

# Refined typechecking with Stardust

Joshua Dunfield

Carnegie Mellon University  
(current affiliation: McGill University)

PLPV 2007  
Freiburg, Germany

5 October 2007

# Conventional Static Typing

- ML, Haskell, ...
- Types express **properties**:
  - *append* : list \* list → list
- Types can be inferred;  
type annotations are documentation  
(except module interfaces)
- “Well typed programs don’t go wrong”

# Conventional Static Typing

- ML, Haskell, ...
- Types express **properties**:
  - $append : list * list \rightarrow list$
- Types can be inferred;  
type annotations are documentation  
(except module interfaces)
- “Well typed programs don’t go wrong”
- We extend conventional static typing  
to verify more precise invariants:
  - $append : -\text{all } a, b : \text{nat} - \underbrace{list(a)}_{\text{length } a} * list(b) \rightarrow list(a+b)$

# Background

- Part of a larger project...
  - Type assignment system\*
  - Tridirectional type system\*
  - Let-normal transform + type system
  - Implementation: Stardust ←— this paper

\* joint work with Frank Pfenning

# Outline

- ☞ **Introduction**
  - Integer Domain—Red-Black Trees
  - Dimension Domain
  - Implementation
  - Related Work
  - Summary and Future Work

# Setting

- Subset of core Standard ML:
  - First class functions
  - Primitive types `int`, `real`, `string`, ...
  - Algebraic datatypes

`datatype list = Nil`

`| Cons of int * list`

`datatype tree = Empty`

`| Tree of key * tree * tree`

✗ syntactic sugar

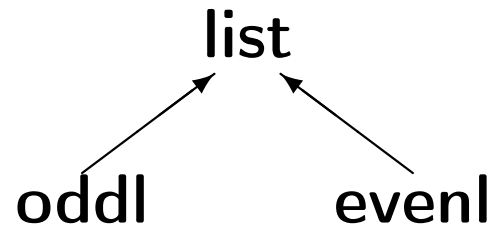
✓ patterns, e.g.

`Black(e, Red lt, Red (rt as (_,_,Red _)))`

# Datasort Refinements

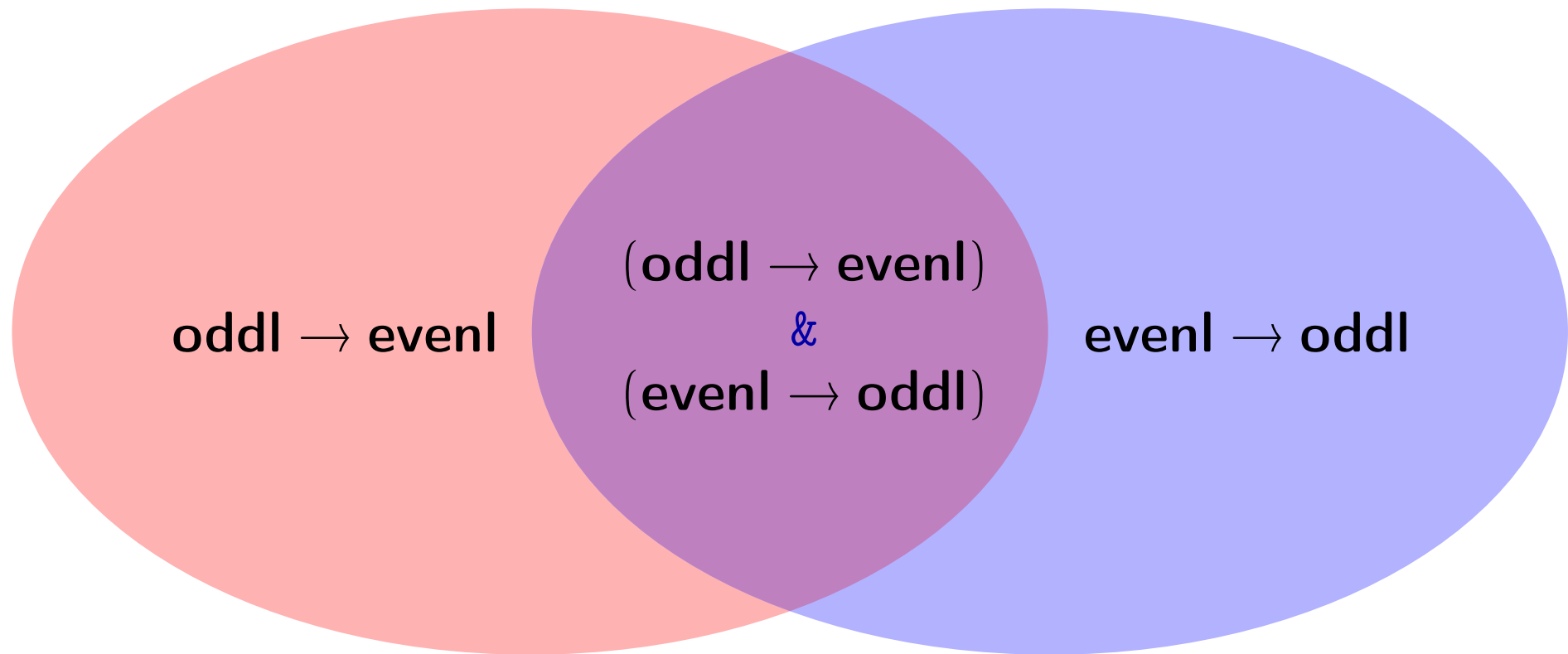
- a.k.a. refinement types [Freeman+Pf. '91, Davies '97,'05]
- Refine an algebraic datatype by a datasort
- **Example:** Lists of integers

$\text{Nil} : \text{list}$	$\text{Nil} : \text{eventl}$
$\text{Cons} : \text{int} * \text{list} \rightarrow \text{list}$	$\text{Cons} : (\text{int} * \text{oddl} \rightarrow \text{eventl})$ $\quad \& (\text{int} * \text{eventl} \rightarrow \text{oddl})$ $\quad \& (\text{int} * \text{list} \rightarrow \text{list})$



- Intersections  $v : A \& B$ :  $v$  has type  $A$  and type  $B$

# Intersection Types ~ Set Intersection



If

$tail : (oddl \rightarrow evenl) \ \& \ (evenl \rightarrow oddl)$

$od : oddl, \quad ev : evenl$

then

$tail(od) : evenl \ \text{and} \ tail(ev) : oddl$



# Index Refinements

- Dependent types with restrictions [Xi & Pfenning '99]
- Refine an algebraic datatype by an index
- Indices drawn from a decidable constraint domain
- **Example:** Lists indexed by their length

$\text{Nil} : \text{list}(0)$

$\text{Cons} : \text{-all } a : \text{nat} - \text{int} * \text{list}(a) \rightarrow \text{list}(a+1)$

$\text{append} : \text{-all } a, b : \text{nat} - \text{list}(a) * \text{list}(b) \rightarrow \text{list}(a+b)$

$\text{filter} : (\text{int} \rightarrow \text{bool})$

$\rightarrow \text{-all } a : \text{nat} - \text{list}(a) \rightarrow (\text{-exists } b : \text{nat} - [b \leq a] \text{list}(b))$

# Index Refinements + Union Types

- **Untagged unions**  $v : A \vee B$ :  $v$  has type  $A$  or type  $B$
- Choice function:  
 $choose : \text{-all } a, b : \text{nat- list}(a) * \text{list}(b) \rightarrow (\text{list}(a) \vee \text{list}(b))$   
 $\text{fun } choose (xs, ys) =$   
    **case**  $random\_bit()$  **of**  $true \Rightarrow xs$   
                                  **|**  $false \Rightarrow ys$
- Assume  $xs : \text{list}(a), ys : \text{list}(b)$
- Show body checks against  $\text{list}(a) \vee \text{list}(b)$ :
  - Show  $xs : \text{list}(a) \vee \text{list}(b)$ — show  $xs : \text{list}(a)$  ✓
  - Show  $ys : \text{list}(a) \vee \text{list}(b)$ — show  $ys : \text{list}(a)$  ✗  
                                  or, show  $ys : \text{list}(b)$  ✓

# Red-black trees

datatype tree = Empty

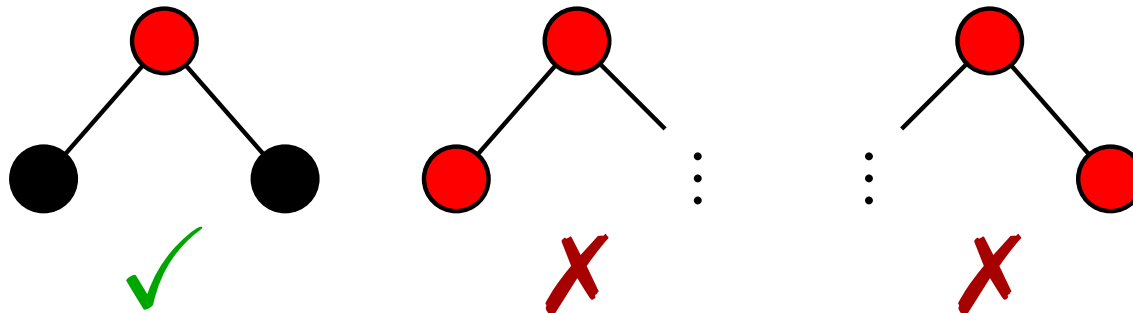
| Red of key \* tree \* tree

| Black of key \* tree \* tree

- Invariants:

(1) Empty nodes are considered black

(2) The children of a red node must be black



(3) For every node  $t$ , there exists  $bh(t)$  s.t. the number of black nodes on **all** paths from  $t$  to its leaves is  $bh(t)$

# Use datasort for invariants (1), (2)

tree

Empty : tree

Red : key \* tree \* tree  $\rightarrow$  tree

Black : key \* tree \* tree  $\rightarrow$  tree

# Use datasort for invariants (1), (2)

tree

↑

rbt

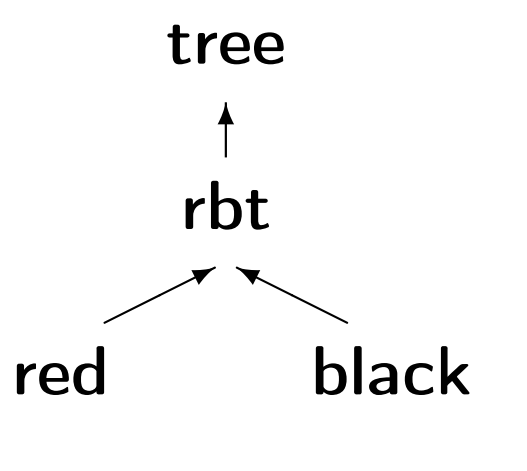
Empty : tree

Red : key \* tree \* tree → tree

Black : key \* tree \* tree → tree

- Add rbt to represent valid red-black trees

# Use datasort for invariants (1), (2)



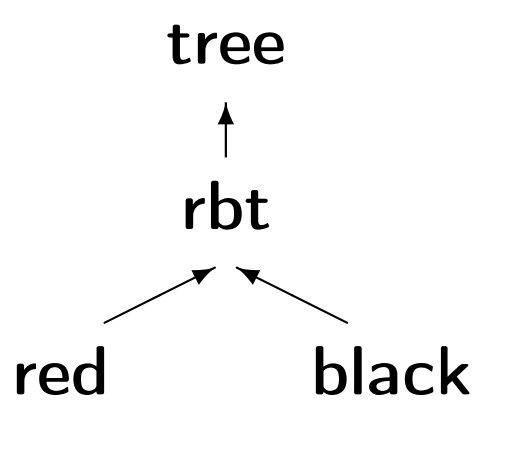
Empty : tree

Red : key \* tree \* tree → tree

Black : key \* tree \* tree → tree

- Add `rbt` to represent valid red-black trees
- Add `red`, `black` to distinguish colors

# Use datasort for invariants (1), (2)



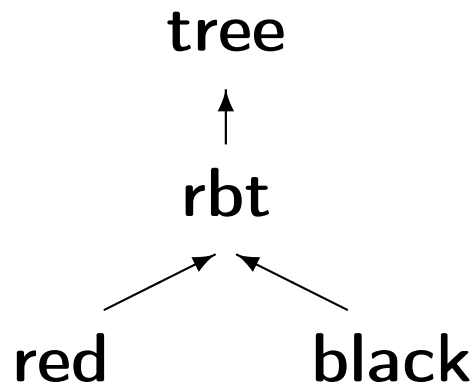
Empty : black

Red :  $\text{key} * \text{tree} * \text{tree} \rightarrow \text{tree}$

Black :  $\text{key} * \text{tree} * \text{tree} \rightarrow \text{tree}$

- Add `rbt` to represent valid red-black trees
- Add `red`, `black` to distinguish colors
- Invariant (1): Empty nodes are considered black

# Use datasort for invariants (1), (2)



Empty : black

Red : key \* tree \* tree  $\rightarrow$  tree

& key \* black \* black  $\rightarrow$  red

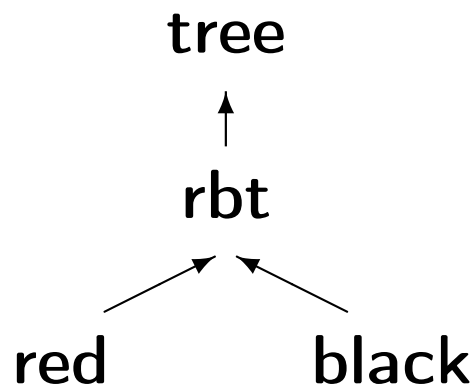
Black : key \* tree \* tree  $\rightarrow$  tree

& key \* rbt \* rbt  $\rightarrow$  black

- Add `rbt` to represent valid red-black trees
- Add `red`, `black` to distinguish colors
- Invariant (1): Empty nodes are considered black
- Invariant (2): The children of a red node must be black



# Use index for invariant (3)



Empty : black

Red :

key \* tree \* tree → tree

& key \* black \* black → red

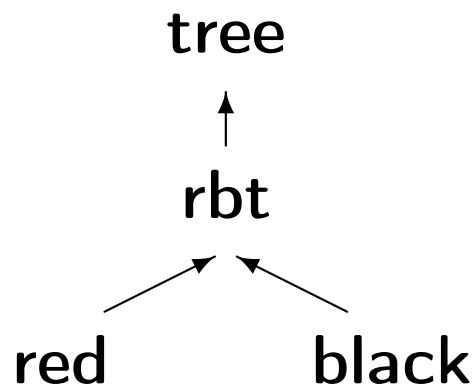
Black :

key \* tree \* tree → tree

& key \* rbt \* rbt → black

- Invariant (3): For every node  $t$ , there exists  $bh(t)$  s.t. the number of black nodes on **all** paths from  $t$  to its leaves is  $bh(t)$

# Use index for invariant (3)



Empty : black

Red :

key \* tree \* tree → tree

& key \* black \* black → red

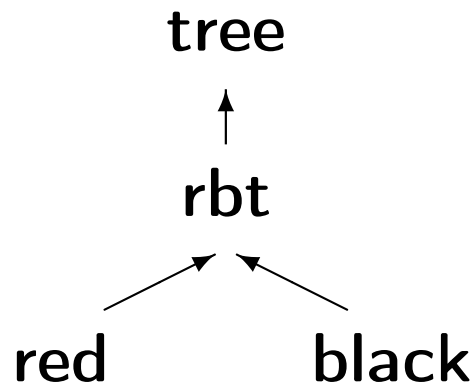
Black :

key \* tree \* tree → tree

& key \* rbt \* rbt → black

- Invariant (3): For every node  $t$ , there exists  $bh(t)$  s.t. the number of black nodes on **all** paths from  $t$  to its leaves is  $bh(t)$
- Index by the black height  $bh(t)$ , a natural number

# Use index for invariant (3)



Empty :  $\text{black}(0)$

Red :  $\text{-all } h:\text{nat-}$

$\text{key} * \text{tree}(h) * \text{tree}(h) \rightarrow \text{tree}(h)$

$\& \text{key} * \text{black}(h) * \text{black}(h) \rightarrow \text{red}(h)$

Black :  $\text{-all } h:\text{nat-}$

$\text{key} * \text{tree}(h) * \text{tree}(h) \rightarrow \text{tree}(h + 1)$

$\& \text{key} * \text{rbt}(h) * \text{rbt}(h) \rightarrow \text{black}(h + 1)$

- Invariant (3): For every node  $t$ , there exists  $bh(t)$  s.t. the number of black nodes on **all** paths from  $t$  to its leaves is  $bh(t)$
- Index by the black height  $bh(t)$ , a natural number

# What the user writes

```
(* [  
  datasort   rbt < tree;  
            red < rbt; black < rbt  
  
  datacon Empty : black(0)  
  datacon Red   : -all h:nat-  
                key * tree(h) * tree(h) → tree(h)  
                & key * black(h) * black(h) → red(h)  
  datacon Black : -all h:nat-  
                key * tree(h) * tree(h) → tree(h + 1)  
                & key * rbt(h) * rbt(h) → black(h + 1)  
]*)
```

# Index Domains

- Integers `int`
  - linear inequalities
- Booleans `bool`
- Dimensions `dim`

- ✓ Introduction
- ☞ **Integer Domain—Red-Black Trees**
  - Dimension Domain
  - Implementation
  - Related Work
  - Summary and Future Work

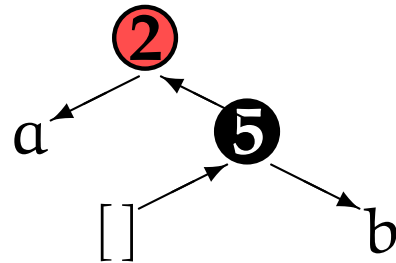
# Integers: Red-Black Trees

- Insertion (with slightly more complex datasorts)
- Deletion: SML/NJ library
  - Found two bugs, each violating the color invariant
  - Buggy module looks correct to clients
  - After bugs fixed, typechecks

# Zipper

$z$

[]

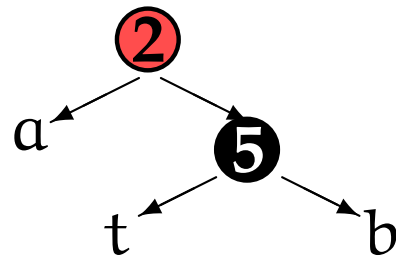


LEFTTB(RIGHTR( $a$ , 2, TOP),  
5,  
 $b$ )

TOP

$zip(z, t)$

$t$

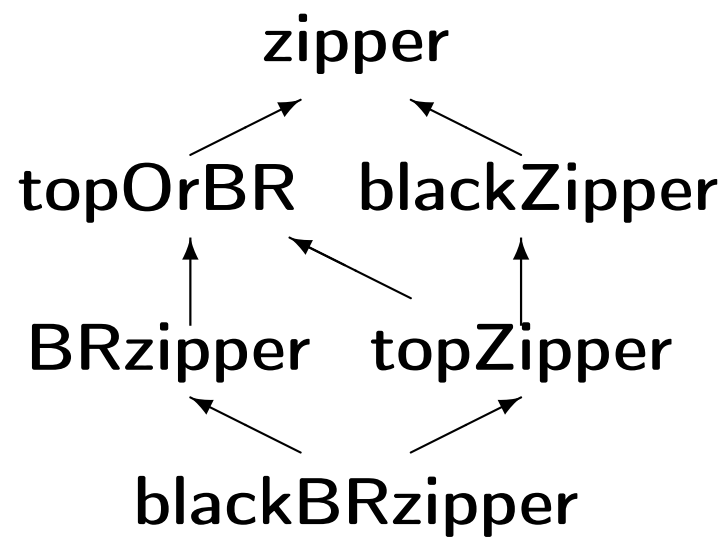


$t$

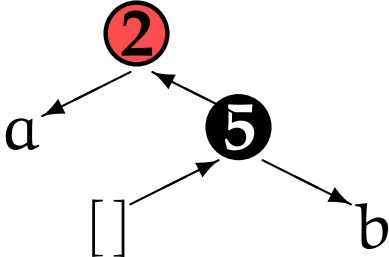
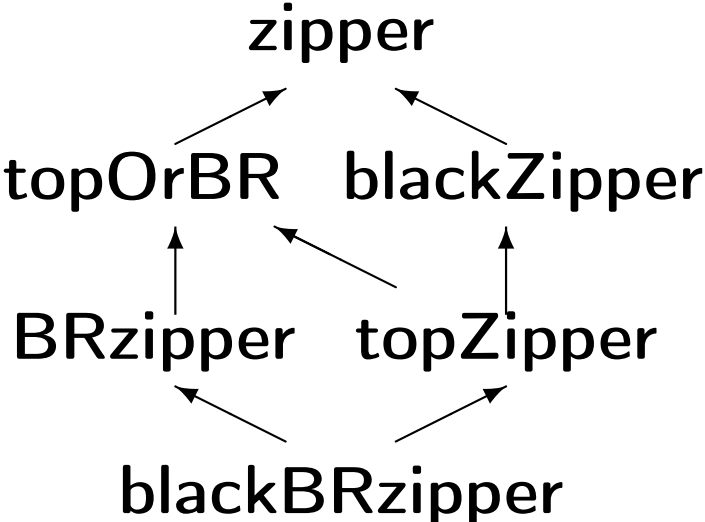
Red(2,  $a$ , Black(5,  $t$ ,  $b$ ))



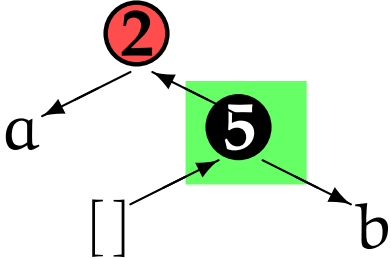
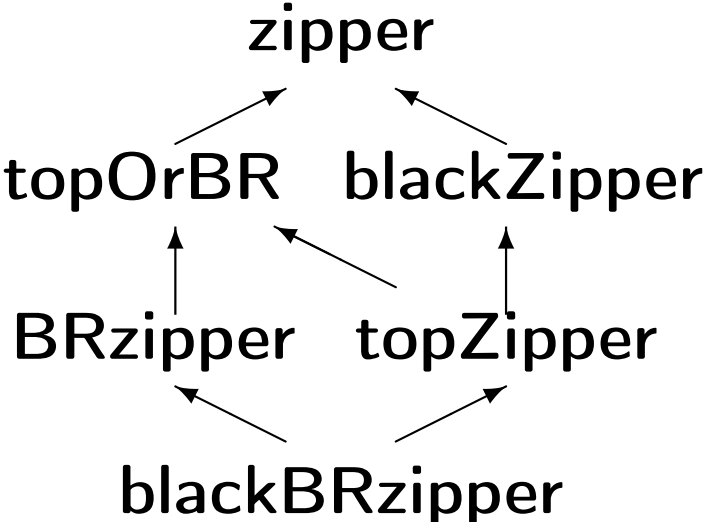
# Datasort Refinement



# Datasort Refinement

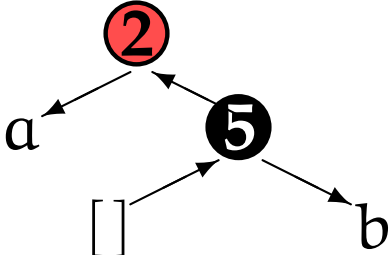
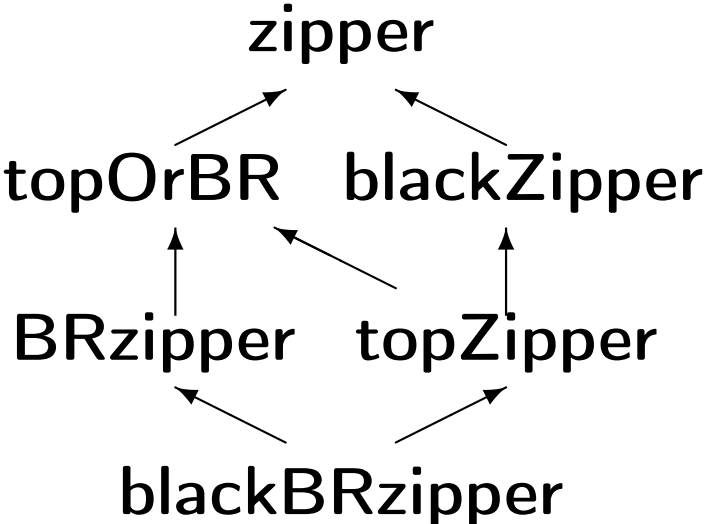


# Datasort Refinement

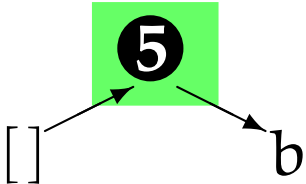


: blackZipper

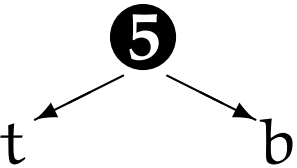
# Datasort Refinement



: blackZipper



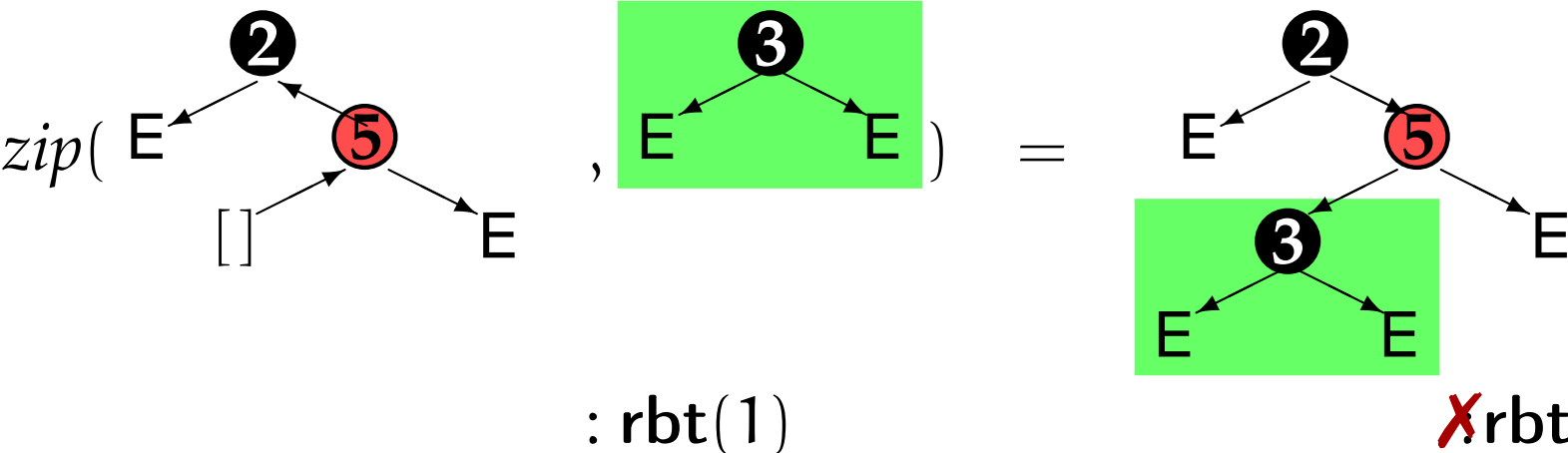
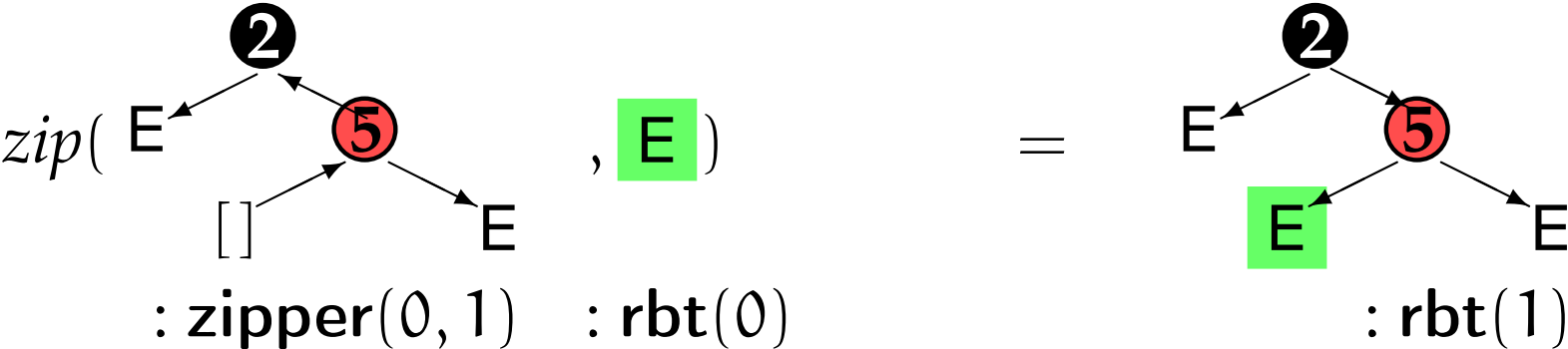
: BRzipper



: black

# Index Refinement: Height When Zipped

$z : \text{zipper}(h, hz)$  when  $\text{zip}(z, t)$  yields  $\text{rbt}(hz)$ ,  
 where  $t : \text{rbt}(h)$



# Zipper Constructor Types

`datacon TOP : -all h:nat- topZipper(h, h)`

...

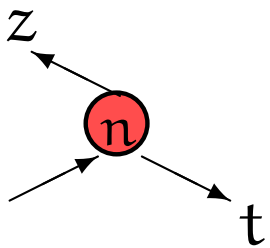
`datacon LEFTR : -all h, hz:nat-`

`int * black(h) * blackZipper(h, hz)`

`→ zipper(h, hz)`

`& int * black(h) * blackBRzipper(h, hz)`

`→ BRzipper(h, hz)`



`LEFTR(n, t, z)` where `t : black(h)`

- ✓ Introduction
- ✓ Integer Domain—Red-Black Trees
- ☞ **Dimension Domain**
  - Implementation
  - Related Work
  - Summary and Future Work

# Dimension Index Sort $\dim$

- Unit confusion a significant factor in the loss of the Mars Climate Orbiter

- Units M, S, KG

- Multiplication and powers:

$$M * KG \quad M^2 \quad S^a$$

- Multiplicative identity 1

- real indexed by  $\dim$



# Dimensions Example

```
(*[ val power : -all d:dim- -all n:int-  
      int(n) → real(d) → real(d^n)  
]*)  
fun power n x =  
  if n = 0 then  
    1.0  
  else if n < 0 then  
    1.0 / power (~n) x  
  else  
    x * power (n-1) x
```

# Dimensions: Invaluable Refinement

Unlike e.g. red-black trees,  
our dimension refinements are **not** tied to values:

- For every  $d$ , each type  $\mathbf{real}(d)$  is inhabited by exactly the same set of floating-point values:

$3.1 : \mathbf{real}(1)$

$3.1 * M : \mathbf{real}(M)$

$3.1 * S : \mathbf{real}(S)$

$3.1 * M * M : \mathbf{real}(M \wedge 2)$

All these expressions yield identical values at runtime!

- Also not based on values: qualified types [Foster], phantom types [Cheney/Hinze, Fluet/Pucella, ...], state refinement [Mandelbaum et al., ...]

# Dimensions: Related Work

- '70s and '80s
- Kennedy '96: Polymorphism, type inference
- Allen '04: abelian classes in Java

- ✓ Introduction
- ✓ Integer Domain—Red-Black Trees
- ✓ Dimension Domain
- ☞ **Implementation**
  - Related Work
  - Summary and Future Work

# Constraint Solving

- Delegate most work to external tools:  
ICS or CVC Lite
  - Accumulate constraints incrementally
- Solve for existential index variables
- Reduce dimension constraints to integer constraints

# (Almost) Ideal Constraint Solver Interface

1. Persistent solver contexts  $\Omega$   
encapsulating index-level assumptions
2.  $\text{ASSERT}(\Omega, P)$  returning
  - Valid if  $\Omega \models P$
  - Invalid if  $\Omega, P \models \perp$
  - Contingent( $\Omega'$ ) if  $\Omega \not\models P$  and  $\Omega, P \not\models \perp$ ,  
yielding a new context  $\Omega' = \Omega, P$
3.  $\text{VALID}(\Omega, P)$  returning
  - Valid if  $\Omega \models P$
  - Invalid otherwise.
4.  $\text{DECLARE}(\Omega, a, \textit{sort})$

# ICS

- Developed at SRI (de Moura, Rueß, Shankar, ...)
- Integers (almost), rationals, bitvectors (not really), functional arrays, ...
- Provides the desired interface
- Fast
- **Future:**
  - Closed source
  - Deprecated, replaced by Yices
  - For our purposes:  
ICS a great improvement on its successor

# ICS

- Developed at SRI (de Moura, Rueß, Shankar, ...)
- Integers (almost), rationals, bitvectors (not really), functional arrays, ...
- Provides the desired interface
- Fast  
except when it hangs
- **Future:**
  - Closed source
  - Deprecated, replaced by Yices
  - For our purposes:  
ICS a great improvement on its successor



# CVC Lite

- Barrett, Berezin, [Dill], [Tinelli], ...
- Integers, rationals, bitvectors, functional arrays, ...
- Does not provide persistent contexts
- Slower than ICS, but more stable
- **Future:**
  - Deprecated, replaced by CVC 3
  - Open source

# A nontraditional application?

- Most satisfiability-modulo-theories systems aimed at large problems
- Stardust needs to solve many small problems
- {Correctness, variety of theories}  $\gg$  speed

# Typechecking Time

Input program		Wall-clock time in seconds		
		ICS	library	stand-alone
redblack-full	$\delta(i)$ &	1.9	8.2	9.2
redblack-full-bug1	$\delta(i)$ &	1.6	6.8	8.1
redblack	$\delta(i)$ &	< 1	< 1	< 1
rbdelete	$\delta(i)$ &✓	*	37.7	31.6
bits	$\delta(i)$ &✓	*	9.5	4.0
bits-un	$\delta(i)$ &✓	33.5	298.5	241.4
others	$\delta(i)$	< 1	< 1	< 1

# Memoization

- Xi's DML:  
index refinements  
No  $\&V$ : typecheck in single pass,  
don't need to memoize
- Davies' SML-CIDRE:  
datasort refinements,  $\&$ : backtracking  
No index refinements—no  $\Sigma$ :  
memoization easy
- Stardust:  
datasort + index,  $\&V$   
memoization ineffective
  - might work with sophisticated irrelevance

- ✓ Introduction
- ✓ Integer Domain—Red-Black Trees
- ✓ Dimension Domain
- ✓ Implementation
- ☞ **Related Work**
  - Summary and Future Work

# Related Work

- [Davies '97/'05, Davies & Pf. '00]: Bidirectional system with  $\delta$ ,  $\&$
- [Xi & Pfenning '99]: Bidirectional, with  $i$ , **-all-**, **-exists-**
- [Coppo et al. '81, Amadio '98]:  $\&$  can characterize normal forms; hence full inference undecidable
- [Reynolds '96]: FORSYTHE with  $\&$  (and type annotations)
- [Pierce '91]: Language with  $\&$ ,  $\vee$ , syntactic markers
- Dependent types ((less) restricted): Cayenne, Epigram, [Chen & Xi '05], [Licata & Harper '05], ...

- ✓ Introduction
- ✓ Integer Domain—Red-Black Trees
- ✓ Dimension Domain
- ✓ Implementation
- ✓ Related Work
- ☞ **Summary and Future Work**

# Summary

- Type system
  - Combination of datasort and index refinements, marker-free intersections, unions, quantifiers
- Data structure refinements, e.g. bitstrings, red-black trees
- Invaluable refinements, e.g. dimensions
- Implementation
  - Delegated constraint solving
  - Speed leaves... room for improvement



# Future Work

- Parametric polymorphism
- Full SML (modules!)
- Index domains: bit vectors, inductive families, finite sets, regular languages, ...
- Memoization and parallelization
- Refinement-based compilation
- Call-by-name (and -need)
- Proof (derivation)  
Counterexample generation
- Suggestion tools

<http://type-refinements.info>

# Extra slides

# Datasort Refinements

*tail* : (eventl  $\rightarrow$  oddl) & (oddl  $\rightarrow$  eventl) & (list  $\rightarrow$  list)

fun *tail* xs =

  case xs of Nil  $\Rightarrow$  raise *Error*

          | Cons(h, t)  $\Rightarrow$  t

- First conjunct (eventl  $\rightarrow$  oddl):
  - Assume xs : eventl
  - Show raise *Error* : oddl and t : oddl
- Second conjunct (oddl  $\rightarrow$  eventl):
  - Assume xs : oddl
  - Show t : eventl

# Parametric Polymorphism: Declarative

- As property type

$$\frac{\Gamma \vdash e \uparrow \Pi a:\gamma. A \quad \bar{\Gamma} \vdash i:\gamma}{\Gamma \vdash e \uparrow [i/a] A} \text{PE}$$
$$\frac{\Gamma \vdash e \uparrow \forall \alpha. B \quad \Gamma \vdash A \text{ [mono]type}}{\Gamma \vdash e \uparrow [A/\alpha] B}$$

# Greed and Failure

- Generate existential variable
- The greedy method [Cardelli '93]
- $choose : \forall \alpha. \alpha \rightarrow \alpha \rightarrow \alpha$

```
fun choose x y =  
  case random_bit() of true  $\Rightarrow$  x  
    | false  $\Rightarrow$  y
```

$x : A, y : B \vdash choose\ x\ y \uparrow \dots$

# Greed and Failure

- Generate existential variable
- The greedy method [Cardelli '93]
- *choose* :  $\forall \alpha. \alpha \rightarrow \alpha \rightarrow \alpha$

```
fun choose x y =  
  case random_bit() of true  $\Rightarrow$  x  
    | false  $\Rightarrow$  y
```

$x : A, y : B \vdash \text{choose } x \ y \uparrow \dots$

- $A \leq B$  and  $B \not\leq A$ :  $y \downarrow A$  **X**

# Greed and Failure

- Generate existential variable
- The greedy method [Cardelli '93]
- *choose* :  $\forall \alpha. \alpha \rightarrow \alpha \rightarrow \alpha$

```
fun choose x y =  
  case random_bit() of true  $\Rightarrow$  x  
    | false  $\Rightarrow$  y
```

$x : A, y : B \vdash \text{choose } x \ y \uparrow \dots$

- $A \leq B$  and  $B \not\leq A$ :  $y \downarrow A$  ✗
- $A \not\leq B$  and  $B \leq A$ :  $y \downarrow A$  ✓



# Greed and Failure

- Generate existential variable
- The greedy method [Cardelli '93]
- *choose* :  $\forall \alpha. \alpha \rightarrow \alpha \rightarrow \alpha$

```
fun choose x y =  
  case random_bit() of true  $\Rightarrow$  x  
    | false  $\Rightarrow$  y
```

$x : A, y : B \vdash \text{choose } x \ y \uparrow \dots$

- $A \leq B$  and  $B \not\leq A$ :  $y \downarrow A$  ✗
- $A \not\leq B$  and  $B \leq A$ :  $y \downarrow A$  ✓
- $A \not\leq B$  and  $B \not\leq A$ :  $y \downarrow A$  ✗

# Unionism

- The greedy method seems to fail...

# Unionism

- The greedy method seems to fail...
- If only we could get an upper bound of A and B!

# Unionism

- The greedy method seems to fail...
- If only we could get an upper bound of A and B!

$A \vee B$

Instantiating  $\alpha = A \vee B$ :

$x : A, y : B \vdash \text{choose } x \ y \uparrow \dots$  ✓

# Typechecker, Phone Home

- Let's not break the subformula property

# Typechecker, Phone Home

- Let's not break the subformula property
- Instead, let the user write the union:

$(*[ \textit{choose} : \forall\alpha. \forall\beta. \alpha \rightarrow \beta \rightarrow (\alpha \vee \beta) ]*)$

# Typechecker, Phone Home

- Let's not break the subformula property
- Instead, let the user write the union:

$(*[ \textit{choose} : \forall\alpha. \forall\beta. \alpha \rightarrow \beta \rightarrow (\alpha \vee \beta) ]*)$

$x : A, y : B \vdash \textit{choose} \ x \ y \ \uparrow \ \dots$

$\alpha = A$

$\beta = B$

# Typechecker, Phone Home

- Let's not break the subformula property
- Instead, let the user write the union:

$(*[ \textit{choose} : \forall\alpha. \forall\beta. \alpha \rightarrow \beta \rightarrow (\alpha \vee \beta) ]*)$

$x : A, y : B \vdash \textit{choose} \ x \ y \ \uparrow \quad A \vee B \quad \checkmark$

$\alpha = A$

$\beta = B$