

CF-GKAT

Efficient Validation of Control-Flow Transformations

Tobias Kappé

NetTCS @ LIACS, October 16, 2024



**Universiteit
Leiden**
The Netherlands



Funded by
the European Union



Joint work with



Cheng Zhang
University College London



David E. Narváez
Virginia Tech



Nico Naus
OU & Virginia Tech

Program equivalence

```
if  $x = 1$  then  
  |  $y++$ ;  
else  
  |  $z := z/2$ ;  
end
```

≡

```
if  $x \neq 1$  then  
  |  $z := z/2$ ;  
else  
  |  $y++$ ;  
end
```

Propositional program equivalence

```
if b then
```

```
| p;
```

```
else
```

```
| q;
```

```
end
```

≡

```
if not b then
```

```
| q;
```

```
else
```

```
| p;
```

```
end
```

Propositional program equivalence

```
s;  
while b do  
  | p;  
end  
while c do  
  | q;  
  | while b do  
  |   | p;  
  |   end  
  end  
end
```

≡

```
s;  
while b or c do  
  | if b then  
  |   | p;  
  |   else  
  |     | q;  
  |   end  
end
```

Kleene Algebra with Tests (KAT)

Fix a set $\{p, q, \dots\}$ of **actions** and a Boolean algebra $\{b, c, \dots\}$ of **tests**

$$b \mid p \mid e + f \mid e; f \mid e^*$$

- ▶ We can embed simple propositional programs:

$$\boxed{\text{if } b \text{ then } e \text{ else } f} = b; e + (\text{not } b); f$$

$$\boxed{\text{while } b \text{ do } e} = (b; e)^* ; (\text{not } b)$$

- ▶ Language semantics in terms of *guarded strings*.
- ▶ Complete and finitary axiomatization
- ▶ **Non-determinism makes equivalence PSPACE-complete**
 - ▶ Basically: translate to automaton and check bisimilarity

(Kozen 1996), (Kozen & Smith 1996)

Guarded Kleene Algebra with Tests (GKAT)

Fix a set $\{p, q, \dots\}$ of **actions** and a Boolean algebra $\{b, c, \dots\}$ of **tests**

$$b \mid p \mid e +_b f \mid e; f \mid e^{(b)}$$

- ▶ Language semantics in terms of *guarded strings*.
- ▶ Complete but infinitary axiomatization (open problem)
- ▶ **Language equivalence is efficiently decidable! (nearly linear)**
 - ▶ Same idea as in KAT, but different composition operators.

(Kozen & Tseng 2008), (Smolka et al. 2020)

GKAT semantics

Guarded strings

Inherited from KAT: sets of *guarded strings*, i.e., words of the form

$$\alpha_1 p_1 \alpha_2 p_2 \alpha_3 p_3 \cdots \alpha_n p_n \alpha_{n+1}$$

where the *atom* $\alpha_i \in At$ tells us the truth value of each test at point i .

GKAT semantics

By example

$$\left[\begin{array}{l} \text{if } b \text{ then} \\ | \quad p; \\ \text{else} \\ | \quad q; \\ \text{end} \end{array} \right] = \left\{ \begin{array}{l} bpb, \\ b\bar{p}\bar{b}, \\ \bar{b}q\bar{b}, \\ \bar{b}qb \end{array} \right\}$$

GKAT semantics

By example

$$\left[\begin{array}{l} \text{while } b \text{ or } c \text{ do} \\ | \text{if } b \text{ then} \\ | | p; \\ | \text{else} \\ | | q; \\ | \text{end} \\ \text{end} \end{array} \right] = \left\{ \begin{array}{l} \bar{b}\bar{c}, \\ \bar{b}c\bar{q}\bar{b}\bar{c}, \\ \bar{b}\bar{c}p\bar{b}c\bar{q}\bar{b}\bar{c}, \\ \dots \end{array} \right\}$$

GKAT decision procedure

GKAT automata

This comes from (Smolka et al. 2020).

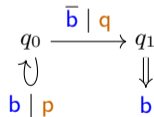
Let Q be a set. The set of *GKAT dynamics* on Q , denoted GQ , consists of functions

$$At \rightarrow \perp + \checkmark + \Sigma \times Q$$

A *GKAT automaton* is a pointed G -coalgebra.

That is, it is a tuple (Q, δ, q_0) where:

- ▶ Q is a (finite) set of states;
- ▶ $q_0 \in Q$ is the *initial state*; and
- ▶ $\delta : Q \rightarrow GQ$ is the *transition function*.



Straightforward semantics in terms of guarded languages.

GKAT decision procedure

Translating to GKAT automata

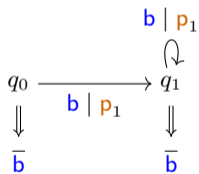
```
while b or c do
  | if b then
  |   | p;
  | else
  |   | q;
  | end
end
```

\mapsto

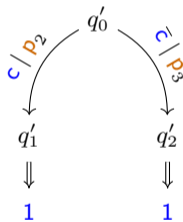
$q_0 \implies \bar{b}\bar{c}$
↻
bc | p
b \bar{c} | p
 $\bar{b}c$ | q

GKAT decision procedure

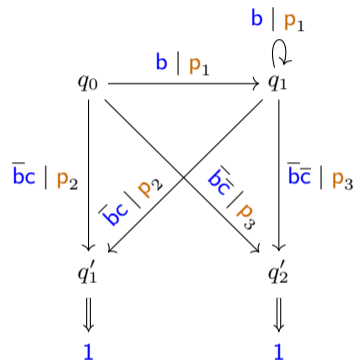
Translating to GKAT automata



$e = \text{while } b \text{ do } p_1$



$f = \text{if } c \text{ then } p_2 \text{ else } p_3$



$g = e \cdot f$

GKAT decision procedure

GKAT automata equivalence

Theorem

For every GKAT expression e , we can construct a GKAT automaton A_e such that $\llbracket e \rrbracket = \llbracket A_e \rrbracket$.

This is a basic syntax-directed translation, à la (Thompson 1968).

Lemma


A_e is at most linear in the size of e .

Equivalence (bisimilarity) of GKAT automata can be easily decided.

Goto removal

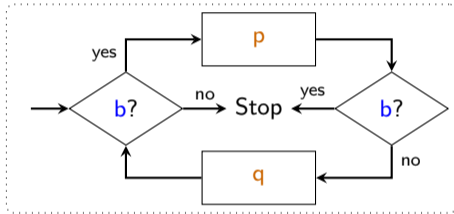
```
L:  
if b then  
  | p;  
  | goto R;  
return;  
R:  
if not b then  
  | q;  
  | goto L;  
return;
```

goto removal
(e.g., calipso)



```
x := 1;  
while x ≠ 0 do  
  | if x = 1 and b then  
  |   | p;  
  |   | x := 2;  
  | else if x = 2 and  
  |   | not b then  
  |   |   | q;  
  |   |   | x := 1;  
  | else  
  |   | x := 0;  
end
```

Decompilation



```
while b do
  | p;
  | if b then
  |   | break;
  | q;
end
```


Map of the Problematique

- ▶ GKAT does not support labels, **goto**, **return**, **break**, variables...
- ▶ It cannot even express the example programs! (Schmid, K. & Silva 2021)
- ▶ They *can* be translated to KAT expressions (Kozen 2008), e.g.:

$$(bpbq)^*(bpb + \bar{b})$$

- ▶ This is a rather non-trivial process...
- ▶ ...and KAT equivalence is PSPACE-hard...
- ▶ ...and really, we just need an automaton.

Enter CF-GKAT

Fix sets of:

- ▶ *actions* $p \in \Sigma$;
- ▶ *tests* $b \in T$;
- ▶ *labels* $\ell \in L$; and
- ▶ *indicator values* $i \in I$.

CF-GKAT is GKAT augmented with non-local flow control:

$$\text{Exp} \ni e, f ::= \mathbf{assert} \ b \mid p \in \Sigma \mid x := i \ (i \in I) \mid e; f \mid \mathbf{if} \ b \ \mathbf{then} \ e \ \mathbf{else} \ f \mid \\ \mathbf{while} \ b \ \mathbf{do} \ e \mid \mathbf{break} \mid \mathbf{return} \mid \mathbf{goto} \ \ell \ (\ell \in L) \mid \mathbf{label} \ \ell \ (\ell \in L)$$

CF-GKAT semantics

Guarded strings with continuations

First step: *guarded strings with continuations*, i.e.,

$$\alpha_1 p_1 \alpha_2 p_2 \alpha_3 p_3 \cdots \alpha_n p_n \alpha_{n+1} \cdot c$$

where c is a *continuation* of the form:

acc i

brk i

ret

jmp (ℓ, i)

A continuation tells us how a partial computation can proceed.

CF-GKAT semantics

By example

$$\left[\begin{array}{l} p; \\ | \\ \text{if } b \text{ then} \\ | \\ \text{break;} \\ | \\ q; \end{array} \right]_i^{\#} = \left\{ \begin{array}{l} bp\bar{b}q\bar{b} \cdot \text{acc } i, \\ bp\bar{b}qb \cdot \text{acc } i, \\ bpb \cdot \text{brk } i \end{array} \right\}$$

$$\left[\begin{array}{l} \text{while } b \text{ do} \\ | \\ p; \\ | \\ \text{if } b \text{ then} \\ | \\ \text{break;} \\ | \\ q; \\ \text{end} \end{array} \right]_i^{\#} = \left\{ \begin{array}{l} bp\bar{b}q\bar{b} \cdot \text{acc } i, \\ bp\bar{b}qb\bar{b}q\bar{b} \cdot \text{acc } i, \\ bpb \cdot \text{acc } i \\ bp\bar{b}qbpb \cdot \text{acc } i \\ \dots \end{array} \right\}$$

CF-GKAT semantics

By example

$$\left[\begin{array}{l} \text{L:} \\ \text{if } b \text{ then} \\ \quad | \quad p; \\ \quad | \quad \text{goto R;} \\ \text{return;} \\ \text{R:} \\ \text{if not } b \text{ then} \\ \quad | \quad q; \\ \quad | \quad \text{goto L;} \\ \text{return;} \end{array} \right]_i^{\sharp} = \left\{ \begin{array}{l} bpb \cdot \text{jmp } (R, i), \\ b\bar{p}\bar{b} \cdot \text{jmp } (R, i), \\ \bar{b} \cdot \text{ret} \end{array} \right\}$$

$$\left[\begin{array}{l} \text{L:} \\ \text{if } b \text{ then} \\ \quad | \quad p; \\ \quad | \quad \text{goto R;} \\ \text{return;} \\ \text{R:} \\ \text{if not } b \text{ then} \\ \quad | \quad q; \\ \quad | \quad \text{goto L;} \\ \text{return;} \end{array} \right]_i^{\text{R}} = \left\{ \begin{array}{l} \bar{b}qb \cdot \text{jmp } (L, i), \\ \bar{b}q\bar{b} \cdot \text{jmp } (L, i), \\ b \cdot \text{ret} \end{array} \right\}$$

Now we have a semantics of the form

$$\llbracket - \rrbracket : L \rightarrow I \rightarrow \text{guarded languages with continuations}$$

where guarded words end at the first **goto** statement.

What we need is a semantics of the form

$$\llbracket - \rrbracket \downarrow : L \rightarrow I \rightarrow \text{guarded languages}$$

where the jumps are “resolved” by stringing them together.

CF-GKAT semantics

Resolution

Resolution, in symbols:

$$\frac{w \cdot c \in \llbracket P \rrbracket_i^\ell \quad c \in \{\mathbf{acc} \ j, \mathbf{ret}\}}{w \in \llbracket P \rrbracket_{\downarrow i}^\ell}$$

$$\frac{w\alpha \cdot \mathbf{jmp}(\ell', j) \in \llbracket P \rrbracket_i^\ell \quad \alpha x \in \llbracket P \rrbracket_{\downarrow j}^{\ell'}}{w\alpha x \in \llbracket P \rrbracket_{\downarrow i}^\ell}$$

$P =$

```
L:
  if b then
    | p;
    | goto R;
  return;
R:
  if not b then
    | q;
    | goto L;
  return;
```

Then we can infer:

$$\frac{b \cdot \mathbf{ret} \in \llbracket P \rrbracket_i^R}{b \in \llbracket P \rrbracket_{\downarrow i}^R}$$
$$\frac{bpb \cdot \mathbf{jmp}(R, i) \in \llbracket P \rrbracket_i^\# \quad b \in \llbracket P \rrbracket_{\downarrow i}^R}{bpb \in \llbracket P \rrbracket_{\downarrow i}^\#}$$

Decision procedure

CF-GKAT automata

Let C be the set of continuations (e.g., `acc i`, `brk i`, etc).

A *CF-GKAT dynamics* on a set X is a function of the form

$$I \rightarrow At \rightarrow \perp + C + \Sigma \times X \times I$$

We write GX for the set of CF-GKAT dynamics on X .

A CF-GKAT automaton is a tuple $A = (Q, \delta, q, \lambda)$ where

- ▶ Q is a set of *states* with $q \in Q$ the *initial state*
- ▶ $\delta : Q \rightarrow GQ$ is the *transition function*
- ▶ $\lambda : L \rightarrow GQ$ is the *jump function*

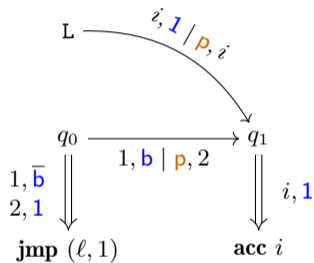
Has a straightforward semantics $\llbracket A \rrbracket : L \rightarrow I \rightarrow$ guarded languages with continuations.

Decision procedure

CF-GKAT automata

```
if b and  $x = 1$  then
|    $x := 2$ ;
|   L:
|   p;
else
|    $x := 1$ ;
|   goto L;
```

\mapsto



Decision procedure

Translating to CF-GKAT automata

Theorem

For every CF-GKAT expression e , we can construct a CF-GKAT automaton A_e such that $\llbracket e \rrbracket = \llbracket A_e \rrbracket$.

Another syntax-directed translation, à la (Thompson 1968).

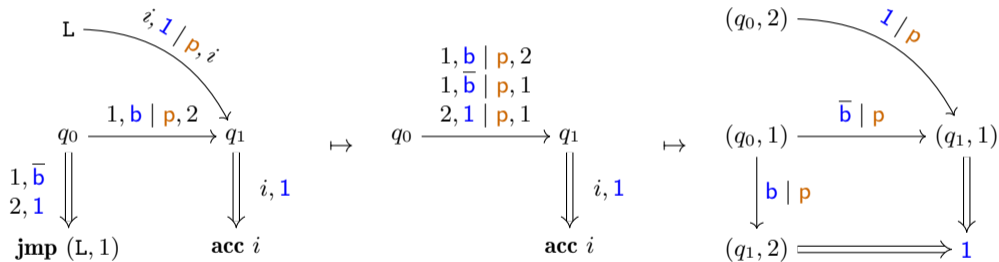
Lemma

A_e is at most linear in the size of e .



Decision procedure

Translating to GKAT automata



Decision procedure

Translating to GKAT automata

Theorem

For every CF-GKAT automaton A , we can construct a GKAT automaton $A\downarrow$ such that $\llbracket A \rrbracket\downarrow = \llbracket A\downarrow \rrbracket$.

Lemma

$A\downarrow$ is at most $|I|$ times as large as A



Validation

Blinding

C programs are still not CF-GKAT programs:

- ▶ Parsing C can be really hard!
- ▶ Same statements mapped to same variable.
- ▶ Which variables are indicators?
- ▶ Macros may hide relevant information.



LLVM/clang to the rescue

Validation

Blinding

```
#define hidden(x) \  
    x=f(x); goto R;
```

```
int foo(int x, int y) {  
    L:  
    if (x == y) {  
        x = f(x);  
        hidden(x);  
    }  
    return;  
    R:  
    if (x != y) {  
        y = g(y);  
        goto L;  
    }  
}
```

↪

```
L:  
if b then  
    | p;  
    | goto R;  
return;  
R:  
if not b then  
    | q;  
    | goto L;
```

Validation

Working on `coreutils`

We chose to work on `mp_factor_using_pollard_rho`:

- ▶ 91 lines of code (before macro expansion)
- ▶ loops nested three levels deep
- ▶ has a **break** statement inside
- ▶ also contains a **goto** for error handling
- ▶ no indicator variables (yet)

Validation

Working on coreutils

First experiment:

- ▶ Blind, compile (clang), and then decompile (Ghidra).
- ▶ Decompiled code has 3 more goto's and labels.
- ▶ Some manual effort to remove decompilation artifacts.
- ▶ Compare original code to decompiled code: OK.

Second experiment:

- ▶ Blind, then eliminate **goto** using indicators (Erosa & Hendren 1994, Cassé et al. 2002)
- ▶ Some manual effort to make indicator detection work.
- ▶ Compare original code to refactored code: OK.

Conclusion & Further Work

- ▶ More experiments are necessary!
 - ▶ Automated testing pipeline
 - ▶ Scale to all of `coreutils`, Linux kernel, ...
 - ▶ Extract control flow from binary (no compilation).
- ▶ Expand support for top-level syntax:
 - ▶ the **continue** statement
 - ▶ multi-way branching using **switch**
 - ▶ **do-while** loops with **break**
- ▶ Integration with existing decompiler
 - ▶ Baked-in translation validation
 - ▶ Continuous integration