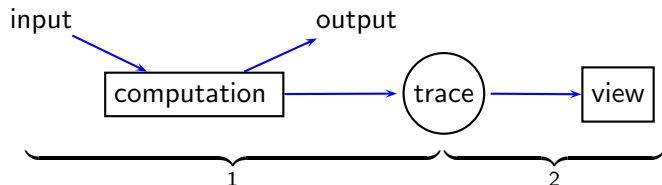# Foundations for Tracing Functional Programs and the Correctness of Algorithmic Debugging

Olaf Chitil and Yong Luo

University of Kent, UK
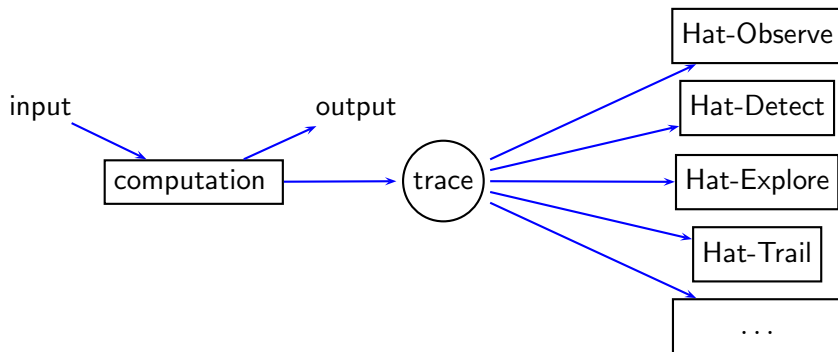Supported by EPSRC grant EP/C516605/1

26th April 2006

# Why Tracing?



```
insert :: Ord a => a -> [a] -> [a]

insert x [] = [x]
insert x (y:ys) =
  if x > y then y : insert x ys
           else x : ys

sort :: Ord a => [a] -> [a]
sort [] = []
sort (x:xs) = insert x (sort xs)

main = getLine >>= putStrLn . sort
```

program

sample text → ▮▮▮▮▮▮▮ → alms

input                computation                output

- Locate a fault (wrong output, run-time error, non-termination).
- Comprehend a program.

# Two-Phase Tracing: A Trace as Data Structure



- Liberates from time arrow of computation.
- Enables views based on different execution models.
  (small-step, big-step, interpreter with environment, denotational)
- Enables compositional views.

- Multi-View Tracer



- Trace = Augmented Redex Trail (ART); distilled as unified trace.

Aim: A theoretical model of this trace and its views.

# Overview

# The Programming Language

## Launchbury's and related semantics

- Subset of $\lambda$-calculus plus `case` for matching.
- Any program can be translated into this core calculus.

## For tracing

- Close relationship between trace and original program essential.
- Language must have most frequently used features:
  - named functions
  - pattern matching

# The Programming Language

## Launchbury's and related semantics

- Subset of $\lambda$-calculus plus `case` for matching.
- Any program can be translated into this core calculus.

## For tracing

- Close relationship between trace and original program essential.
- Language must have most frequently used features:
  - named functions
  - pattern matching

$\Rightarrow$ Higher-order term rewriting system

```
sort [] = []                          or  sort = foldr insert []
sort (x:xs) = insert x (sort xs)

insert x [] = [x]
insert x (y:ys) = if x > y then y : (insert x ys) else x : ys
```

# What is a Good Trace?

Program + input determine every detail of computation.

# What is a Good Trace?

Program + input determine every detail of computation.

⇒ Trace gives efficient access to certain details of computation.

# What is a Good Trace?

Program $+$ input determine every detail of computation.
$\Rightarrow$ Trace gives <span style="color:red">efficient</span> access to certain details of computation.

What is a computation? Semantics answers:

- Term rewriting: A sequence of expressions.
  $$t_1 \rightarrow t_2 \rightarrow t_3 \rightarrow t_4 \rightarrow t_5 \rightarrow \ldots \rightarrow t_n$$
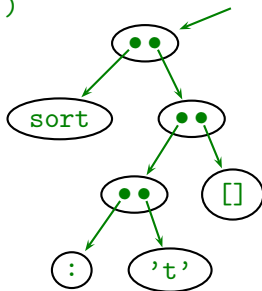- Natural semantics: A proof tree.

# What is a Good Trace?

Program $+$ input determine every detail of computation.

$\Rightarrow$ Trace gives efficient access to certain details of computation.

What is a computation? Semantics answers:

- Term rewriting: A sequence of expressions.

  $t_1 \rightarrow t_2 \rightarrow t_3 \rightarrow t_4 \rightarrow t_5 \rightarrow \ldots \rightarrow t_n$

- Natural semantics: A proof tree.

But

- Lots of redundancy.
- Much structure already lost.
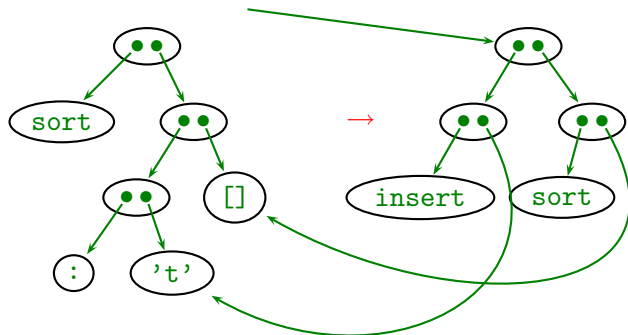
# Graph Rewriting I

`sort ('t':[])`



```
sort [] = []
sort (x:xs) = insert x (sort xs)
```
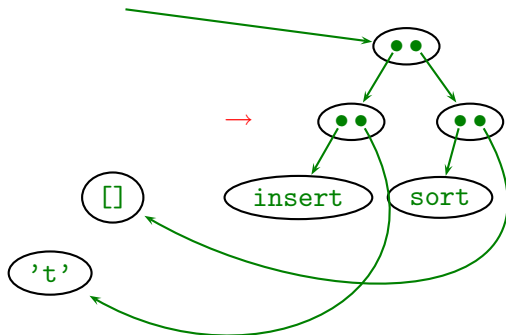
# Graph Rewriting I



```
sort [] = []
sort (x:xs) = insert x (sort xs)
```

- Create new nodes for right-hand-side.
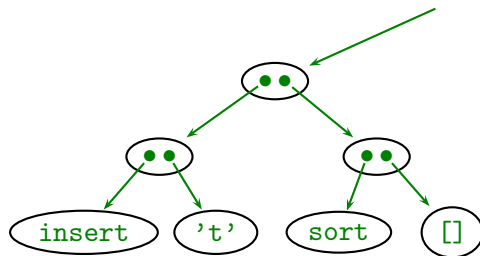- Nodes of subexpressions are shared.

# Graph Rewriting I



```
sort [] = []
sort (x:xs) = insert x (sort xs)
```

- Create new nodes for right-hand-side.
- Nodes of subexpressions are shared.
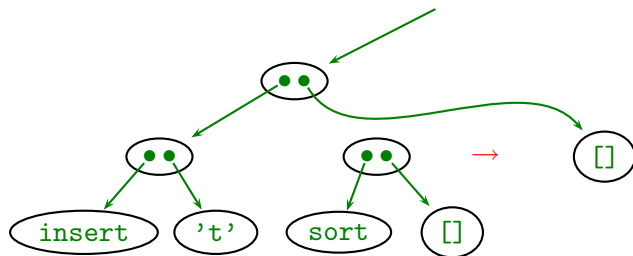- Some old nodes become garbage.

# Graph Rewriting II



```
sort [] = []
sort (x:xs) = insert x (sort xs)

insert x [] = [x]
insert x (y:ys) = if x > y then y : (insert x ys) else x : ys
```

# Graph Rewriting II



```
sort [] = []
sort (x:xs) = insert x (sort xs)

insert x [] = [x]
insert x (y:ys) = if x > y then y : (insert x ys) else x : ys
```

- Application node of redex replaced by new node.

```
sort [] = []
sort (x:xs) = insert x (sort xs)

insert x [] = [x]
insert x (y:ys) = if x > y then y : (insert x ys) else x : ys
```
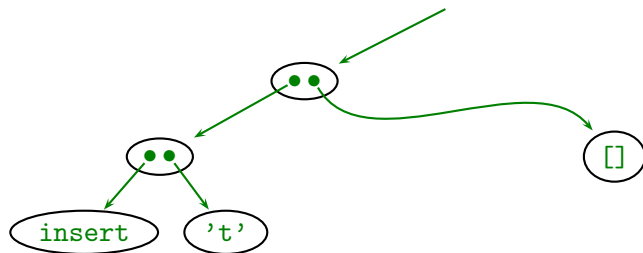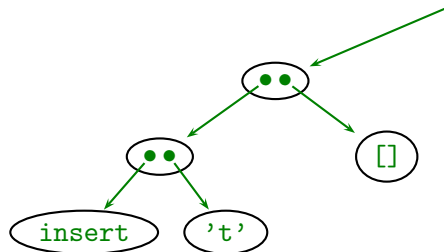
- Application node of redex replaced by new node.

```
sort [] = []
sort (x:xs) = insert x (sort xs)

insert x [] = [x]
insert x (y:ys) = if x > y then y : (insert x ys) else x : ys
```

```
sort [] = []
sort (x:xs) = insert x (sort xs)

insert x [] = [x]
insert x (y:ys) = if x > y then y : (insert x ys) else x : ys
```
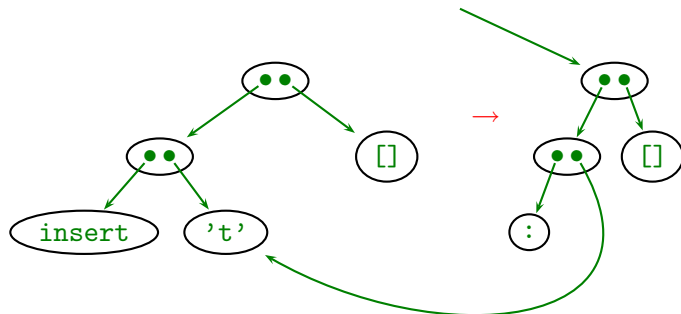
```
sort [] = []
sort (x:xs) = insert x (sort xs)

insert x [] = [x]
insert x (y:ys) = if x > y then y : (insert x ys) else x : ys
```

- New nodes for right-hand-side, connected via result pointer.
- Only add to graph, never remove.
- Sharing ensures compact representation.

- New nodes for right-hand-side, connected via result pointer.
- Only add to graph, never remove.
- Sharing ensures compact representation.

- New nodes for right-hand-side, connected via result pointer.
- Only add to graph, never remove.
- Sharing ensures compact representation.

# The Node Labels



$$
\begin{array}{llll}
\text{term constructor} & T & := & a & \text{atom} \\
& & | & n\,m & \text{application of nodes} \\[1em]
\text{atom} & a & := & f \mid C \mid 42 \mid \dots & \text{defined variable, data constructor} \\
& & & & \text{atomic literal, } \dots
\end{array}
$$

- pointers instead of edges

# The Node Naming Scheme



## Aim

- not distinguish isomorphic graphs
- avoid inconvenience of isomorphism classes

# The Node Naming Scheme



### Aim

- not distinguish isomorphic graphs
- avoid inconvenience of isomorphism classes

### Solution

- standard representation with node describing path from root
- path at creation time (sharing later)
- path independent of evaluation order

# The Node Naming Scheme II



- Reduction edge implicitly given through existence of node.
- Node encodes parent = top node of redex causing its creation:

$$\text{parent}(n\text{t}) = n$$
$$\text{parent}(n\text{l}) = \text{parent}(n)$$
$$\text{parent}(n\text{r}) = \text{parent}(n)$$
$$\text{parent}(\varepsilon) = \text{undefined}$$

- Easy to identify right-hand-side of rule: same parent.

# Projections

- Reduction edge implicitly given through existence of node.
- Every redex should be parent of at least one node.
  (otherwise reduction unreachable from computation result)

```
True && x = x
not True = False
```

# Projections

- Reduction edge implicitly given through existence of node.
- Every redex should be parent of at least one node.
  (otherwise reduction unreachable from computation result)

$\Rightarrow$ A projection requires an indirection as result.



```
True && x = x
not True = False
```

| term constructor | $T$ | := | $a$ | atom |
| | | $\mid$ | $n\,m$ | application of nodes |
| | | $\mid$ | $n$ | indirection |

| atom | $a$ | := | $x \mid C \mid 42 \mid \ldots$ | variable, data constructor, ... |

# Projections

- Reduction edge implicitly given through existence of node.
- Every redex should be parent of at least one node.
  (otherwise reduction unreachable from computation result)

⇒ A projection requires an indirection as result.

```
True && x = x
not True = False
```



| term constructor | $T$ | := | $a$ | atom |
| | | | $n\,m$ | application of nodes |
| | | | $n$ | indirection |
| atom | $a$ | := | $x \mid C \mid 42 \mid \ldots$ | variable, data constructor, . . . |

A trace $G$ for initial term $M$ and program $P$ is a partial function from nodes to term constructors, $G : n \mapsto T$, defined by

- The unshared graph representation of $M$, $\text{graph}_G(\varepsilon, M)$, is a trace.
- If $G$ is a trace and
    - $L = R$ an equation of the program $P$,
    - $\sigma$ a substitution replacing argument variables by nodes,
    - $\text{match}_G(n, L\sigma)$,
    - $n\text{t} \notin \text{dom}(G)$,

  then $G \cup \text{graph}_G(n\text{t}, R\sigma)$ is a trace.

No evaluation order is fixed.

# Unshared Graph Representation

For the initial term and right-hand-sides of equation.

$\text{graph}(t, \text{insert rlr}(\text{sort rr})) =$



## Definition

$$\text{graph}(n, a) = \{(n, a)\}$$

$$\text{graph}(n, m) = \{(n, m)\}$$

$$\text{graph}(n, M\,N) = \begin{cases} \{(n, M\,N)\} & \text{, if } M, N \text{ are nodes} \\ \{(n, M\,n\text{r})\} \cup \text{graph}(n\text{r}, N) & \text{, if only } M \text{ is a node} \\ \{(n, n\text{l}\,N)\} \cup \text{graph}(n\text{l}, M) & \text{, if only } N \text{ is a node} \\ \{(n, n\text{l}\,n\text{r})\} \cup \text{graph}(n\text{l}, M) \cup \text{graph}(n\text{r}, N), & \text{otherwise} \end{cases}$$

# Matching

Matching a node with an instance of the left-hand-side of an equation.



$\mathrm{match}_G(\varepsilon, \mathtt{sort}\ (\mathrm{rlr:rr}))$

## Definition

$\mathrm{match}_G(n, M) = \text{if } M \text{ is a node then } n = M \text{ else } \mathrm{matchT}_G(\mathrm{last}_G(n), M)$

$\mathrm{matchT}_G(a, M) = (a = M)$

$\mathrm{matchT}_G(n, M) = \mathrm{matchT}_G(\mathrm{last}_G(n), M)$

$\mathrm{matchT}_G(n\,o, M) = \exists N, O.\,(N\,O = M) \wedge \mathrm{match}_G(n, N) \wedge \mathrm{match}_G(o, O)$

$\mathrm{last}_G(n) = \text{if } n\mathrm{t} \in \mathrm{dom}(G) \text{ then } \mathrm{last}_G(n\mathrm{t}) \text{ else } G(n)$

# The Most Evaluated Form of a Node

A node represents many terms, in particular a most evaluated one.



$\mathsf{mef}_G(\mathsf{tr}) = \texttt{[]}$

$\mathsf{mef}_G(\varepsilon) = \texttt{(:)} \quad \texttt{'t'} \quad \texttt{[]}$

## Definition

$$\mathsf{mef}_G(n) = \mathsf{mefT}_G(\mathsf{last}_G(n))$$

$$\mathsf{mefT}_G(a) = a$$

$$\mathsf{mefT}_G(n) = \mathsf{mef}_G(n)$$

$$\mathsf{mefT}_G(n\,m) = \mathsf{mef}_G(n)\,\mathsf{mef}_G(m)$$

# Redexes and Big-Step Reductions



$$\mathrm{redex}_G(\mathrm{t}) = \texttt{insert 't' []}$$
$$\mathrm{bigstep}_G(\mathrm{t}) = \texttt{insert 't' [] = (:) 't' []}$$

## Definition

For any redex node $n$, i.e., $n\mathrm{t} \in \mathrm{dom}(G)$

$$\mathrm{redex}_G(n) = \begin{cases} \mathrm{mef}_G(m)\,\mathrm{mef}_G(o) & \text{, if } G(n) = m\,o \\ a & \text{, if } G(n) = a \end{cases}$$

$$\mathrm{bigstep}_G(n) = \mathrm{redex}_G(n) = \mathrm{mef}_G(n)$$

## Properties of the ART

- closed (no dangling nodes)
- domain prefix-closed
- acyclic
- strongly confluent
- no application contains a node ending in t
- only a node ending in t can be an indirection
- if $nl \in \text{dom}(G)$, then $G(n) = nl\, m$
- if $nr \in \text{dom}(G)$, then $G(n) = m\, nr$
- if $nt \in \text{dom}(G)$, then $\text{redex}_G(n) = L\sigma$ and $\text{reduct}_G(n) = R\sigma$
  for some program equation $L = R$ and substitution $\sigma$

Give non-inductive definition of ART based on properties?

# Reduct of a Small Step Reduction



$reduct_G(\varepsilon) = $ insert 't' (sort [])

## Definition

$reduct_G(n) = reductB_G(nt)$

$$reductB_G(n) = \begin{cases} a & \text{, if } G(n) = a \\ mef_G(m) & \text{, if } G(n) = m \\ reductB_G(nl)\, reductB_G(nr) & \text{, if } G(n) = nl\, nr \\ reductB_G(nl)\, mef_G(o) & \text{, if } G(n) = nl\, o \text{ and } o \neq nr \\ mef_G(m)\, reductB_G(nr) & \text{, if } G(n) = m\, nr \text{ and } m \neq nl \\ mef_G(m)\, mef_G(o) & \text{, if } G(n) = m\, o, \, m \neq nl \text{ and } o \neq nr \end{cases}$$

# Views of the Trace

- Observation of Expressions and Functions
- Following Redex Trails
- Algorithmic Debugging

# Observation of Expressions and Functions

Observation of function sort:

```
sort "sort" = "os"
sort "ort" = "o"
sort "rt" = "r"
sort "t" = "t"
sort "" = ""
```

Observation of function insert:

```
insert 's' "o" = "os"
insert 's' "" = "s"
insert 'o' "r" = "o"
insert 'r' "t" = "r"
insert 't' "" = "t"
```

Big step reductions of redex nodes.

# Following Redex Trails

```
Output: ------------------------------------------------------
os\n

Trail: ------- Insert.hs line: 10 col: 25 ------------------
<- putStrLn "os"
<- insert 's' "o" | if True
<- insert 'o' "r" | if False
<- insert 'r' "t" | if False
<- insert 't' []
<- sort []
```

- Go backwards from observed failure to fault.
- Which redex created this expression?
- To prove: every reduction step reachable from final result.

# Algorithmic Debugging

```
sort "sort" = "os"?   n

insert 's' "o" = "os"?   y

sort "ort" = "o"?   n

insert 'o' "r" = "o"?   n

Bug identified:
  "Insert.hs":8-9:
  insert x [] = [x]
  insert x (y:ys) = if x > y then y:(insert x ys) else x:ys
```

# The Evaluation Dependency Tree



```
                        main = {IO}

    sort "sort" = "os"              putStrLn "os" = {IO}

  sort "ort" = "o"              insert 's' "o" = "os"

                        's' > 'o' = True    insert 's' "" = "s"

sort "rt" = "r"        insert 'o' "r" = "o"

sort "t" = "t"    insert 'r' "t" = "r"    'o' > 'r' = False

sort "" = ""    insert 't' "" = "t"    'r' > 't' = False
```

# The Evaluation Dependency Tree



```
                        main = {IO}
                 /                        \
    sort "sort" = "os"              putStrLn "os" = {IO}
       /          \
sort "ort" = "o"     insert 's' "o" = "os"
                        /              \
           's' > 'o' = True      insert 's' "" = "s"
   /                \
sort "rt" = "r"    insert 'o' "r" = "o"
   |                 /            \
sort "t" = "t"   insert 'r' "t" = "r"    'o' > 'r' = False
   /        \
sort "" = ""  insert 't' "" = "t"   'r' > 't' = False
```

# The Evaluation Dependency Tree

# The Evaluation Dependency Tree

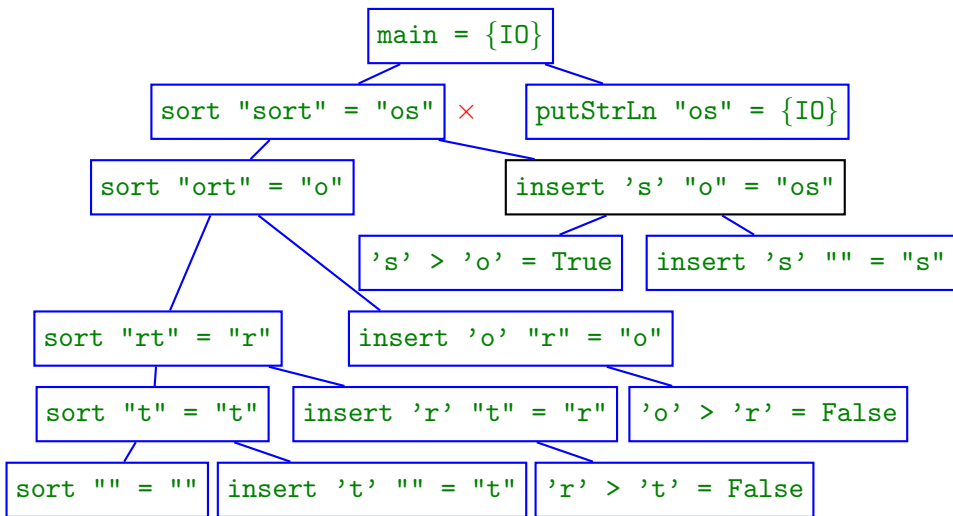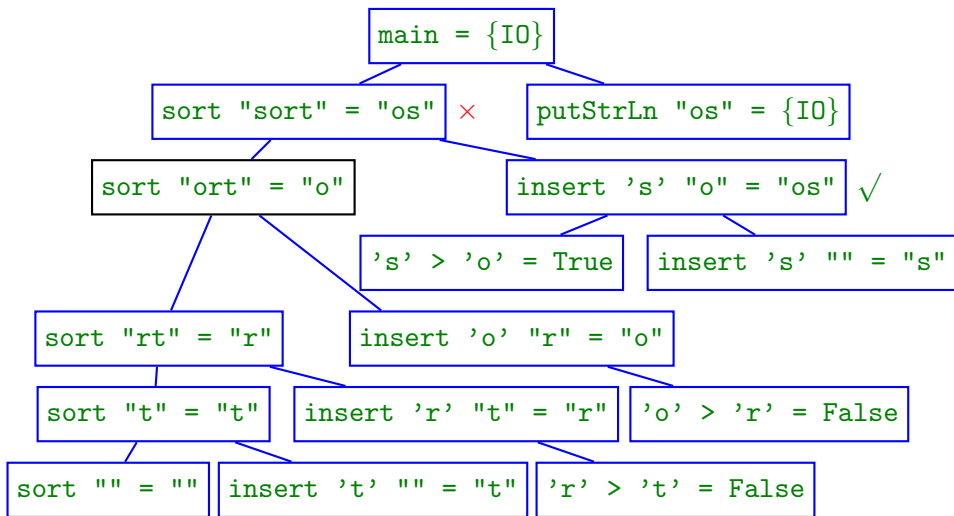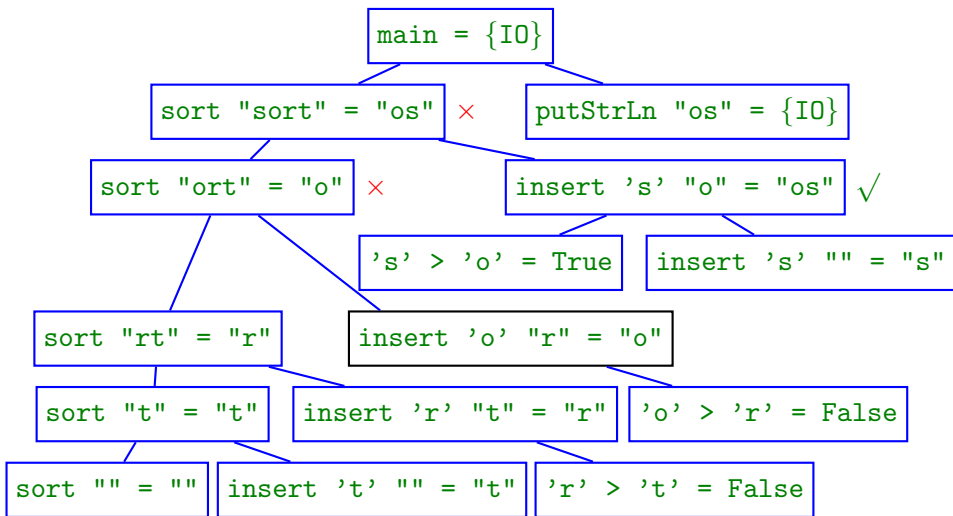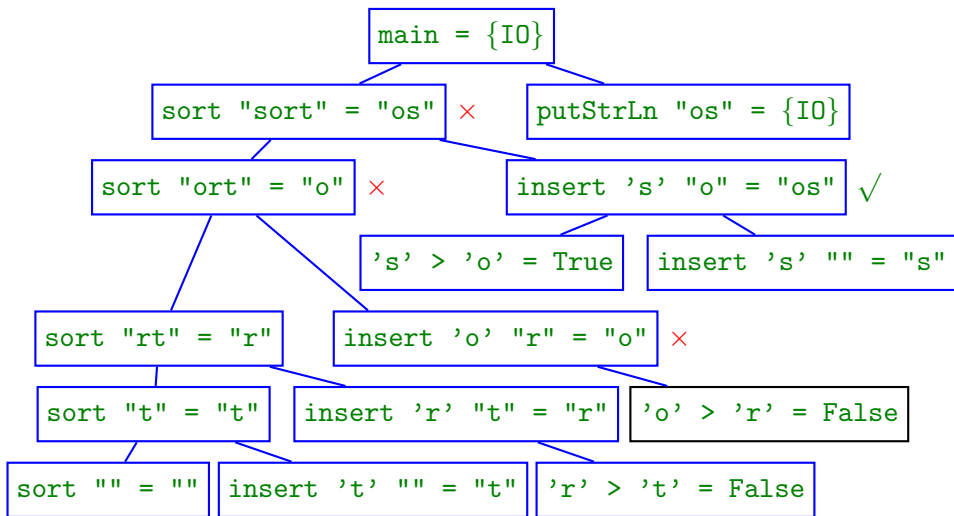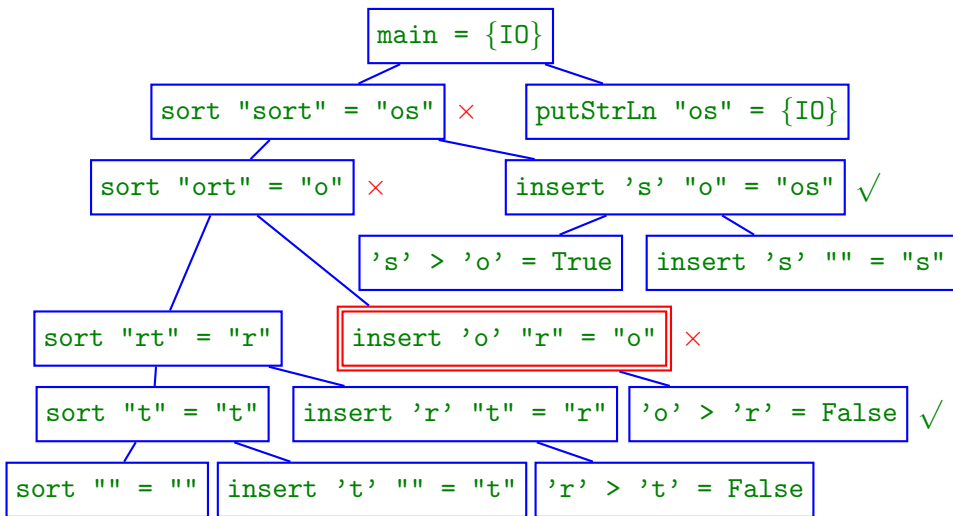# The Evaluation Dependency Tree

# The Evaluation Dependency Tree



```
                    ┌─────────────────┐
                    │ main = {IO}     │
                    └─────────────────┘
          ┌──────────────────────┐    ┌──────────────────────────┐
          │ sort "sort" = "os"   │ ×  │ putStrLn "os" = {IO}     │
          └──────────────────────┘    └──────────────────────────┘
      ┌──────────────────────┐        ┌──────────────────────────┐
      │ sort "ort" = "o"     │ ×      │ insert 's' "o" = "os"    │ √
      └──────────────────────┘        └──────────────────────────┘
                              ┌──────────────────┐  ┌──────────────────────────┐
                              │ 's' > 'o' = True │  │ insert 's' "" = "s"      │
                              └──────────────────┘  └──────────────────────────┘
  ┌──────────────────────┐    ┌──────────────────────────┐
  │ sort "rt" = "r"      │    │ insert 'o' "r" = "o"     │ ×
  └──────────────────────┘    └──────────────────────────┘
  ┌──────────────────┐  ┌──────────────────────────┐  ┌──────────────────────┐
  │ sort "t" = "t"   │  │ insert 'r' "t" = "r"     │  │ 'o' > 'r' = False    │
  └──────────────────┘  └──────────────────────────┘  └──────────────────────┘
┌──────────────┐ ┌──────────────────────────┐ ┌──────────────────────┐
│ sort "" = "" │ │ insert 't' "" = "t"      │ │ 'r' > 't' = False    │
└──────────────┘ └──────────────────────────┘ └──────────────────────┘
```
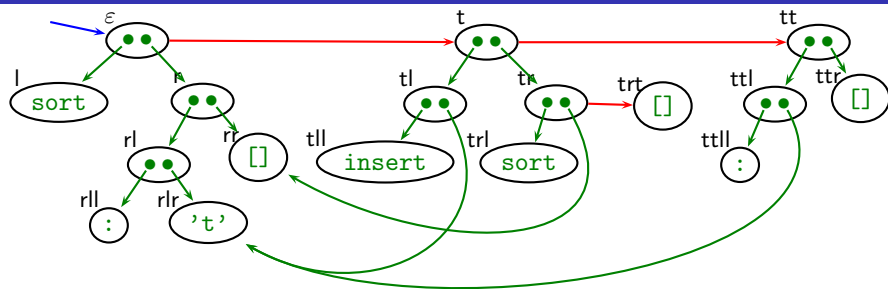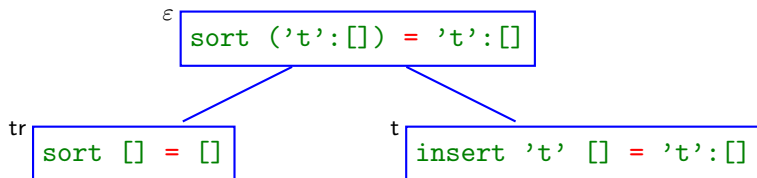
# The Evaluation Dependency Tree

- Every redex node $n$ yields a tree node $n$ labelled $\text{bigstep}_G(n)$.
- Tree node $n$ is child of tree node $\text{parent}(n)$.
- Usually root label $\text{bigstep}_G(\varepsilon) = \text{main} = \ldots$

# Correctness of Algorithmic Debugging: The Property

If node $n$ incorrect and all its children correct, then node $n$ faulty, i.e., its equation is faulty.

$\varepsilon$
```
sort ('t':[]) = 't':[]
```

tr
```
sort [] = []
```

t
```
insert 't' [] = 't':[]
```

### Definition

Tree node $n$ incorrect $\quad\Leftrightarrow\quad \text{redex}_G(n) \ncong_I \text{mef}_G(n)$.

Tree node $n$ faulty $\quad\Leftrightarrow\quad \text{redex}_G(n) \ncong_I \text{reduct}_G(n)$.
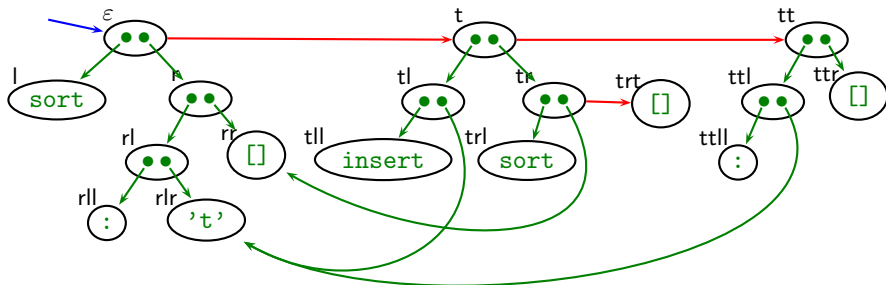
If tree node $n$ faulty, then for its program equation $L = R$ exists substitution $\sigma$ such that $L\sigma \ncong_I R\sigma$.

# Correctness of Algorithmic Debugging: Main Theorem

## Theorem

Let $n$ be a redex node. If for all redex nodes $m$ with $parent(m) = n$ we have $redex_G(m) \cong_I mef_G(m)$, then $reduct_G(n) \cong_I mef_G(n)$.

With $redex_G(n) \ncong_I mef_G(n)$ follows $redex_G(n) \ncong_I reduct_G(n)$.

**Proof.**

Generalise property: Let $n \in \mathrm{dom}(G)$. If for all redex nodes $m$ with $\mathrm{parent}(m) = \mathrm{parent}(n)$ we have $\mathrm{redex}_G(m) \cong_{\mathsf{I}} \mathrm{mef}_G(m)$, then $\mathrm{reductB}_G(n) \cong_{\mathsf{I}} \mathrm{mef}_G(n)$.

Induction over $\mathrm{hight}_G(n) = \max\{|o| \mid o \in \{\mathsf{l}, \mathsf{r}\}^* \wedge no \in \mathrm{dom}(G)\}$. $\qquad\square$

# Future Work

- Still play with definitions.
- Extend model further:
  - Drop non-needed nodes from ART (unevaluated expressions).
  - Model run-time error with error value.
  - Allow local function definitions ($\Rightarrow$ free variables).
  - Share reductions of constants ($\Rightarrow$ cycles in graph).
  - Describe strict and mixed semantics.
- Prove further properties.

# Conclusions

- Simple model amenable to proof.
- Contains a wealth of information about computation.
- Models real-world trace of Haskell tracer Hat.
- Proved correctness of algorithmic debugging.