underneath the operational semantics is to have the trace and the execution path being constructed in a similar way as in the loop language. Take the query request as an example, the argument \mathbf{e} which contains ssa variables will be evaluated to a value v first before the request is sent to the database in rule **ssa-query-arg**. The trace expands in the rule **ssa-query** likewise in the loop language. The query q , a primitive symbol representing the abstract query in both the ssa language and the loop language, makes no difference in the two languages. Since we add the extra part $[\bar{x}, \bar{x}_1, \bar{x}_2], [\bar{y}, \bar{y}_1, \bar{y}_2], [\bar{z}, \bar{z}_1, \bar{z}_2]$ in the `if` command compared to its counterpart in the while language, the rules relevant to the `if` condition (**ssa-if-t** and **ssa-if-f**) use the command `ifvar(\bar{x}, \bar{x}')` to update the ssa memory \mathbf{m} with the mapping from the new generated variable \mathbf{x} in \bar{x} to the appropriate value $\mathbf{m}(\mathbf{x}')$ where \mathbf{x}' is the corresponding variable w.r.t \mathbf{x} in \bar{x}' . The rule **ssa-ifvar** reflects the usage of `ifvar(\bar{x}, \bar{x}')`. It is easier to understand the usage of `ifvar(\bar{x}, \bar{x}')` in the rule **ssa-if-t** when we think about how ssa works: in the ssa form, when a variable to be used may come from two sources (e.g. \mathbf{x}_1 and \mathbf{x}_2 in the rule), it generates a new variable \mathbf{x} , assigning it with $\phi(\mathbf{x}_1, \mathbf{x}_2)$, and replaces the variable to be used with newly assigned \mathbf{x} . We know that in the future program after the `if` command, only the variables \bar{x} will be available instead of \bar{x}_1, \bar{x}_2 from two branches. For the evaluation of the program after the `if` command, we need to tell the memory the exact value of the newly generated variable \mathbf{x} , which is the value stored in \mathbf{x}_1 when the conditional \mathbf{b} is true, or the value in \mathbf{x}_2 when \mathbf{b} is false. To this end, the internal command `ifvar(\bar{x}, \bar{x}')` plays its role. For the `if` rule, we need to instantiate the variables from \bar{x} whose values come from two branches, \bar{y} whose values from then branch or assignment before the `if` command, and \bar{z} whose values from else branch or before the `if` command. Correspondingly, we need to have three `ifvar` commands.

The evaluation of while depends on the while iteration counter \mathbf{n} and the guard \mathbf{b} . When \mathbf{b} is evaluated to `true`, the while is still executing, and all the variables \mathbf{x} in \bar{x} of the loop body \mathbf{c} are replaced as the corresponding variables in \bar{x}_1 in the first iteration($n = 0$), or \bar{x}_2 in other iterations($n > 0$). The while turns to an exit when $\mathbf{n} > 0$, and the memory \mathbf{m} updates the mapping of variables in \bar{x} with \bar{x}_1 if the guard \mathbf{b} evaluates to `false`, which means the while body is not executed once. When the while enters the exit after executing the body a few times(n), the variables in \bar{x} is instantiated with the value from the body $\mathbf{m}(\bar{x}_2)$.

3.4 Event and Trace

[Events preserve the order of execution. The order relation is defined in Definition 16.

Definition 16 (Order of Annotated Variables / Events). .

[JL:

$$\boxed{Memory \times Command \times VTrace \times VCounter \rightarrow Memory \times Command \times VTrace \times VCounter}$$

$$\boxed{\langle \mathbf{m}, \mathbf{c}, \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}', \mathbf{c}', \tau', \text{vcnt}' \rangle}$$

$$\frac{\langle \mathbf{m}, \mathbf{e} \rangle \rightarrow \mathbf{e}'}{\langle \mathbf{m}, [\mathbf{x} \leftarrow \mathbf{e}]^l, \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}, [\mathbf{x} \leftarrow \mathbf{e}']^l, \tau, \text{vcnt} \rangle} \text{ assn-e}$$

$$\frac{\text{vcnt}'(\mathbf{x}) = \text{vcnt}(\mathbf{x}) + 1 \quad \text{av} = (\mathbf{x}, l, \text{vcnt}'(\mathbf{x}), v)}{\langle \mathbf{m}, [\mathbf{x} \leftarrow v]^l, \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}[\nu/\mathbf{x}], [\text{skip}]^l, \tau + +[\text{av}], \text{vcnt}' \rangle} \text{ assn-v}$$

$$\frac{}{\langle \mathbf{m}, \text{while } [\mathbf{b}]^l, \mathbf{n}, [\bar{\mathbf{x}}', \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] \text{ do } \mathbf{c}, \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}, \text{if}_w(\mathbf{b}, \mathbf{n}, [\bar{\mathbf{x}}_1/\bar{\mathbf{x}}'], [\bar{\mathbf{x}}', \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], \mathbf{c}[\bar{\mathbf{x}}_1/\bar{\mathbf{x}}']; \text{while } [\mathbf{b}]^l, (\mathbf{n} + 1), [\bar{\mathbf{x}}', \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] \text{ do } \mathbf{c}, \text{skip}), \tau, \text{vcnt} \rangle} \text{ ssa-while-b}$$

$$\frac{\langle \mathbf{m}, \mathbf{b} \rangle \rightarrow \mathbf{b}'}{\langle \mathbf{m}, \text{if}_w(\mathbf{b}, [\bar{\mathbf{x}}', \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], \mathbf{n}, \mathbf{c}_1, \mathbf{c}_2), \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}, \text{if}_w(\mathbf{b}', [\bar{\mathbf{x}}', \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], \mathbf{n}, \mathbf{c}_1, \mathbf{c}_2), \tau, \text{vcnt} \rangle} \text{ ssa-ifw-b}$$

$$\frac{}{\langle \mathbf{m}, \text{if}_w(\text{true}, [\bar{\mathbf{x}}', \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], \mathbf{n}, \mathbf{c}; \text{while } [\mathbf{b}]^l, \mathbf{n}, [\bar{\mathbf{x}}', \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] \text{ do } \mathbf{c}, \text{skip}), \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}, \mathbf{c}; \text{while } [\mathbf{b}]^l, \mathbf{n}, [\bar{\mathbf{x}}', \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] \text{ do } \mathbf{c}, \tau, \text{vcnt} \rangle} \text{ ssa-ifw-true}$$

$$\frac{n = 0 \rightarrow i = 1 \quad n > 0 \rightarrow i = 2}{\langle \mathbf{m}, \text{if}_w(\text{false}, [\bar{\mathbf{x}}', \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], \mathbf{n}, \mathbf{c}; \text{while } [\mathbf{b}]^l, \mathbf{n}, [\bar{\mathbf{x}}', \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] \text{ do } \mathbf{c}, \text{skip}), \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}, \text{skip}; \text{ifvar}(\bar{\mathbf{x}}', \bar{\mathbf{x}}_i), \tau, \text{vcnt} \rangle} \text{ ssa-ifw-false}$$

$$\frac{\langle \mathbf{m}, \psi \rangle \rightarrow_q \psi'}{\langle \mathbf{m}, [\mathbf{x} \leftarrow \text{query}(\psi)]^l, \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}, [\mathbf{x} \leftarrow \text{query}(\psi')]^l, \tau, \text{vcnt} \rangle} \text{ ssa-query-e}$$

$$\frac{\text{query}(\alpha) = v \quad \text{vcnt}'[\mathbf{x}] = \text{vcnt}[\mathbf{x}] + 1 \quad \text{av} = (\mathbf{x}, \alpha, l, \text{vcnt}'[\mathbf{x}])}{\langle \mathbf{m}, [\mathbf{x} \leftarrow \text{query}(\alpha)]^l, \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}[\nu/\mathbf{x}], \text{skip}, \tau + +[\text{av}], \text{vcnt}' \rangle} \text{ ssa-query-v}$$

$$\frac{\langle \mathbf{m}, \mathbf{c}_1, \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}', \mathbf{c}'_1, \tau', \text{vcnt}' \rangle}{\langle \mathbf{m}, \mathbf{c}_1; \mathbf{c}_2, \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}', \mathbf{c}'_1; \mathbf{c}_2, \tau', \text{vcnt}' \rangle} \text{ ssa-seq1} \quad \frac{}{\langle \mathbf{m}, [\text{skip}]^l; \mathbf{c}_2, \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}, \mathbf{c}_2, \tau, \text{vcnt} \rangle} \text{ ssa-seq2}$$

$$\frac{\langle \mathbf{m}, \mathbf{b} \rangle \rightarrow_b \mathbf{b}'}{\langle \mathbf{m}, \text{if } ([\mathbf{b}]^l, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2], \mathbf{c}_1, \mathbf{c}_2), \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}, [\text{if } ([\mathbf{b}']^l, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2], \mathbf{c}_1, \mathbf{c}_2)]^l, \tau, \text{vcnt} \rangle} \text{ ssa-if-b}$$

$$\frac{}{\langle \mathbf{m}, \text{if } ([\text{true}]^l, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2], \mathbf{c}_1, \mathbf{c}_2), \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}, \mathbf{c}_1; \text{ifvar}(\bar{\mathbf{x}}, \bar{\mathbf{x}}_2); \text{ifvar}(\bar{\mathbf{y}}, \bar{\mathbf{y}}_1); \text{ifvar}(\bar{\mathbf{z}}, \bar{\mathbf{z}}_2), \tau, \text{vcnt} \rangle} \text{ ssa-if-t}$$

$$\frac{}{\langle \mathbf{m}, \text{if } ([\text{false}]^l, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2], \mathbf{c}_1, \mathbf{c}_2), \tau, \text{vcnt} \rangle \rightarrow \langle \mathbf{m}, \mathbf{c}_2; \text{ifvar}(\bar{\mathbf{x}}, \bar{\mathbf{x}}_2); \text{ifvar}(\bar{\mathbf{y}}, \bar{\mathbf{y}}_1); \text{ifvar}(\bar{\mathbf{z}}, \bar{\mathbf{z}}_2), \tau, \text{vcnt} \rangle} \text{ ssa-if-f}$$

$$\frac{}{\langle \mathbf{m}, \text{ifvar}(\bar{\mathbf{x}}, \bar{\mathbf{x}}'), \tau, \text{vcnt} \rangle \rightarrow \langle (\bar{\mathbf{x}} \rightarrow \mathbf{m}(\bar{\mathbf{x}}')) :: \mathbf{m}, \text{skip}, \tau, \text{vcnt} \rangle} \text{ ssa-ifvar}$$

]

Given 2 annotated queries $\text{av}_1 = (x_1, v_1, l_1, n_1), \text{av}_2 = (x_2, v_2, l_2, n_2) :$

$$\text{av}_1 <_{\text{av}} \text{av}_2 \triangleq \begin{cases} n_1 < n_2 & l_1 = l_2 \\ w_1 <_w w_2 & \text{o.w.} \end{cases}$$

$\text{av}_1 \geq_{\text{av}} \text{av}_2$ is defined vice versa.

]] Given the evaluation rules for query expression, we define its equivalence relation in Definition 17.

Definition 17 (Equivalence of Query). . Given a memory m and 2 query expressions ψ_1, ψ_2 s.t., $FV(\psi_1) \in \text{dom}(m)$ and $FV(\psi_2) \in \text{dom}(m)$:

$$\psi_1 =_q^m \psi_2 \triangleq \begin{cases} \text{true} & \exists \alpha_1, \alpha_2. (\langle m, \psi_1 \rangle \rightarrow_q \alpha_1 \wedge \langle m, \psi_2 \rangle \rightarrow_q \alpha_2) \\ & \wedge (\forall r \in \mathcal{QD}. \exists v. \text{ s.t., } \langle m, \alpha_1[r/\chi] \rangle \rightarrow_a v \wedge \langle m, \alpha_2[r/\chi] \rangle \rightarrow_a v) \\ \text{false} & \text{o.w.} \end{cases}$$

, where $FV(\psi)$ is the set of free variables in the query expression ψ . $\psi_1 \neq_q^m \psi_2$ is defined vice versa. We use $=_q$ and \neq_q as the shorthands for $=_q^\square$ and \neq_q^\square .

Then, we have the corresponding equivalence relation between 2 annotated queries defined in Definition 18:]]

Definition 18 (Equivalence of Annotated Variables / Events). Given 2 annotated queries $\text{av}_1 = (x_1, v_1, l_1, n_1), \text{av}_2 = (x_2, v_2, l_2, n_2) :$

$$\text{av}_1 =_{\text{av}} \text{av}_2 \triangleq (l_1 = l_2 \wedge w_1 =_w w_2 \wedge \alpha_1 =_q \alpha_2)$$

$\text{av}_1 \neq_{\text{av}} \text{av}_2$ is defined vice versa.

]] [[Given an annotated variable av and a trace t , the appending operation $\text{av} :: t$ is the standard list appending operation, appends av to the head of trace t . The concatenation operation between 2 traces t_1 and t_2 , i.e., $t_1 ++ t_2$ is the standard list concatenation operation as follows:

$$t_1 ++ t_2 \triangleq \begin{cases} t_2 & t_1 = [] \\ \text{av} :: (t'_1 ++ t_2) & t_1 = \text{av} :: t'_1 \end{cases} \quad (8)$$

The subtraction operation between 2 traces t_1 and t_2 , i.e., $t_1 - t_2$ is defined as follows:

$$t_1 - t_2 \triangleq t_3 \text{ s.t., } t_2 ++ t_3 = t_1 \quad (9)$$

Given an annotated query av , av belongs to a trace t , i.e., $\text{av} \in_{\text{av}} t$ are defined as follows:

$$\text{av} \in_{\text{av}} t \triangleq \begin{cases} \text{false} & t = [] \\ \text{true} & t = \text{av}' :: t' \quad \text{av} =_{\text{av}} \text{av}' \\ \text{av} \in t' & t = \text{av}' :: t' \quad \text{av} \neq_{\text{av}} \text{av}' \end{cases} \quad (10)$$

]] [[

Definition 19 (Equivalence of Program). Given 2 programs c_1 and c_2 :

$$c_1 =_c c_2 \triangleq \begin{cases} \text{true} & c_1 = \text{skip} \wedge c_2 = \text{skip} \\ \forall m. \exists v. \langle m, e_1 \rangle \rightarrow_a^* v \wedge \langle m, e_1 \rangle \rightarrow_a^* v & c_1 = x \leftarrow e_1 \wedge c_2 = x \leftarrow e_2 \\ \psi_1 =_q \psi_2 & c_1 = x \leftarrow \text{query}(\psi_1) \wedge c_2 = x \leftarrow \text{query}(\psi_2) \\ c_1^f =_c c_2^f \wedge c_1^t =_c c_2^t & c_1 = \text{if } (b, c_1^f, c_1^t) \wedge c_2 = \text{if } (b, c_2^f, c_2^t) \\ c_1^i =_c c_2^i & c_1 = \text{while } b \text{ do } c_1^i \wedge c_2 = \text{while } b \text{ do } c_2^i \\ c_1^h =_c c_2^h \wedge c_1^t =_c c_2^t & c_1 = c_1^h; c_1^t \wedge c_2 = c_2^h; c_2^t \end{cases}$$

$c_1 \neq_c c_2$ is defined vice versa.

Given 2 programs c and c' , c' is a sub-program of c , i.e., $c' \in_c c$ is defined as:

$$c' \in_c c \triangleq \exists c_1, c_2, c''. \text{ s.t., } c =_c c_1; c''; c_2 \wedge c' =_c c'' \quad (11)$$

]] [[

Definition 20 (Assigned Variables (aVar)). Given a program \mathbf{c} , its assigned variables aVar is a vector containing all variables newly assigned in the program preserving the order. It is defined as follows:

$$\text{aVar}_{\mathbf{c}} \triangleq \begin{cases} [\mathbf{x}] & \mathbf{c} = [\mathbf{x} \leftarrow \mathbf{e}]^{(l,w)} \\ [\mathbf{x}] & \mathbf{c} = [\mathbf{x} \leftarrow \text{query}(\boldsymbol{\psi})]^{(l,w)} \\ \text{aVar}_{\mathbf{c}_1} ++ \text{aVar}_{\mathbf{c}_2} & \mathbf{c} = \mathbf{c}_1; \mathbf{c}_2 \\ \text{aVar}_{\mathbf{c}_1} ++ \text{aVar}_{\mathbf{c}_2} ++ [\bar{\mathbf{x}}, \bar{\mathbf{y}}, \bar{\mathbf{z}}] & \mathbf{c} = \text{if}([\mathbf{b}]^{(l,w)}, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_2, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_2, \bar{\mathbf{y}}_3], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_2, \bar{\mathbf{z}}_3], \mathbf{c}_1, \mathbf{c}_2) \\ \text{aVar}_{\mathbf{c}'} ++ [\bar{\mathbf{x}}] & \mathbf{c} = \text{while}([\mathbf{b}]^{(l,w)}, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_2, \bar{\mathbf{x}}_2], \mathbf{c}') \end{cases}$$

Definition 21 (Query Variables (qVar)). .

Given a program c , its query variables qVar is a vector containing all variables newly assigned by a query in the program, $\text{qVar} \subset \mathcal{VAR}$. It is defined as follows:

$$\text{qVar}_{\mathbf{c}} \triangleq \begin{cases} [] & \mathbf{c} = [\mathbf{x} \leftarrow \mathbf{e}]^{(l,w)} \\ [\mathbf{x}] & \mathbf{c} = [\mathbf{x} \leftarrow \text{query}(\boldsymbol{\psi})]^{(l,w)} \\ \text{aVar}_{\mathbf{c}_1} ++ \text{aVar}_{\mathbf{c}_2} & \mathbf{c} = \mathbf{c}_1; \mathbf{c}_2 \\ \text{aVar}_{\mathbf{c}_1} ++ \text{aVar}_{\mathbf{c}_2} ++ [\bar{\mathbf{x}}, \bar{\mathbf{y}}, \bar{\mathbf{z}}] & \mathbf{c} = \text{if}([\mathbf{b}]^{(l,w)}, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_2, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_2, \bar{\mathbf{y}}_3], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_2, \bar{\mathbf{z}}_3], \mathbf{c}_1, \mathbf{c}_2) \\ \text{aVar}_{\mathbf{c}'} ++ [\bar{\mathbf{x}}] & \mathbf{c} = \text{while}([\mathbf{b}]^{(l,w)}, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_2, \bar{\mathbf{x}}_2], \mathbf{c}') \end{cases}$$

We are abusing the notations and operators from list here. The notation $[]$ represents an empty vector and $x :: A$ represents add an element x to the head of the vector A . The concatenation operation between 2 vectors A_1 and A_2 , i.e., $A_1 ++ A_2$ is mimic the standard list concatenation operations as follows:

$$A_1 ++ A_2 \triangleq \begin{cases} A_2 & A_1 = [] \\ x :: (A'_1 ++ A_2) & A_1 = x :: A'_1 \end{cases} \quad (12)$$

We use index within parenthesis to denote the access to the element of corresponding location, $A(i)$ denotes the element at location i in the vector A and $M(i, j)$ denotes the element at location i -th row, j -th column in the matrix M .

Definition 22 (Initial Variable Counter vcnt_c^0). Given a program c with its assigned variables aVar_c of length N , its initial variable counter vcnt_c^0 maps all the variable to 0, i.e.:

$$\text{vcnt}_c^0(x) = 0, x = \text{aVar}_c(i) \forall i = 1, \dots, N$$

We define some properties and prove lemmas for trace and event w.r.t. the operational semantics as follows. [[

Definition 23 (Well-formed Trace). A trace t is well formed if and only if it preserves the following two properties:

- (Uniqueness) $\forall \text{av}_1, \text{av}_2 \in_{\text{av}} t. (\text{av}_1 \neq_{\text{av}} \text{av}_2)$
- (Ordering) $\forall \text{av}_1, \text{av}_2 \in_{\text{av}} t. (\text{av}_1 <_{\text{av}} \text{av}_2) \iff \exists t_1, t_2, t_3, \text{av}'_1, \text{av}'_2. \text{ s.t., } (\text{av}_1 =_{\text{av}} \text{av}'_1) \wedge (\text{av}_2 =_{\text{av}} \text{av}'_2) \wedge t_1 ++ [\text{av}'_1] ++ t_2 ++ [\text{av}'_2] ++ t_3 = t$

]]
[[

Theorem 3.1 (Variable Trace Generated from Operational Semantics is Well-formed). .

Given a program c , with arbitrary starting memory m , trace τ and variable counter vcnt if $\langle m, c, \tau, \text{vcnt} \rangle \rightarrow^* \langle m', \text{skip}, \tau', \text{vcnt}' \rangle$, then $(\tau' - \tau)$ is a well formed trace with respect to program c , m and w , denoted as $m, c \models \tau' - \tau$.

Proof. Proof in File: "thm_os_wf_trace.tex". □

]] [[

Lemma 8 (While Map Remains Unchanged (Invariant)). Given a program c with a starting memory m , trace t and while map w , s.t., $\langle m, c, t, w \rangle \rightarrow^* \langle m', \text{skip}, t', w' \rangle$ and $\text{Labels}(c) \cap \text{Keys}(w) = \emptyset$, then

$$w = w'$$

Proof of Lemma 8. Proof in File: "lem_wunchange.tex" ■

]] [[

Lemma 9 (Trace is Written Only). Given a program c with starting trace t_1 and t_2 , for arbitrary starting memory m and while map w , if there exist evaluations

$$\langle m, c, t_1, w \rangle \rightarrow^* \langle m'_1, \text{skip}, t'_1, w'_1 \rangle$$

$$\langle m, c, t_2, w \rangle \rightarrow^* \langle m'_2, \text{skip}, t'_2, w'_2 \rangle$$

then:

$$m'_1 = m'_2 \wedge w'_1 = w'_2$$

Proof of Lemma 9. Proof in File: "lem_twriteonly.tex" ■

]] [[

Lemma 10 (Trace Uniqueness). Given a program c with a starting memory m , [WQ: a while map w ,] for any starting trace t_1 and t_2 , if there exist evaluations

$$\langle m, c, t_1, w \rangle \rightarrow^* \langle m'_1, \text{skip}, t'_1, w'_1 \rangle$$

$$\langle m, c, t_2, w \rangle \rightarrow^* \langle m'_2, \text{skip}, t'_2, w'_2 \rangle$$

then:

$$t'_1 - t_1 = t'_2 - t_2$$

Proof of Lemma 10. Proof in File: "lem_tunique.tex" ■

]] [[

Corollary 3.1.1.

$$\text{av} \in_{\text{av}} t \implies \exists t_1, t_2, \text{av}'. \text{ s.t., } (\text{av} =_{\text{av}} \text{av}') \wedge t_1 ++ [\text{av}'] ++ t_2 = t$$

Proof. Proof in File: "coro_aqintrace.tex" ■

]]

Lemma 11 (Trace Non-Decreasing). .

[JL: For any program c with a starting memory m , trace t and while map w :

$$\langle m, c, t, w \rangle \rightarrow \langle m, c', t', w' \rangle \implies \exists t'', \text{ s.t., } t + t'' = t'$$

/

Proof. Proof is obvious by induction on the operational semantic rules applied in the transition .

By induction on the operational semantic rules applied in the transition $\langle m, c, t, w \rangle \rightarrow \langle m, c', t', w' \rangle$, we have cases for each rule. By observation on the rules, the trace t remains unchanged in all the rules except the only one **query-v**. So, the rule **query-v** is the only interesting case to be discussed as following.

• **case:**

$$\frac{\text{query}(\alpha) = v}{\langle m, [x \leftarrow \text{query}(\alpha)]^l, t, w \rangle \rightarrow \langle m, \text{skip}, t + +[(\alpha, l, w)], w \rangle} \text{query-v}$$

In this case, we have $c' = \text{skip}$, $t' = t + +[(\alpha, l, w)]$, $m' = m[v/x]$ and $w' = w$.

Let $t'' = [(\alpha, l, w)]$, we have $t + +[(\alpha, l, w)] = t'$, i.e., $t + t'' = t'$. This case is proved.

□

]]

3.5 Trace-based Adaptivity

We define adaptivity through a query-based dependency graph. In our model, an *analyst* asks a sequence of queries to the mechanism, and the analyst receives the answers to these queries from the mechanism. A query is adaptively chosen by the analyst when the choice of this query is affected by answers from previous queries. In this model, the adaptivity we are interested in is the length of the longest sequence of such adaptively chosen queries, among all the queries the data analyst asks to the mechanism. Also, when the analyst asks a query, the only information the analyst will have will be the answers to previous queries and the state of the program. It means that when we want to know if this query is adaptively chosen, we only need to check whether the choice of this query will be affected by changes of answers to previous queries. There are two possible situations that can affect the choice of a query, either the query argument directly uses the results of previous queries (data dependency), or the control flow of the program with respect to a query (whether to ask this query or not) depends on the results of previous queries (control flow dependency).

As a first step, we give a definition of when one query may depend on a previous query, which is supposed to consider both control dependency and data dependency. We first look at two possible candidates:

1. One query may depend on a previous query if and only if a change of the answer to the previous query may also change the result of the query.
2. One query may depend on a previous query if and only if a change of the answer to the previous query may also change the appearance of the query.

The first candidate works well by witnessing the result of one query according to the change of the answer of another query. We can easily find that the two queries have nothing to do with each other in a simple example $c = x \leftarrow \text{query}(\chi(1)); y \leftarrow \text{query}(\chi(2))$. This candidate definition works well with respect to data dependency. However, it fails to handle control dependency since it just monitors the changes to the answer of a query when the answer of previous queries returned change. The key point is that this query may also not be asked because of an analyst decision which depend on the answers of previous queries. An example of this situation is shown in program c_1 as follows.

$$c_1 = x \leftarrow \text{query}(\chi(1)); \text{if } (x > 2, y \leftarrow \text{query}(\chi(2)), \text{skip})$$

We choose the second candidate, which performs well by witnessing the appearance of one query $\text{query}(\chi(2))$ upon the change of the result of one previous query $\text{query}(\chi(1))$ in c_1 . It considers the control dependency, and at the same, does not miss the data dependency. In particular, the arguments of a query characterizes it. In this sense, if the data used in the arguments changes due to a different answer to a certain previous query, the appearance of the query may change as well. This situation is also captured by our definition. Let us look at another variant of program c , p_2 , in which the queries equipped with functions using previously assigned variables storing answer of its previous query.

$$c_2 = x \leftarrow \text{query}(\chi(2)); y \leftarrow \text{query}(x + \chi(3))$$

As a reminder, in the `While` language, the query request is composed by two components: a symbol query representing a linear query type and the argument e , which represents the function specifying what the query asks. So we do think $\text{query}(\chi(1))$ is different from $\text{query}(\chi(2))$. Informally, we think $\text{query}(x + \chi(3))$ may depend on the query $\text{query}(\chi(2))$, because equipped function of the former $x + \chi(3)$ depend on the data assigned with $\text{query}(\chi(2))$. We can see the appearance definition catches data dependency in such a way, since $\text{query}(x + \chi(2))$ will not be the same query if the value of x is changed.

We give a formal definition of variable may dependency based on the trace-based operational semantics as follows. [\[JL\]](#):

Definition 24 (Annotated Variables / Events May Dependency). .

One event $\text{av}_2 = (x_2, v_2, l_2, n_2)$ may depend on another one $\text{av}_1 = (x_1, v_1, l_1, n_1)$ in a program c , with a starting memory m and hidden database D , denoted as $\text{DEP}_{\text{av}}(\text{av}_1, \text{av}_2, c, m, D)$ is defined below.

$$\begin{aligned} & \exists \mathbf{m}_1, \mathbf{m}_3, \tau_1, \tau_3, \mathbf{c}_2, v_1, (\alpha_1 \vee \mathbf{e}_1). \\ & \left(\begin{aligned} & \langle \mathbf{m}, \mathbf{c}, [], [] \rangle \rightarrow^* \langle \mathbf{m}_1, [\mathbf{x}_1 \leftarrow \text{query}(\alpha_1) / (v_1)]^{l_1}; \mathbf{c}_1, \text{qt}_1, \tau_1, w_1 \rangle \xrightarrow{\text{ssa-query-v} / (\text{asn-v})} \\ & \langle \mathbf{m}_1[v_1/\mathbf{x}], \mathbf{c}_2, \text{qt}'_1, \tau_1 + +[\text{av}_1], w_1 \rangle \rightarrow^* \langle \mathbf{m}_3, \text{skip}, \text{qt}_3, \tau_3, w_3 \rangle \\ & \wedge \left(\begin{aligned} & \text{av}_2 \in (\tau'_3 - (\tau_1 + +[\text{av}_1])) \\ & \Rightarrow \exists v \in \mathcal{QD}, v \neq v_1, \mathbf{m}'_3, \text{qt}'_3, \tau'_3, w'_3. \langle \mathbf{m}_1[v/\mathbf{x}], \mathbf{c}_2, \text{qt}'_1, \tau_1 + +[\text{av}_1], w_1 \rangle \\ & \rightarrow^* (\langle \mathbf{m}'_3, \text{skip}, \text{qt}'_3, \tau'_3, w'_3 \rangle \\ & \wedge \text{av}_2 \notin (\tau'_3 - (\tau_1 + +[\text{av}_1]))) \end{aligned} \right) \\ & \wedge \left(\begin{aligned} & \text{av}_2 \notin (\tau_3 - (\tau_1 + +[\text{av}_1])) \\ & \Rightarrow \exists v \in \mathcal{QD}, v \neq v_1, \mathbf{m}'_3, \text{qt}'_3, \tau'_3, w'_3. \langle \mathbf{m}_1[v/\mathbf{x}], \mathbf{c}_2, \text{qt}'_1, \tau_1 + +[\text{av}_1], w_1 \rangle \\ & \rightarrow^* (\langle \mathbf{m}'_3, \text{skip}, \text{qt}'_3, \tau'_3, w'_3 \rangle \\ & \wedge \text{av}_2 \in (\tau'_3 - (\tau_1 + +[\text{av}_1]))) \end{aligned} \right) \end{aligned} \right) \end{aligned}$$

Definition 25 (Annotated Variables May Dependency – Version 2). .

One event av_2 may depend on another one av_1 in a program \mathbf{c} , with a starting memory \mathbf{m} and hidden

database D , denoted as $\text{DEP}_{\text{av}}(\text{av}_1, \text{av}_2, c, \mathbf{m}, D)$ is defined below.

$$\begin{aligned}
& \exists \mathbf{m}, \mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_3, \mathbf{m}'_2, \mathbf{m}'_3, \tau_1, \tau_2, \tau'_2, t_1, t_2, t'_2, \mathbf{c}_1, \mathbf{c}_2, v'_1, \mathbf{e}_2. \\
& \left(\begin{array}{l} \langle \mathbf{m}, \mathbf{c}, [] \rangle \rightarrow^* \langle \mathbf{m}_1, [\mathbf{x}_1 \leftarrow v_1]^{l_1}; \mathbf{c}_1, \tau_1, t_1 \rangle \\ \wedge \langle \mathbf{m}_1[v_1/\mathbf{x}_1], \mathbf{c}_1, \tau_1 ++ [\text{av}_1], t_1[\mathbf{x}_1] ++ \rangle \\ \rightarrow^* \langle \mathbf{m}_2, [\mathbf{x}_2 \leftarrow \mathbf{e}_2]^{l_2}; \mathbf{c}_2, \tau_2, t_2 \rangle \\ \rightarrow^* \langle \mathbf{m}_3, \mathbf{c}_2, \tau_2 ++ [\text{av}_2], t_2[\mathbf{x}_2] ++ \rangle \\ \wedge \langle \mathbf{m}_1[v'_1/\mathbf{x}_1], \mathbf{c}_1, \tau_1, t_1 \rangle \rightarrow^* \langle \mathbf{m}'_2, \mathbf{c}_2, \tau'_2, t'_2 \rangle \\ \wedge \text{av}_2 \notin \tau'_2 \end{array} \right) \quad \begin{array}{l} \text{av}_1 = (\mathbf{x}_1, v_1, l_1, n_1) \\ \text{av}_2 = (\mathbf{x}_2, v_2, l_2, n_2) \end{array} \\
& \exists \mathbf{m}, \mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_3, \mathbf{m}'_2, \mathbf{m}'_3, \tau_1, \tau_2, \tau'_2, t_1, t_2, t'_2, \mathbf{c}_1, \mathbf{c}_2, v'_1, \psi_2. \\
& \left(\begin{array}{l} \langle \mathbf{m}, \mathbf{c}, [] \rangle \rightarrow^* \langle \mathbf{m}_1, [\mathbf{x}_1 \leftarrow \text{query}(\alpha_1)]^{l_1}; \mathbf{c}_1, \tau_1, t_1 \rangle \\ \wedge \langle \mathbf{m}_1[v_1/\mathbf{x}_1], \mathbf{c}_1, \tau_1 ++ [\text{av}_1], t_1[\mathbf{x}_1] ++ \rangle \\ \rightarrow^* \langle \mathbf{m}_2, [\mathbf{x}_2 \leftarrow \text{query}(\psi_2)]^{l_2}; \mathbf{c}_2, \tau_2, t_2 \rangle \\ \rightarrow^* \langle \mathbf{m}_3, \mathbf{c}_2, \tau_2 ++ [\text{av}_2], t_2[\mathbf{x}_2] ++ \rangle \\ \wedge \langle \mathbf{m}_1[v'_1/\mathbf{x}_1], \mathbf{c}_1, \tau_1, t_1 \rangle \rightarrow^* \langle \mathbf{m}'_2, \mathbf{c}_2, \tau'_2, t'_2 \rangle \\ \wedge \text{av}_2 \notin \tau'_2 \end{array} \right) \quad \begin{array}{l} \text{av}_1 = (\mathbf{x}_1, \alpha_1, l_1, n_1) \\ \text{av}_2 = (\mathbf{x}_2, \alpha_2, l_2, n_2) \end{array} \\
& \exists \mathbf{m}, \mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_3, \mathbf{m}'_2, \mathbf{m}'_3, \tau_1, \tau_2, \tau'_2, t_1, t_2, t'_2, \mathbf{c}_1, \mathbf{c}_2, v'_1, \psi_2. \\
& \left(\begin{array}{l} \langle \mathbf{m}, \mathbf{c}, [] \rangle \rightarrow^* \langle \mathbf{m}_1, [\mathbf{x}_1 \leftarrow v_1]^{l_1}; \mathbf{c}_1, \tau_1, t_1 \rangle \\ \wedge \langle \mathbf{m}_1[v_1/\mathbf{x}_1], \mathbf{c}_1, \tau_1 ++ [\text{av}_1], t_1[\mathbf{x}_1] ++ \rangle \\ \rightarrow^* \langle \mathbf{m}_2, [\mathbf{x}_2 \leftarrow \text{query}(\psi_2)]^{l_2}; \mathbf{c}_2, \tau_2, t_2 \rangle \\ \rightarrow^* \langle \mathbf{m}_3, \mathbf{c}_2, \tau_2 ++ [\text{av}_2], t_2[\mathbf{x}_2] ++ \rangle \\ \wedge \langle \mathbf{m}_1[v'_1/\mathbf{x}_1], \mathbf{c}_1, \tau_1, t_1 \rangle \rightarrow^* \langle \mathbf{m}'_2, \mathbf{c}_2, \tau'_2, t'_2 \rangle \\ \wedge \text{av}_2 \notin \tau'_2 \end{array} \right) \quad \begin{array}{l} \text{av}_1 = (\mathbf{x}_1, v_1, l_1, n_1) \\ \text{av}_2 = (\mathbf{x}_2, \alpha_2, l_2, n_2) \end{array} \\
& \exists \mathbf{m}, \mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_3, \mathbf{m}'_2, \mathbf{m}'_3, \tau_1, \tau_2, \tau'_2, t_1, t_2, t'_2, \mathbf{c}_1, \mathbf{c}_2, v'_1, \mathbf{e}_2. \\
& \left(\begin{array}{l} \langle \mathbf{m}, \mathbf{c}, [] \rangle \rightarrow^* \langle \mathbf{m}_1, [\mathbf{x}_1 \leftarrow v_1]^{l_1}; \mathbf{c}_1, \tau_1, t_1 \rangle \\ \wedge \langle \mathbf{m}_1[v_1/\mathbf{x}_1], \mathbf{c}_1, \tau_1 ++ [\text{av}_1], t_1[\mathbf{x}_1] ++ \rangle \\ \rightarrow^* \langle \mathbf{m}_2, [\mathbf{x}_2 \leftarrow \mathbf{e}_2]^{l_2}; \mathbf{c}_2, \tau_2, t_2 \rangle \\ \rightarrow^* \langle \mathbf{m}_3, \mathbf{c}_2, \tau_2 ++ [\text{av}_2], t_2[\mathbf{x}_2] ++ \rangle \\ \wedge \langle \mathbf{m}_1[v'_1/\mathbf{x}_1], \mathbf{c}_1, \tau_1, t_1 \rangle \rightarrow^* \langle \mathbf{m}'_2, \mathbf{c}_2, \tau'_2, t'_2 \rangle \\ \wedge \text{av}_2 \notin \tau'_2 \end{array} \right) \quad \begin{array}{l} \text{av}_1 = (\mathbf{x}_1, \alpha_1, l_1, n_1) \\ \text{av}_2 = (\mathbf{x}_2, v_2, l_2, n_2) \end{array}
\end{aligned}$$

Definition 26 (Variable May Dependency). .

Given a program \mathbf{c} with its assigned variables $\text{aVar}_{\mathbf{c}}$, one variable $\mathbf{x}_2 \in \text{aVar}_{\mathbf{c}}$ may depend on another variable $\mathbf{x}_1 \in \text{aVar}_{\mathbf{c}}$ in \mathbf{c} denoted as $\text{DEP}_{\text{var}}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{c})$ is defined below.

$$\exists \text{av}_1, \text{av}_2, \mathbf{m}, D. \pi_1(\text{av}_1) = \mathbf{x}_1 \wedge \pi_1(\text{av}_2) = \mathbf{x}_2 \wedge \text{DEP}_{\text{av}}(\text{av}_1, \text{av}_2, c, \mathbf{m}, D)$$

Definition 27 (Execution Based Dependency Graph). .

Given a program \mathbf{c} , a database D , a starting memory \mathbf{m} with its assigned variables $\text{aVar}_{\mathbf{c}}$ and initial variable counter $\text{vcnt}_{\mathbf{c}}^0$ with its corresponding execution: $\langle \mathbf{m}, \mathbf{c}, [], \text{vcnt}_{\mathbf{c}}^0 \rangle \rightarrow^* \langle \mathbf{m}', \text{skip}, \tau, \text{vcnt} \rangle$, the dependency graph $\mathbf{G}_{\text{trace}}(\mathbf{c}, \mathbf{m}, D) = (V, E, W, \text{qF})$ is defined as:

$$\begin{array}{lll}
\text{Vertices} & V & := \{x \in \mathcal{VAR} \mid x = \text{aVar}_{\mathbf{c}}(i); i = 0, \dots, |\text{aVar}_{\mathbf{c}}|\} \\
\text{Directed Edges} & E & := \{(x, x') \in \mathcal{VAR} \times \mathcal{VAR} \mid \text{DEP}_{\text{var}}(x, x', c); x = \text{aVar}_{\mathbf{c}}(i); x' = \text{aVar}_{\mathbf{c}}(j); i, j = 0, \dots, |\text{aVar}_{\mathbf{c}}|\} \\
\text{Weights} & W & := \{(x, n) \in \mathcal{VAR} \times \mathbb{N} \mid n = \text{vcnt}(x); x = \text{aVar}_{\mathbf{c}}(i); i = 0, \dots, |\text{aVar}_{\mathbf{c}}|\} \\
\text{Query Flags} & \text{qF} & := \left\{ (x, n) \in V \times \{0, 1\} \mid \left\{ \begin{array}{ll} n = 1 & x \in \text{qVar}_{\mathbf{c}} \\ n = 0 & o.w. \end{array} \right\}; x = \text{aVar}_{\mathbf{c}}(i); i = 0, \dots, |\text{aVar}_{\mathbf{c}}| \right\}
\end{array}$$

Definition 28 (Finite Walk (k)). .

Given a labeled weighted graph $G = (V, E, W, qF)$, a finite walk k in G is a sequence of edges $(e_1 \dots e_{n-1})$ for which there is a sequence of vertices (v_1, \dots, v_n) such that:

- $e_i = (v_i, v_{i+1})$ for every $1 \leq i < n$.
- every vertex $v \in V$ appears in this vertices sequence (v_1, \dots, v_n) of k at most $W(v)$ times.

(v_1, \dots, v_n) is the vertex sequence of this walk.

Length of this finite walk k is the number of vertices in its vertex sequence, i.e., $\text{len}(k) = n$.

Given a labeled weighted graph $G = (V, E, W, qF)$, we use $\mathcal{WALK}(G)$ to denote a set containing all finite walks k in G ; and $k_{v_1 \rightarrow v_2} \in \mathcal{WALK}(G)$ where $v_1, v_2 \in V$ denotes the walk from vertex v_1 to v_2 .

Definition 29 (Length of Finite Walk w.r.t. Query (len_q)). .

Given a labeled weighted graph $G = (V, E, W, qF)$ and a finite walk k in G with its vertex sequence (v_1, \dots, v_n) , the length of k w.r.t query is defined as:

$$\text{len}_q(k) = \text{len}(v \mid v \in (v_1, \dots, v_n) \wedge F(v) = 1)$$

, where $(v \mid v \in (v_1, \dots, v_n) \wedge F(v) = 1)$ is a subsequence of k 's vertex sequence.

Given a program c with a starting memory m and database D , we generate its program-based graph $\mathbf{G}_{\text{trace}}(\mathbf{c}, \mathbf{m}, \mathbf{D}) = (V, E, W, qF)$. Then the adaptivity bound based on program analysis for \mathbf{c} is the number of query vertices on a finite walk in $\mathbf{G}_{\text{prog}}(\mathbf{c})$. This finite walk satisfies:

- the number of query vertices on this walk is maximum
- the visiting times of each vertex v on this walk is bound by its weight $W(v)$.

It is formally defined in 30.

Definition 30 (Adaptivity of A Program). .

Given a program \mathbf{c} in SSA language, its adaptivity is defined for all possible starting SSA memory \mathbf{m} and database D as follows:

$$A(c) = \max \{ \text{len}_q(k) \mid \mathbf{m} \in \mathcal{SM}, D \in \mathcal{DB}, k \in \mathcal{WALK}(\mathbf{G}_{\text{trace}}(\mathbf{c}, \mathbf{m}, \mathbf{D})) \}$$

]

[[The following lemma describes a property of the trace-based dependency graph. For any program c with a database D and a starting memory m , the directed edges in its trace-based dependency graph can only be constructed from nodes representing smaller annotated queries to annotated queries of greater order. There doesn't exist backward edges with direction from greater annotated queries to smaller ones.]]

Lemma 12. [Edges are Forwarding Only].

Given a program c , a database D , a starting memory m and the corresponding trace-based dependency graph $G(c, D, m) = (V, E)$, for any directed edge $(av', av) \in E$, this is not the case that:

$$av' \geq_{av} av$$

Proof. Proof in File: "edge_forward.tex".

□

Lemma 13. *[Trace-based Dependency Graph is Directed Acyclic].
Every trace-based dependency graph is a directed acyclic graph.*

Proof. Proof is obvious based on the Lemma 12. □

Lemma 14 (Adaptivity is Bounded). .

Given the program c with a certain database D and starting memory m , the $A(c)$ w.r.t. the D and m is bounded, i.e.,:

$$\langle m, c, [], [] \rangle \rightarrow^* \langle m', \text{skip}, t', w' \rangle \implies A_{D,m}(c) \leq |t'|$$

Proof. Proof is obvious based on the Lemma 13. □

3.6 SSA Transformation and Soundness of Transformation

in File "ssa_transform_sound.tex"

4 AdaptFun

There are four steps to get the adaptivity of a program c based on analyzing the program.

1. Collecting the variables that are newly assigned in the program (via assignment expressions). These variables are stored in an assigned variable vector $aVar$. We also track extra information of each assigned variable (whether it is assigned by a query result, or showing up in loop, or showing up in `if` expression or o.w.) and store it in a vector F of the same size as $aVar$.
2. Tracking the data flow relations between all these assigned variables. These informations are stored in a matrix M , whose size is $|aVar| \times |aVar|$.
3. Estimating the reachability bound of each variable in $aVar$.
4. With all these informations from previous steps, generating a program-based dependency graph G_{prog} and compute the adaptivity bound.

In the following subsections, we first define the notations and symbols being used in **AdaptFun** with a simple example for understanding these definitions. Then we present the algorithmic analysis rules, which is the core of the **AdaptFun**, with 3 examples illustrating how **AdaptFun** works. In the following subsections, we present the adaptivity analysis based on the **AdaptFun**'s analyzing results, and the soundness w.r.t. the trace-based analyzing results in previous sections.

4.1 Notations

Definition 31 (Assigned Variables ($aVar$)). *Given a program c , its assigned variables $aVar$ is a vector containing all variables newly assigned in the program preserving the order. It is defined as follows:*

$$aVar_c \triangleq \begin{cases} [x] & c = [x \leftarrow e]^{(l,w)} \\ [x] & c = [x \leftarrow \text{query}(\psi)]^{(l,w)} \\ aVar_{c_1} ++ aVar_{c_2} & c = c_1; c_2 \\ aVar_{c_1} ++ aVar_{c_2} ++ [\bar{x}, \bar{y}, \bar{z}] & c = \text{if} ([b]^{(l,w)}, [\bar{x}, \bar{x}_2, \bar{x}_2], [\bar{y}, \bar{y}_2, \bar{y}_3], [\bar{z}, \bar{z}_2, \bar{z}_3], c_1, c_2) \\ aVar_{c'} ++ [\bar{x}] & c = \text{while} ([b]^{(l,w)}, [\bar{x}, \bar{x}_2, \bar{x}_2], c') \end{cases}$$

[JL: We are abusing the notations and operators from list here. The notation $[]$ represents an empty vector and $x :: A$ represents add an element x to the head of the vector A . The concatenation operation between 2 vectors A_1 and A_2 , i.e., $A_1 ++ A_2$ is mimic the standard list concatenation operations as follows:

$$A_1 ++ A_2 \triangleq \begin{cases} A_2 & A_1 = [] \\ x :: (A'_1 ++ A_2) & A_1 = x :: A'_1 \end{cases} \quad (13)$$

We use index within parenthesis to denote the access to the element of corresponding location, $A(i)$ denotes the element at location i in the vector A and $M(i, j)$ denotes the element at location i -th row, j -th column in the matrix M .]

Consider the program c below in the left hand side as an example, its assigned variables $aVar$ (short for $aVar(c)$) is as in the right hand side is shown as follows:

$$c = \begin{bmatrix} [x_1 \leftarrow \text{query}(0)]^1; \\ [x_2 \leftarrow x_1 + 1]^2; \\ [x_3 \leftarrow x_2 + 2]^3 \end{bmatrix}; \quad aVar = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

Lemma 15. For any program \mathbf{c} , every variable in $\mathbf{aVar}(\mathbf{c})$ is distinct

Proof. It is due to the SSA nature. We can prove it by induction on \mathbf{c} . □

[JL:

Definition 32 (Variable Flags (F)). .

Given a program \mathbf{c} with its assigned variables \mathbf{aVar} , the \mathbf{F} is a vector of the same length as \mathbf{aVar} , s.t. for each variable \mathbf{x} showing up as the i -th element in \mathbf{aVar} (i.e., $\mathbf{x} = \mathbf{aVar}(i)$), $\mathbf{F}(i) \in \{0, 1, 2\}$ is defined as follows:

$$\mathbf{F}(i) := \left\{ \begin{array}{ll} 2 & \mathbf{x} = \mathbf{aVar}(i) \wedge (\exists \psi. \text{ s.t., } [\mathbf{x} \leftarrow \text{query}(\psi)]^l \in_c \mathbf{c}) \\ & \mathbf{x} = \mathbf{aVar}(i) \wedge \\ 1 & \left((\exists \mathbf{c}', \mathbf{e}, \mathbf{b}, l, l'. \text{ while } [\mathbf{b}]^l \text{ do } \mathbf{c}' \in_c \mathbf{c} \wedge [\mathbf{x} \leftarrow \mathbf{e}]^{l'} \in_c \mathbf{c}') \vee \right. \\ & \left. (\exists \mathbf{b}, l, l_1, l_2, \mathbf{c}_1, \mathbf{c}_2, \mathbf{e}_1, \mathbf{e}_2. \text{ if } ([\mathbf{b}]^l, \mathbf{c}_1, \mathbf{c}_2) \in_c \mathbf{c} \wedge ([\mathbf{x} \leftarrow \mathbf{e}_1]^{l_1} \in_c \mathbf{c}_1 \vee [\mathbf{x} \leftarrow \mathbf{e}_2]^{l_2} \in_c \mathbf{c}_2)) \right) \\ 0 & \text{o.w.} \end{array} \right\}.$$

Operations on \mathbf{F} are defined as follows:

$$\begin{aligned} \mathbf{F}_1 \uplus \mathbf{F}_2(i) &:= \begin{cases} k & k = \max\{\mathbf{F}_1(i), \mathbf{F}_2(i)\} \wedge |\mathbf{F}_1| = |\mathbf{F}_2| \\ 0 & \text{o.w.} \end{cases} & i = 1, \dots, |\mathbf{F}_1| \\ \mathbf{F} \uplus n(i) &:= \max\{\mathbf{F}(i), n\} & i = 1, \dots, |\mathbf{F}| \\ [n]^k(i) &:= n & i = 1, \dots, k \wedge |[n]^k| = k \end{aligned} \quad (14)$$

[[Given a program \mathbf{c} with its assigned variables \mathbf{aVar} , and two variables \mathbf{x}, \mathbf{y} showing up as i -th, j -th elements in \mathbf{aVar} (i.e., $\mathbf{x} = \mathbf{aVar}(i)$ and $\mathbf{y} = \mathbf{aVar}(j)$), we say \mathbf{y} flows to \mathbf{x} in \mathbf{c} if and only if $j < i$ and the value of \mathbf{y} directly or indirectly influence the evaluation of the value of \mathbf{x} as follows:

- **(Directly Influence)** The program \mathbf{c} contains either a command $\mathbf{x} \leftarrow \mathbf{e}$ or $\mathbf{x} \leftarrow \text{query}(\psi)$, such that \mathbf{y} shows up as a free variable in \mathbf{e} or ψ . We use $\text{flowsTo}(\mathbf{x}, \mathbf{y}, \mathbf{c})$ to denote \mathbf{y} flows to \mathbf{x} in \mathbf{c} .
- **(Indirectly Influence)** The program \mathbf{c} contains either a while loop command or if condition command, such that \mathbf{y} shows up in the guard and \mathbf{x} shows up in the left hand of an assignment command in the body.

This is formally defined in 33. We use $FV(e)$, $FV(\mathbf{b})$ and $FV(\psi)$ denote the set of free variables in expression e , boolean expression \mathbf{b} and query expression ψ respectively.

Definition 33 (Data Flows between Assigned Variables (flowsTo)). .

Given a program \mathbf{c} with its assigned variables \mathbf{aVar} , and two variables \mathbf{x}, \mathbf{y} s.t., $\mathbf{x} = \mathbf{aVar}(i)$ and $\mathbf{y} = \mathbf{aVar}(j)$, \mathbf{y} flows to \mathbf{x} in \mathbf{c} , i.e., $\text{flowsTo}(\mathbf{x}, \mathbf{y}, \mathbf{c})$ is defined as:

$$\text{flowsTo}(\mathbf{x}, \mathbf{y}, \mathbf{c}) \triangleq (j < i) \wedge \left(\begin{array}{l} (\exists \mathbf{e}, l. [\mathbf{x} \leftarrow \mathbf{e}]^l \in_c \mathbf{c} \wedge \mathbf{y} \in FV(\mathbf{e})) \\ \vee (\exists \psi, l. [\mathbf{x} \leftarrow \text{query}(\psi)]^l \in_c \mathbf{c} \wedge \mathbf{y} \in FV(\psi)) \\ \vee (\exists \mathbf{c}', \mathbf{e}, \mathbf{b}, l, l'. \text{ while } [\mathbf{b}]^l \text{ do } \mathbf{c}' \in_c \mathbf{c} \wedge [\mathbf{x} \leftarrow \mathbf{e}]^{l'} \in_c \mathbf{c}' \wedge \mathbf{y} \in FV(\mathbf{b})) \\ \vee (\exists \mathbf{b}, l, l_1, l_2, \mathbf{c}_1, \mathbf{c}_2, \mathbf{e}_1, \mathbf{e}_2. \text{ if } ([\mathbf{b}]^l, \mathbf{c}_1, \mathbf{c}_2) \in_c \mathbf{c} \wedge ([\mathbf{x} \leftarrow \mathbf{e}_1]^{l_1} \in_c \mathbf{c}_1 \vee [\mathbf{x} \leftarrow \mathbf{e}_2]^{l_2} \in_c \mathbf{c}_2) \wedge \mathbf{y} \in FV(\mathbf{b})) \end{array} \right).$$

]]

Definition 34 (Data Flow Matrix (M)). *Given a program \mathbf{c} with its assigned variables \mathbf{aVar} of length N , its data flow matrix \mathbf{M} is a matrix of size $N \times N$ s.t. $\forall \mathbf{x}, \mathbf{y} \in \mathbf{aVar}. \mathbf{x} = \mathbf{aVar}(i), \mathbf{y} = \mathbf{aVar}(j)$:*

$$\mathbf{M}(i, j) \triangleq \begin{cases} 1 & \text{flowsTo}(\mathbf{x}, \mathbf{y}, \mathbf{c}) \\ 0 & \text{o.w.} \end{cases}, \mathbf{x} = \mathbf{aVar}(i); \mathbf{y} = \mathbf{aVar}(j); i, j = 1, \dots, N.$$

Operations on the data flow matrices are defined as follows:

$$\mathbf{M}_1; \mathbf{M}_2 := \mathbf{M}_2 \cdot \mathbf{M}_1 + \mathbf{M}_1 + \mathbf{M}_2 \quad (15)$$

Consider the same program c as above, its data flow matrix \mathbf{M} and \mathbf{F} for the program c is as follows:

$$\mathbf{c} = \begin{array}{l} [\mathbf{x}_1 \leftarrow \text{query}(\mathbf{0})]^1; \\ [\mathbf{x}_2 \leftarrow \mathbf{x}_1 + \mathbf{1}]^2; \\ [\mathbf{x}_3 \leftarrow \mathbf{x}_2 + \mathbf{2}]^3 \end{array} \quad \mathbf{M} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \end{bmatrix}, \mathbf{F} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

There are two special matrices used for generating the data flow matrix \mathbf{M} in the analysis algorithm. They are the left matrix \mathbf{LM}_i and right matrix $\mathbf{R}_{(e,i)}$.

Given a program \mathbf{c} with its assigned variables \mathbf{aVar} of length N , the left matrix \mathbf{LM}_i generates a matrix of 1 column, N rows, where the i -th row is 1 and all the other rows are 0.

Definition 35 (Left Matrix (\mathbf{LM}_i)). .

Given a program \mathbf{c} with its assigned variables \mathbf{aVar} of length N , the left matrix \mathbf{LM}_i is defined as follows:

$$\mathbf{LM}_i(j) := \begin{cases} 1 & j = i \\ 0 & \text{o.w.} \end{cases}, j = 1, \dots, N.$$

Given a program \mathbf{c} with its assigned variables \mathbf{aVar} of length N , the right matrix $\mathbf{RM}_{e,i}$ generates a matrix of one row and N columns, where the locations of free variables in e is marked as 1.

Definition 36 (Right Matrix (\mathbf{RM}_e)). .

Given a program \mathbf{c} with its assigned variables \mathbf{aVar} of length N , the right matrix \mathbf{RM}_e is defined as follows:

$$\mathbf{RM}_e(j) := \begin{cases} 1 & \mathbf{x} \in \mathbf{FV}(e) \\ 0 & \text{o.w.} \end{cases}, \mathbf{x} = \mathbf{aVar}(j), j = 1, \dots, N.$$

Using the same example program \mathbf{c} as above with assigned variables $\mathbf{aVar} = [\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3]$, the left and right matrices w.r.t. its 2-nd command $[\mathbf{x}_2 \leftarrow \mathbf{x}_1 + \mathbf{1}]^2$ are as follows:

$$\mathbf{LM}_1 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \mathbf{RM}_{\mathbf{x}_1+1} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$$

4.2 Algorithmic Analysis Rules

Variable Collection Algorithm, VetxCol The **VetxCol** algorithm shows how the assigned variables \mathbf{aVar} are collected (via the command $\mathbf{x} \leftarrow \mathbf{e}$ or $\mathbf{x} \leftarrow \text{query}(\psi)$) from the program \mathbf{c} in the first step, along with constructing the flag for each variable, i.e., \mathbf{F} . The algorithmic rules for **VetxCol** algorithm is defined in Figure 3. It has the form: $[\mathbf{JL}: \text{VetxCol}(\mathbf{aVar}; \mathbf{F}; \mathbf{c}) \rightarrow (\mathbf{aVar}'; \mathbf{F}')]$. The input of **VetxCol** is a program \mathbf{c} , the assigned variables \mathbf{aVar} collected before the program \mathbf{c} as well as the flags \mathbf{F} for every corresponding variable. The output of the algorithm is the updated assigned variables \mathbf{aVar}' and flags

[JL:

$$\begin{array}{c}
\frac{}{\text{VetxCol}(\text{aVar}; F; [\mathbf{x} \leftarrow \mathbf{e}]^l) \rightarrow (\text{aVar} ++ [\mathbf{x}]; F ++ [0])} \text{VetxCol-asgn} \\
\\
\frac{}{\text{VetxCol}(\text{aVar}; F; [\mathbf{x} \leftarrow \text{query}(\boldsymbol{\psi})]^l) \rightarrow (\text{aVar} ++ [\mathbf{x}]; F ++ [2])} \text{VetxCol-query} \\
\\
\frac{\begin{array}{l} \text{VetxCol}(\text{aVar}; []; \mathbf{c}_1) \rightarrow (\text{aVar}_1; F_1) \quad \text{VetxCol}(\text{aVar}_1; []; \mathbf{c}_2) \rightarrow (\text{aVar}_2; F_2) \quad \text{aVar}' = [\bar{\mathbf{x}}] ++ [\bar{\mathbf{y}}] ++ [\bar{\mathbf{z}}] \\ k = \text{len}(\text{aVar}') \quad \text{aVar}_3 = \text{aVar}_2 ++ \text{aVar}' \quad F_3 = F ++ ((F_1 ++ F_2) \uplus 1) ++ ([1]^k) \end{array}}{\text{VetxCol}(\text{aVar}; F; [\text{if } (\mathbf{b}, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2], \mathbf{c}_1, \mathbf{c}_2)]^l) \rightarrow (\text{aVar}_3; F_3)} \text{VetxCol-if} \\
\\
\frac{\text{VetxCol}(\text{aVar}; F\mathbf{c}_1) \rightarrow (\text{aVar}_1; F_1) \quad \text{VetxCol}(\text{aVar}_1; F_1; \mathbf{c}_2) \rightarrow (\text{aVar}_2; F_2)}{\text{VetxCol}(\text{aVar}; F; (\mathbf{c}_1; \mathbf{c}_2)) \rightarrow (\text{aVar}_2; F_2)} \text{VetxCol-seq} \\
\\
\frac{\text{VetxCol}(\text{aVar}; []; \mathbf{c}) \rightarrow (\text{aVar}'; F') \quad \text{aVar}'' = \text{aVar}' ++ [\bar{\mathbf{x}}] \quad F'' = F ++ (F' \uplus 1) ++ ([1]^{\text{len}([\bar{\mathbf{x}}])})}{\text{VetxCol}(\text{aVar}; F; \text{while } [\mathbf{b}]^l, \mathbf{n}, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] \text{ do } \mathbf{c}) \rightarrow (\text{aVar}''; F'')} \text{VetxCol-while} \\
]
\end{array}$$

Figure 3: The Algorithmic Rules of **VetxCol**

F' thorough the program \mathbf{c} The assignment commands are the source of variables **VetxCol** collecting, in the case **VetxCol-asgn** and **VetxCol-query**, the output assigned variables are extended by \mathbf{x} .

When it comes to the `if ... then ... else` command in the rule **VetxCol-if**, variables assigned in the then branch \mathbf{c}_1 , as well as the variables assigned in the else branch \mathbf{c}_2 , and the new generated variables $\bar{\mathbf{x}}, \bar{\mathbf{y}}, \bar{\mathbf{z}}$ in $[\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2]$.

The sequence command $\mathbf{c}_1; \mathbf{c}_2$ is standard by accumulating the predicted variables in the two commands \mathbf{c}_1 and \mathbf{c}_2 preserving their order.

The while command `while $\mathbf{b}, [\bar{\mathbf{x}}] \dots$ do \mathbf{c}` considers the newly generated variables by SSA transformation $\bar{\mathbf{x}}$ as well and the newly assigned variables in its body \mathbf{c} .

Below we present the definition for a valid index, to have a clear understanding on the variable collecting algorithm: [JL:

Definition 37 (Valid Index (Remove?)). *Given an assigned variable list aVar , $\text{aVar} \models (\mathbf{c}, i_1, i_2)$ iff $\text{aVar}' = \text{aVar}[0, \dots, i_1 - 1], \text{aVar}'; \mathbf{c} \rightarrow \text{aVar}'' \wedge \text{aVar}'' = \text{aVar}[0, \dots, i_2 - 1]$.*

]

Data Flow Matrix Generating Algorithm In this data flow matrix generating algorithm, we analyze the data flow information among all assigned variables aVar collected via the the **VetxCol** algorithm of length N . We track the data flow relations between all these assigned variables. These informations are stored in a matrix M , whose size is $N \times N$. The algorithm to fill in the matrix is of the form: [JL: $\text{FlowGen}(\Gamma; \mathbf{c}; \text{aVar}) \rightarrow (M)$] $\text{FlowGen}(\Gamma; \mathbf{c}; i_1, i_2) \rightarrow (M; F)$. Γ is a vector records the variables the current program \mathbf{c} depends on, the index i_1 is a pointer which refers to the position of the first new-generated variable in \mathbf{c} in the assigned variables aVar , and i_2 points to the first new variable that is not in \mathbf{c} (if exists). [JL:

Definition 38 (Valid Gamma (Remove?)). $\Gamma \models i_1$ iff $\forall i \geq i_1, \Gamma(i_1) = 0$.

$$\begin{array}{c}
\boxed{\Gamma \vdash_{M,F}^{i_1, i_2} c} \\
\\
\frac{M = LM_i * (RM_{e,i} + \Gamma)}{\text{FlowGen}(\Gamma; [\mathbf{x} \leftarrow \mathbf{e}]^l; i) \rightarrow (M; F_0; i + 1)} \text{FlowGen-assign} \\
\\
\frac{M = LM_i * (RM_{e,i} + \Gamma) \quad F = LM_i \quad F(i) = 1}{\text{FlowGen}(\Gamma; [\mathbf{x} \leftarrow \text{query}(\mathbf{e})]^l; i) \rightarrow (M; F; i + 1)} \text{FlowGen-query} \\
\\
\frac{\begin{array}{c} \text{FlowGen}(\Gamma + RM_{b,i_1}; \mathbf{c}_1; i_1) \rightarrow (M_1; F_1; i_2) \quad \text{FlowGen}(\Gamma + RM_{b,i_1}; \mathbf{c}_2; i_2) \rightarrow (M_2; F_2; i_3) \\ \text{FlowGen}(\Gamma; [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]; i_3) \rightarrow (M_x; F_\emptyset; i_3 + |\bar{\mathbf{x}}|) \\ \text{FlowGen}(\Gamma; [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2]; i_3 + |\bar{\mathbf{x}}|) \rightarrow (M_y; F_\emptyset; i_3 + |\bar{\mathbf{x}}| + |\bar{\mathbf{y}}|) \\ \text{FlowGen}(\Gamma; [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2]; i_3 + |\bar{\mathbf{x}}| + |\bar{\mathbf{y}}|) \rightarrow (M_z; F_\emptyset; i_3 + |\bar{\mathbf{x}}| + |\bar{\mathbf{y}}| + |\bar{\mathbf{z}}|) \quad M = (M_1 + M_2) + M_x + M_y + M_z \end{array}}{\text{FlowGen}(\Gamma; \text{if } ([\mathbf{b}]^l, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2], \mathbf{c}_1, \mathbf{c}_2); i_1) \rightarrow (M; F_1 \uplus F_2 \uplus 2; i_3 + |\bar{\mathbf{x}}| + |\bar{\mathbf{y}}| + |\bar{\mathbf{z}}|)} \text{FlowGen-if} \\
\\
\frac{\begin{array}{c} \text{FlowGen}(\Gamma; \mathbf{c}_1; i_1) \rightarrow (M_1; F_1; i_2) \quad \text{FlowGen}(\Gamma; \mathbf{c}_2; i_2) \rightarrow (M_2; F_2; i_3) \end{array}}{\text{FlowGen}(\Gamma; (\mathbf{c}_1; \mathbf{c}_2); i_1) \rightarrow ((M_1; M_2); F_1 \uplus V_2; i_3)} \text{FlowGen-seq} \\
\\
\frac{\begin{array}{c} B = |\bar{\mathbf{x}}| \quad A = |\mathbf{c}| \quad \text{FlowGen}(\Gamma; [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]; i + (B + A)) \rightarrow (M_1; V_1; i + B + (B + A)) \\ \text{FlowGen}(\Gamma; \mathbf{c}; i + B + (B + A)) \rightarrow (M_2; F_2; i + B + A + (B + A)) \\ \text{FlowGen}(\Gamma; [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]; i + (B + A)) \rightarrow (M; F; i + (B + A) + B) \\ M' = M + (M_1 + M_2) \quad F' = F \uplus ((F_1 \uplus F_2) \uplus 2) \end{array}}{\text{FlowGen}(\Gamma; \text{while } [\mathbf{b}]^l \mathbf{n} [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] \text{ do } \mathbf{c}; i) \rightarrow (M'; F'; i + (B + A) + B)} \text{FlowGen-while} \\
\\
\text{[JL: Updated Flow Generation Algorithm] [JL: } \boxed{\Gamma \vdash_{M, aVar} \mathbf{c}} \\
\\
\frac{\mathbf{x} = aVar(i) \quad M = LM_i * (RM_{e,i} + \Gamma)}{\text{FlowGen}(\Gamma; [\mathbf{x} \leftarrow \mathbf{e}]^l; aVar) \rightarrow (M)} \text{FlowGen-assign} \\
\\
\frac{\mathbf{x} = aVar(i) \quad M = LM_i * (RM_{e,i} + \Gamma)}{\text{FlowGen}(\Gamma; [\mathbf{x} \leftarrow \text{query}(\boldsymbol{\psi})]^l; aVar) \rightarrow (M)} \text{FlowGen-query} \\
\\
\frac{\begin{array}{c} \text{FlowGen}(\Gamma + RM_{b,i_1}; \mathbf{c}_1; aVar) \rightarrow (M_1) \quad \text{FlowGen}(\Gamma + RM_{b,i_1}; \mathbf{c}_2; aVar) \rightarrow (M_2) \\ \text{FlowGen}(\Gamma; [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]; aVar) \rightarrow (M_x) \quad \text{FlowGen}(\Gamma; [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2]; aVar) \rightarrow (M_y) \\ \text{FlowGen}(\Gamma; [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2]; aVar) \rightarrow (M_z) \quad M = (M_1 + M_2) + M_x + M_y + M_z \end{array}}{\text{FlowGen}(\Gamma; \text{if } ([\mathbf{b}]^l, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2], \mathbf{c}_1, \mathbf{c}_2)) \rightarrow (M)} \text{FlowGen-if} \\
\\
\frac{\begin{array}{c} \text{FlowGen}(\Gamma; \mathbf{c}_1; aVar) \rightarrow (M_1) \quad \text{FlowGen}(\Gamma; \mathbf{c}_2; aVar) \rightarrow (M_2) \end{array}}{\text{FlowGen}(\Gamma; (\mathbf{c}_1; \mathbf{c}_2); aVar) \rightarrow ((M_1; M_2))} \text{FlowGen-seq} \\
\\
\frac{\begin{array}{c} \text{FlowGen}(\Gamma + RM_{b,i_1}; \mathbf{c}; aVar) \rightarrow (M_1) \quad \text{FlowGen}(\Gamma; [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]; aVar) \rightarrow (M_2) \end{array}}{\text{FlowGen}(\Gamma; \text{while } [\mathbf{b}]^l \mathbf{n} [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] \text{ do } \mathbf{c}; aVar) \rightarrow (M_1 + M_2)} \text{FlowGen-while}
\end{array}$$

] Below we define the valid data flow matrix, to have a clear understanding on the data flow generating algorithm:

Definition 39 (Valid Matrix). For a assigned variables aVar , $\text{aVar} \models (\mathbf{M}, \mathbf{F})$ iff the cardinality of aVar equals to the one of \mathbf{F} , $|\text{aVar}| = |\mathbf{F}|$ and the matrix \mathbf{M} is of size $|\mathbf{F}| \times |\mathbf{F}|$.

[JL:

Definition 40 (Valid Matrix). Given a program \mathbf{c} with its assigned variables aVar , $\text{aVar} \models \mathbf{M}$ iff the cardinality of \mathbf{M} equals to the product of aVar 's cardinality, i.e., $|\mathbf{M}| = |\text{aVar}| \times |\text{aVar}|$.

]

Reachability Bounds Given a program \mathbf{c} with its assigned variables aVar , we use the $\text{RechBound}(\mathbf{x}, \mathbf{c})$ algorithm, from paper [2], to estimate the reachability bound for each variable $\mathbf{x} \in \text{aVar}$. The input of RechBound is a program \mathbf{c} in SSA language and a variable \mathbf{x} from \mathbf{c} . The output of $\text{RechBound}(\mathbf{x}, \mathbf{c})$ is an integer representing the reachability bound of \mathbf{x} in \mathbf{c} .

The following example programs $\mathbf{c2}$ and $\mathbf{c3}$ with while loop illustrate how the algorithm works. The collected assigned variables, $\text{aVar}_{\mathbf{c2}}$ and $\text{aVar}_{\mathbf{c3}}$, data flow matrix $\mathbf{M}_{\mathbf{c2}}$ and $\mathbf{M}_{\mathbf{c3}}$ and variable flags $\mathbf{F}_{\mathbf{c2}}$ and $\mathbf{F}_{\mathbf{c3}}$ for program $\mathbf{c2}$ and $\mathbf{c3}$ are presented in the right hand side.

$$\mathbf{c2} \triangleq \begin{array}{l} [\mathbf{x}_1 \leftarrow \text{query}(1)]^1; \\ [\mathbf{i}_1 \leftarrow 0]^2; \\ \text{while } [\mathbf{i}_1 < 2]^3 \\ \quad [\mathbf{x}_3, \mathbf{x}_1, \mathbf{x}_2], [\mathbf{i}_3, \mathbf{i}_1, \mathbf{i}_2] \text{ do} \\ \quad \quad ([\mathbf{y}_1 \leftarrow \text{query}(2)]^4; \\ \quad \quad [\mathbf{x}_2 \leftarrow \mathbf{y}_1 + \mathbf{x}_3]^5; \\ \quad \quad [\mathbf{i}_2 \leftarrow 1 + \mathbf{i}_3]^6); \\ [\mathbf{z}_1 \leftarrow \mathbf{x}_3 + 2]^7 \end{array}, \quad \text{aVar}_{\mathbf{c2}} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_3 \\ \mathbf{y}_1 \\ \mathbf{x}_2 \\ \mathbf{z}_1 \\ \mathbf{i}_1 \\ \mathbf{i}_2 \\ \mathbf{i}_3 \end{bmatrix}, \quad \mathbf{M}_{\mathbf{c2}} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}, \quad \mathbf{F}_{\mathbf{c2}} = \begin{bmatrix} 1 \\ 2 \\ 1 \\ 2 \\ 0 \\ 0 \\ 2 \\ 1 \end{bmatrix}$$

$$\mathbf{c3} \triangleq \begin{array}{l} [\mathbf{x}_1 \leftarrow \text{query}(1)]^1; \\ [\mathbf{i}_1 \leftarrow 1]^2; \\ \text{while } [\mathbf{i}_1 < 0]^3, \\ \quad [\mathbf{x}_3, \mathbf{x}_1, \mathbf{x}_2], [\mathbf{i}_3, \mathbf{i}_1, \mathbf{i}_2] \text{ do} \\ \quad \quad ([\mathbf{y}_1 \leftarrow \text{query}(2)]^3; \\ \quad \quad [\mathbf{x}_2 \leftarrow \mathbf{y}_1 + \mathbf{x}_3]^5); \\ [\mathbf{z}_1 \leftarrow \mathbf{x}_3 + 2]^6 \end{array}, \quad \text{aVar}_{\mathbf{c3}} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{i}_1 \\ \mathbf{x}_3 \\ \mathbf{i}_3 \\ \mathbf{z}_1 \end{bmatrix}, \quad \mathbf{M}_{\mathbf{c3}} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix}, \quad \mathbf{F}_{\mathbf{c3}} = \begin{bmatrix} 1 \\ 0 \\ 2 \\ 2 \\ 0 \end{bmatrix}$$

We can now look at the if statement.

$$\mathbf{c4} \triangleq \begin{array}{l} [\mathbf{x}_1 \leftarrow \text{query}(1)]^1; \\ [\mathbf{y}_1 \leftarrow \text{query}(2)]^2; \\ \text{if } (\mathbf{x}_1 + \mathbf{y}_1 == 5)^3, \\ \quad [\mathbf{x}_4, \mathbf{x}_2, \mathbf{x}_3], [], [\mathbf{y}_3, \mathbf{y}_1, \mathbf{y}_2] \\ \quad \text{then } [\mathbf{x}_2 \leftarrow \text{query}(3)]^4 \\ \quad \text{else } [\mathbf{x}_3 \leftarrow \text{query}(4)]^5; \\ \mathbf{y}_2 \leftarrow 2) \\ [\mathbf{z}_1 \leftarrow \mathbf{x}_4 + \mathbf{y}_3]^6 \end{array}, \quad \text{aVar}_{\mathbf{c4}} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{y}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \mathbf{y}_2 \\ \mathbf{x}_4 \\ \mathbf{y}_3 \\ \mathbf{z}_1 \end{bmatrix}, \quad \mathbf{M}_{\mathbf{c4}} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 \end{bmatrix}, \quad \mathbf{F}_{\mathbf{c4}} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

4.3 Adaptivity Based on Program Analysis in AdaptFun

Definition 41 (Program-Based Dependency Graph). .

Given a program \mathbf{c} with its assigned variables aVar of length N , s.t., $\Gamma \vdash_{M_c, F_c} \mathbf{c}$, its program-based graph $G(\mathbf{c}) = (V, E, W, F)$ is defined as:

$$\begin{aligned} \text{Vertices } V &:= \{\mathbf{x} \mid \mathbf{x} = \text{aVar}(i); i = 1, \dots, N\} \\ \text{Directed Edges } E &:= \{(\mathbf{x}_1, \mathbf{x}_2) \mid (\mathbf{x}_1 = \text{aVar}(i) \wedge \mathbf{x}_2 = \text{aVar}(j) \wedge M_c(i, j) \geq 1); i, j = 1, \dots, N\} \\ \text{Weights } W &:= \bigcup \left\{ (v, w) \in V \times (\mathbb{N} \cup e) \mid v \in V \wedge F(v) > 0 \wedge w = \text{RechBound}(v, c) \right\} \\ &\quad \bigcup \left\{ (v, 1) \in V \times \{1\} \mid v \in V \wedge F(v) = 0 \right\} \\ \text{Query Flags } qF &:= \{(\mathbf{x}, n) \in V \times \{0, 1\} \mid \left\{ \begin{array}{ll} n = 1 & F_c(i) = 2 \\ n = 0 & o.w. \end{array} \right\}; \mathbf{x} = \text{aVar}(i); i = 1, \dots, N\} \end{aligned}$$

Definition 42 (Finite Walk (k)). .

Given a labeled weighted graph $G = (V, E, W, qF)$, a finite walk k in G is a sequence of edges $(e_1 \dots e_{n-1})$ for which there is a sequence of vertices (v_1, \dots, v_n) such that:

- $e_i = (v_i, v_{i+1})$ for every $1 \leq i < n$.
- every vertex $v \in V$ appears in this vertices sequence (v_1, \dots, v_n) of k at most $W(v)$ times.

(v_1, \dots, v_n) is the vertex sequence of this walk.

[JL: Length of this finite walk k is the number of vertices in its vertex sequence, i.e., $\text{len}(k) = n$.]

[JL: Given a labeled weighted graph $G = (V, E, W, qF)$, we use $\mathcal{WALK}(G)$ to denote a set containing all finite walks k in G ; and $k_{v_1 \rightarrow v_2} \in \mathcal{WALK}(G)$ where $v_1, v_2 \in V$ denotes the walk from vertex v_1 to v_2 .]

Definition 43 (Length of Finite Walk w.r.t. Query (len_q)). .

Given a labeled weighted graph $G = (V, E, W, qF)$ and a finite walk k in G with its vertex sequence (v_1, \dots, v_n) , the length of k w.r.t query is defined as:

$$\text{len}_q(k) = \text{len}(v \mid v \in (v_1, \dots, v_n) \wedge qF(v) = 2)$$

, where $(v \mid v \in (v_1, \dots, v_n) \wedge qF(v) = 2)$ is a subsequence of k 's vertex sequence.

Given a program \mathbf{c} , we generate its program-based graph $\mathbf{G}_{\text{prog}}(\mathbf{c}) = (V, E, W, qF)$. Then the adaptivity bound based on program analysis for \mathbf{c} is the number of query vertices on a finite walk in $\mathbf{G}_{\text{prog}}(\mathbf{c})$. This finite walk satisfies:

- the number of query vertices on this walk is maximum
- the visiting times of each vertex v on this walk is bound by its reachability bound $W(v)$.

It is formally defined in 44.

Definition 44 (Program-Based Adaptivity). .

Given a program \mathbf{c} and its program-based graph $\mathbf{G}_{\text{prog}}(\mathbf{c}) = (V, E, W, qF)$, the program-based adaptivity for c is defined as

$$A_{\text{prog}}(\mathbf{c}) := \max\{\text{len}_q(k) \mid k \in \mathcal{WALK}(\mathbf{G}_{\text{prog}}(\mathbf{c}))\}.$$

4.4 [[Soundness of the AdaptFun]]

[JL:

Theorem 4.1 (Soundness of the AdaptFun). . *Given a program \mathbf{c} , we have:*

$$A_{\text{prog}}(\mathbf{c}) \geq A(\mathbf{c}).$$

] [JL:

Proof. Given a program \mathbf{c} , we construct its program-based graph $\mathbf{G}_{\text{prog}}(\mathbf{c}) = (V, E, W, q_F)$ by Definition 41. According to the Definition 44, we have:

$$A_{\text{prog}}(\mathbf{c}) := \max\{\text{len}_q(k) \mid k \in \mathcal{WALK}(\mathbf{G}_{\text{prog}}(\mathbf{c}))\}.$$

According to the Definition ??, we have the trace-based adaptivity as follows:

$$A(\mathbf{c}) = \max\{\text{len}(p) \mid \mathbf{m} \in \mathcal{SM}, D \in \mathcal{DB}, p \in \mathcal{PATH}(\mathbf{G}_{\text{trace}}(\mathbf{c}, D, \mathbf{m}))\}$$

Then, we need to show:

$$\max\{\text{len}(p) \mid \mathbf{m} \in \mathcal{SM}, D \in \mathcal{DB}, p \in \mathcal{PATH}(\mathbf{G}_{\text{trace}}(\mathbf{c}, D, \mathbf{m}))\} \leq \max\{\text{len}_q(\mathbf{k}) \mid \mathbf{k} \in \mathcal{WALK}(\mathbf{G}_{\text{prog}}(\mathbf{c}))\}$$

It is sufficient to show that:

$$\forall p, \mathbf{m}, D, \text{ s.t., } p \in \mathcal{PATH}(\mathbf{G}_{\text{trace}}(\mathbf{c}, D, \mathbf{m})), \exists \mathbf{k} \in \mathcal{WALK}(\mathbf{G}_{\text{prog}}(\mathbf{c})) \wedge \text{len}(\mathbf{p}) \leq \text{len}_q(\mathbf{k})$$

Taking an arbitrary starting memory m and an arbitrary underlying database D , we construct a trace-based graph $\mathbf{G}_{\text{trace}}(\mathbf{c}, D, \mathbf{m}) = (V, E)$ by the definition ??.

Let $\mathbf{G}_{\text{mid}}(\mathbf{c}, \mathbf{m}, D) = \{\mathbf{V}_{\text{mid}}, \mathbf{E}_{\text{mid}}, \mathbf{qF}_{\text{mid}}\}$ be the intermediate graph by Definition 45.

By Lemma 21, we know:

$$\forall p, \mathbf{m}, D, \text{ s.t., } p \in \mathcal{PATH}(\mathbf{G}_{\text{trace}}(\mathbf{c}, D, \mathbf{m})), \exists \mathbf{p}' \in \mathcal{PATH}(\mathbf{G}_{\text{mid}}(\mathbf{c}, \mathbf{m}, D)) \wedge \text{len}(\mathbf{p}) = \text{len}_q(\mathbf{p}')$$

Then it is sufficient to show that:

$$\forall p, \mathbf{m}, D, \text{ s.t., } p \in \mathcal{PATH}(\mathbf{G}_{\text{mid}}(\mathbf{c}, D, \mathbf{m})), \exists \mathbf{k} \in \mathcal{WALK}(\mathbf{G}_{\text{prog}}(\mathbf{c})) \wedge \text{len}_q(p) \leq \text{len}_q(\mathbf{k})$$

We prove a stronger statement instead:

$$\forall p, \mathbf{m}, D, \text{ s.t., } p \in \mathcal{PATH}(\mathbf{G}_{\text{mid}}(\mathbf{c}, D, \mathbf{m})), \exists \mathbf{k} \in \mathcal{WALK}(\mathbf{G}_{\text{prog}}(\mathbf{c})) \wedge \text{len}_q(p) = \text{len}_q(\mathbf{k})$$

By Lemma 22, let g be the surjective function $g : \mathbf{V}_{\text{prog}} \rightarrow \mathbf{V}_{\text{mid}}$ s.t.:

$$\forall \text{av} \in \mathbf{V}_{\text{mid}}. \mathbf{qF}_{\text{prog}}(f(\text{av})) = \mathbf{qF}_{\text{mid}}(\text{av}) \wedge |\text{image}(f(\text{av}))| \leq W(f(\text{av})).$$

Let \mathbf{m} and D be an arbitrary memory and database D , taking an arbitrary path $p_{\text{av}_1 \rightarrow \text{av}_n} \in \mathcal{PATH}(\mathbf{G}_{\text{mid}}(\mathbf{c}, D, \mathbf{m}))$ with:

Edge sequence: (e, \dots, e_{n-1})

Vertices sequence: (av_1, \dots, av_n) .

By Lemma 24, let $h : \mathcal{PATH}(\mathbf{G}_{\text{mid}}(\mathbf{c}, D, \mathbf{m})) \rightarrow \mathcal{WALK}(\mathbf{G}_{\text{prog}}(\mathbf{c}))$ be the surjective function satisfies:

$$\forall p_{av_1 \rightarrow av_n} \in \mathcal{PATH}(\mathbf{G}_{\text{mid}}(\mathbf{c}, D, \mathbf{m})) \text{ with } \begin{cases} \text{edge sequence:} & (e, \dots, e_{n-1}) \\ \text{vertices sequence:} & (av_1, \dots, av_n) \end{cases}$$

$$\exists k_{f(av_1) \rightarrow f(av_n)} \in \mathcal{WALK}(\mathbf{G}_{\text{prog}}(\mathbf{c})) \text{ with } \begin{cases} \text{edge sequence:} & (g(e), \dots, g(e_{n-1})) \\ \text{vertices sequence:} & (f(av_1), \dots, f(av_n)) \end{cases}$$

We have the walk: $k_{f(av_1) \rightarrow f(av_n)} \in \mathcal{WALK}(\mathbf{G}_{\text{prog}}(\mathbf{c}))$ with:

Edges sequence: $(g(e), \dots, g(e_{n-1}))$

Vertices sequence: $(f(av_1), \dots, f(av_n))$.

It is sufficient to show

$$\text{len}_q(p_{av_1 \rightarrow av_n}) = \text{len}_q(k_{f(av_1) \rightarrow f(av_n)})$$

Unfold the definition of len_q , it is suffice to show:

$$\text{len}(av \mid av \in (av_1, \dots, av_n) \wedge \mathbf{qF}_{\text{mid}}(av) = 2) = \text{len}(f(av) \mid f(av) \in (f(av_1), \dots, f(av_n)) \wedge \mathbf{qF}_{\text{prog}}(f(av)) = 2) \quad (a)$$

By Lemma 22, we know:

$$\forall av \in V_{\text{mid}}. \mathbf{qF}_{\text{mid}}(av) = \mathbf{qF}_{\text{prog}}(f(av)) \quad (b)$$

By rewriting (b) in (a), we have this case proved.

[[

Definition 45 (Intermediate Graph \mathbf{G}_{mid}). .

\mathcal{AV} : Annotated Variables based on program execution

Given a program \mathbf{c} with its assigned variables \mathbf{aVar} of length N , a database D , a starting memory \mathbf{m} , s.t., $\Gamma \vdash_{\mathbf{M}_c, \mathbf{F}_c} \mathbf{c}$, the intermediate graph $\mathbf{G}_{\text{mid}}(\mathbf{c}, \mathbf{m}, D) = (V, E, F)$ is defined as:

$$\begin{aligned} \text{Vertices } V &:= \{av \in \mathcal{AV} \mid \exists \mathbf{m}', w', qt, \tau. \text{ s.t., } \langle \mathbf{m}, \mathbf{c}, [], [], [] \rangle \rightarrow^* \langle \mathbf{m}', \text{skip}, qt, \tau, w' \rangle \wedge av \in \tau\} \\ \text{Directed Edges } E &:= \{(av, av') \in \mathcal{AV} \times \mathcal{AV} \mid \text{flowsTo}(av, av', \mathbf{c}, \mathbf{m}, D)\} \\ \text{Flags } F &:= \{(av, n) \in V \times \{0, 1, 2\} \mid (\pi_1(av) = \mathbf{aVar}(i) \wedge n = F_c(i)); i = 1, \dots, N\} \end{aligned}$$

]]

[[

Lemma 16 (DEP_{var} is Transitive). .

Given a program \mathbf{c} , with a starting memory \mathbf{m} and a hidden database D , s.t., $\langle \mathbf{m}, \mathbf{c}, [], [], [] \rangle \rightarrow^* \langle \mathbf{m}', \text{skip}, qt, \tau, w \rangle$. Then, $\forall av_1, av_2, av_3 \in \tau$:

$$\left(\text{DEP}_{\text{var}}(av_1, av_2, \mathbf{c}, \mathbf{m}, D) \wedge \text{DEP}_{\text{var}}(av_2, av_3, \mathbf{c}, \mathbf{m}, D) \right) \implies \text{DEP}_{\text{var}}(av_1, av_3, \mathbf{c}, \mathbf{m}, D)$$

of Lemma 16. Proof by unfolding and rewriting the Definition 26. ■

]]

[[

Lemma 17 (flowsTo is Transitive ??). .

Given a program \mathbf{c} with its assigned variables \mathbf{aVar} of length N . Then $\forall x_1, x_2, x_3 \in \mathbf{aVar}$

$$\left(\text{flowsTo}(x_1, x_2) \wedge \text{flowsTo}(x_2, x_3) \right) \implies \text{flowsTo}(x_1, x_3)$$

of Lemma 17. Proof by unfolding the Definition 33. ■

]]

[[

Lemma 18 (DEP_q Implies DEP_{var}). .

Given a program \mathbf{c} , with a starting memory \mathbf{m} and a hidden database D , s.t., $\langle \mathbf{m}, \mathbf{c}, [], [], [] \rangle \rightarrow^* \langle \mathbf{m}', \text{skip}, \text{qt}, \tau, w \rangle$. Then, $\forall \text{av}_1, \text{av}_2 \in \text{qt}$

$$\text{DEP}_q(\text{av}_1, \text{av}_2, \mathbf{c}, \mathbf{m}, D) \implies \text{DEP}_{\text{var}}(\pi_2(\text{av}_1), \pi_2(\text{av}_2), \mathbf{c}, \mathbf{m}, D)$$

of Lemma 18. Proof by unfolding the Definition 26 and Definition ?? ■

]]

[[

Lemma 19 (DEP_{var} Implies flowsTo). .

Given a program \mathbf{c} , with a starting memory \mathbf{m} and a hidden database D , s.t., $\langle \mathbf{m}, \mathbf{c}, [], [], [] \rangle \rightarrow^* \langle \mathbf{m}', \text{skip}, \text{qt}, \tau, w \rangle$. Then, $\forall \text{av}_1, \text{av}_2 \in \tau$

$$\text{DEP}_{\text{var}}(\text{av}_1, \text{av}_2, \mathbf{c}, \mathbf{m}, D) \implies \text{flowsTo}(\pi_1(\text{av}_1), \pi_1(\text{av}_2))$$

of Lemma 18. Proof by showing contradiction based on the Definition 26 and Definition 33. Let $\text{av}_1, \text{av}_2 \in \tau$ be 2 arbitrary annotated variables in the variable trace τ , s.t., $\text{DEP}_{\text{var}}(\text{av}_1, \text{av}_2, \mathbf{c}, \mathbf{m}, D)$. Unfolding the DEP_{var} definition, we have: ■

]]

[[

Lemma 20 (Injective Mapping of vertices from $\mathbf{G}_{\text{trace}}$ to \mathbf{G}_{mid}). .

$$\mathbf{G}_{\text{trace}}(\mathbf{c}) = \{\mathbf{V}_{\text{trace}}, \mathbf{E}_{\text{trace}}\}$$

$$\mathbf{G}_{\text{mid}}(\mathbf{c}, \mathbf{m}, D) = \{\mathbf{V}_{\text{mid}}, \mathbf{E}_{\text{mid}}, \mathbf{qF}_{\text{mid}}\}$$

$$\exists \text{ injective } f: \mathcal{AQ} \rightarrow \mathcal{AV}. \forall \text{av} \in \mathbf{V}_{\text{trace}}. f(\text{av}) \in \mathbf{V}_{\text{mid}} \wedge \mathbf{qF}_{\text{mid}}(f(\text{av})) = 2$$

Proof. Proving by Definition 45 and Definition 44. ■

]]

[[

Lemma 21 (One-on-One Mapping from E of $\mathbf{G}_{\text{trace}}$ to $\mathcal{PATH}(\mathbf{G}_{\text{mid}})$). .

$$\mathbf{G}_{\text{trace}}(\mathbf{c}) = \{\mathbf{V}_{\text{trace}}, \mathbf{E}_{\text{trace}}\}$$

$$\mathbf{G}_{\text{mid}}(\mathbf{c}, \mathbf{m}, D) = \{\mathbf{V}_{\text{mid}}, \mathbf{E}_{\text{mid}}, \mathbf{qF}_{\text{mid}}\}$$

An injective function $f: \mathbf{V}_{\text{trace}} \rightarrow \mathbf{V}_{\text{mid}}$ s.t., $\forall \text{av} \in \mathbf{V}_{\text{trace}}. \mathbf{qF}_{\text{mid}}(f(\text{av})) = 2$

$$\forall e = (\text{av}_1, \text{av}_2) \in \mathbf{E}_{\text{trace}}. \exists p_{f(\text{av}_1) \rightarrow f(\text{av}_2)} \in \mathcal{PATH}(\mathbf{G}_{\text{mid}}(\mathbf{c}, D, \mathbf{m}))$$

Proof. Proving by Lemma 20 and Definition 45 and acyclic property of $\mathbf{G}_{\text{trace}}$ and \mathbf{G}_{mid} . ■

]]

[[

Lemma 22 (Surjective Mapping of Vertices from \mathbf{G}_{mid} to \mathbf{G}_{prog}). .

$$\mathbf{G}_{\text{mid}}(\mathbf{c}, \mathbf{m}, D) = \{\mathbf{V}_{\text{mid}}, \mathbf{E}_{\text{mid}}, \mathbf{qF}_{\text{mid}}\}$$

$$\mathbf{G}_{\text{prog}}(\mathbf{c}) = \{\mathbf{V}_{\text{prog}}, \mathbf{E}_{\text{prog}}, \mathbf{qF}_{\text{prog}}, \mathbf{W}_{\text{prog}}\}$$

\exists surjective $f: \mathcal{AV} \rightarrow \mathcal{SVA}\mathcal{R}$.

$$\forall \text{av} \in \mathbf{V}_{\text{mid}}. f(\text{av}) \in \mathbf{V}_{\text{prog}} \wedge \mathbf{qF}_{\text{prog}}(f(\text{av})) = \mathbf{qF}_{\text{mid}}(\text{av}) \wedge |\text{image}(f(\text{av}))| \leq W(f(\text{av}))$$

Proof. Proving by Definition 45. ■

]]

[[

Lemma 23 (Surjective Mapping from E of \mathbf{G}_{mid} to E of \mathbf{G}_{prog}). .

$$\mathbf{G}_{\text{mid}}(\mathbf{c}, \mathbf{m}, D) = \{\mathbf{V}_{\text{mid}}, \mathbf{E}_{\text{mid}}, \mathbf{qF}_{\text{mid}}\}$$

$$\mathbf{G}_{\text{prog}}(\mathbf{c}) = \{\mathbf{V}_{\text{prog}}, \mathbf{E}_{\text{prog}}, \mathbf{qF}_{\text{prog}}, \mathbf{W}_{\text{prog}}\}$$

A surjective function $f: \mathbf{V}_{\text{prog}} \rightarrow \mathbf{V}_{\text{mid}}$ s.t., $\forall \text{av} \in \mathbf{V}_{\text{mid}}. \mathbf{qF}_{\text{prog}}(f(\text{av})) = \mathbf{qF}_{\text{mid}}(\text{av}) \wedge |\text{image}(f(\text{av}))| \leq W(f(\text{av}))$

$$\exists \text{ surjective } g: \mathbf{E}_{\text{mid}} \rightarrow \mathbf{E}_{\text{prog}}. \forall e_{\text{mid}} = (\text{av}_1, \text{av}_2) \in \mathbf{E}_{\text{mid}}. \exists e_{\text{prog}} = (f(\text{av}_1), f(\text{av}_2)) \in \mathbf{E}_{\text{prog}}$$

Proof. Proving by Lemma 22. ■

]]

[[

Lemma 24 (Surjective Mapping from $\mathcal{PATH}(\mathbf{G}_{\text{mid}})$ to $\mathcal{WALK}(\mathbf{G}_{\text{prog}})$). .

$$\mathbf{G}_{\text{mid}}(\mathbf{c}, \mathbf{m}, D) = \{\mathbf{V}_{\text{mid}}, \mathbf{E}_{\text{mid}}, \mathbf{qF}_{\text{mid}}\}$$

$$\mathbf{G}_{\text{prog}}(\mathbf{c}) = \{\mathbf{V}_{\text{prog}}, \mathbf{E}_{\text{prog}}, \mathbf{qF}_{\text{prog}}, \mathbf{W}_{\text{prog}}\}$$

A surjective function $f: \mathbf{V}_{\text{prog}} \rightarrow \mathbf{V}_{\text{mid}}$ s.t., $\forall \text{av} \in \mathbf{V}_{\text{mid}}. \mathbf{qF}_{\text{prog}}(f(\text{av})) = \mathbf{qF}_{\text{mid}}(\text{av}) \wedge |\text{image}(f(\text{av}))| \leq W(f(\text{av}))$

A surjective function $g: \mathbf{E}_{\text{mid}} \rightarrow \mathbf{E}_{\text{prog}}$ s.t., $\forall e_{\text{mid}} = (\text{av}_1, \text{av}_2) \in \mathbf{E}_{\text{mid}}. \exists e_{\text{prog}} = (f(\text{av}_1) \rightarrow f(\text{av}_2)) \in \mathbf{E}_{\text{prog}}$

\exists surjective $h: \mathcal{PATH}(\mathbf{G}_{\text{mid}}(\mathbf{c}, \mathbf{m}, D)) \rightarrow \mathcal{WALK}(\mathbf{G}_{\text{prog}}(\mathbf{c}))$ s.t.:

$$\forall p_{\text{av}_1 \rightarrow \text{av}_2} \in \mathcal{PATH}(\mathbf{G}_{\text{mid}}(\mathbf{c}, \mathbf{m}, D)) \text{ with } \begin{cases} \text{edge sequence:} & (e, \dots, e_{n-1}) \\ \text{vertices sequence:} & (\text{av}_1, \dots, \text{av}_n) \end{cases}$$

$$\exists k_{f(\text{av}_1) \rightarrow f(\text{av}_2)} \in \mathcal{WALK}(\mathbf{G}_{\text{prog}}(\mathbf{c})) \text{ with } \begin{cases} \text{edge sequence:} & (g(e), \dots, g(e_{n-1})) \\ \text{vertices sequence:} & (f(\text{av}_1), \dots, f(\text{av}_n)) \end{cases}$$

Proof. Proving by induction on the length of $l = p_{\text{av}_1 \rightarrow \text{av}_2} \in \mathcal{PATH}(\mathbf{G}_{\text{mid}}(\mathbf{c}, \mathbf{m}, D))$, and Lemma 23 and Lemma 22.

case: $l = 1$:

case: $l = l' + 1, l' \geq 1$:

]]

]

□

5 [[Examples]]

Example 5.1 (TwoRound Algorithm).

$$\begin{array}{ll}
 \begin{array}{l}
 [i \leftarrow 1]^1; \\
 [a_1 \leftarrow []]^2; \\
 \text{while } [i < k]^3, 0, \text{ do} \\
 \quad [x \leftarrow \text{query}()^4]; \\
 \quad [a \leftarrow x :: a]^5 [i_2 \leftarrow i_3 + 1]^6; \\
 \quad [l \leftarrow q_{k+1}(a)]^7
 \end{array}
 & \Rightarrow \quad
 \begin{array}{l}
 [i \leftarrow 1]^1; \\
 [a_1 \leftarrow []]^2; \\
 \text{while } [i < k]^3, 0, [a_3, a_1, a_2][i_3, i_1, i_2] \text{ do} \\
 \quad [x_1 \leftarrow q]^4; \\
 \quad [a_2 \leftarrow x_1 :: a_3]^5 [i_2 \leftarrow i_3 + 1]^6; \\
 \quad [l \leftarrow q_{k+1}(a_3)]^7
 \end{array}
 \end{array}
 \triangleq TR(k) \quad \Rightarrow \quad TR^{ssa}$$

$\text{Adapt}(TR) = 2$

Using **AdaptFun**, we first generate a assigned variables G from an empty list $[]$ and empty while map $[]$.

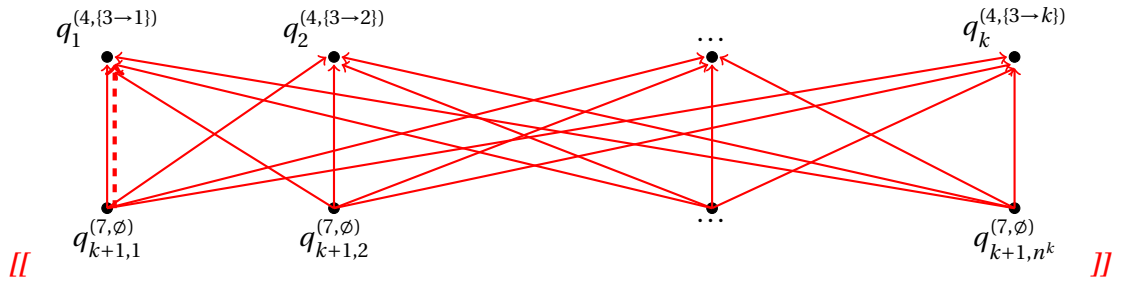
$$[]; []; TR^{ssa} \rightarrow G; w \wedge w = []$$

$$G_{k=2} = [a_1^2, a_3^{(2,[2:1])}, x_1^4, a_2^5, i_3^3, i_2^6, i_1^2, l_1^7, l_1^{(5,[1])}]$$

We denote a_1^1 short for $a_1^{(1,[1])}$ and $a_3^{(2,1)}$ short for $a_3^{(2,[2:1])}$, where the label $(2,1)$ represents at line number 2 and in the 1 st iteration.

$$M = \begin{bmatrix} a_1^2 & a_3^3 & x_1^4 & a_2^5 & i_3^3 & i_2^6 & i_1^2 & l_1^7 \\ a_1^2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ a_3^3 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ x_1^4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ a_2^5 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ i_3^3 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ i_2^6 & 0 & 0 & 0 & 0 & 1 & 1 & 0 \\ i_1^2 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ l_1^7 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}, V = \begin{bmatrix} a_1^2 & 2 \\ a_3^3 & 2 \\ x_1^4 & 1 \\ a_2^5 & 2 \\ i_3^3 & 2 \\ i_2^6 & 2 \\ i_1^2 & 2 \\ l_1^7 & 1 \end{bmatrix}$$

$$M = \begin{bmatrix} a_1^1 & a_3^{(2,1)} & x_1^{(3,1)} & a_2^{(4,1)} & a_3^{(2,2)} & x_1^{(3,2)} & a_2^{(4,2)} & a_3^2 & l_1^5 \\ a_1^1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ a_3^{(2,1)} & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ x_1^{(3,1)} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ a_2^{(4,1)} & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ a_3^{(2,2)} & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ x_1^{(3,2)} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ a_2^{(4,2)} & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 \\ a_3^2 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ l_1^5 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}, V = \begin{bmatrix} a_1^1 & 0 \\ a_3^{(2,1)} & 0 \\ x_1^{(3,1)} & 1 \\ a_2^{(4,1)} & 0 \\ a_3^{(2,2)} & 0 \\ x_1^{(3,2)} & 1 \\ a_2^{(4,2)} & 0 \\ a_3^2 & 0 \\ l_1^5 & 1 \end{bmatrix}$$



Example 5.2 (Multi-Round Algorithm).

$$\begin{aligned}
 MR \triangleq & \begin{pmatrix} [i \leftarrow 1]^1; \\ [I \leftarrow []]^2; \\ \text{while } [i < k]^3 \text{ do} \\ \quad [p \leftarrow c]^4; \\ \quad [a \leftarrow \text{query}(p, I)]^5; \\ \quad [I \leftarrow \text{update}(I, (a, p))]^6; \\ \quad [i \leftarrow i + 1]^7 \\ \end{pmatrix} \Rightarrow MR^{ssa} \triangleq \begin{pmatrix} [i \leftarrow 1]^1; \\ [I \leftarrow []]^2; \\ \text{while } [i < k]^3 0, [I_3, I_1, I_2] \\ \quad \text{do} \\ \quad \quad [p_1 \leftarrow c]^4; \\ \quad \quad [a \leftarrow \text{query}(p_1, I_2)]^5; \\ \quad \quad [I_2 \leftarrow \text{update}(I_3, (a_1, p_1))]^6; \\ \quad \quad [i \leftarrow i + 1]^7 \\ \end{pmatrix}
 \end{aligned}$$

$\text{Adapt}(MR) = k$.

Using **AdaptFun**, we first generate a assigned variables G from an empty list $[]$ and empty whlemap \emptyset .

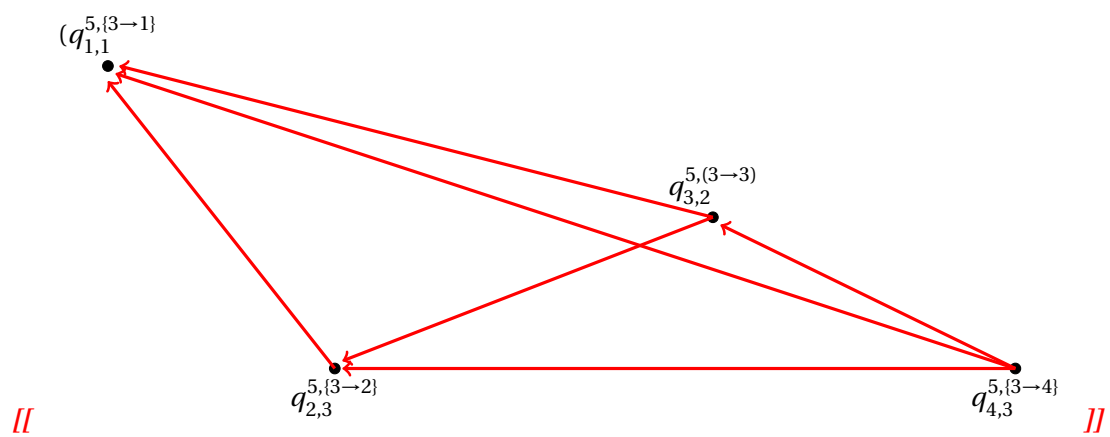
$$[]; \emptyset; MR^{ssa} \rightarrow G; w \wedge w = \emptyset$$

.

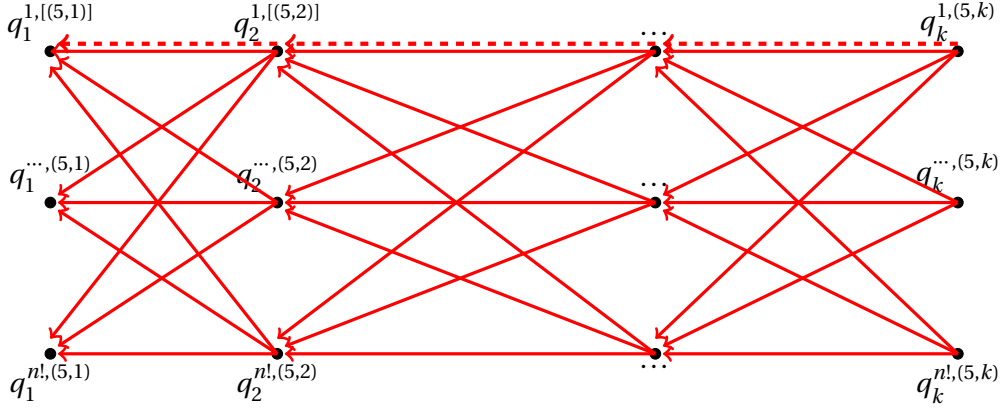
$$G_{k=2} = [i_1^1, I_1^2, i_3^3, I_3^3, p_1^4, a_1^5, I_2^6, i_2^7]$$

We denote I_1^1 short for $I_1^{(1, \emptyset)}$ and $I_3^{(2, 1)}$ short for $I_3^{(2, [2:1])}$, where the label $(2, 1)$ represents at line number 2 and in the 1 st iteration.

$$M = \begin{bmatrix} i_1^1 & I_1^2 & i_3^3 & I_3^3 & p_1^4 & a_1^5 & I_2^6 & i_2^7 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}, V = \begin{bmatrix} i_1^1 & 0 \\ I_1^2 & 0 \\ i_3^3 & 2 \\ I_3^3 & 2 \\ p_1^4 & 2 \\ a_1^5 & 1 \\ I_2^6 & 2 \\ i_2^7 & 2 \end{bmatrix}$$



$\forall k. \forall D$, we have $A(TR^L) = (k-1)$ given all possible execution traces. //



//

6 Non Determinism

Non-Determinism of queries. When evaluating a query $\text{query}(\alpha)$ on a given database D , in addition to obtain a result v from the database $v = \text{query}(\alpha)(D)$, we assume there is an underlying mechanism that will perform extra manipulations on v . The mechanism is considered as primitive operations in our language, behaving as black box to programmers. There are different kinds of mechanisms, such as adding noise sampled from certain probabilistic distribution to the result [1]. Because of the randomness of the underlying mechanism, the evaluation of a query $\text{query}(\alpha)$ is non-deterministic. That's the reason, in the Definition ??, given a fixed database D , there will be a query domain \mathcal{QD} where $\text{query}(\alpha)(D)$ can be evaluated to different values $v \in \mathcal{QD}$.

On the other hand, in the operational semantics rule **query-v**:

$$\frac{\text{query}(\alpha) = v}{\langle m, [x \leftarrow \text{query}(\alpha)]^l, t, w \rangle \rightarrow \langle m[v/x], \text{skip}, (t + +[(\alpha, l, w)]), w \rangle} \text{query-v}$$

, we evaluate the query given database D based on an assumption that the underlying mechanism is fixed. This fixed mechanism only adds constant 0 to the original result v returned from the database, i.e., $v = \text{query}(\alpha)(D)$.

The Lemma 25 and 26 formalize this property.

Lemma 25 (Semi-Determinism). .

for any program c with a starting memory m , trace t and while label w , if program c contains neither $[x \leftarrow \text{query}(\psi)]^l$ nor $[x \leftarrow \text{query}(\alpha)]^l$ for any ψ and α , then

$$\bigwedge \left\{ \begin{array}{l} \langle m, c, t, w \rangle \rightarrow^* \langle m_1, \text{skip}, t_1, w_1 \rangle \\ \langle m, c, t, w \rangle \rightarrow^* \langle m_2, \text{skip}, t_2, w_2 \rangle \end{array} \right\} \implies (m_1 = m_2 \wedge t_1 = t_2 \wedge w_1 = w_2)$$

Proof. Proof is obvious by induction on the operational semantics rules. □

Lemma 26 (Query Semi-Determinism). .

Given a program $c; x \leftarrow \text{query}(\psi); c'$ with a starting memory m , trace t and while label w , s.t. c contains neither $[x \leftarrow \text{query}(\psi)]^l$ nor $[x \leftarrow \text{query}(\alpha)]^l$ for any ψ and α , then:

$$\bigwedge \left\{ \begin{array}{l} \langle m, c; x \leftarrow \text{query}(\psi); c', t, w \rangle \rightarrow^* \langle m_1, x \leftarrow \text{query}(\alpha_1); c', t_1, w_1 \rangle \\ \langle m, c; x \leftarrow \text{query}(\psi); c', t, w \rangle \rightarrow^* \langle m_2, x \leftarrow \text{query}(\alpha_2); c', t_2, w_2 \rangle \end{array} \right\} \implies (\alpha_1 = \alpha_2 \wedge m_1 = m_2 \wedge t_1 = t_2 \wedge w_1 = w_2)$$

Proof. Proof is obvious by induction on the operational semantics rules. \square

7 Analysis of Generalization Error

Example 7.1 (Two Round Algorithm).

$$\begin{aligned}
 & [a_1 \leftarrow []]^1; \\
 & \text{loop } [k]^2 \ (a_2 \leftarrow f(1, a_1, a_3)) \\
 & \quad \text{do} \\
 TR^H(k) \triangleq & \left([x_1 \leftarrow \text{query}()]^3; \right. \\
 & [a_3 \leftarrow x_1 :: a_2]^4; \\
 & \left. [l \leftarrow q_{k+1}(a_3)]^5 \right)
 \end{aligned}$$

Example 7.2 (Multi-Round Algorithm).

$$\begin{aligned}
 & [I_2 \leftarrow []]^1; \\
 & \text{loop } [k]^2 \ (I_2 \leftarrow f(2, I_1, I_3)) \\
 & \quad \text{do} \\
 MR^H \triangleq & \left([p_1 \leftarrow c]^3; \right. \\
 & [a_1 \leftarrow \delta(\text{query}(p, I_2))]^4; \\
 & \left. [I_3 \leftarrow \text{update}(I_2, (a_1, p))]^5 \right)
 \end{aligned}$$

By applying different mechanisms $\delta()$ over the queries $\text{query}(\cdot)$, we have different error bounds.

Gaussian Mechanism: $N(0, \sigma)$ [1]:

Adaptivity $r = 2$: $\sigma = O\left(\frac{\sqrt{r \log(k)}}{\sqrt{n}}\right)$ (also known as expected error);

Adaptivity unknown: $\sigma = O\left(\frac{\sqrt[4]{k}}{\sqrt{n}}\right)$;

Mean Squared Error Bound: $\frac{1}{2n} \min_{\lambda \in [0,1]} \left(\frac{2\rho k n - \ln(1-\lambda)}{\lambda} \right) + 2\mathbb{E}_{Z_i \sim N(0, \frac{1}{2n^2\rho})} \left[\max_{i \in [k]} (Z_i^2) \right]$

Confidence Bounds: minimize τ where $\tau \geq \sqrt{\frac{2}{n\beta} \min_{\lambda \in [0,1]} \left(\frac{2\rho k n - \ln(1-\lambda)}{\lambda} \right)}$ and $\tau \geq \frac{2}{n} \sqrt{\frac{\ln(4n/\beta)}{\rho'}}$ with confidence level $1 - \beta$.

(ϵ, δ) - DP mechanism:

Confidence Bounds: $\tau \geq \sqrt{\frac{48}{n} \ln(4/\beta)}$ with $\epsilon \leq \frac{\tau}{4}$ and $\delta = \exp\left(\frac{-4 \ln(8/\beta)}{\tau}\right)$

Sample Splitting:

Expected Error: $O\left(\frac{\sqrt{k \log(k)}}{\sqrt{n}}\right)$

Thresholdout: B, σ, T, h

Confidence bounds: $\tau = \max \left\{ \sqrt{\frac{2\zeta}{h\beta}}, 2\sigma \ln\left(\frac{\beta}{2}\right), \sqrt{\frac{1}{\beta}} \cdot \left(\sqrt{T^2 + 56\sigma^2} + \sqrt{\frac{\zeta}{4h}} \right) \right\}$, for $\zeta = \min_{\lambda \in [0,1]} \left(\frac{2B(\sigma^2 h) - \ln(1-\lambda)}{\lambda} \right)$

References

- [1] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Leon Roth. Preserving statistical validity in adaptive data analysis. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pages 117–126, 2015.
- [2] Sumit Gulwani and Florian Zuleger. The reachability-bound problem. In *Proceedings of the 31st ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '10*, page 292–304, New York, NY, USA, 2010. Association for Computing Machinery.