

Technical Report COMP-001-2009
Computing Department
Lancaster University

From stack traces to call-trees: outline of a proof

François Taïani

January 22, 2009

Abstract

This report outlines a proof of the theorem described in COSMOPEN: *Dynamic reverse engineering on a budget, How cheap observation techniques can be used to reconstruct complex multi-level behaviour* [1] on the construction of a minimal call-trees from stack traces. Please refer to original publication for details regarding the context of the work, the notation used, and definitions. We repeat below the pseudo-code of the call-tree construction algorithm, and present some key definitions in more detail. We then outline a proof of the characterisation of the call-tree constructed by the algorithm.

Computing Department
Infolab21, South Drive
Lancaster University
LANCASTER LA1 4WA
United Kingdom
Tel.: +44 (0) 1524 51 03 38
mail: francois.taiani@comp.lancs.ac.uk

1 Introduction

This short document outlines a proof of the theorem described in [1] on the construction of a minimal call-trees from stack traces. Please refer to [1] for details regarding the context of the work, the notation used, and definitions. We repeat below the pseudo-code of the call-tree construction algorithm, and present some key definitions in more detail. We then outline a proof of the characterisation of the call-tree constructed by the algorithm.

2 Call-tree construction algorithm

Algorithm 1 Transforming a stack trace sequence into a call-tree

```
1: ActiveTreePath  $\leftarrow$  (); V  $\leftarrow$  E  $\leftarrow$   $\emptyset$ 
2: traceSequence = (trace1, trace2, ..., traceN)
3: for  $i = 1$  to |traceSequence| do
4:   AddTraceToGraph(traceSequence[i], V, E, ActiveTreePath)
5: end for
6: procedure AddTraceToGraph(trace, V, E, ActiveTreePath)
7:   for  $j = 1$  to  $\min(|ActiveTreePath|, |trace| - 1)$  do
8:     if  $symbol(ActiveTreePath[j]) \neq trace[j]$  then break
9:   end for
10:  truncate(ActiveTreePath,  $j$ )
11:  previousNode  $\leftarrow last(ActiveTreePath)$ 
12:  for  $k = j$  to |trace| do
13:    newNode  $\leftarrow new Node(trace[k])$ 
14:    ActiveTreePath  $\leftarrow ActiveTreePath + (newNode)$ 
15:    V  $\leftarrow V \cup \{newNode\}$ 
16:    if previousNode  $\neq NIL$  then E  $\leftarrow E \cup \{(previousNode, newNode)\}$ 
17:    previousNode  $\leftarrow newNode$ 
18:  end for
19: end procedure
```

In algorithm 1, time stamps (also known as *sequence numbers*) are allocated by the constructor **new** *Node*(symbol) (line 13) in the order in which nodes are created. The same timestamp is never allocated more than once.

3 Definitions and theorem

Call tree A *call-tree* is a directed tree (V, E) that is associated with two labelling functions: *symbol* and *timeStamp*.

symbol associates each node $x \in V$ of the tree with the name of a function (its *symbol*). Multiple nodes might be associated with the same symbol.

timeStamp associates each node $x \in V$ with an integer (its *time stamp*). Time stamps are unique (no two nodes have the same time stamp) and the time stamp of a child node is strictly greater than that of its parent:

$$\forall x, y \in V : \quad timeStamp(x) = timeStamp(y) \Rightarrow x = y \quad (1)$$

$$\forall x, y \in V : \quad (x, y) \in E \Rightarrow timeStamp(x) < timeStamp(y) \quad (2)$$

Trace compatibility We say that a call-tree (V, E) is *trace-compatible* (compatible for short) with a sequence of observed stack traces $traceSequence = (trace_1, trace_2, \dots, trace_N)$ if and only if there exists a surjective

mapping Φ from the frames $\langle \text{trace}_i, j \rangle$ of the observed traces onto the nodes of \mathbb{V} , such that Φ fulfils the following properties P0, P1, P2 and P3.

PROPERTY P0: SYMBOL CONSISTENCY Stack frames are mapped onto nodes with the same symbol:

$$\forall \text{trace}_i \in \text{traceSequence}, 1 \leq k \leq |\text{trace}_i| : \text{symbol}(\Phi(\langle \text{trace}_i, k \rangle)) = \text{trace}_i[k] \quad (3)$$

PROPERTY P1: TRACE INCLUSION Each observed trace is mapped by Φ onto a path of the tree that starts from the tree's root:

$$\begin{aligned} \forall \mathbf{t}_i \in \text{traceSequence} : \\ \Phi(\langle \mathbf{t}_i, 1 \rangle) \text{ is the root of } (\mathbb{V}, \mathbb{E}) \wedge \\ \forall 1 \leq j < |\mathbf{t}_i| : \left(\Phi(\langle \mathbf{t}_i, j \rangle), \Phi(\langle \mathbf{t}_i, j+1 \rangle) \right) \in \mathbb{E} \end{aligned} \quad (4)$$

PROPERTY P2: BREAKPOINT DISCRIMINATION The last frame of every stack trace (e.g. the frame that corresponds to a breakpoint activation) is mapped to a node that is distinct from those of the stack traces that precede it:

$$\forall \mathbf{t}_i, \mathbf{t}_j \in \text{traceSequence} : i < j \Rightarrow \forall 1 \leq k \leq |\mathbf{t}_i| : \Phi(\langle \mathbf{t}_j, |\mathbf{t}_j| \rangle) \neq \Phi(\langle \mathbf{t}_i, k \rangle) \quad (5)$$

PROPERTY P3: ORDER CONSERVATION The sequence numbers of the nodes of the tree reflect the order in which traces are observed. More precisely: if a trace trace_j is observed after another trace trace_i , then the tail of trace_j that does not overlap with trace_i is mapped onto nodes with higher sequence numbers than the nodes of trace_i .

$$\begin{aligned} \forall \text{trace}_i, \text{trace}_j \in \text{traceSequence}, \forall k \in [1, \min(|\text{trace}_i|, |\text{trace}_j|)] : \\ (j > i) \wedge \Phi(\langle \text{trace}_i, k \rangle) \neq \Phi(\langle \text{trace}_j, k \rangle) \Rightarrow \\ \forall l \in [k, |\text{trace}_j|], \forall m \in [1, |\text{trace}_i|] : \\ \text{timeStamp}(\Phi(\langle \text{trace}_j, l \rangle)) > \text{timeStamp}(\Phi(\langle \text{trace}_i, m \rangle)) \end{aligned} \quad (6)$$

Homomorphism between call-trees A *call-tree homomorphism* h from a call-tree (\mathbb{V}, \mathbb{E}) onto another call-tree (\mathbb{W}, \mathbb{F}) is a graph homomorphism from (\mathbb{V}, \mathbb{E}) to (\mathbb{W}, \mathbb{F}) that fulfils the following additional properties:

$$\forall x \in \mathbb{V} : \text{symbol}(h(x)) = \text{symbol}(x) \quad (7)$$

$$\begin{aligned} \forall r, s \in h[\mathbb{V}] : \left(\forall x \in h^{-1}[\{r\}], \forall y \in h^{-1}[\{s\}] : \text{timeStamp}(x) \leq \text{timeStamp}(y) \right) \\ \Rightarrow \text{timeStamp}(r) \leq \text{timeStamp}(s) \end{aligned} \quad (8)$$

where $f[A]$ (resp. $f^{-1}[A]$) denotes the image (resp. the pre-image) of the set A by the function f .

A *call-tree isomorphism* is a bijective call-tree homomorphism. Note that for a call-tree isomorphism (8) is equivalent to

$$\forall x, y \in \mathbb{V} : \text{timeStamp}(x) \leq \text{timeStamp}(y) \Rightarrow \text{timeStamp}(h(x)) \leq \text{timeStamp}(h(y)) \quad (9)$$

Theorem 3.1. *The set of call-trees that are trace-compatible with a sequence of observed stack traces is non-empty. The call-tree of this set with the smallest number of nodes is unique (modulo a call-tree isomorphism) and is the call-tree constructed by algorithm 1.*

4 Proof

Let's consider a trace sequence $\text{traceSequence} = (\text{trace}_1, \text{trace}_2, \dots, \text{trace}_N)$, and let's call $\mathbf{G}_0 = (\mathbb{V}_0, \mathbb{E}_0)$ the call-tree constructed by algorithm 1 from this trace sequence.

Lemma 4.1. \mathbf{G}_0 is trace compatible with traceSequence .

Proof. Let's note ActiveTreePath_i the value of ActiveTreePath just after trace_i has been processed by $\text{AddTraceToGraph}()$ at line 5. A mapping Φ_0 from the frames of traceSequence onto \mathbf{G}_0 can be defined as:

$$\forall \text{trace}_i \in \text{traceSequence}, \forall j \in [1, |\text{trace}_i|] : \Phi_0(\langle \text{trace}_i, j \rangle) \equiv \text{ActiveTreePath}_i[j] \quad (10)$$

By construction Φ_0 is a function (i.e. it associates each frame $\langle \text{trace}_i, j \rangle$ with exactly one node of \mathbf{G}_0), and by construction of algorithm 1, Φ_0 is surjective.

Also because of lines 8, 13, and 14 and the constructor property of the **new** $\text{Node}(\text{trace}[k])$, we have

$$\text{symbol}(\Phi_0(\langle \text{trace}_i, j \rangle)) = \text{symbol}(\text{ActiveTreePath}_i[j]) = \text{trace}_i[j]$$

and consequently Φ_0 fulfils property P0.

Property P1 follows from the creation of new edges at line 16 and the fact that all traces start with the same symbol (the program's entry point): $\forall t_i \in \text{traceSequence} : t_i[1] = \text{trace}_1[1]$.

Property P2 follows from the bound $|\text{trace}| - 1$ put on the loop that compares ActiveTreePath with the current trace at line 7. '-1' forces the last frame of the trace (and hence the symbol corresponding to a breakpoint activation) to be included as a new node in the call-tree.

The proof of property P3 is slightly more technical. In this proof we will use the following sublemma:

Sublemma 4.1.1.

$$\begin{aligned} & \forall \text{trace}_i \in \text{traceSequence}, \forall k, l \in [1, |\text{trace}_i|] : \\ & k < l \Rightarrow \text{timeStamp}(\Phi_0(\langle \text{trace}_i, k \rangle)) < \text{timeStamp}(\Phi_0(\langle \text{trace}_i, l \rangle)) \end{aligned} \quad (11)$$

Proof. This sublemma directly follows from

$$\begin{aligned} & \forall i \in [1, |\text{traceSequence}|], \forall k, l \in [1, |\text{traceSequence}[i]|] : \\ & k < l \Rightarrow \text{timeStamp}(\text{ActiveTreePath}_i[k]) < \text{timeStamp}(\text{ActiveTreePath}_i[l]) \end{aligned} \quad (12)$$

which is itself an invariant of the loop at line 3. This invariant can be proved by recursion using the properties of time stamps and node creation. (The details of this short proof are left out for place reason.) \square

To prove P3, let's now assume $\Phi_0(\langle \text{trace}_i, k \rangle) \neq \Phi_0(\langle \text{trace}_j, k \rangle)$ for some $i > j$ and $k \in [1, \min(|\text{trace}_i|, |\text{trace}_j|)]$. The proof needs to consider two cases depending whether all traces between i and j are at least of length k or not.

Case 1 $\forall r \in [i, j] : |\text{trace}_r| \geq k$

In this case we can consider the largest $l \in [i, j - 1]$ such that $\Phi_0(\langle \text{trace}_l, k \rangle) \neq \Phi_0(\langle \text{trace}_{l+1}, k \rangle)$. Let's note this value l_0 :

$$l_0 \equiv \max\left(\left\{l \in [i, j - 1] \mid \Phi_0(\langle \text{trace}_l, k \rangle) \neq \Phi_0(\langle \text{trace}_{l+1}, k \rangle)\right\}\right) \quad (13)$$

By construction of l_0 we have

$$\Phi_0(\langle \text{trace}_{l_0}, k \rangle) \neq \Phi_0(\langle \text{trace}_{l_0+1}, k \rangle) = \Phi_0(\langle \text{trace}_{l_0+2}, k \rangle) = \dots = \Phi_0(\langle \text{trace}_j, k \rangle) \quad (14)$$

This implies (by definition of Φ_0 in equation (10))

$$\text{ActiveTreePath}_{l_0}[k] \neq \text{ActiveTreePath}_{l_0+1}[k] \quad (15)$$

This again implies that a new node was constructed for $\text{trace}_{l_0+1}[k]$ at line 13. Because timestamps are allocated in the order in which nodes are created, this node's timestamp is higher than those of any previously created nodes, and in particular than the timestamp of $\text{ActiveTreePath}_i[m] = \Phi_0(\langle \text{trace}_i, m \rangle)$ for all $m \in [1, |\text{trace}_i|]$.

$$\forall m \in [1, |\text{trace}_i|] : \text{timeStamp}(\Phi_0(\langle \text{trace}_i, m \rangle)) < \text{timeStamp}(\Phi_0(\langle \text{trace}_{l_0+1}, k \rangle)) \quad (16)$$

From (14) and (16) we thus get (since $\Phi_0(\langle \text{trace}_{l_0+1}, k \rangle) = \Phi_0(\langle \text{trace}_j, k \rangle)$):

$$\forall m \in [1, |\text{trace}_i|] : \text{timeStamp}(\Phi_0(\langle \text{trace}_i, m \rangle)) < \text{timeStamp}(\Phi_0(\langle \text{trace}_j, k \rangle)) \quad (17)$$

Combined with sublemma 4.1.1, equation (17) gives

$$\begin{aligned} \forall m \in [1, |\text{trace}_i|], l \in [k, |\text{trace}_j|] : \\ \text{timeStamp}(\Phi_0(\langle \text{trace}_i, m \rangle)) < \text{timeStamp}(\Phi_0(\langle \text{trace}_j, l \rangle)) \end{aligned} \quad (18)$$

which concludes the proof of P3 in this case.

Case 2 $\exists r \in [i, j] : |\text{trace}_r| < k$

In that case we can consider

$$r_0 \equiv \max\left(\left\{r \in [i+1, j-1] \mid |\text{trace}_r| < k\right\}\right) \quad (19)$$

By construction, note that we have:

$$\forall r \in [r_0+1, j] : |\text{trace}_r| \geq k \quad (20)$$

Two subcases follow.

Case 2.1

If $\exists l \in [r_0+1, j-1] : \Phi_0(\langle \text{trace}_l, k \rangle) \neq \Phi_0(\langle \text{trace}_{l+1}, k \rangle)$ then we can define

$$l_0 \equiv \max\left(\left\{l \in [r_0+1, j-1] \mid \Phi_0(\langle \text{trace}_l, k \rangle) \neq \Phi_0(\langle \text{trace}_{l+1}, k \rangle)\right\}\right) \quad (21)$$

and we are back to case 1.

Case 2.2

If $\forall l \in [r_0+1, j-1] : \Phi_0(\langle \text{trace}_l, k \rangle) = \Phi_0(\langle \text{trace}_{l+1}, k \rangle)$ then we note that because $|\text{trace}_{r_0}| < k$ and $|\text{trace}_{r_0+1}| \geq k$ then a new node was constructed for $\text{trace}_{r_0+1}[k]$ at line 13. Because $\Phi_0(\langle \text{trace}_{r_0+1}, k \rangle) = \Phi_0(\langle \text{trace}_j, k \rangle)$, we are again back to case 1.

This concludes the proof of P3, and more generally the proof of lemma 4.1. □

Lemma 4.2. *If a call-tree $\mathbf{H} = (W, F)$ is trace compatible with traceSequence through a mapping $\Phi_{\mathbf{H}}$, then*

$$\forall w \in W : \Phi_0\left[\Phi_{\mathbf{H}}^{-1}\left[\{w\}\right]\right] \text{ is a singleton}$$

where $f[A]$ (resp. $f^{-1}[A]$) denotes the image (resp. the pre-image) of the set A by the function f .

Proof. Let's first note that because $\Phi_{\mathbf{H}}$ is by definition surjective, $\Phi_{\mathbf{H}}^{-1}[\{x\}]$ is non-empty, and $\Phi_0[\Phi_{\mathbf{H}}^{-1}[\{x\}]]$ therefore contains at least one element.

Let's now consider

$$x, y \in \Phi_0 \left[\Phi_{\mathbf{H}}^{-1}[\{w\}] \right] \quad (22)$$

Our goal is to prove that $x = y$

(22) means that there exists $\text{trace}_i, \text{trace}_j \in \text{traceSequence}$, $k \in [1, |\text{trace}_i|]$ and $l \in [1, |\text{trace}_j|]$ such that

$$\langle \text{trace}_i, k \rangle, \langle \text{trace}_j, l \rangle \in \Phi_{\mathbf{H}}^{-1}[\{w\}] \quad (23)$$

$$\Phi_0(\langle \text{trace}_i, k \rangle) = x \quad (24)$$

$$\Phi_0(\langle \text{trace}_j, l \rangle) = y \quad (25)$$

(23) means in turn that:

$$\Phi_{\mathbf{H}}(\langle \text{trace}_i, k \rangle) = \Phi_{\mathbf{H}}(\langle \text{trace}_j, l \rangle) = w \quad (26)$$

Because \mathbf{H} is a tree, Property P1 of \mathbf{H} and (26) imply that $l = k$ and that

$$\forall m \in [1, k] : \Phi_{\mathbf{H}}(\langle \text{trace}_i, m \rangle) = \Phi_{\mathbf{H}}(\langle \text{trace}_j, m \rangle) \quad (27)$$

If $i = j$, (24) and (25), and $l = k$ imply trivially that $x = y$. Without loss of generality, we will therefore assume in the remaining that $i < j$.

We will now prove by contradiction that

$$\forall r \in [i + 1, j] : |\text{trace}_r| > k \quad (28)$$

$$\forall r \in [i, j], \forall m \in [1, k] : \Phi_{\mathbf{H}}(\langle \text{trace}_r, m \rangle) = \Phi_{\mathbf{H}}(\langle \text{trace}_i, m \rangle) \quad (29)$$

Let's first consider (28) and assume there exists $r \in [i + 1, j]$ such that $|\text{trace}_r| \leq k$. Let's note $m_r \equiv |\text{trace}_r|$. By applying property P2 of \mathbf{H} on trace_i and trace_r we get

$$\Phi_{\mathbf{H}}(\langle \text{trace}_r, m_r \rangle) \neq \Phi_{\mathbf{H}}(\langle \text{trace}_i, m_r \rangle) \quad (30)$$

Because of (27) and $m_r \equiv |\text{trace}_r| \leq k$ we also have

$$\Phi_{\mathbf{H}}(\langle \text{trace}_r, m_r \rangle) \neq \Phi_{\mathbf{H}}(\langle \text{trace}_j, m_r \rangle) \quad (31)$$

From this we also conclude that $r < j$. Since $i < r < j$, we can apply the property P3 of \mathbf{H} twice, first on trace_i and trace_r and then on trace_r and trace_j , and we eventually get:

$$\text{timeStamp} \left(\Phi_{\mathbf{H}}(\langle \text{trace}_i, m_r \rangle) \right) < \text{timeStamp} \left(\Phi_{\mathbf{H}}(\langle \text{trace}_r, m_r \rangle) \right) < \text{timeStamp} \left(\Phi_{\mathbf{H}}(\langle \text{trace}_j, m_r \rangle) \right) \quad (32)$$

(32) implies that $\Phi_{\mathbf{H}}(\langle \text{trace}_i, m_r \rangle) \neq \Phi_{\mathbf{H}}(\langle \text{trace}_j, m_r \rangle)$, which contradicts (27), thus proving (28).

Turning to (29), let's assume there exists $r \in [i, j]$ and $n \in [1, k]$ such that

$$\Phi_{\mathbf{H}}(\langle \text{trace}_r, n \rangle) \neq \Phi_{\mathbf{H}}(\langle \text{trace}_i, n \rangle) \quad (33)$$

We are back to equation (30) with m_r replaced by n , and we derive the same contradiction, thus proving (29).

Using the property P0 of \mathbf{H} , we derive from (29)

$$\forall r \in [i, j], \forall m \in [1, k] : \text{trace}_r[m] = \text{trace}_i[m] \quad (34)$$

By construction of the loop at lines 3- 5 of algorithm 1, we derive from (34) and (28) that

$$\forall r \in [i, j], \forall m \in [1, k] : \text{ActiveTreePath}_r[m] = \text{ActiveTreePath}_i[m] \quad (35)$$

and by definition of Φ_0 (equation (10)) that

$$\Phi_0(\langle \text{trace}_i, k \rangle) = \Phi_0(\langle \text{trace}_j, k \rangle) \quad (36)$$

which with (24) and (25) and $l = k$ leads us to $x = y$ and concludes the proof of lemma 4.2. \square

Lemma 4.3. *If a call-tree $\mathbf{H} = (W, F)$ is trace compatible with traceSequence, then there exists a surjective call-tree homomorphism from \mathbf{H} onto $\mathbf{G}_0 = (V_0, E_0)$.*

Proof. Consider the following mapping h from \mathbf{H} onto \mathbf{G}_0 :

$$\forall w \in W : h(w) \equiv \text{the unique element of } \Phi_0 \left[\Phi_{\mathbf{H}}^{-1} \left[\{w\} \right] \right] \quad (37)$$

Because of lemma 4.2, h is well-defined. Furthermore, because of the properties of set images and pre-images by a function, and because Φ_0 is surjective, we have

$$\forall x \in V_0 : \exists w \in \Phi_{\mathbf{H}} \left[\Phi_0^{-1}[x] \right] : h(w) = x \quad (38)$$

h is therefore surjective.

In the following we first prove that h is a graph homomorphism, and then that h fulfils (7) and (8), and is thus a call tree homomorphism.

Let's consider $v, w \in W$ such that $(v, w) \in F$. We need to prove that $(h(v), h(w)) \in E_0$.

Because $\Phi_{\mathbf{H}}$ is surjective, there exists a stack frame $\langle t, k \rangle$ such that

$$\Phi_{\mathbf{H}}(\langle t, k \rangle) = w \quad (39)$$

Because $(v, w) \in F$, w is not the root of \mathbf{H} , and by applying property P1 if \mathbf{H} on trace t , we conclude that $k > 1$. Again from property P1 we derive

$$\left(\Phi_{\mathbf{H}}(\langle t, k-1 \rangle), \Phi_{\mathbf{H}}(\langle t, k \rangle) \right) \in F \quad (40)$$

Because \mathbf{H} is a tree, w has at most one parent. $(v, w) \in F$ and (40) thus imply:

$$\Phi_{\mathbf{H}}(\langle t, k-1 \rangle) = v \quad (41)$$

From (39) and (41) we derive

$$\langle t, k \rangle \in \Phi_{\mathbf{H}}^{-1}[w] \quad (42)$$

$$\langle t, k-1 \rangle \in \Phi_{\mathbf{H}}^{-1}[v] \quad (43)$$

And thus after applying the mapping Φ_0 and using the definition of h (37):

$$\Phi_0(\langle t, k \rangle) \in \Phi_0 \left[\Phi_{\mathbf{H}}^{-1}[w] \right] = \{h(w)\} \quad (44)$$

$$\Phi_0(\langle t, k-1 \rangle) \in \Phi_0 \left[\Phi_{\mathbf{H}}^{-1}[v] \right] = \{h(v)\} \quad (45)$$

We finally get:

$$\Phi_0(\langle t, k \rangle) = h(w) \quad (46)$$

$$\Phi_0(\langle t, k-1 \rangle) = h(v) \quad (47)$$

By applying property P1 of \mathbf{G}_0 on trace \mathfrak{t} we conclude:

$$(h(v), h(w)) = \left(\Phi_0(\langle \mathfrak{t}, k-1 \rangle), \Phi_0(\langle \mathfrak{t}, k \rangle) \right) \in E_0 \quad (48)$$

which shows that h is a graph homomorphism from \mathbf{H} onto \mathbf{G}_0 .

We will now turn to equations (7) and (8) to prove that h is a *call-tree homomorphism* and complete the lemma.

Let's consider $w \in W$. Using property P0 of \mathbf{H} we have

$$\forall \langle \mathfrak{t}, k \rangle \in \Phi_{\mathbf{H}}^{-1}[w] : \text{symbol}(w) = \text{symbol}\left(\Phi_{\mathbf{H}}(\langle \mathfrak{t}, k \rangle)\right) = \mathfrak{t}[k] \quad (49)$$

Using now the definition of h (37) and the property P0 of \mathbf{G}_0 we have

$$\forall \langle \mathfrak{t}, k \rangle \in \Phi_{\mathbf{H}}^{-1}[w] : \text{symbol}(h(w)) = \text{symbol}\left(\Phi_0(\langle \mathfrak{t}, k \rangle)\right) = \mathfrak{t}[k] \quad (50)$$

Because $\Phi_{\mathbf{H}}$ is surjective, and therefore $\Phi_{\mathbf{H}}^{-1}[w]$ is non-empty, we derive from (49) and (50) that:

$$\text{symbol}(w) = \text{symbol}(h(w)) \quad (51)$$

which proves (7) for h .

Turning now to the proof of property (8), let's consider $r, s \in h(W) = V_0$, and let's assume

$$\forall u \in h^{-1}[r], \forall v \in h^{-1}[s] : \text{timeStamp}(u) \leq \text{timeStamp}(v) \quad (52)$$

Let's note R and S the sets of stack frames such that

$$R \equiv \Phi_0^{-1}[r] \wedge S \equiv \Phi_0^{-1}[s] \quad (53)$$

By definition of h (37) and lemma 4.2 we have

$$h^{-1}[r] = \Phi_{\mathbf{H}}[R] \quad (54)$$

$$h^{-1}[s] = \Phi_{\mathbf{H}}[S] \quad (55)$$

Because by definition of R we have $\forall \langle \mathfrak{t}_i, k \rangle \in R : \Phi_0(\langle \mathfrak{t}_i, k \rangle) = r$, property P1 of \mathbf{G}_0 and the fact that \mathbf{G}_0 is a tree means all frames of R are at the same stack depth k_r :

$$\forall \langle \mathfrak{t}_i, k \rangle \in R : k = k_r \quad (56)$$

The same is true of S . There exists a stack depth m_s such that:

$$\forall \langle \mathfrak{t}_j, m \rangle \in S : m = m_s \quad (57)$$

Let's consider two stack frames from R and S :

$$\langle \mathfrak{t}_i, k_r \rangle \in R \wedge \langle \mathfrak{t}_j, m_s \rangle \in S \quad (58)$$

Note that by definition of R and S (53), we have

$$\begin{aligned} \Phi_0(\langle \mathfrak{t}_i, k_r \rangle) &= r \\ \Phi_0(\langle \mathfrak{t}_j, m_s \rangle) &= s \end{aligned} \quad (59)$$

We first prove by contradiction that the following holds:

$$(j \geq i) \vee \left(\forall k \leq k_r : \Phi_{\mathbf{H}}(\langle \mathfrak{t}_j, k \rangle) = \Phi_{\mathbf{H}}(\langle \mathfrak{t}_i, k \rangle) \right) \quad (60)$$

Let's assume the negation of (60):

$$(j < i) \wedge \left(\exists k \leq k_r : \Phi_{\mathbf{H}}(\langle t_j, k \rangle) \neq \Phi_{\mathbf{H}}(\langle t_i, k \rangle) \right) \quad (61)$$

We can apply property P3 of \mathbf{H} on t_j and t_i , and we get

$$\text{timeStamp}(\Phi_{\mathbf{H}}(\langle t_j, m_s \rangle)) < \text{timeStamp}(\Phi_{\mathbf{H}}(\langle t_i, k_r \rangle)) \quad (62)$$

Because $\Phi_{\mathbf{H}}(\langle t_j, m_s \rangle) \in \Phi_{\mathbf{H}}[S] = h^{-1}[s]$ and $\Phi_{\mathbf{H}}(\langle t_i, k_r \rangle) \in \Phi_{\mathbf{H}}[R] = h^{-1}[r]$, (62) contradicts (52). We therefore conclude that (60) must be true.

Because

$$(i = j) \Rightarrow \left(\Phi_{\mathbf{H}}(\langle t_j, k_r \rangle) = \Phi_{\mathbf{H}}(\langle t_i, k_r \rangle) \right)$$

and

$$\left(\Phi_{\mathbf{H}}(\langle t_j, k_r \rangle) = \Phi_{\mathbf{H}}(\langle t_i, k_r \rangle) \right) \Rightarrow \left(\forall k \leq k_r : \Phi_{\mathbf{H}}(\langle t_j, k \rangle) = \Phi_{\mathbf{H}}(\langle t_i, k \rangle) \right)$$

(using the fact that \mathbf{H} is a tree and property P1), (60) is furthermore equivalent to

$$(j > i) \vee \left(\forall k \leq k_r : \Phi_{\mathbf{H}}(\langle t_j, k \rangle) = \Phi_{\mathbf{H}}(\langle t_i, k \rangle) \right) \quad (63)$$

The remaining of the proof considers two cases, depending whether $m_s < k_r$.

Case 1 $m_s < k_r$

If we assume $\Phi_{\mathbf{H}}(\langle t_j, m_s \rangle) = \Phi_{\mathbf{H}}(\langle t_i, m_s \rangle)$, then $m_s < k_r$, property P1 of \mathbf{H} , and (2) give us

$$\text{timeStamp}(\Phi_{\mathbf{H}}(\langle t_j, m_s \rangle)) = \text{timeStamp}(\Phi_{\mathbf{H}}(\langle t_i, m_s \rangle)) \quad (64)$$

$$< \text{timeStamp}(\Phi_{\mathbf{H}}(\langle t_i, k_r \rangle)) \quad (65)$$

The same equation as (62), which yields the same contradiction with (52). From this we conclude that:

$$\Phi_{\mathbf{H}}(\langle t_j, m_s \rangle) \neq \Phi_{\mathbf{H}}(\langle t_i, m_s \rangle) \quad (66)$$

With (63) and $m_s < k_r$ (case assumption), this implies that

$$j > i \quad (67)$$

(66) and (67) mean we can apply property P3 of Φ_0 to t_i and t_j . We get

$$\text{timeStamp}(\Phi_0(\langle t_j, m_s \rangle)) > \text{timeStamp}(\Phi_0(\langle t_i, k_r \rangle)) \quad (68)$$

Using (59), this yields

$$\text{timeStamp}(s) > \text{timeStamp}(r) \quad (69)$$

and proves (8) in this case.

Case 2 $m_s \geq k_r$

We consider in turn the two alternatives of (63).

If $j > i$, then by construction of algorithm 1 and the property of the node constructor `new Node(trace[k])` at line 13, we have

$$\text{timeStamp}(\Phi_0(\langle t_i, k_r \rangle)) \leq \text{timeStamp}(\Phi_0(\langle t_j, k_r \rangle)) \quad (70)$$

Because of sublemma 4.1.1 on Φ_0 and $m_s \geq k_r$ we also have

$$timeStamp(\Phi_0(\langle t_j, k_r \rangle)) \leq timeStamp(\Phi_0(\langle t_j, m_s \rangle)) \quad (71)$$

From (70) and (71) we derive

$$timeStamp(\Phi_0(\langle t_i, k_r \rangle)) \leq timeStamp(\Phi_0(\langle t_j, m_s \rangle)) \quad (72)$$

which using (59), gives us

$$timeStamp(r) \leq timeStamp(s) \quad (73)$$

and proves (8) when $j > i$.

If we turn to the second term of (63) and assume

$$\Phi_{\mathbf{H}}(\langle t_j, k_r \rangle) = \Phi_{\mathbf{H}}(\langle t_i, k_r \rangle) \quad (74)$$

Because of (58) and (54) we have

$$\Phi_{\mathbf{H}}(\langle t_i, k_r \rangle) \in \Phi_{\mathbf{H}}[R] = h^{-1}[r] \quad (75)$$

With (74), we thus have

$$\Phi_{\mathbf{H}}(\langle t_j, k_r \rangle) \in h^{-1}[r] \quad (76)$$

From which we derive

$$\langle t_j, k_r \rangle \in \Phi_{\mathbf{H}}^{-1}[h^{-1}[r]] \quad (77)$$

$$\Phi_0(\langle t_j, k_r \rangle) \in \Phi_0[\Phi_{\mathbf{H}}^{-1}[h^{-1}[r]]] \quad (78)$$

By definition of h (37) and lemma 4.2 this yields

$$\Phi_0(\langle t_j, k_r \rangle) = r \quad (79)$$

Because $m_s \geq k_r$, sublemma 4.1.1 on t_j gives us

$$timeStamp(\Phi_0(\langle t_j, m_s \rangle)) \geq timeStamp(\Phi_0(\langle t_j, k_r \rangle)) \quad (80)$$

This, with (59) and (79) leads to

$$timeStamp(s) \geq timeStamp(r) \quad (81)$$

and concludes the proof (8) for h .

h is thus a surjective graph homomorphism that fulfils (7) and (8), which proves lemma 4.3. \square

Theorem 3.1. *The set of call-trees that are trace-compatible with a sequence of observed stack traces is non-empty. The call-tree of this set with the smallest number of nodes is unique (modulo a call-tree isomorphism) and is the call-tree constructed by algorithm 1.*

Proof. As before, let traceSequence be a sequence of observed stack traces, and $\mathbf{G}_0 = (E_0, V_0)$ the call-tree constructed by algorithm 1 for traceSequence.

Because of lemma 4.1, the set of call-trees that are trace-compatible with traceSequence contains \mathbf{G}_0 and is therefore non-empty.

Let's consider $\mathbf{H} = (W, F)$ a call tree that is trace-compatible with traceSequence. Because of lemma 4.3, there exists a surjective call-tree homomorphism h from \mathbf{H} onto \mathbf{G}_0 . Because h is surjective we have

$$|E_0| \leq |W| \quad (82)$$

i.e., \mathbf{G}_0 has the smallest number of nodes of all call-trees compatible with traceSequence.

Let's now consider $\mathbf{H}_0 = (W_0, F_0)$, a call tree that is trace-compatible with traceSequence and has the same number of nodes as \mathbf{G}_0 .

$$|E_0| = |W_0| \tag{83}$$

Because of lemma 4.3, there exists a surjective call-tree homomorphism h_0 from \mathbf{H}_0 onto \mathbf{G}_0 .

Because h_0 is surjective, and its domain and co-domain are finite and contain the same number of elements, then h_0 is a bijection. \mathbf{H}_0 onto \mathbf{G}_0 are therefore identical modulo a call-tree isomorphism.

\mathbf{G}_0 is therefore unique (modulo a call-tree isomorphism), which concludes the proof of theorem 3.1.

□

References

- [1] François Taïani, Marc-Olivier Killijian, and Jean-Charles Fabre. Cosmopen: Dynamic reverse engineering on a budget, how cheap observation techniques can be used to reconstruct complex multi-level behaviour. Technical report, Lancaster University, Computing Departement, 2008. submitted for publication.