# Functional Data Structures

## with Isabelle/HOL

Tobias Nipkow

Fakultät für Informatik
Technische Universität München

2023-8-7

# Part II

# Functional Data Structures

# Chapter 6

# Sorting

$sorted :: ('a::linorder)\ list \Rightarrow bool$

$sorted\ [] = True$
$sorted\ (x\ \#\ ys) = ((\forall y \in set\ ys.\ x \leq y) \wedge sorted\ ys)$

# Correctness of sorting

Specification of $sort :: ('a::linorder)\ list \Rightarrow 'a\ list$:

$$sorted\ (sort\ xs)$$

Is that it? How about

$$set\ (sort\ xs) = set\ xs$$

Better: every $x$ occurs as often in $sort\ xs$ as in $xs$.

More succinctly:

$$mset\ (sort\ xs) = mset\ xs$$

where $mset :: 'a\ list \Rightarrow 'a\ multiset$

# What are multisets?

### Sets with (possibly) repeated elements

Some operations:

$$
\begin{array}{rcl}
\{\#\} & :: & {}'a \; multiset \\
add\_mset & :: & {}'a \Rightarrow {}'a \; multiset \Rightarrow {}'a \; multiset \\
+ & :: & {}'a \; multiset \Rightarrow {}'a \; multiset \Rightarrow {}'a \; multiset \\
mset & :: & {}'a \; list \Rightarrow {}'a \; multiset \\
set\_mset & :: & {}'a \; multiset \Rightarrow {}'a \; set
\end{array}
$$

Import $HOL-Library.Multiset$

# HOL/Data_Structures/Sorting.thy

Insertion Sort Correctness

# Principle: Count function calls

For every function $f :: \tau_1 \Rightarrow ... \Rightarrow \tau_n \Rightarrow \tau$
define a *timing function* $T_f :: \tau_1 \Rightarrow ... \Rightarrow \tau_n \Rightarrow nat$:

Translation of defining equations:
$$\mathcal{E}[\![f\ p_1\ ...\ p_n = e]\!]\ =\ (T_f\ p_1\ ...\ p_n = \mathcal{T}[\![e]\!] + 1)$$

Translation of expressions:
$$\mathcal{T}[\![g\ e_1\ ...\ e_k]\!]\ =\ \mathcal{T}[\![e_1]\!] + ... + \mathcal{T}[\![e_k]\!] + T_g\ e_1\ ...\ e_k$$

All other operations (variable access, constants, constructors, primitive operations on $bool$ and numbers) cost 1 time unit

# Example: @

$\mathcal{E}[\![\ [\ ]\ @\ ys = ys\ ]\!]$
$=\ (T_@\ [\ ]\ ys = \mathcal{T}[\![ys]\!]\ +\ 1)$
$=\ (T_@\ [\ ]\ ys = 2)$

$\mathcal{E}[\![\ (x\ \#\ xs)\ @\ ys = x\ \#\ (xs\ @\ ys)\ ]\!]$
$=\ (T_@\ (x\ \#\ xs)\ ys = \mathcal{T}[\![x\ \#\ (xs\ @\ ys)]\!]\ +\ 1)$
$=\ (T_@\ (x\ \#\ xs)\ ys = T_@\ xs\ ys\ +\ 5)$

$\mathcal{T}[\![x\ \#\ (xs\ @\ ys)]\!]$
$=\ \mathcal{T}[\![x]\!]\ +\ \mathcal{T}[\![xs\ @\ ys]\!]\ +\ T_\#\ x\ (xs\ @\ ys)$
$=\ 1\ +\ (\mathcal{T}[\![xs]\!]\ +\ \mathcal{T}[\![ys]\!]\ +\ T_@\ xs\ ys)\ +\ 1$
$=\ 1\ +\ (1\ +\ 1\ +\ T_@\ xs\ ys)\ +\ 1$

# *if* and *case*

So far we model a call-by-value semantics

Conditionals and case expressions are evaluated lazily.

$\mathcal{T}[\![\text{if } b \text{ then } e_1 \text{ else } e_2]\!]$
$= \mathcal{T}[\![b]\!] + (\text{if } b \text{ then } \mathcal{T}[\![e_1]\!] \text{ else } \mathcal{T}[\![e_2]\!])$

$\mathcal{T}[\![\text{case } e \text{ of } p_1 \Rightarrow e_1 \mid \ldots \mid p_k \Rightarrow e_k]\!]$
$= \mathcal{T}[\![e]\!] + (\text{case } e \text{ of } p_1 \Rightarrow \mathcal{T}[\![e_1]\!] \mid \ldots \mid p_k \Rightarrow \mathcal{T}[\![e_k]\!])$

Also special: let $x = t_1$ in $t_2$

# $O(.)$ is enough

$\implies$ Reduce all additive constants to 1

## Example

$T_@ \ (x \ \# \ xs) \ ys = \ T_@ \ xs \ ys + 5 \rightsquigarrow$
$T_@ \ (x \ \# \ xs) \ ys = \ T_@ \ xs \ ys + 1$

This means we count only

- the defined functions via $T_f$ and
- $+1$ for the function call itself.

All other operations (variables etc) cost 0, not 1.

# Discussion

- The definition of $T_f$ from $f$ can be automated.
- The correctness of $T_f$ could be proved w.r.t. a semantics that counts computation steps.
- Precise complexity bounds (as opposed to $O(.)$) would require a formal model of (at least) the compiler and the hardware.

# HOL/Data_Structures/Sorting.thy

Insertion sort complexity

```
merge :: 'a list ⇒ 'a list ⇒ 'a list
merge [] ys = ys
merge xs [] = xs
merge (x # xs) (y # ys) =
(if x ≤ y then x # merge xs (y # ys)
 else y # merge (x # xs) ys)

msort :: 'a list ⇒ 'a list
msort xs =
(let n = length xs
 in if n ≤ 1 then xs
    else merge (msort (take (n div 2) xs))
           (msort (drop (n div 2) xs)))
```

# Number of comparisons

$C\_merge :: {}'a\ list \Rightarrow {}'a\ list \Rightarrow nat$
$C\_msort :: {}'a\ list \Rightarrow nat$

## Lemma
$C\_merge\ xs\ ys$

## Theorem
$length\ xs = 2^k \implies C\_msort\ xs \leq k * 2^k$

# HOL/Data_Structures/Sorting.thy

Merge Sort

$msort\_bu :: \ 'a \ list \Rightarrow \ 'a \ list$

$msort\_bu \ xs = \ merge\_all \ (map \ (\lambda x. \ [x]) \ xs)$

$merge\_all :: \ 'a \ list \ list \Rightarrow \ 'a \ list$

$merge\_all \ [] = []$

$merge\_all \ [xs] = xs$

$merge\_all \ xss = \ merge\_all \ (merge\_adj \ xss)$

$merge\_adj :: \ 'a \ list \ list \Rightarrow \ 'a \ list \ list$

$merge\_adj \ [] = []$

$merge\_adj \ [xs] = [xs]$

$merge\_adj \ (xs \ \# \ ys \ \# \ zss) =$

$merge \ xs \ ys \ \# \ merge\_adj \ zss$

# Number of comparisons

$C\_merge\_adj :: {'}a\ list\ list \Rightarrow nat$

$C\_merge\_all :: {'}a\ list\ list \Rightarrow nat$

$C\_msort\_bu :: {'}a\ list \Rightarrow nat$

## Theorem

$length\ xs = 2^k \implies C\_msort\_bu\ xs \leq k * 2^k$

# HOL/Data_Structures/Sorting.thy

Bottom-Up Merge Sort

# Even better

Make use of already sorted subsequences

## Example

Sorting $[7, 3, 1, 2, 5]$:

do not start with $[[7], [3], [1], [2], [5]]$

but with $[[1, 3, 7], [2, 5]]$

# Archive of Formal Proofs

```
https://www.isa-afp.org/entries/
    Efficient-Mergesort.shtml
```

# Chapter 7

## Binary Trees

**5** Binary Trees

**6** Basic Functions

**7** (Almost) Complete Trees

**5** Binary Trees

**6** Basic Functions

**7** (Almost) Complete Trees

HOL/Library/Tree.thy

# Binary trees

**datatype** $'a\ tree = Leaf\ |\ Node\ ('a\ tree)\ 'a\ ('a\ tree)$

Abbreviations:
$$\langle\rangle \equiv Leaf$$
$$\langle l,\ a,\ r\rangle \equiv Node\ l\ a\ r$$

Most of the time: tree = binary tree

# Tree traversal

$inorder :: {}'a\ tree \Rightarrow {}'a\ list$

$inorder\ \langle\rangle = []$
$inorder\ \langle l,\ x,\ r\rangle = inorder\ l\ @\ [x]\ @\ inorder\ r$

$preorder :: {}'a\ tree \Rightarrow {}'a\ list$

$preorder\ \langle\rangle = []$
$preorder\ \langle l,\ x,\ r\rangle = x\ \#\ preorder\ l\ @\ preorder\ r$

$postorder :: {}'a\ tree \Rightarrow {}'a\ list$

$postorder\ \langle\rangle = []$
$postorder\ \langle l,\ x,\ r\rangle = postorder\ l\ @\ postorder\ r\ @\ [x]$

# Size

$size :: {}'a\ tree \Rightarrow nat$
$|\langle\rangle| = 0$
$|\langle l,\ {}_-,\ r\rangle| = |l| + |r| + 1$

$size1 :: {}'a\ tree \Rightarrow nat$
$|\langle\rangle|_1 = 1$
$|\langle l,\ {}_-,\ r\rangle|_1 = |l|_1 + |r|_1$

**Lemma** $|t|_1 = |t| + 1$

Warning: $|.|$ and $|.|_1$ only on slides

# Height

$height :: {'}a\ tree \Rightarrow nat$
$h(\langle\rangle) = 0$
$h(\langle l,\ \_,\ r\rangle) = max\ (h(l))\ (h(r)) + 1$

Warning: $h(.)$ only on slides

**Lemma** $h(t) \leq |t|$

**Lemma** $|t|_1 \leq 2^{h(t)}$

# Minimal height

$min\_height :: \ 'a \ tree \Rightarrow nat$
$mh(\langle\rangle) = 0$
$mh(\langle l, \ \_, \ r\rangle) = min \ (mh(l)) \ (mh(r)) + 1$

Warning: $mh(.)$ only on slides

**Lemma** $mh(t) \leq h(t)$

**Lemma** $2^{mh(t)} \leq |t|_1$

# Complete tree

$complete :: {'}a\ tree \Rightarrow bool$
$complete\ \langle\rangle\ =\ True$
$complete\ \langle l,\ \_,\ r\rangle\ =$
$(h(l)\ =\ h(r)\ \wedge\ complete\ l\ \wedge\ complete\ r)$

**Lemma** $complete\ t\ =\ (mh(t)\ =\ h(t))$

**Lemma** $complete\ t\ \Longrightarrow\ |t|_1\ =\ 2^{h(t)}$

**Lemma** $\neg\ complete\ t\ \Longrightarrow\ |t|_1\ <\ 2^{h(t)}$
**Lemma** $\neg\ complete\ t\ \Longrightarrow\ 2^{mh(t)}\ <\ |t|_1$

**Corollary** $|t|_1\ =\ 2^{h(t)}\ \Longrightarrow\ complete\ t$
**Corollary** $|t|_1\ =\ 2^{mh(t)}\ \Longrightarrow\ complete\ t$

# Almost complete tree

$acomplete :: {'}a\ tree \Rightarrow bool$
$acomplete\ t = (h(t) - mh(t) \leq 1)$

Almost complete trees have optimal height:
**Lemma** If $acomplete\ t$ and $|t| \leq |t'|$ then $h(t) \leq h(t')$.

# Warning

- The terms *complete* and *almost complete* are not defined uniquely in the literature.

- For example,
  Knuth calls *complete* what we call *almost complete*.

# Chapter 8

# Search Trees

Most of the material focuses on
BSTs = binary search trees

# BSTs represent sets

Any tree represents a set:

$set\_tree :: {}'a\ tree \Rightarrow {}'a\ set$

$set\_tree\ \langle\rangle = \{\}$
$set\_tree\ \langle l,\ x,\ r\rangle = set\_tree\ l \cup \{x\} \cup set\_tree\ r$

A BST represents a set that can be searched in time $O(h(t))$

Function $set\_tree$ is called an *abstraction function*
because it maps the implementation
to the abstract mathematical object

# $bst$

$bst :: {}'a\ tree \Rightarrow bool$

$bst\ \langle\rangle\ =\ True$
$bst\ \langle l,\ a,\ r\rangle\ =$
$((\forall\, x {\in} set\_tree\ l.\ x\ <\ a)\ \wedge$
$\ (\forall\, x {\in} set\_tree\ r.\ a\ <\ x)\ \wedge\ bst\ l\ \wedge\ bst\ r)$

Type ${}'a$ must be in class $linorder$ (${}'a :: linorder$) where $linorder$ are *linear orders* (also called *total orders*).

Note: $nat$, $int$ and $real$ are in class $linorder$

# Set interface

An implementation of sets of elements of type $\,'a\,$ must provide

- An implementation type $\,'s$
- $empty :: {}'s$
- $insert :: {}'a \Rightarrow {}'s \Rightarrow {}'s$
- $delete :: {}'a \Rightarrow {}'s \Rightarrow {}'s$
- $isin :: \quad {}'s \Rightarrow {}'a \Rightarrow bool$

# Map interface

Instead of a set, a search tree can also implement a map from $'a$ to $'b$:

- An implementation type $\ 'm$
- $empty :: 'm$
- $update :: 'a \Rightarrow 'b \Rightarrow 'm \Rightarrow 'm$
- $delete :: 'a \Rightarrow 'm \Rightarrow 'm$
- $lookup :: 'm \Rightarrow 'a \Rightarrow 'b\ option$

Sets are a special case of maps

# Comparison of elements

We assume that the element type $'a$ is a linear order

Instead of using $<$ and $\leq$ directly:

**datatype** $cmp\_val = LT \mid EQ \mid GT$

$cmp\ x\ y =$
(if $x < y$ then $LT$ else if $x = y$ then $EQ$ else $GT$)

Implementation type: $'a\ tree$

$empty = Leaf$

$insert\ x\ \langle\rangle = \langle\langle\rangle,\ x,\ \langle\rangle\rangle$
$insert\ x\ \langle l,\ a,\ r\rangle = ($case $cmp\ x\ a$ of
$\qquad\qquad\qquad LT \Rightarrow \langle insert\ x\ l,\ a,\ r\rangle$
$\qquad\qquad\ \ |\ EQ \Rightarrow \langle l,\ a,\ r\rangle$
$\qquad\qquad\ \ |\ GT \Rightarrow \langle l,\ a,\ insert\ x\ r\rangle)$

$$isin \; \langle \rangle \; x = \; False$$

$$isin \; \langle l, \; a, \; r \rangle \; x = (\mathsf{case} \; cmp \; x \; a \; \mathsf{of}$$

$$LT \Rightarrow isin \; l \; x$$

$$| \; EQ \Rightarrow True$$

$$| \; GT \Rightarrow isin \; r \; x)$$

$delete\ x\ \langle \rangle = \langle \rangle$
$delete\ x\ \langle l,\ a,\ r \rangle =$
(case $cmp\ x\ a$ of
  $LT \Rightarrow \langle delete\ x\ l,\ a,\ r \rangle$
 | $EQ \Rightarrow$ if $r = \langle \rangle$ then $l$
        else let $(a',\ r') = split\_min\ r$ in $\langle l,\ a',\ r' \rangle$
 | $GT \Rightarrow \langle l,\ a,\ delete\ x\ r \rangle$)

$split\_min\ \langle l,\ a,\ r \rangle =$
(if $l = \langle \rangle$ then $(a,\ r)$
 else let $(x,\ l') = split\_min\ l$ in $(x,\ \langle l',\ a,\ r \rangle))$

# Why is this implementation correct?

Because   $empty$   $insert$   $delete$   $isin$
simulate    $\{\}$    $\cup\{.\}$   $-\{.\}$   $\in$

$$set\_tree\ empty = \{\}$$
$$set\_tree\ (insert\ x\ t) = set\_tree\ t \cup \{x\}$$
$$set\_tree\ (delete\ x\ t) = set\_tree\ t - \{x\}$$
$$isin\ t\ x = (x \in set\_tree\ t)$$

Under the assumption $bst\ t$

# Also: *bst* must be invariant

$$bst\ empty$$
$$bst\ t \implies bst\ (insert\ x\ t)$$
$$bst\ t \implies bst\ (delete\ x\ t)$$

**8** Unbalanced BST
  Implementation
  Correctness
  Correctness Proof Method Based on Sorted Lists

# Key idea

Local definition:

$sorted$ means sorted w.r.t. $<$

No duplicates!

$\implies$   $bst\ t$ can be expressed as $sorted(inorder\ t)$

Conduct proofs on sorted lists, not sets

# Two kinds of invariants

- Unbalanced trees only need the invariant $bst$
- More efficient search trees come with additional *structural invariants* $=$ balance criteria.

# Correctness via sorted lists

Correctness proofs of (almost) all search trees
covered in this course
can be automated.

Except for the structural invariants.

Therefore we concentrate on the latter.

For details see file See `HOL/Data_Structures/Set_Specs.thy` and
T. Nipkow. *Automatic Functional Correctness Proofs for Functional Search Trees.* Interactive Theorem Proving, LNCS, 2016.

A methodological interlude:

> A closer look at ADT principles
> and their realization in Isabelle

> Set and binary search tree as examples
> (ignoring *delete*)

**9** Abstract Data Types
  Defining ADTs
  Using ADTs
  Implementing ADTs

ADT $=$ *interface* $+$ *specification*

## Example (Set interface)

$empty :: \ 's$
$insert :: \ 'a \Rightarrow \ 's \Rightarrow \ 's$
$isin :: \quad 's \Rightarrow \ 'a \Rightarrow \ bool$

We assume that each ADT describes one

*Type of Interest* $T$

Above: $T = \ 's$

# Model-oriented specification

Specify type $T$ via a model $=$ existing HOL type $A$
Motto: $T$ should behave like $A$

Specification of "behaves like" via an

- *abstraction function* $\alpha :: T \Rightarrow A$

Only some elements of $T$ represent elements of $A$:

- *invariant* $invar :: T \Rightarrow bool$

$\alpha$ and $invar$ are part of the interface,
but only for specification and verification purposes

# Example (Set ADT)

$empty :: \ldots$
$insert :: \ldots$
$isin :: \ldots$
$set :: \quad 's \Rightarrow 'a \; set$   (name arbitrary)
$invar :: \; 's \Rightarrow bool$    (name arbitrary)

$$set \; empty = \{\}$$

$invar \; s \Longrightarrow \quad set(insert \; x \; s) = set \; s \cup \{x\}$

$invar \; s \Longrightarrow \quad isin \; s \; x = (x \in set \; s)$

$$invar \; empty$$

$invar \; s \Longrightarrow \quad invar(insert \; x \; s)$

# In Isabelle: **locale**

**locale** $Set =$

**fixes** $empty :: \ 's$

**fixes** $insert :: \ 'a \Rightarrow \ 's \Rightarrow \ 's$

**fixes** $isin :: \ 's \Rightarrow \ 'a \Rightarrow bool$

**fixes** $set :: \ 's \Rightarrow \ 'a \ set$

**fixes** $invar :: \ 's \Rightarrow bool$

**assumes** $set \ empty = \{\}$

**assumes** $invar \ s \implies isin \ s \ x = (x \in set \ s)$

**assumes** $invar \ s \implies set(insert \ x \ s) = set \ s \cup \{x\}$

**assumes** $invar \ empty$

**assumes** $invar \ s \implies invar(insert \ x \ s)$

See HOL/Data_Structures/Set_Specs.thy

# Formally, in general

To ease notation, generalize $\alpha$ and $invar$ (conceptually):
$\alpha$ is the identity and $invar$ is $True$
on types other than $T$

Specification of each interface function $f$ (on $T$):

- $f$ must behave like some function $f_A$ (on $A$):
  $$invar\ t_1 \wedge ... \wedge invar\ t_n \implies$$
  $$\alpha(f\ t_1\ ...\ t_n) = f_A\ (\alpha\ t_1)\ ...\ (\alpha\ t_n)$$
  ($\alpha$ is a homomorphism)

- $f$ must preserve the invariant:
  $$invar\ t_1 \wedge ... \wedge invar\ t_n \implies invar(f\ t_1\ ...\ t_n)$$

The purpose of an ADT is to provide a context
for implementing generic algorithms
parameterized with the interface functions of the ADT.

## Example

**locale** $Set =$
**fixes** ...
**assumes** ...
**begin**

**fun** $set\_of\_list$ **where**
$set\_of\_list\ [] = empty\ |$
$set\_of\_list\ (x\ \#\ xs) = insert\ x\ (set\_of\_list\ xs)$

**lemma** $invar(set\_of\_list\ xs)$
**by**$(induction\ xs)$
  $(auto\ simp{:}\ invar\_empty\ invar\_insert)$

**end**

1. Implement interface
2. Prove specification

## Example

Define functions *isin* and *insert* on type $'a\ tree$ with invariant *bst*.

Now implement locale *Set*:

# In Isabelle: **interpretation**

**interpretation** *Set*
**where** *empty = Leaf* **and** *isin = isin*
**and** *insert = insert* **and** *set = set_tree* **and** *invar = bst*
**proof**
  **show** *set_tree Leaf = {}* ⟨*proof*⟩
**next**
  **fix** *s* **assume** *bst s*
  **show** *set_tree (insert x s) = set_tree s ∪ {x}*
  ⟨*proof*⟩
**next**
⋮
**qed**

Interpretation of $Set$ also yields

- function $set\_of\_list :: {}'a\ list \Rightarrow {}'a\ tree$
- theorem $bst\ (set\_of\_list\ xs)$

Now back to search trees …

HOL/Data_Structures/
Tree23_Set.thy

# 2-3 Trees

**datatype** $'a\ tree23 = \langle\rangle$
  | $Node2\ ('a\ tree23)\ 'a\ ('a\ tree23)$
  | $Node3\ ('a\ tree23)\ 'a\ ('a\ tree23)\ 'a\ ('a\ tree23)$

Abbreviations:

$$\langle l,\ a,\ r\rangle \quad\equiv\quad Node2\ l\ a\ r$$
$$\langle l,\ a,\ m,\ b,\ r\rangle \quad\equiv\quad Node3\ l\ a\ m\ b\ r$$

# *isin*

*isin* ⟨*l, a, m, b, r*⟩ *x* =
(case *cmp x a* of
    *LT* ⇒ *isin l x*
 | *EQ* ⇒ *True*
 | *GT* ⇒ case *cmp x b* of
                *LT* ⇒ *isin m x*
            | *EQ* ⇒ *True*
            | *GT* ⇒ *isin r x*)

Assumes the usual ordering invariant

# Structural invariant *complete*

All leaves are at the same level:

$complete \; \langle \rangle = True$

$complete \; \langle l, \; \llcorner, \; r \rangle =$
$(h(l) = h(r) \wedge complete \; l \wedge complete \; r)$

$complete \; \langle l, \; \llcorner, \; m, \; \llcorner, \; r \rangle =$
$(h(l) = h(m) \wedge h(m) = h(r) \; \wedge$
$\; complete \; l \wedge complete \; m \wedge complete \; r)$

## Lemma
$complete \; t \Longrightarrow 2^{h(t)} \leq |t| + 1$

# Insertion

The idea:

$$
\begin{array}{rcl}
Leaf & \rightsquigarrow & Node2 \\
Node2 & \rightsquigarrow & Node3 \\
Node3 & \rightsquigarrow & \text{overflow, pass 1 element back up}
\end{array}
$$

# Insertion

Two possible return values:

- tree accommodates new element without increasing height: $TI\ t$
- tree overflows: $OF\ l\ x\ r$

**datatype** $'a\ upI = TI\ ('a\ tree23)$
$|\ OF\ ('a\ tree23)\ 'a\ ('a\ tree23)$

$treeI :: 'a\ upI \Rightarrow 'a\ tree23$
$treeI\ (TI\ t) = t$
$treeI\ (OF\ l\ a\ r) = \langle l,\ a,\ r \rangle$

# Insertion

$insert :: {}'a \Rightarrow {}'a\ tree23 \Rightarrow {}'a\ tree23$

$insert\ x\ t = treeI\ (ins\ x\ t)$

$ins :: {}'a \Rightarrow {}'a\ tree23 \Rightarrow {}'a\ upI$

# Insertion

$ins\ x\ \langle\rangle\ =\ OF\ \langle\rangle\ \ x\ \langle\rangle$

$ins\ x\ \langle l,\ a,\ r\rangle\ =$

case $cmp\ x\ a$ of

  $LT \Rightarrow$ case $ins\ x\ l$ of

        $TI\ l' \Rightarrow\ TI\ \langle l',\ a,\ r\rangle$

       $|\ OF\ l_1\ \ b\ \ l_2 \Rightarrow\ TI\ \langle l_1,\ b,\ l_2,\ a,\ r\rangle$

$|\ EQ \Rightarrow\ TI\ \langle l,\ a,\ r\rangle$

$|\ GT \Rightarrow$ case $ins\ x\ r$ of

        $TI\ r' \Rightarrow\ TI\ \langle l,\ a,\ r'\rangle$

       $|\ OF\ r_1\ \ b\ \ r_2 \Rightarrow\ TI\ \langle l,\ a,\ r_1,\ b,\ r_2\rangle$

# Insertion

$ins\ x\ \langle l,\ a,\ m,\ b,\ r\rangle\ =$

case $cmp\ x\ a$ of

$\quad LT \Rightarrow$ case $ins\ x\ l$ of

$\qquad\qquad TI\ l' \Rightarrow TI\ \langle l',\ a,\ m,\ b,\ r\rangle$

$\qquad\qquad |\ OF\ l_1\ c\ l_2 \Rightarrow OF\ \langle l_1,\ c,\ l_2\rangle\ a\ \langle m,\ b,\ r\rangle$

$|\ EQ \Rightarrow TI\ \langle l,\ a,\ m,\ b,\ r\rangle$

$|\ GT \Rightarrow$

$\quad$ case $cmp\ x\ b$ of

$\qquad LT \Rightarrow$ case $ins\ x\ m$ of

$\qquad\qquad\ TI\ m' \Rightarrow TI\ \langle l,\ a,\ m',\ b,\ r\rangle$

$\qquad\qquad |\ OF\ m_1\ c\ m_2 \Rightarrow OF\ \langle l,\ a,\ m_1\rangle\ c\ \langle m_2,\ b,\ r\rangle$

$\quad |\ EQ \Rightarrow TI\ \langle l,\ a,\ m,\ b,\ r\rangle$

$\quad |\ GT \Rightarrow$ case $ins\ x\ r$ of

$\qquad\qquad TI\ r' \Rightarrow TI\ \langle l,\ a,\ m,\ b,\ r'\rangle$

# Insertion preserves *complete*

Lemma
*complete t $\Longrightarrow$*
*complete (treeI (ins a t)) $\wedge$ hI (ins a t) = h(t)*
*where hI :: 'a upI $\Rightarrow$ nat*
*hI (TI t) = h(t)*
*hI (OF l a r) = h(l)*
**Proof** by induction on $t$. Base and step automatic.

Corollary
*complete t $\Longrightarrow$ complete (insert a t)*

# Deletion

The idea:

$$Node3 \quad \rightsquigarrow \quad Node2$$
$$Node2 \quad \rightsquigarrow \quad \textcolor{red}{\text{underflow}}, \text{ height decreases by 1}$$

Underflow: merge with siblings on the way up

# Deletion

Two possible return values:

- height unchanged: $TD\ t$
- height decreased by 1: $UF\ t$

**datatype** $'a\ upD\ =\ TD\ ('a\ tree23)\ |\ UF\ ('a\ tree23)$

$treeD\ (TD\ t)\ =\ t$
$treeD\ (UF\ t)\ =\ t$

# Deletion

$delete :: {}'a \Rightarrow {}'a\ tree23 \Rightarrow {}'a\ tree23$

$delete\ x\ t = treeD\ (del\ x\ t)$

$del :: {}'a \Rightarrow {}'a\ tree23 \Rightarrow {}'a\ upD$

# Deletion

$del\ x\ \langle\rangle\ =\ TD\ \langle\rangle$

$del\ x\ \langle\langle\rangle,\ a,\ \langle\rangle\rangle\ =$
(if $x = a$ then $UF\ \langle\rangle$ else $TD\ \langle\langle\rangle,\ a,\ \langle\rangle\rangle$)

$del\ x\ \langle\langle\rangle,\ a,\ \langle\rangle,\ b,\ \langle\rangle\rangle\ =\ ...$

$del\ x\ \langle l,\ a,\ r \rangle =$
$(\mathsf{case}\ cmp\ x\ a\ \mathsf{of}$
  $LT \Rightarrow node21\ (del\ x\ l)\ a\ r$
$|\ EQ \Rightarrow \mathsf{let}\ (a',\ t) = split\_min\ r\ \mathsf{in}\ node22\ l\ a'\ t$
$|\ GT \Rightarrow node22\ l\ a\ (del\ x\ r))$

$node21\ (TD\ t_1)\ a\ t_2 = TD\ \langle t_1,\ a,\ t_2 \rangle$
$node21\ (UF\ t_1)\ a\ \langle t_2,\ b,\ t_3 \rangle = UF\ \langle t_1,\ a,\ t_2,\ b,\ t_3 \rangle$
$node21\ (UF\ t_1)\ a\ \langle t_2,\ b,\ t_3,\ c,\ t_4 \rangle =$
$TD\ \langle\langle t_1,\ a,\ t_2 \rangle,\ b,\ \langle t_3,\ c,\ t_4 \rangle\rangle$

Analogous: $node22$

# Deletion preserves *complete*

After 13 simple lemmas:

## Lemma
$complete\ t \implies complete\ (treeD\ (del\ x\ t))$

## Corollary
$complete\ t \implies complete\ (delete\ x\ t)$

# Beyond 2-3 trees

**datatype** $'a\ tree234 =$
   $Leaf\ |\ Node2\ ...\ |\ Node3\ ...\ |\ Node4\ ...$

Like 2-3 trees, but with many more cases

The general case:

B-trees and $(a, b)$-trees

HOL/Data_Structures/
RBT_Set.thy

# Relationship to 2-3-4 trees

Idea:   encode 2-3-4 trees as binary trees;
        use color to express grouping

$$
\begin{aligned}
\langle\rangle &\approx \langle\rangle \\
\langle t_1, a, t_2 \rangle &\approx \langle t_1, a, t_2 \rangle \\
\langle t_1, a, t_2, b, t_3 \rangle &\approx \langle \langle t_1, a, t_2 \rangle, b, t_3 \rangle \text{ or } \langle t_1, a, \langle t_2, b, t_3 \rangle \rangle \\
\langle t_1, a, t_2, b, t_3, c, t_4 \rangle &\approx \langle \langle t_1, a, t_2 \rangle, b, \langle t_3, c, t_4 \rangle \rangle
\end{aligned}
$$

Red means "I am part of a bigger node"

# Structural invariants

- The root is
- Every $\langle \rangle$ is considered Black.
- If a node is <span style="color:red">Red</span>,
- All paths from a node to a leaf have the same number of

# Red-black trees

**datatype** $color = Red \mid Black$

**type_synonym** $'a\ rbt = ('a \times color)\ tree$

Abbreviations:

$$
\begin{array}{rcl}
R\ l\ a\ r & \equiv & Node\ l\ (a,\ Red)\ r \\
B\ l\ a\ r & \equiv & Node\ l\ (a,\ Black)\ r
\end{array}
$$

# Color

$color :: {}'a\ rbt \Rightarrow color$
$color\ \langle\rangle\ =\ Black$
$color\ \langle\_,\ (\_,\ c),\ \_\rangle\ =\ c$

$paint :: color \Rightarrow {}'a\ rbt \Rightarrow {}'a\ rbt$
$paint\ c\ \langle\rangle\ =\ \langle\rangle$
$paint\ c\ \langle l,\ (a,\ \_),\ r\rangle\ =\ \langle l,\ (a,\ c),\ r\rangle$

# Structural invariants

$rbt :: {}'a\ rbt \Rightarrow bool$
$rbt\ t = (invc\ t \wedge invh\ t \wedge color\ t = Black)$

$invc :: {}'a\ rbt \Rightarrow bool$
$invc\ \langle\rangle = True$
$invc\ \langle l,\ (\_,\ c),\ r\rangle =$
$((c = Red \longrightarrow color\ l = Black \wedge color\ r = Black) \wedge$
$\ invc\ l \wedge invc\ r)$

# Structural invariants

$invh :: {}'a\ rbt \Rightarrow bool$
$invh\ \langle\rangle\ =\ True$
$invh\ \langle l,\ (\_,\ \_),\ r\rangle\ =\ (bh(l)\ =\ bh(r)\ \wedge\ invh\ l\ \wedge\ invh\ r)$

$bheight :: {}'a\ rbt \Rightarrow nat$
$bh(\langle\rangle)\ =\ 0$
$bh(\langle l,\ (\_,\ c),\ \_\rangle)\ =$
$(\text{if}\ c\ =\ Black\ \text{then}\ bh(l)\ +\ 1\ \text{else}\ bh(l))$

# Logarithmic height

## Lemma

$rbt\ t \implies h(t) \leq 2 * \log_2 |t|_1$

Intuition: $h(t)\ /\ 2 \leq bh(t) \leq mh(t) \leq \log_2 |t|_1$

# Insertion

$insert :: \, 'a \Rightarrow \, 'a \; rbt \Rightarrow \, 'a \; rbt$

$insert \; x \; t = paint \; Black \; (ins \; x \; t)$

$ins :: \, 'a \Rightarrow \, 'a \; rbt \Rightarrow \, 'a \; rbt$

$ins \; x \; \langle \rangle = R \; \langle \rangle \; x \; \langle \rangle$

$ins \; x \; (B \; l \; a \; r) = (\text{case} \; cmp \; x \; a \; \text{of}$

$\qquad\qquad\qquad LT \Rightarrow baliL \; (ins \; x \; l) \; a \; r$

$\qquad\qquad\qquad | \; EQ \Rightarrow B \; l \; a \; r$

$\qquad\qquad\qquad | \; GT \Rightarrow baliR \; l \; a \; (ins \; x \; r))$

$ins \; x \; (R \; l \; a \; r) = (\text{case} \; cmp \; x \; a \; \text{of}$

$\qquad\qquad\qquad LT \Rightarrow R \; (ins \; x \; l) \; a \; r$

$\qquad\qquad\qquad | \; EQ \Rightarrow R \; l \; a \; r$

$\qquad\qquad\qquad | \; GT \Rightarrow R \; l \; a \; (ins \; x \; r))$

# Adjusting colors

$baliL$, $baliR$ :: $'a\ rbt \Rightarrow\ 'a \Rightarrow\ 'a\ rbt \Rightarrow\ 'a\ rbt$

- Combine arguments $l\ a\ r$ into tree, ideally $\langle l,\ a,\ r \rangle$
- Treat invariant violation <span style="color:red">Red-Red</span> in $l/r$

  $baliL\ (R\ (R\ t_1\ a_1\ t_2)\ a_2\ t_3)\ a_3\ t_4$
  $\quad = R\ (B\ t_1\ a_1\ t_2)\ a_2\ (B\ t_3\ a_3\ t_4)$
  $baliL\ (R\ t_1\ a_1\ (R\ t_2\ a_2\ t_3))\ a_3\ t_4$
  $\quad = R\ (B\ t_1\ a_1\ t_2)\ a_2\ (B\ t_3\ a_3\ t_4)$

- Principle: replace <span style="color:red">Red-Red</span> by <span style="color:blue">Red-Black</span>
- Final equation:

  $baliL\ l\ a\ r = B\ l\ a\ r$

- Symmetric: $baliR$

# Preservation of invariant

After 14 simple lemmas:

## Theorem
$rbt\ t \implies rbt\ (insert\ x\ t)$

# Proof in CLRS

The **while** loop in lines 1–15 maintains the following three-part invariant at the start of each iteration of the loop:

a. Node $z$ is red.

b. If $z.p$ is the root, then $z.p$ is black.

c. If the tree violates any of the red-black properties, then it violates at most one of them, and the violation is of either property 2 or property 4. If the tree violates property 2, it is because $z$ is the root and is red. If the tree violates property 4, it is because both $z$ and $z.p$ are red.

Part (c), which deals with violations of red-black properties, is more central to showing that RB-INSERT-FIXUP restores the red-black properties than parts (a) and (b), which we use along the way to understand situations in the code. Because we'll be focusing on node $z$ and nodes near it in the tree, it helps to know from part (a) that $z$ is red. We shall use part (b) to show that the node $z.p.p$ exists when we reference it in lines 2, 3, 7, 8, 13, and 14.

Recall that we need to show that a loop invariant is true prior to the first iteration of the loop, that each iteration maintains the loop invariant, and that the loop invariant gives us a useful property at loop termination.

We start with the initialization and termination arguments. Then, as we examine how the body of the loop works in more detail, we shall argue that the loop maintains the invariant upon each iteration. Along the way, we shall also demonstrate that each iteration of the loop has two possible outcomes: either the pointer $z$ moves up the tree, or we perform some rotations and then the loop terminates.

**Initialization:** Prior to the first iteration of the loop, we started with a red-black tree with no violations, and we added a red node $z$. We show that each part of the invariant holds at the time RB-INSERT-FIXUP is called:

a. When RB-INSERT-FIXUP is called, $z$ is the red node that was added.

b. If $z.p$ is the root, then $z.p$ started out black and did not change prior to the call of RB-INSERT-FIXUP.

c. We have already seen that properties 1, 3, and 5 hold when RB-INSERT-FIXUP is called.

If the tree violates property 2, then the red root must be the newly added node $z$, which is the only internal node in the tree. Because the parent and both children of $z$ are the sentinel, which is black, the tree does not also violate property 4. Thus, this violation of property 2 is the only violation of red-black properties in the entire tree.

If the tree violates property 4, then, because the children of node $z$ are black sentinels and the tree had no other violations prior to $z$ being added, the

**Figure 13.5** Case 1 of the procedure RB-INSERT-FIXUP. Property 4 is violated, since $z$ and its parent $z.p$ are both red. We take the same action (whether (a) $z$ is a right child or (b) $z$ is a left child. Each of the subtrees $\alpha$, $\beta$, $\gamma$, $\delta$, and $\varepsilon$ has a black root, and each has the same black-height. The code for case 1 changes the colors of some nodes, preserving property 5: all downward simple paths from a node to a leaf have the same number of blacks. The **while** loop continues with node $z$'s grandparent $z.p.p$ as the new node $z$. Any violation of property 4 can now occur only between the new $z$, which is red, and its parent, if it is red as well.

If node $z$'s uncle is black and $z$ is a right child ...

**Case 2:** $z$'s uncle $y$ is black and $z$ is a right child
**Case 3:** $z$'s uncle $y$ is black and $z$ is a left child

In cases 2 and 3, the color of $z$'s uncle $y$ is black. We distinguish the two cases according to whether $z$ is a right or left child of $z.p$. Lines 10–11 constitute case 2, which is shown in Figure 13.6 together with case 3. In case 2, node $z$ is a right child of its parent. We immediately use a left rotation to transform the situation into case 3 (lines 12–14), in which node $z$ is a left child. Because

violation must be because both $z$ and $z.p$ are red. Moreover, the tree violates no other red-black properties.

**Termination:** When the loop terminates, it does so because $z.p$ is black. (If $z$ is the root, then $z.p$ is the sentinel $T.nil$, which is black.) Thus, the tree does not violate property 4 at loop termination. By the loop invariant, the only property that might fail to hold is property 2. Line 16 restores this property, too, so that when RB-INSERT-FIXUP terminates, all the red-black properties hold.

**Maintenance:** We actually need to consider six cases in the **while** loop, but three of them are symmetric to the other three, depending on whether line 2 determines $z$'s parent $z.p$ to be a left child or a right child of $z$'s grandparent $z.p.p$. We have given the code only for the situation in which $z.p$ is a left child. The node $z.p.p$ exists, since by part (b) of the loop invariant, if $z.p$ is the root, then $z.p$ is black. But since we enter a loop iteration only if $z.p$ is red, we know that $z.p.p$ cannot be the root. Hence, $z.p.p$ exists.

We distinguish case 1 from cases 2 and 3 by the color of $z$'s parent's sibling, or "uncle." Line 3 makes $y$ point to $z$'s uncle $z.p.p.right$, and line 4 tests $y$'s color. If $y$ is red, then we execute case 1. Otherwise, control passes to cases 2 and 3. In all three cases, $z$'s grandparent $z.p.p$ is black, since its parent $z.p$ is red, and property 4 is violated only between $z$ and $z.p$.

**Case 1:** $z$'s uncle $y$ is red

Figure 13.5 shows the situation for case 1 (lines 5–8), which occurs when both $z.p$ and $y$ are red. Because $z.p.p$ is black, we can color both $z.p$ and $y$ black, thereby fixing the problem of $z$ and $z.p$ both being red, and we can color $z.p.p$ red, thereby maintaining property 5. We then repeat the **while** loop with $z.p.p$ as the new node $z$. The pointer $z$ moves up two levels in the tree.

Now, we show that case 1 maintains the loop invariant at the start of the next iteration. We use $z$ to denote node $z$ in the current iteration, and $z' = z.p.p$ to denote the node that will be called node $z$ at the test in line 1 upon the next iteration.

a. Because this iteration colors $z.p.p$ red, node $z'$ is red at the start of the next iteration.

b. The node $z'.p$ is $z.p.p.p$ in this iteration, and the color of this node does not change. If this node is the root, it was black prior to this iteration, and it remains black at the start of the next iteration.

c. We have already argued that case 1 maintains property 5, and it does not introduce a violation of properties 1 or 3.

**Figure 13.6** Cases 2 and 3 of the procedure RB-INSERT-FIXUP. As in case 1, property 4 is violated in either case 2 or case 3 because $z$ and its parent $z.p$ are both red. Each of the subtrees $\alpha$, $\beta$, $\gamma$, and $\delta$ has a black root (in $z$, and $z$ from property 4, and $\delta$ because otherwise we would be in case 1), and each has the same black-height. We transform case 2 into case 3 by a left rotation, which preserves property 5: all downward simple paths from a node to a leaf have the same number of blacks. Case 3 causes some color changes and a right rotation, which also preserves property 5. The **while** loop then terminates, because property 4 is satisfied: there are no longer two red nodes in a row.

both $z$ and $z.p$ are red, the rotation affects neither the black-height of nodes nor property 5. Whether we enter case 3 directly or through case 2, $z$'s uncle $y$ is black, since otherwise we would have executed case 1. Additionally, the node $z.p.p$ exists, since we have argued that this node existed at the time that lines 2 and 3 were executed, and after moving $z$ up one level in line 10 and then down one level in line 11, the identity of $z.p.p$ remains unchanged. In case 3, we execute some color changes and a right rotation, which preserve property 5, and then, since we no longer have two red nodes in a row, we are done. The **while** loop does not iterate another time, since $z.p$ is now black.

We now show that cases 2 and 3 maintain the loop invariant. (As we have just argued, $z.p$ will be black upon the next test in line 1, and the loop body will not execute again.)

a. Case 2 makes $z$ point to $z.p$, which is red. No further changes to $z$ or its color occurs in cases 2 and 3.

b. Case 3 makes $z.p$ black, so that if $z.p$ is the root at the start of the next iteration, it is black.

c. As in case 1, properties 1, 3, and 5 are maintained in cases 2 and 3.

Since node $z$ is not the root in cases 2 and 3, we know that there is no violation of property 2. Cases 2 and 3 do not introduce a violation of property 2, since the only node that is made red becomes a child of a black node by the rotation in case 3.

Cases 2 and 3 correct the loop violation of property 4, and they do not introduce another violation.

# Deletion code

$delete\ x\ t\ =\ paint\ Black\ (del\ x\ t)$

$del\ \_\ \langle\rangle\ =\ \langle\rangle$
$del\ x\ \langle l,\ (a,\ \_),\ r\rangle\ =$
(case $cmp\ x\ a$ of
   $LT\ \Rightarrow$
     if $l \neq \langle\rangle\ \wedge\ color\ l\ =\ Black$
     then $baldL\ (del\ x\ l)\ a\ r$ else $R\ (del\ x\ l)\ a\ r$
| $EQ\ \Rightarrow$
     if $r\ =\ \langle\rangle$ then $l$
     else let $(a',\ r')\ =\ split\_min\ r$
        in if $color\ r\ =\ Black$ then $baldR\ l\ a'\ r'$
          else $R\ l\ a'\ r'$
| $GT\ \Rightarrow$

# Deletion code

$split\_min\ \langle l,\ (a,\ \_),\ r \rangle =$
(if $l = \langle \rangle$ then $(a,\ r)$
 else let $(x,\ l') = split\_min\ l$
        in $(x,$ if $color\ l = Black$ then $baldL\ l'\ a\ r$
                else $R\ l'\ a\ r))$

$baldL\ (R\ t_1\ a\ t_2)\ b\ t_3 = R\ (B\ t_1\ a\ t_2)\ b\ t_3$
$baldL\ t_1\ a\ (B\ t_2\ b\ t_3) = baliR\ t_1\ a\ (R\ t_2\ b\ t_3)$
$baldL\ t_1\ a\ (R\ (B\ t_2\ b\ t_3)\ c\ t_4) =$
$R\ (B\ t_1\ a\ t_2)\ b\ (baliR\ t_3\ c\ (paint\ Red\ t_4))$
$baldL\ t_1\ a\ t_2 = R\ t_1\ a\ t_2$

# Deletion proof

After a number of lemmas:

$\llbracket invh\ t;\ invc\ t \rrbracket$
$\implies invh\ (del\ x\ t)\ \wedge$
  $(color\ t\ =\ Red\ \longrightarrow$
   $bh(del\ x\ t)\ =\ bh(t)\ \wedge\ invc\ (del\ x\ t))\ \wedge$
  $(color\ t\ =\ Black\ \longrightarrow$
   $bh(del\ x\ t)\ =\ bh(t)\ -\ 1\ \wedge\ invc2\ (del\ x\ t))$

$rbt\ t\ \implies rbt\ (delete\ x\ t)$

# Code and proof in CLRS

# Source of code

Insertion:

Okasaki's *Purely Functional Data Structures*

Deletion partly based on:

Stefan Kahrs. *Red Black Trees with Types*.
J. Functional Programming. 1996.

**⑫ More Search Trees**
   AVL Trees
   Weight-Balanced Trees
   AA Trees
   Scapegoat Trees

# AVL Trees

[Adelson-Velskii & Landis 62]

- Every node $\langle l, \_, r \rangle$ must be balanced:
  $|h(l) - h(r)| \leq 1$
- Verified Isabelle implementation:
  `HOL/Data_Structures/AVL_Set.thy`

# Weight-Balanced Trees

[Nievergelt & Reingold 72,73]

- Parameter: balance factor $0 < \alpha \leq 0.5$
- Every node $\langle l, {-}, r \rangle$ must be balanced:
  $\alpha \leq |l|_1 / (|l|_1 + |r|_1) \leq 1 - \alpha$
- Insertion and deletion: single and double rotations depending on subtle numeric conditions
- Nievergelt and Reingold incorrect
- Mistakes discovered and corrected by [Blum & Mehlhorn 80] and [Hirai & Yamamoto 2011]
- Verified implementation
  in Isabelle's *Archive of Formal Proofs*.

# AA trees

[Arne Andersson 93, Ragde 14]

- Simulation of 2-3 trees by binary trees
  $\langle t_1, a, t_2, b, t_3 \rangle \rightsquigarrow \langle t_1, a, \langle t_2, b, t_3 \rangle \rangle$
- Height field (or single bit) to distinguish single from double node
- Code short but opaque
- 4 bugs in $delete$ in [Ragde 14]: non-linear pattern; going down wrong subtree; missing function call; off by 1

# AA trees

After corrections, the proofs:

- Code relies on tricky pre- and post-conditions that need to be found
- Structural invariant preservation requires most of the work

# Scapegoat trees

[Anderson 89, Igal & Rivest 93]

Central idea:

Don't rebalance every time,
Rebuild when the tree gets "too unbalanced"

- Tricky: amortized logarithmic complexity analysis
- Verified implementation
  in Isabelle's *Archive of Formal Proofs*.

# One by one (Union)

Let $c(x) =$ cost of adding 1 element to set of size $x$

Cost of adding $m$ elements to a set of $n$ elements:

$$c(n) + \cdots + c(n + m - 1)$$

$\implies$ choose $m \le n \implies$ smaller into bigger

If $c(x) = \log_2 x \implies$
Cost $= O(m * \log_2(n + m)) = O(m * \log_2 n)$

Similar for intersection and difference.

- We can do better than $O(m * \log_2 n)$
- This section:

    *A parallel divide and conquer approach*
- Cost: $\Theta(m * \log_2(\frac{n}{m} + 1))$
- Works for many kinds of balanced trees
- For ease of presentation: use concrete type $tree$

# Uniform *tree* type

Red-Black trees, AVL trees, weight-balanced trees, etc can all be implemented with $'b$-augmented trees:

$('a \times 'b)$ *tree*

We work with this type of trees without committing to any particular kind of balancing schema.

# Just $join$

Can synthesize all BST interface functions from just one function:

$$join\ l\ a\ r\ \approx\ Node\ l\ (a,\ \_)\ r + \text{rebalance}$$

Thus $join$ determines the balancing schema

# Just $join$

Given $join :: tree \Rightarrow {}'a \Rightarrow tree \Rightarrow tree$
(where $tree$ abbreviates $({}'a, {}'b)\ tree$), implement

$union :: tree \Rightarrow tree \Rightarrow tree$
$inter :: tree \Rightarrow tree \Rightarrow tree$
$diff :: tree \Rightarrow tree \Rightarrow tree$

$union\ t_1\ t_2 =$
$(\text{if}\ t_1 = \langle\rangle\ \text{then}\ t_2$
$\ \text{else if}\ t_2 = \langle\rangle\ \text{then}\ t_1$
$\qquad \text{else case}\ t_1\ \text{of}$
$\qquad\qquad \langle l_1,\ (a,\ b),\ r_1\rangle \Rightarrow$
$\qquad\qquad\quad \text{let}\ (l_2,\ x,\ r_2) = split\ t_2\ a;$
$\qquad\qquad\qquad l' = union\ l_1\ l_2;$
$\qquad\qquad\qquad r' = union\ r_1\ r_2$
$\qquad\qquad\quad \text{in}\ join\ l'\ a\ r')$

$split :: tree \Rightarrow {}'a \Rightarrow tree \times bool \times tree$

$split \langle\rangle \, \_ = (\langle\rangle, \, False, \, \langle\rangle)$

$split \langle l, \, (a, \, \_), \, r\rangle \, x =$

(case $cmp \, x \, a$ of

   $LT \Rightarrow$

     let $(l_1, \, b, \, l_2) = split \, l \, x$

     in $(l_1, \, b, \, join \, l_2 \, a \, r)$

 | $EQ \Rightarrow (l, \, True, \, r)$

 | $GT \Rightarrow$

     let $(r_1, \, b, \, r_2) = split \, r \, x$

     in $(join \, l \, a \, r_1, \, b, \, r_2))$

$inter\ t_1\ t_2 =$
(if $t_1 = \langle\rangle$ then $\langle\rangle$
 else if $t_2 = \langle\rangle$ then $\langle\rangle$
    else case $t_1$ of
          $\langle l_1,\ (a,\ x),\ r_1 \rangle \Rightarrow$
            let $(l_2,\ b,\ r_2) = split\ t_2\ a;$
               $l' = inter\ l_1\ l_2;$
               $r' = inter\ r_1\ r_2$
           in if $b$ then $join\ l'\ a\ r'$
              else $join2\ l'\ r')$

$join2 :: tree \Rightarrow tree \Rightarrow tree$

$join2\ l\ \langle\rangle = l$

$join2\ l\ \langle v,\ va,\ vb\rangle =$

$(\textsf{let}\ (x,\ y) = split\_min\ \langle v,\ va,\ vb\rangle\ \textsf{in}\ join\ l\ x\ y)$

$split\_min :: tree \Rightarrow\ 'a \times\ tree$

$split\_min\ \langle l,\ (a,\ \_),\ r\rangle =$

$(\textsf{if}\ l = \langle\rangle\ \textsf{then}\ (a,\ r)$

$\ \textsf{else let}\ (m,\ l') = split\_min\ l\ \textsf{in}\ (m,\ join\ l'\ a\ r))$

$$diff\ t_1\ t_2 =$$

$$(\text{if } t_1 = \langle\rangle \text{ then } \langle\rangle$$

$$\text{else if } t_2 = \langle\rangle \text{ then } t_1$$

$$\text{else case } t_2 \text{ of}$$

$$\langle l_2,\ (a,\ b),\ r_2 \rangle \Rightarrow$$

$$\text{let } (l_1,\ x,\ r_1) = split\ t_1\ a;$$

$$l' = diff\ l_1\ l_2;$$

$$r' = diff\ r_1\ r_2$$

$$\text{in } join2\ l'\ r')$$

# *insert* and *delete*

*insert x t* = (let *(l, b, r)* = *split t x* in *join l x r*)

*delete x t* = (let *(l, b, r)* = *split t x* in *join2 l r*)

# Specification of $join$ and $inv$

- $set\_tree\ (join\ l\ a\ r) = set\_tree\ l \cup \{a\} \cup set\_tree\ r$
- $bst\ \langle l,\ (a,\ b),\ r \rangle \implies bst\ (join\ l\ a\ r)$

Also required: structural invariant $inv$:

- $inv\ \langle \rangle$
- $inv\ \langle l,\ (a,\ b),\ r \rangle \implies inv\ l \wedge inv\ r$
- $[\![ inv\ l;\ inv\ r ]\!] \implies inv\ (join\ l\ a\ r)$

Locale context for def of $union$ etc

# Specification of *union, inter, diff*

ADT/Locale $Set2$ = extension of locale $Set$ with

- *union, inter, diff* $:: \, 's \Rightarrow \, 's \Rightarrow \, 's$
- $\llbracket invar \; s_1; \; invar \; s_2 \rrbracket$
  $\implies set \; (union \; s_1 \; s_2) = set \; s_1 \cup set \; s_2$
- $\llbracket invar \; s_1; \; invar \; s_2 \rrbracket \implies invar \; (union \; s_1 \; s_2)$
- ...*inter* ...
- ...*diff* ...

We focus on *union*.

See HOL/Data_Structures/Set_Specs.thy

# Correctness lemmas
# for $union$ etc code

In the context of $join$ specification:

- $bst\ t_2 \implies$
  $set\_tree\ (union\ t_1\ t_2) = set\_tree\ t_1 \cup set\_tree\ t_2$
- $[\![bst\ t_1;\ bst\ t_2]\!] \implies bst\ (union\ t_1\ t_2)$
- $[\![inv\ t_1;\ inv\ t_2]\!] \implies inv\ (union\ t_1\ t_2)$

Proofs automatic (more complex for $inter$ and $diff$)

Implementation of locale $Set2$:

**interpretation** $Set2$ **where** $union = union$ …
**and** $set = set\_tree$ **and** $invar = (\lambda t.\ bst\ t \wedge inv\ t)$

HOL/Data_Structures/
Set2_Join.thy

**13** Union, Intersection, Difference on BSTs
   Correctness
   Join for Red-Black Trees

# $join\ l\ a\ r$ — The idea

Assume $l$ is "smaller" than $r$ :

- Descend along the left spine of $r$
  until you find a subtree $t$ of the same "size" as $l$.
- Replace $t$ by $\langle l,a,t \rangle$.
- Rebalance on the way up.

$join\ l\ x\ r =$
(if $bheight\ r < bheight\ l$
 then $paint\ Black\ (joinR\ l\ x\ r)$
 else if $bheight\ l < bheight\ r$
     then $paint\ Black\ (joinL\ l\ x\ r)$ else $B\ l\ x\ r)$

$joinL\ l\ x\ r =$
(if $bheight\ r \leq bheight\ l$ then $R\ l\ x\ r$
 else case $r$ of
     $\langle l',\ (x',\ Red),\ r' \rangle \Rightarrow R\ (joinL\ l\ x\ l')\ x'\ r'$
   $|\ \langle l',\ (x',\ Black),\ r' \rangle \Rightarrow baliL\ (joinL\ l\ x\ l')\ x'\ r')$

Need to store black height in each node
for logarithmic complexity

Thys/Set2_Join_RBT.thy

# Literature

The idea of "just $join$":

Stephen Adams. *Efficient Sets — A Balancing Act.*
J. Functional Programming, volume 3, number 4, 1993.

The precise analysis:

Guy E. Blelloch, D. Ferizovic, Y. Sun.
*Just Join for Parallel Ordered Sets.*
ACM Symposium on Parallelism in Algorithms and
Architectures 2016.

# Trie
### [Fredkin, CACM 1960]

Name: *reTRIEval*

- Tries are search trees indexed by lists
- Tries are tree-shaped DFAs

# Example Trie

$\{$ a, an, can, car, cat $\}$

**14 Tries and Patricia Tries**

Tries via Functions

Binary Tries and Patricia Tries

HOL/Data_Structures/
Trie_Fun.thy

# Trie

**datatype** $'a\ trie = Nd\ bool\ ('a \Rightarrow 'a\ trie\ option)$

Function update notation:

$f(a := b) = (\lambda x.\ \text{if } x = a \text{ then } b \text{ else } f\ x)$

$f(a \mapsto b) = f(a := Some\ b)$

Next: Implementation of ADT $Set$

$empty = Nd\ False\ (\lambda_-.\ None)$

# *isin*

$isin\ (Nd\ b\ m)\ [] = b$

$isin\ (Nd\ b\ m)\ (k\ \#\ xs) = (\textsf{case}\ m\ k\ \textsf{of}$
$\qquad\qquad\qquad\qquad None \Rightarrow False$
$\qquad\qquad\qquad\ |\ Some\ t \Rightarrow isin\ t\ xs)$

$insert\ []\ (Nd\ b\ m)\ =\ Nd\ True\ m$

$insert\ (x\ \#\ xs)\ (Nd\ b\ m)\ =$
let $s =$ case $m\ x$ of
$\qquad\quad None \Rightarrow empty$
$\qquad\ |\ Some\ t \Rightarrow t$
in $Nd\ b\ (m(x \mapsto insert\ xs\ s))$

# *delete*

$delete\ []\ (Nd\ b\ m) = Nd\ False\ m$

$delete\ (x\ \#\ xs)\ (Nd\ b\ m) =$
$Nd\ b\ (\text{case}\ m\ x\ \text{of}$
$\qquad None \Rightarrow m$
$\quad |\ Some\ t \Rightarrow m(x \mapsto delete\ xs\ t))$

Does not shrink trie — exercise!

# Correctness:
# Abstraction function

$set :: {}'a\ trie \Rightarrow {}'a\ list\ set$

$set\ t = \{xs.\ isin\ t\ xs\}$

Invariant is $True$

# Correctness theorems

- $set\ empty = \{\}$
- $isin\ t\ xs = (xs \in set\ t)$
- $set\ (insert\ xs\ t) = set\ t \cup \{xs\}$
- $set\ (delete\ xs\ t) = set\ t - \{xs\}$

No lemmas required

# Abstraction function via $isin$

$$set\ t = \{xs.\ isin\ t\ xs\}$$

- Trivial definition
- Reusing code ($isin$) *may* complicate proofs.
- Separate abstract mathematical definition *may* simplify proofs

Also possible for some other ADTs, e.g. for Map:
$lookup :: {'}t \Rightarrow ({'}a \Rightarrow {'}b\ option)$

HOL/Data_Structures/
Tries_Binary.thy

# Trie

**datatype** $trie = Lf \mid Nd\ bool\ (trie \times trie)$

Auxiliary functions on pairs:

$sel2 :: bool \Rightarrow {}'a \times {}'a \Rightarrow {}'a$
$sel2\ b\ (a_1,\ a_2) = (\text{if } b \text{ then } a_2 \text{ else } a_1)$

$mod2 :: ({}'a \Rightarrow {}'a) \Rightarrow bool \Rightarrow {}'a \times {}'a \Rightarrow {}'a \times {}'a$
$mod2\ f\ b\ (a_1,\ a_2) = (\text{if } b \text{ then } (a_1,\ f\ a_2) \text{ else } (f\ a_1,\ a_2))$

$empty = Lf$

# *isin*

*isin Lf ks = False*

*isin (Nd b lr) ks =* (case *ks* of
$\qquad\qquad\quad$ [] $\Rightarrow$ *b*
$\qquad\qquad$ | *k # x* $\Rightarrow$ *isin (sel2 k lr) x*)

# insert

*insert* [] *Lf = Nd True (Lf, Lf)*

*insert* [] *(Nd b lr) = Nd True lr*

*insert (k # ks) Lf =*
*Nd False (mod2 (insert ks) k (Lf, Lf))*

*insert (k # ks) (Nd b lr) =*
*Nd b (mod2 (insert ks) k lr)*

# *delete*

$delete\ ks\ Lf = Lf$

$delete\ ks\ (Nd\ b\ lr) =$
case $ks$ of
$\quad [] \Rightarrow node\ False\ lr$
$\mid k\ \#\ ks' \Rightarrow node\ b\ (mod2\ (delete\ ks')\ k\ lr)$

Shrink trie if possible:

$node\ b\ lr = ($if $\neg\ b \wedge lr = (Lf,\ Lf)$ then $Lf$ else $Nd\ b\ lr)$

# Correctness of implementation

Abstraction function:

$$set\_trie\ t = \{xs.\ isin\ t\ xs\}$$

- $isin\ (insert\ xs\ t)\ ys = (xs = ys \lor isin\ t\ ys)$
  $\implies set\_trie\ (insert\ xs\ t) = set\_trie\ t \cup \{xs\}$
- $isin\ (delete\ xs\ t)\ ys = (xs \neq ys \land isin\ t\ ys)$
  $\implies set\_trie\ (delete\ xs\ t) = set\_trie\ t - \{xs\}$

# From tries to Patricia tries

# Patricia trie

**datatype** *trieP = LfP*
  | *NdP (bool list) bool (trieP × trieP)*

# $isinP$

$isinP\ LfP\ ks = False$

$isinP\ (NdP\ ps\ b\ lr)\ ks =$
(let $n = length\ ps$
 in if $ps = take\ n\ ks$
   then case $drop\ n\ ks$ of
         $[] \Rightarrow b$
         $\mid k\ \#\ ks' \Rightarrow isinP\ (sel2\ k\ lr)\ ks'$
   else $False$)

# Splitting lists

$split\ xs\ ys = (zs,\ xs',\ ys')$
iff $zs$ is the longest common prefix of $xs$ and $ys$
and $xs'/ys'$ is the remainder of $xs/ys$

# insertP

$insertP\ ks\ LfP = NdP\ ks\ True\ (LfP,\ LfP)$

$insertP\ ks\ (NdP\ ps\ b\ lr) =$
case $split\ ks\ ps$ of
  $(qs,\ [],\ []) \Rightarrow NdP\ ps\ True\ lr$
$|\ (qs,\ [],\ p\ \#\ ps') \Rightarrow$
    let $t = NdP\ ps'\ b\ lr$
    in $NdP\ qs\ True$ (if $p$ then $(LfP,\ t)$ else $(t,\ LfP)$)
$|\ (qs,\ k\ \#\ ks',\ []) \Rightarrow NdP\ ps\ b\ (mod2\ (insertP\ ks')\ k\ lr)$
$|\ (qs,\ k\ \#\ ks',\ p\ \#\ ps') \Rightarrow$
    let $tp = NdP\ ps'\ b\ lr;\ tk = NdP\ ks'\ True\ (LfP,\ LfP)$
    in $NdP\ qs\ False$ (if $k$ then $(tp,\ tk)$ else $(tk,\ tp)$)

# deleteP

$deleteP\ ks\ LfP = LfP$

$deleteP\ ks\ (NdP\ ps\ b\ lr) =$
$(case\ split\ ks\ ps\ of$
$\quad (qs,\ ks',\ p\#ps') \Rightarrow NdP\ ps\ b\ lr\ |$
$\quad (qs,\ k\#ks',\ []) \Rightarrow$
$\qquad nodeP\ ps\ b\ (mod2\ (deleteP\ ks')\ k\ lr)\ |$
$\quad (qs,\ [],\ []) \Rightarrow nodeP\ ps\ False\ lr)$

# Stepwise data refinement

View $trieP$ as an implementation ("refinement") of $trie$

|  | Type | Abstraction function |
|---|---|---|

$bool\ list\ set$

$\uparrow$       $set\_trie$

$trie$

$\uparrow$       $abs\_trieP$

$trieP$

$\Longrightarrow$ Modular correctness proof of $trieP$

# $abs\_trieP :: trieP \Rightarrow trie$

$abs\_trieP\ LfP = Lf$

$abs\_trieP\ (NdP\ ps\ b\ (l,\ r)) =$
$prefix\_trie\ ps\ (Nd\ b\ (abs\_trieP\ l,\ abs\_trieP\ r))$

$prefix\_trie :: bool\ list \Rightarrow trie \Rightarrow trie$

# Correctness of *trieP* w.r.t. *trie*

- *isinP t ks = isin (abs_trieP t) ks*
- *abs_trieP (insertP ks t) = insert ks (abs_trieP t)*
- *abs_trieP (deleteP ks t) = delete ks (abs_trieP t)*

*isin (prefix_trie ps t) ks =*
*(ps = take (length ps) ks ∧ isin t (drop (length ps) ks))*
*prefix_trie ks (Nd True (Lf, Lf)) = insert ks Lf*
*insert ps (prefix_trie ps (Nd b lr)) = prefix_trie ps (Nd True lr)*
*insert (ks @ ks') (prefix_trie ks t) = prefix_trie ks (insert ks' t)*
*prefix_trie (ps @ qs) t = prefix_trie ps (prefix_trie qs t)*
*split ks ps = (qs, ks', ps') ⟹*
*ks = qs @ ks' ∧ ps = qs @ ps' ∧ (ks' ≠ [] ∧ ps' ≠ [] ⟶ hd ks' ≠ hd ps')*
*(prefix_trie xs t = Lf) = (xs = [] ∧ t = Lf)*
*(abs_trieP t = Lf) = (t = LfP)*
*delete xs (prefix_trie xs (Nd b (l, r))) =*
(if *(l, r) = (Lf, Lf)* then *Lf* else *prefix_trie xs (Nd False (l, r))*)
*delete (xs @ ys) (prefix_trie xs t) =*
(if *delete ys t = Lf* then *Lf* else *prefix_trie xs (delete ys t)*)

# Correctness of $trieP$ w.r.t. $bool\ list\ set$

Define $set\_trieP = set\_trie \circ abs\_trieP$

$\implies$ Overall correctness by trivial composition of correctness theorems for $trie$ and $trieP$

Example:

$set\_trieP\ (insertP\ xs\ t) = set\_trieP\ t \cup \{xs\}$

follows directly from

$abs\_trieP\ (insertP\ ks\ t) = insert\ ks\ (abs\_trieP\ t)$
$set\_trie\ (insert\ xs\ t) = set\_trie\ t \cup \{xs\}$

# Chapter 9

# Priority Queues

# Priority queue informally

Collection of elements with priorities

Operations:

- empty
- emptiness test
- insert
- get element with minimal priority
- delete element with minimal priority

We focus on the priorities:
element = priority

# Priority queues are multisets

The same element can be contained multiple times
in a priority queue

$$\implies$$

The abstract view of a priority queue is a multiset

# Interface of implementation

The type of elements (= priorities) $'a$ is a linear order

An implementation of a priority queue of elements of type $'a$ must provide

- An implementation type $'q$
- $empty :: 'q$
- $is\_empty :: 'q \Rightarrow bool$
- $insert :: 'a \Rightarrow 'q \Rightarrow 'q$
- $get\_min :: 'q \Rightarrow 'a$
- $del\_min :: 'q \Rightarrow 'q$

# More operations

- $merge :: {}'q \Rightarrow {}'q \Rightarrow {}'q$
  Often provided
- decrease key/priority
  A bit tricky in functional setting

# Correctness of implementation

A priority queue represents a multiset of priorities.
Correctness proof requires:

Abstraction function:     $mset :: {}'q \Rightarrow {}'a\ multiset$
Invariant:                $invar :: {}'q \Rightarrow bool$

# Correctness of implementation

Must prove $\quad invar\ q \Longrightarrow$

$mset\ empty = \{\#\}$

$is\_empty\ q = (mset\ q = \{\#\})$

$mset\ (insert\ x\ q) = mset\ q + \{\#x\#\}$

$mset\ q \neq \{\#\} \Longrightarrow get\_min\ q = Min\_mset\ (mset\ q)$

$mset\ q \neq \{\#\} \Longrightarrow$
$mset\ (del\_min\ q) = mset\ q - \{\#get\_min\ q\#\}$

$invar\ empty$
$invar\ (insert\ x\ q)$
$invar\ (del\_min\ q)$

# Terminology

A binary tree is a *heap* if for every subtree the root is $\leq$ all elements in that subtree.

$heap \; \langle \rangle = True$
$heap \; \langle l, \; m, \; r \rangle =$
$((\forall \, x \in set\_tree \; l \cup set\_tree \; r. \; m \leq x) \; \wedge$
$\; heap \; l \wedge heap \; r)$

The term "heap" is frequently used synonymously with "priority queue".

# Priority queue via heap

- $empty = \langle\rangle$
- $is\_empty\ h = (h = \langle\rangle)$
- $get\_min\ \langle\_,\ a,\ \_\rangle = a$
- Assume we have $merge$
- $insert\ a\ t = merge\ \langle\langle\rangle,\ a,\ \langle\rangle\rangle\ t$
- $del\_min\ \langle l,\ a,\ r\rangle = merge\ l\ r$

# Priority queue via heap

A naive merge:

$merge\ t_1\ t_2 = ($ case $(t_1, t_2)$ of
  $(\langle\rangle, \_) \Rightarrow t_2\ |$
  $(\_, \langle\rangle) \Rightarrow t_1\ |$
  $(\langle l_1, a_1, r_1\rangle, \langle l_2, a_2, r_2\rangle) \Rightarrow$
    if $a_1 \leq a_2$ then $\langle merge\ l_1\ r_1,\ a_1,\ t_2\rangle$
    else $\langle t_1,\ a_2,\ merge\ l_2\ r_2\rangle$

Challenge: how to maintain some kind of balance

HOL/Data_Structures/
Leftist_Heap.thy

# Leftist tree informally

In a *leftist tree*, the minimum height of every left child is $\geq$ the minimum height of its right sibling.

$\implies$ m.h. = length of right spine



Merge descends along the right spine.
Thus m.h. bounds number of steps.

If m.h. of right child gets too large: swap with left child.

# Implementation type

**type_synonym** $'a\ lheap = ('a \times nat)\ tree$

Abstraction function:

$mset\_tree :: 'a\ lheap \Rightarrow 'a\ multiset$

$mset\_tree\ \langle\rangle = \{\#\}$

$mset\_tree\ \langle l,\ (a,\ \_),\ r\rangle =$
$\{\#a\#\} + mset\_tree\ l + mset\_tree\ r$

# Leftist tree

$ltree :: {}'a\ lheap \Rightarrow bool$

$ltree\ \langle\rangle\ =\ True$

$ltree\ \langle l,\ (\_,\ n),\ r\rangle\ =$

$(mh(r) \leq mh(l) \wedge n = mh(r) + 1 \wedge ltree\ l \wedge ltree\ r)$

$mht :: {}'a\ lheap \Rightarrow nat$

$mht\ \langle\rangle\ =\ 0$

$mht\ \langle\_,\ (\_,\ n),\ \_\rangle\ =\ n$

# Leftist heap invariant

$$invar\ h = (heap\ h \land ltree\ h)$$

# *merge*

Principle: descend on the right

$merge \langle\rangle \ t = t$

$merge \ t \ \langle\rangle = t$

$merge \ (\langle l_1, (a_1, \_), r_1\rangle =: t_1) \ (\langle l_2, (a_2, \_), r_2\rangle =: t_2) =$
(if $a_1 \leq a_2$ then $node \ l_1 \ a_1 \ (merge \ r_1 \ t_2)$
 else $node \ l_2 \ a_2 \ (merge \ t_1 \ r_2))$

$node :: \ 'a \ lheap \Rightarrow \ 'a \Rightarrow \ 'a \ lheap \Rightarrow \ 'a \ lheap$

$node \ l \ a \ r =$
(let $mhl = mht \ l; \ mhr = mht \ r$
 in if $mhr \leq mhl$ then $\langle l, (a, mhr + 1), r\rangle$
    else $\langle r, (a, mhl + 1), l\rangle)$

## *merge*

$merge\ (\langle l_1,\ (a_1,\ n_1),\ r_1 \rangle =:\ t_1)$
$(\langle l_2,\ (a_2,\ n_2),\ r_2 \rangle =:\ t_2) =$
$(\text{if}\ a_1 \leq a_2\ \text{then}\ node\ l_1\ a_1\ (merge\ r_1\ t_2)$
$\text{else}\ node\ l_2\ a_2\ (merge\ t_1\ r_2))$

Function *merge* terminates because
decreases with every recursive call.

# Functional correctness proofs

### including preservation of *invar*

Straightforward

# Logarithmic complexity

Correlation of rank and size:

**Lemma** $2^{mh(t)} \leq |t|_1$

Complexity measures $T\_merge$, $T\_insert$ $T\_del\_min$:
count calls of $merge$.

**Lemma** $[\![ ltree\ l;\ ltree\ r ]\!]$
$\implies T\_merge\ l\ r \leq mh(l) + mh(r) + 1$

**Corollary** $[\![ ltree\ l;\ ltree\ r ]\!]$
$\implies T\_merge\ l\ r \leq \log_2 |l|_1 + \log_2 |r|_1 + 1$

**Corollary**
$ltree\ t \implies T\_insert\ x\ t \leq \log_2 |t|_1 + 3$

**Corollary**
$ltree\ t \implies T\_del\_min\ t \leq 2 * \log_2 |t|_1 + 1$

Can we avoid the height info in each node?

# Archive of Formal Proofs

```
https://www.isa-afp.org/entries/Priority_
          Queue_Braun.shtml
```

# What is a Braun tree?

$braun :: \ 'a \ tree \Rightarrow bool$

$braun \ \langle \rangle \ = \ True$
$braun \ \langle l, \ x, \ r \rangle \ =$
$((|l| = |r| \ \vee \ |l| = |r| + 1) \ \wedge \ braun \ l \ \wedge \ braun \ r)$

1

**Lemma** $braun \ t \Longrightarrow 2^{h(t)} \leq 2 * |t| + 1$

# Idea of invariant maintenance

$braun \; \langle \rangle \; = \; True$
$braun \; \langle l, \; x, \; r \rangle \; =$
$((|l| = |r| \; \lor \; |l| = |r| + 1) \; \land \; braun \; l \; \land \; braun \; r)$

Let $t = \langle l, \; x, \; r \rangle$. Assume $braun \; t$

Add element: to $r$, then swap subtrees: $t' = \langle r', \; x, \; l \rangle$
To prove $braun \; t'$: $|l| \leq |r'| \; \land \; |r'| \leq |l| + 1$ $\qquad \square$

Delete element: from $l$, then swap subtrees: $t' = \langle r, \; x, \; l' \rangle$
To prove $braun \; t'$: $|l'| \leq |r| \; \land \; |r| \leq |l'| + 1$ $\qquad \square$

# Priority queue implementation

Implementation type: $'a\ tree$

Invariants: $heap$ and $braun$

No $merge$ — $insert$ and $del\_min$ defined explicitly

# *insert*

*insert* :: $'a \Rightarrow 'a\ tree \Rightarrow 'a\ tree$
*insert* $a\ \langle\rangle = \langle\langle\rangle,\ a,\ \langle\rangle\rangle$
*insert* $a\ \langle l,\ x,\ r\rangle =$
(if $a < x$ then $\langle insert\ x\ r,\ a,\ l\rangle$ else $\langle insert\ a\ r,\ x,\ l\rangle$)

Correctness and preservation of invariant straightforward.

# del_min

$del\_min :: {}'a\ tree \Rightarrow {}'a\ tree$
$del\_min\ \langle\rangle = \langle\rangle$
$del\_min\ \langle\langle\rangle,\ x,\ r\rangle = \langle\rangle$
$del\_min\ \langle l,\ x,\ r\rangle =$
$(\mathsf{let}\ (y,\ l') = del\_left\ l\ \mathsf{in}\ sift\_down\ r\ y\ l')$

- ❶ Delete leftmost element $y$
- ❷ Sift $y$ from the root down

Reminiscent of heapsort, but not quite …

# del_left

$del\_left :: {}'a\ tree \Rightarrow {}'a \times {}'a\ tree$

$del\_left\ \langle\langle\rangle,\ x,\ r\rangle = (x,\ r)$

$del\_left\ \langle l,\ x,\ r\rangle =$
($\mathsf{let}\ (y,\ l') = del\_left\ l\ \mathsf{in}\ (y,\ \langle r,\ x,\ l'\rangle))$

# sift_down

$sift\_down :: \, 'a \; tree \Rightarrow \, 'a \Rightarrow \, 'a \; tree \Rightarrow \, 'a \; tree$

$sift\_down \; \langle\rangle \; a \; \_ = \langle\langle\rangle, \; a, \; \langle\rangle\rangle$

$sift\_down \; \langle\langle\rangle, \; x, \; \_\rangle \; a \; \langle\rangle =$
(if $a \leq x$ then $\langle\langle\langle\rangle, \; x, \; \langle\rangle\rangle, \; a, \; \langle\rangle\rangle$
 else $\langle\langle\langle\rangle, \; a, \; \langle\rangle\rangle, \; x, \; \langle\rangle\rangle$)

$sift\_down \; (\langle l_1, \; x_1, \; r_1\rangle =: t_1) \; a \; (\langle l_2, \; x_2, \; r_2\rangle =: t_2) =$
if $a \leq x_1 \wedge a \leq x_2$ then $\langle t_1, \; a, \; t_2\rangle$
else if $x_1 \leq x_2$ then $\langle sift\_down \; l_1 \; a \; r_1, \; x_1, \; t_2\rangle$
     else $\langle t_1, \; x_2, \; sift\_down \; l_2 \; a \; r_2\rangle$

Maintains $braun$

# Functional correctness proofs
## for $del\_min$

Many lemmas, mostly straightforward

# Logarithmic complexity

Running time of $insert$, $del\_left$ and $sift\_down$ (and therefore $del\_min$) bounded by height

Remember: $braun\ t \implies 2^{h(t)} \leq 2 * |t| + 1$

$\implies$

Above running times logarithmic in size

# Source of code

Based on code from

L.C. Paulson. *ML for the Working Programmer*. 1996

based on code from Chris Okasaki.

# Sorting with priority queue

$pq\ [] = empty$
$pq\ (x\#xs) = insert\ x\ (pq\ xs)$

$mins\ q =$
(if $is\_empty\ q$ then $[]$
 else $get\_min\ h\ \#\ mins\ (del\_min\ h))$

$sort\_pq = mins \circ pq$

Complexity of $sort$: $O(n \log n)$
if all priority queue functions have complexity $O(\log n)$

HOL/Data_Structures/
Binomial_Heap.thy

# Numerical method

Idea: only use trees $t_i$ of size $2^i$

## Example

To store (in binary) $11001$ elements: $[t_0, 0, 0, t_3, t_4]$

Merge $\approx$ addition with carry
Needs function to combine   two trees of size $2^i$
$\qquad\qquad\qquad\qquad\qquad\qquad$ into one tree of size $2^{i+1}$

# Binomial tree

**datatype** $'a\ tree =$
   $Node\ (rank\colon nat)\ (root\colon 'a)\ ('a\ tree\ list)$

Invariant: Node of rank $r$  has children $[t_{r-1}, \ldots, t_0]$
                                   of ranks $[r-1, \ldots, 0]$

$btree\ (Node\ r\ x\ ts) =$
$((\forall\ t \in set\ ts.\ btree\ t) \land map\ rank\ ts = rev\ [0..<r])$

Lemma
$btree\ t \implies |t| = 2^{rank\ t}$

# Combining two trees

How to combine   two trees of rank $i$
                  into one tree of rank $i{+}1$

$link\ (Node\ r\ x_1\ ts_1 =:\ t_1)\ (Node\ r'\ x_2\ ts_2 =:\ t_2) =$
$(\text{if}\ x_1 \leq x_2\ \text{then}\ Node\ (r + 1)\ x_1\ (t_2\ \#\ ts_1)$
$\ \text{else}\ Node\ (r + 1)\ x_2\ (t_1\ \#\ ts_2))$

# Binomial heap

Use sparse representation for binary numbers:
$[t_0, 0, 0, t_3, t_4]$ represented as $[ (0, t_0), (3, t_3), (4, t_4) ]$

**type_synonym** $'a\ heap = 'a\ tree\ list$

Remember: $tree$ contains rank

Invariant:
$invar\ ts =$
$((\forall\ t \in set\ ts.\ bheap\ t) \land sorted\_wrt\ (<)\ (map\ rank\ ts))$
$bheap\ t = (btree\ t \land heap\ t)$
$heap\ (Node\ \_\ x\ ts) = (\forall\ t \in set\ ts.\ heap\ t \land x \leq root\ t)$

# Inserting a tree into a heap

Intuition: propagate a carry

Precondition:
Rank of inserted tree $\leq$ ranks of trees in heap

$ins\_tree\ t\ [] = [t]$
$ins\_tree\ t_1\ (t_2\ \#\ ts) =$
$(\text{if}\ rank\ t_1 < rank\ t_2\ \text{then}\ t_1\ \#\ t_2\ \#\ ts$
$\ \text{else}\ ins\_tree\ (link\ t_1\ t_2)\ ts)$

## *merge*

$merge\ ts_1\ [] = ts_1$
$merge\ []\ ts_2 = ts_2$
$merge\ (t_1\ \#\ ts_1 =:\ h_1)\ (t_2\ \#\ ts_2 =:\ h_2) =$
(if $rank\ t_1 < rank\ t_2$ then $t_1\ \#\ merge\ ts_1\ h_2$
 else if $rank\ t_2 < rank\ t_1$ then $t_2\ \#\ merge\ h_1\ ts_2$
     else $ins\_tree\ (link\ t_1\ t_2)\ (merge\ ts_1\ ts_2))$

Intuition: Addition of binary numbers

Note: Handling of carry *after* recursive call

# Get/delete minimum element

All trees are min-heaps.

Smallest element may be any root node:

$$ts \neq [] \implies get\_min \; ts = Min \; (set \; (map \; root \; ts))$$

Similar:

$get\_min\_rest :: \; 'a \; tree \; list \Rightarrow \; 'a \; tree \times \; 'a \; tree \; list$

Returns tree with minimal root, and remaining trees

$del\_min \; ts =$
(case $get\_min\_rest \; ts$ of
  $(Node \; r \; x \; ts_1, \; ts_2) \Rightarrow merge \; (rev \; ts_1) \; ts_2)$

Why $rev$? Rank decreasing in $ts_1$ but increasing in $ts_2$

# Complexity

Recall: $btree\ t \Longrightarrow |t| = 2^{rank\ t}$

$\Longrightarrow$ length of heap logarithmic in number of elements:
$invar\ ts \Longrightarrow length\ ts \leq \log_2 (|ts| + 1)$

Complexity of operations: linear in length of heap

Proofs straightforward?

# Complexity of $merge$

$merge\ (t_1\ \#\ ts_1\ =:\ h_1)\ (t_2\ \#\ ts_2\ =:\ h_2)\ =$
$(\text{if } rank\ t_1 < rank\ t_2 \text{ then } t_1\ \#\ merge\ ts_1\ h_2$
$\quad \text{else if } rank\ t_2 < rank\ t_1 \text{ then } t_2\ \#\ merge\ h_1\ ts_2$
$\qquad \text{else } ins\_tree\ (link\ t_1\ t_2)\ (merge\ ts_1\ ts_2))$

Complexity of $ins\_tree$: $T\_ins\_tree\ t\ ts \leq length\ ts + 1$

A call $merge\ t_1\ t_2$ (where $length\ ts_1 = length\ ts_2 = n$)
can lead to calls of $ins\_tree$ on lists of length 1, ..., $n$.

$\sum \in O(n^2)$

# Complexity of $merge$

$merge\ (t_1\ \#\ ts_1 =:\ h_1)\ (t_2\ \#\ ts_2 =:\ h_2) =$
$(\text{if } rank\ t_1 < rank\ t_2 \text{ then } t_1\ \#\ merge\ ts_1\ h_2$
$\quad \text{else if } rank\ t_2 < rank\ t_1 \text{ then } t_2\ \#\ merge\ h_1\ ts_2$
$\qquad\quad \text{else } ins\_tree\ (link\ t_1\ t_2)\ (merge\ ts_1\ ts_2))$

Relate time and length of input/output:

$T\_ins\_tree\ t\ ts + length\ (ins\_tree\ t\ ts) = 2 + length\ ts$

$T\_merge\ ts_1\ ts_2 + length\ (merge\ ts_1\ ts_2)$
$\leq 2 * (length\ ts_1 + length\ ts_2) + 1$

Yields desired linear bound!

# Sources

The inventor of the binomial heap:
Jean Vuillemin.
A Data Structure for Manipulating Priority Queues.
*CACM*, 1978.

The functional version:
Chris Okasaki. *Purely Functional Data Structures.*
Cambridge University Press, 1998.

# Priority queues so far

*insert*, *del_min* (and *merge*)
have logarithmic complexity

# Skew Binomial Heap

Similar to binomial heap, but involving also
*skew binary numbers*:

$d_1 \ldots d_n$ represents $\sum_{i=1}^{n} d_i * (2^{i+1} - 1)$
where $d_i \in \{0, 1, 2\}$

# Complexity

Skew binomial heap:

$$insert \text{ in time } O(1)$$
$$del\_min \text{ and } merge \text{ still } O(\log n)$$

Fibonacci heap (imperative!):

$$insert \text{ and } merge \text{ in time } O(1)$$
$$del\_min \text{ still } O(\log n)$$

Every operation in time $O(1)$?

# Puzzle

Design a functional queue
with (worst case) constant time $enq$ and $deq$ functions

# Chapter 10

## Amortized Complexity

# 20 Amortized Complexity
## Motivation
Formalization
Simple Classical Examples

# Example

$n$ increments of a binary counter starting with $0$

- WCC of one increment? $O(\log_2 n)$
- WCC of $n$ increments? $O(n * \log_2 n)$
- $O(n * \log_2 n)$ is too pessimistic!
- Every second increment is cheap and compensates for the more expensive increments
- Fact: WCC of $n$ increments is $O(n)$

WCC = worst case complexity

# The problem

WCC of individual operations
may lead to overestimation of
WCC of sequences of operations

# Amortized analysis

Idea:

Try to determine the average cost of each operation
(in the worst case!)

Use cheap operations to pay for expensive ones

Method:
- Cheap operations pay extra (into a "bank account"), making them more expensive
- Expensive operations withdraw money from the account, making them cheaper

# Bank account $=$ *Potential*

- The potential ("credit") is implicitly "stored" in the data structure.

- Potential $\Phi :: data\text{-}structure \Rightarrow non\text{-}neg.\ number$
  tells us how much credit is stored in a data structure

- Increase in potential $=$
  deposit to pay for *later* expensive operation

- Decrease in potential $=$
  withdrawal to pay for expensive operation

# Back to example: counter

Increment:

- Actual cost: 1 for each bit flip
- Bank transaction:
    - pay in 1 for final $0 \to 1$ flip
    - take out 1 for each $1 \to 0$ flip

    $\implies$ increment has amortized cost $2 = 1{+}1$

Formalization via potential:

$\Phi \; counter =$ the number of 1's in $counter$

# Data structure

Given an implementation:

- Type $\tau$
- Operation(s) $f :: \tau \Rightarrow \tau$
  (may have additional parameters)
- Initial value: $init :: \tau$
  (function "empty")

Needed for complexity analysis:

- Time/cost: $T\_f :: \tau \Rightarrow num$
  ($num =$ some numeric type
  $nat$ may be inconvenient)
- Potential $\Phi :: \tau \Rightarrow num$   (creative spark!)

Need to prove: $\Phi\ s \geq 0$ and $\Phi\ init = 0$

# Amortized and real cost

Sequence of operations: $f_1, \ldots, f_n$

Sequence of states:

$$s_0 := init, \; s_1 := f_1 \; s_0, \ldots, s_n := f_n \; s_{n-1}$$

Amortized cost := real cost + potential difference

$$A_{i+1} := T\_f_{i+1} \; s_i + \Phi \; s_{i+1} - \Phi \; s_i$$

$$\implies$$

Sum of amortized costs $\geq$ sum of real costs

$$
\begin{aligned}
\sum_{i=1}^{n} A_i \;&=\; \sum_{i=1}^{n} \left( T\_f_i \; s_{i-1} + \Phi \; s_i - \Phi \; s_{i-1} \right) \\
&=\; \left( \sum_{i=1}^{n} T\_f_i \; s_{i-1} \right) + \Phi \; s_n - \Phi \; init \\
&\geq\; \sum_{i=1}^{n} T\_f_i \; s_{i-1}
\end{aligned}
$$

# Verification of amortized cost

For each operation $f$:
provide an upper bound for its amortized cost

$$A\_f :: \tau \Rightarrow num$$

and prove

$$T\_f\, s + \Phi(f\, s) - \Phi\, s \leq A\_f\, s$$

# Back to example: counter

$incr :: bool\ list \Rightarrow bool\ list$
$incr\ [] = [True]$
$incr\ (False\ \#\ bs) = True\ \#\ bs$
$incr\ (True\ \#\ bs) = False\ \#\ incr\ bs$

$init = []$

$\Phi\ bs = length\ (filter\ id\ bs)$

## Lemma
$T\_incr\ bs + \Phi\ (incr\ bs) - \Phi\ bs = 2$
**Proof** by induction

# Proof obligation summary

- $\Phi \ s \geq 0$
- $\Phi \ init = 0$
- For every operation $f :: \tau \Rightarrow ... \Rightarrow \tau$:
  $T\_f \ s \ \overline{x} + \Phi(f \ s \ \overline{x}) - \Phi \ s \leq A\_f \ s \ \overline{x}$

If the data structure has an invariant $invar$:
assume precondition $invar \ s$

If $f$ takes 2 arguments of type $\tau$:
$T\_f \ s_1 \ s_2 \ \overline{x} + \Phi(f \ s_1 \ s_2 \ \overline{x}) - \Phi \ s_1 - \Phi \ s_2 \leq A\_f \ s_1 \ s_2 \ \overline{x}$

# Warning: real time

Amortized analysis unsuitable for real time applications:

Real running time for individual calls
may be much worse than amortized time

# Warning: single threaded

Amortized analysis is only correct for single threaded uses of the data structure.

Single threaded = no value is used more than once

Otherwise:

$$
\begin{aligned}
&\textit{let} \quad counter = 0; \\
&\qquad bad = \text{increment } counter \; 2^n - 1 \text{ times}; \\
&\qquad \_ = incr \; bad; \\
&\qquad \_ = incr \; bad; \\
&\qquad \_ = incr \; bad; \\
&\qquad \vdots
\end{aligned}
$$

# Warning: observer functions

*Observer function*: does not modify data structure

$\implies$ Potential difference $= 0$

$\implies$ amortized cost $=$ real cost

$\implies$ Must analyze WCC of observer functions

This makes sense because

  Observer functions do not consume their arguments!

Legal:    *let* $bad$ $=$    create unbalanced data structure
                       with high potential;
           _ $=$    $observer\ bad$;
           _ $=$    $observer\ bad$;
           $\vdots$

# Archive of Formal Proofs

```
https://www.isa-afp.org/entries/Amortized_
              Complexity.shtml
```

# Fact

Can reverse $[x_1, \ldots, x_n]$ onto $ys$ in $n$ steps:

$\quad ([x_1, \ x_2, \ x_3, \ \ldots, \ x_n], \ ys)$
$\rightarrow ([x_2, \ x_3, \ \ldots, \ x_n], \ x_1 \ \# \ ys)$
$\rightarrow ([x_3, \ \ldots, \ x_n], \ x_2 \ \# \ x_1 \ \# \ ys)$
$\vdots$
$\rightarrow ([], \ x_n \ \# \ \ldots \ \# \ x_1 \ \# \ ys)$

# The problem with $(front,\ rear)$ queues

- Only <span style="color:red">amortized</span> linear complexity of $enq$ and $deq$
- Problem: $([],\ rear)$ requires reversal of $rear$

# Solution

- Do not wait for $([], rear)$
- Compute new front $front @ rev\ rear$
  early and slowly
- In parallel with $enq$ and $deq$ calls
- Using a 'copy' of $front$ and $rear$
  "shadow queue"

# Solution

When to start? When $|front| = n$ and $|rear| = n+1$

Two phases:

$$front \quad \xrightarrow{n} \quad rev\ front$$

$$\searrow n$$

$$front\ @\ rev\ rear$$

$$\nearrow$$

$$rear \quad \xrightarrow{n+1} \quad rev\ rear$$

Must finish before original $front$ is empty.
$\Rightarrow$ Must take two steps in every $enq$ and $deq$ call

# Complication

Calls of $deq$ remove elements from the original $front$

Cannot easily remove them from the modified copy of $front$

Solution:
- Remember how many elements have been removed
- Better: how many elements are still valid

# Example

| | | |
|---|---|---|
| $enq$: | $([1..5],\ [11..6],$ | $Idle)$ |
| $\rightarrow$ | $([1..5],\ [],$ | $R\ (0,\ [1..5],[],\quad [11..6],[])$ |
| $\rightarrow^2$ | | $R\ (2,\ [3..5],[2..1],\ [9..6],[10..11])$ |
| $deq$: | $([2..5],\ [],$ | $R\ (1,\ [3..5],[2..1],\ [9..6],[10..11])$ |
| $\rightarrow^2$ | | $R\ (3,\ [5],[4..1],\quad [7..6],[8..11])$ |
| $enq$: | $([2..5],\ [12],$ | $R\ (3,\ [5],[4..1],\quad [7..6],[8..11])$ |
| $\rightarrow$ | | $R\ (4,\ [],[5..1],\quad [6],[7..11])$ |
| $\rightarrow$ | | $A\ (4,\ [5..1],\ [6..11])$ |
| $deq$: | $([3..5],\ [12],$ | $A\ (3,\ [5..1],\ [6..11])$ |
| $\rightarrow^2$ | | $A\ (1,\ [3..1],\ [4..11])$ |
| $deq$: | $([4..5],\ [12],$ | $A\ (0,\ [3..1],\ [4..11])$ |
| $\rightarrow$ | | $Done\ [4..11])$ |
| $\rightarrow$ | $([4..11],\ [12],$ | $Idle)$ |

# The shadow queue

**datatype** $'a$ $status$ =
    $Idle$ |
    $Rev$ ($nat$) ($'a$ $list$) ($'a$ $list$) ($'a$ $list$) ($'a$ $list$) |
    $App$ ($nat$) ($'a$ $list$) ($'a$ $list$) |
    $Done$ ($'a$ $list$)

# Shadow step

$exec :: {}'a\ status \Rightarrow {}'a\ status$

$exec\ Idle = Idle$

$exec\ (Rev\ ok\ (x\ \#\ f)\ f'\ (y\ \#\ r)\ r')$
$= Rev\ (ok + 1)\ f\ (x\ \#\ f')\ r\ (y\ \#\ r')$

$exec\ (Rev\ ok\ [\,]\ f'\ [y]\ r') = App\ ok\ f'\ (y\ \#\ r')$

$exec\ (App\ (ok + 1)\ (x\ \#\ f')\ r') = App\ ok\ f'\ (x\ \#\ r')$

$exec\ (App\ 0\ f'\ r') = Done\ r'$

$exec\ (Done\ v) = Done\ v$

# Dequeue from shadow queue

$invalidate :: {}'a\ status \Rightarrow {}'a\ status$

$invalidate\ Idle = Idle$
$invalidate\ (Rev\ ok\ f\ f'\ r\ r') = Rev\ (ok - 1)\ f\ f'\ r\ r'$
$invalidate\ (App\ (ok + 1)\ f'\ r') = App\ ok\ f'\ r'$
$invalidate\ (App\ 0\ f'\ (x\ \#\ r')) = Done\ r'$
$invalidate\ (Done\ v) = Done\ v$

# The whole queue

**record** $'a\ queue =$    $front$  :: $'a\ list$
$lenf$  :: $nat$
$rear$  :: $'a\ list$
$lenr$  :: $nat$
$status$ :: $'a\ status$

# $enq$ and $deq$

$enq\ x\ q =$
$check\ (q(\!|rear := x\ \#\ rear\ q,\ lenr := lenr\ q + 1|\!))$

$deq\ q =$
$check$
$\ (q(\!|lenf := lenf\ q - 1,\ front := tl\ (front\ q),$
$\qquad status := invalidate\ (status\ q)|\!))$

$check\ q =$
$(\text{if } lenr\ q \leq lenf\ q \text{ then } exec2\ q$
$\text{ else let } newstate =$
$\qquad\qquad Rev\ 0\ (front\ q)\ []\ (rear\ q)\ []$
$\quad \text{ in } exec2$
$\qquad (q(\!|lenf := lenf\ q + lenr\ q,$
$\qquad\qquad status := newstate,$
$\qquad\qquad rear := [],\ lenr := 0|\!)))$

$exec2\ q = (case\ exec\ (exec\ q)\ of$
$\qquad\qquad Done\ fr \Rightarrow q(\!|status = Idle,\ front = fr|\!)\ |$
$\qquad\qquad newstatus \Rightarrow q(\!|status = newstatus|\!))$

# Correctness

The proof is
- easy because all functions are non-recursive
  ($\implies$ constant running time!)
- tricky because of invariant

# *status* invariant

$inv\_st \ (Rev \ ok \ f \ f' \ r \ r') =$
$(|f| + 1 = |r| \ \wedge \ |f'| = |r'| \ \wedge \ ok \leq |f'|)$
$inv\_st \ (App \ ok \ f' \ r') = (ok \leq |f'| \ \wedge \ |f'| < |r'|)$
$inv\_st \ Idle = True$
$inv\_st \ (Done \ \_) = True$

# Queue invariant

$invar\ q =$
$(lenf\ q = |front\_list\ q| \land$
$\quad lenr\ q = |rev\ (rear\ q)| \land$
$\quad lenr\ q \le lenf\ q \land$
$\quad (\textsf{case}\ status\ q\ \textsf{of}$
$\qquad Rev\ ok\ f\ f'\ r\ r' \Rightarrow$
$\qquad\quad 2 * lenr\ q \le |f'| \land$
$\qquad\quad ok \ne 0 \land 2 * |f| + ok + 2 \le 2 * |front\ q|$
$\quad |\ App\ ok\ f\ r \Rightarrow$
$\qquad\quad 2 * lenr\ q \le |r| \land ok + 1 \le 2 * |front\ q|$
$\quad |\ \_ \Rightarrow True) \land$
$(\exists\ rest.\ front\_list\ q = front\ q\ @\ rest) \land$
$(\nexists fr.\ status\ q = Done\ fr) \land inv\_st\ (status\ q))$

# Queue invariant

$front\_list\ q =$
(case $status\ q$ of
  $Idle \Rightarrow front\ q$
| $Rev\ ok\ f\ f'\ r\ r' \Rightarrow rev\ (take\ ok\ f')\ @\ f\ @\ rev\ r\ @\ r'$
| $App\ ok\ f'\ x \Rightarrow rev\ (take\ ok\ f')\ @\ x$
| $Done\ f \Rightarrow f)$

# Archive of Formal Proofs

```
https://www.isa-afp.org/entries/Hood_
      Melville_Queue.shtml
```

# Inventors

Robert Hood and Robert Melville.
Real-Time Queue Operation in Pure LISP.
*Information Processing Letters*, 1981.

# Generalization

Real-time *double-ended* queue

Inventors: Hood (1982), Chuang and Goldberg (1993)
Verifiers: Toth and Nipkow (2023)
4500 lines of Isabelle (Hood-Melville queue: 800)

# Archive of Formal Proofs

```
https:
//www.isa-afp.org/entries/Skew_Heap.shtml
```

A *skew heap* is a self-adjusting heap (priority queue)

Functions $insert$, $merge$ and $del\_min$
have amortized logarithmic complexity.

Functions $insert$ and $del\_min$ are defined via $merge$

# Implementation type

Ordinary binary trees

Invariant: $heap$

# $merge$

$merge \; \langle \rangle \;\; t = t$
$merge \; h \; \langle \rangle = h$

Swap subtrees when descending:

$merge \; (\langle l_1, \; a_1, \; r_1 \rangle =: t_1) \; (\langle l_2, \; a_2, \; r_2 \rangle =: t_2) =$
$(\text{if } a_1 \leq a_2 \text{ then } \langle merge \; t_2 \; r_1, \; a_1, \; l_1 \rangle$
$\;\; \text{else } \langle merge \; t_1 \; r_2, \; a_2, \; l_2 \rangle)$

Function $merge$ terminates because …?

# *merge*

Very similar to leftist heap but

- subtrees are *always* swapped
- no size information needed

# Functional correctness proofs

Straightforward

# 22 Skew Heap

Amortized Analysis

# Archive of Formal Proofs

```
https://www.isa-afp.org/theories/amortized_
        complexity/#Skew_Heap_Analysis
```

# Logarithmic amortized complexity

**Theorem**
$T\_merge\ t_1\ t_2 + \Phi\ (merge\ t_1\ t_2) - \Phi\ t_1 - \Phi\ t_2$
$\leq 3 * \log_2\ (|t_1|_1 + |t_2|_1) + 1$

# Towards the proof

Right heavy:
$rh\ l\ r = (\text{if } |l| < |r| \text{ then } 1 \text{ else } 0)$

Number of right heavy nodes on left spine:
$lrh\ \langle\rangle = 0$
$lrh\ \langle l,\ \_,\ r\rangle = rh\ l\ r + lrh\ l$

## Lemma
$2^{lrh\ t} \leq |t| + 1$

## Corollary
$lrh\ t \leq \log_2 |t|_1$

# Towards the proof

Right heavy:
$rh\ l\ r = ($if $|l| < |r|$ then $1$ else $0)$

Number of not right heavy nodes on right spine:
$rlh\ \langle\rangle = 0$
$rlh\ \langle l,\ _-,\ r\rangle = 1 -\ rh\ l\ r +\ rlh\ r$

## Lemma
$2^{rlh\ t} \leq |t| + 1$

## Corollary
$rlh\ t \leq \log_2 |t|_1$

# Potential

The potential is the number of right heavy nodes:

$\Phi \langle\rangle = 0$

$\Phi \langle l, \_, r \rangle = \Phi\ l + \Phi\ r + rh\ l\ r$

## Lemma

$T\_merge\ t_1\ t_2 + \Phi\ (merge\ t_1\ t_2) - \Phi\ t_1 - \Phi\ t_2$
$\leq lrh\ (merge\ t_1\ t_2) + rlh\ t_1 + rlh\ t_2 + 1$

```
by(induction t1 t2 rule: merge.induct)(auto)
```

# Node-Node case

Let $t_1 = \langle l_1,\ a_1,\ r_1 \rangle$, $t_2 = \langle l_2,\ a_2,\ r_2 \rangle$.
Case $a_1 \leq a_2$. Let $m = merge\ t_2\ r_1$

$T\_merge\ t_1\ t_2 + \Phi\ (merge\ t_1\ t_2) - \Phi\ t_1 - \Phi\ t_2$
$= T\_merge\ t_2\ r_1 + 1 + \Phi\ m + \Phi\ l_1 + rh\ m\ l_1$
$\quad - \Phi\ t_1 - \Phi\ t_2$
$= T\_merge\ t_2\ r_1 + 1 + \Phi\ m + rh\ m\ l_1$
$\quad - \Phi\ r_1 - rh\ l_1\ r_1 - \Phi\ t_2$
$\leq lrh\ m + rlh\ t_2 + rlh\ r_1 + rh\ m\ l_1 + 2 - rh\ l_1\ r_1$
$\quad$ by IH
$= lrh\ m + rlh\ t_2 + rlh\ t_1 + rh\ m\ l_1 + 1$
$= lrh\ (merge\ t_1\ t_2) + rlh\ t_1 + rlh\ t_2 + 1$

# Main proof

$T\_merge\ t_1\ t_2\ +\ \Phi\ (merge\ t_1\ t_2)\ -\ \Phi\ t_1\ -\ \Phi\ t_2$
$\leq lrh\ (merge\ t_1\ t_2)\ +\ rlh\ t_1\ +\ rlh\ t_2\ +\ 1$
$\leq \log_2 |merge\ t_1\ t_2|_1\ +\ \log_2 |t_1|_1\ +\ \log_2 |t_2|_1\ +\ 1$
$=\ \log_2 (|t_1|_1\ +\ |t_2|_1\ -\ 1)\ +\ \log_2 |t_1|_1\ +\ \log_2 |t_2|_1\ +\ 1$
$\leq \log_2 (|t_1|_1\ +\ |t_2|_1)\ +\ \log_2 |t_1|_1\ +\ \log_2 |t_2|_1\ +\ 1$
$\leq \log_2 (|t_1|_1\ +\ |t_2|_1)\ +\ 2 * \log_2 (|t_1|_1\ +\ |t_2|_1)\ +\ 1$
  because $\log_2 x\ +\ \log_2 y\ \leq 2 * \log_2 (x\ +\ y)$ if $x,y > 0$
$=\ 3 * \log_2 (|t_1|_1\ +\ |t_2|_1)\ +\ 1$

# *insert* and *del_min*

Easy consequences:

## Lemma
$T\_insert\ a\ t\ +\ \Phi\ (insert\ a\ t)\ -\ \Phi\ t$
$\leq\ 3 * \log_2 (|t|_1 + 2) + 2$

## Lemma
$T\_del\_min\ t\ +\ \Phi\ (del\_min\ t)\ -\ \Phi\ t$
$\leq\ 3 * \log_2 (|t|_1 + 2) + 2$

# Sources

The inventors of skew heaps:

Daniel Sleator and Robert Tarjan.
Self-adjusting Heaps.
*SIAM J. Computing*, 1986.

The formalization is based on

Anne Kaldewaij and Berry Schoenmakers.
The Derivation of a Tighter Bound for Top-down Skew
Heaps. *Information Processing Letters*, 1991.

# Archive of Formal Proofs

```
https:
//www.isa-afp.org/entries/Splay_Tree.shtml
```

A *splay tree* is a self-adjusting binary search tree.

Functions $isin$, $insert$ and $delete$
have amortized logarithmic complexity.

## Definition (splay)

Become wider or more separated.

## Example

The river splayed out into a delta.

# Splay tree

Implementation type $=$ binary tree

Key operation $splay\ a$:
1. Search for $a$ ending up at $x$
   where $x = a$ or $x$ is a leaf node.
2. Move $x$ to the root of the tree by rotations.
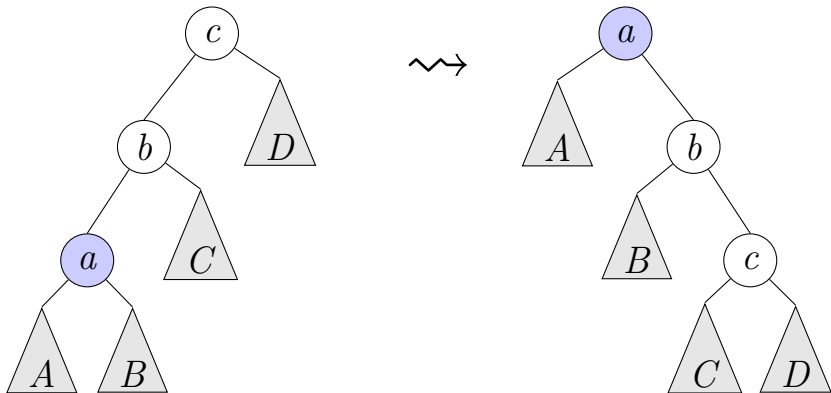
Derived operations $isin/insert/delete\ a$ :
1. $splay\ a$
2. Perform $isin/insert/delete$ action
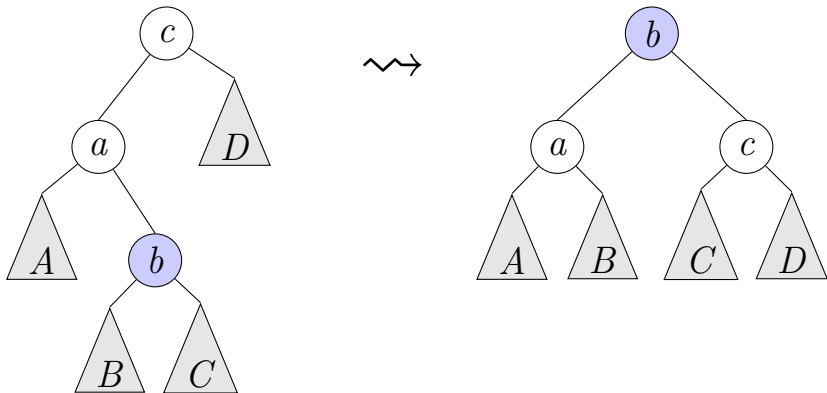
# Key ideas

Move to root

Double rotations

# Zig-zag

# Zig-zig and zig-zag

Zig-zig $\neq$ two single rotations

Zig-zag $=$ two single rotations

# Functional definition

$splay :: {}'a \Rightarrow {}'a\ tree \Rightarrow {}'a\ tree$

# Zig-zig and zig-zag

$[\![ x < b;\ x < c;\ AB \neq \langle\rangle ]\!]$
$\implies splay\ x\ \langle\langle AB,\ b,\ C\rangle,\ c,\ D\rangle =$
$\quad$ (case $splay\ x\ AB$ of
$\quad\quad \langle A,\ a,\ B\rangle \implies \langle A,\ a,\ \langle B,\ b,\ \langle C,\ c,\ D\rangle\rangle\rangle)$

$[\![ x < c;\ c < a;\ BC \neq \langle\rangle ]\!]$
$\implies splay\ c\ \langle\langle A,\ x,\ BC\rangle,\ a,\ D\rangle =$
$\quad$ (case $splay\ c\ BC$ of
$\quad\quad \langle B,\ b,\ C\rangle \implies \langle\langle A,\ x,\ B\rangle,\ b,\ \langle C,\ a,\ D\rangle\rangle)$

# Some base cases

$x < b \implies splay\ x\ \langle\langle A,\ x,\ B\rangle,\ b,\ C\rangle = \langle A,\ x,\ \langle B,\ b,\ C\rangle\rangle$

$x < a \implies$
$splay\ x\ \langle\langle\langle\rangle,\ a,\ A\rangle,\ b,\ B\rangle = \langle\langle\rangle,\ a,\ \langle A,\ b,\ B\rangle\rangle$

# Functional correctness proofs

Automatic

# Archive of Formal Proofs

```
https://www.isa-afp.org/theories/amortized_
        complexity/#Splay_Tree_Analysis
```

# Potential

Sum of logarithms of the size of all nodes:

$\Phi \langle \rangle = 0$

$\Phi \langle l,\ a,\ r \rangle = \varphi \langle l,\ a,\ r \rangle + \Phi\ l + \Phi\ r$

where $\varphi\ t = \log_2\ (|t|\ +\ 1)$

Amortized complexity of $splay$:

$A\_splay\ a\ t\ =\ T\_splay\ a\ t\ +\ \Phi\ (splay\ a\ t)\ -\ \Phi\ t$

# Analysis of *splay*

**Theorem**
$[\![ bst\ t;\ \langle l,\ a,\ r \rangle \in subtrees\ t ]\!]$
$\implies A\_splay\ a\ t \leq 3 * (\varphi\ t - \varphi\ \langle l,\ a,\ r \rangle) + 1$

**Corollary**

$[\![ bst\ t;\ x \in set\_tree\ t ]\!]$
$\implies A\_splay\ x\ t \leq 3 * (\varphi\ t - 1) + 1$

**Corollary** $bst\ t \implies A\_splay\ x\ t \leq 3 * \varphi\ t + 1$

**Lemma**
$[\![ t \neq \langle \rangle;\ bst\ t ]\!]$
$\implies \exists\ x' \in set\_tree\ t.$
$\qquad splay\ x'\ t = splay\ x\ t\ \wedge$
$\qquad T\_splay\ x'\ t = T\_splay\ x\ t$

$$insert$$

## Definition

$insert \; x \; t =$
(if $t = \langle \rangle$ then $\langle \langle \rangle, \, x, \, \langle \rangle \rangle$
 else case $splay \; x \; t$ of
        $\langle l, \, a, \, r \rangle \Rightarrow$ case $cmp \; x \; a$ of
                $LT \Rightarrow \langle l, \, x, \, \langle \langle \rangle, \, a, \, r \rangle \rangle$
                $| \; EQ \Rightarrow \langle l, \, a, \, r \rangle$
                $| \; GT \Rightarrow \langle \langle l, \, a, \, \langle \rangle \rangle, \, x, \, r \rangle )$

Counting only the cost of $splay$:

## Lemma

$bst \; t \Longrightarrow$
$T\_insert \; x \; t + \Phi \; (insert \; x \; t) - \Phi \; t \leq 4 * \varphi \; t + 3$

$$delete$$

## Definition

$delete\ x\ t =$
(if $t = \langle\rangle$ then $\langle\rangle$
 else case $splay\ x\ t$ of
        $\langle l,\ a,\ r\rangle \Rightarrow$
          if $x \neq a$ then $\langle l,\ a,\ r\rangle$
          else if $l = \langle\rangle$ then $r$
              else case $splay\_max\ l$ of
                      $\langle l',\ m,\ r'\rangle \Rightarrow \langle l',\ m,\ r\rangle)$

## Lemma

$bst\ t \Longrightarrow$
$T\_delete\ a\ t + \Phi\ (delete\ a\ t) - \Phi\ t \leq 6 * \varphi\ t + 3$

# Remember

Amortized analysis is only correct for single threaded
uses of a data structure.

Otherwise:

> *let* $counter = 0;$
> $bad =$ increment $counter$ $2^n - 1$ times;
> $\_ = incr\ bad;$
> $\_ = incr\ bad;$
> $\_ = incr\ bad;$
> $\vdots$

$$isin :: {}'a\ tree \Rightarrow {}'a \Rightarrow bool$$

Single threaded $\implies$ $isin\ t\ a$ eats up $t$

Otherwise:

      *let*   $bad =$ build unbalanced splay tree;
          $\_ = isin\ bad\ a;$
          $\_ = isin\ bad\ a;$
          $\_ = isin\ bad\ a;$
          $\vdots$

# Solution 1:

$$isin :: {}'a\ tree \Rightarrow {}'a \Rightarrow bool \times {}'a\ tree$$

Observer function returns new data structure:

Definition
$isin\ t\ a =$
(let $t' = splay\ a\ t$ in (case $t'$ of
$$\langle\rangle \Rightarrow False$$
$$|\ \langle l,\ x,\ r \rangle \Rightarrow a = x,$$
$$t'))$$

# Solution 2:
$$isin = splay;\ is\_root$$

Client uses $splay$ before calling $is\_root$:

## Definition
$is\_root :: {}'a \Rightarrow {}'a\ tree \Rightarrow bool$
$is\_root\ x\ t = ($case $t$ of
$\qquad\qquad \langle\rangle \Rightarrow False$
$\qquad\quad |\ \langle l,\ a,\ r\rangle \Rightarrow x = a)$

May call $is\_root\ \_\ t$ multiple times (with the same $t$!)
because $is\_root$ takes constant time

$\implies$ $is\_root\ \_\ t$ does not eat up $t$

Splay trees have an imperative flavour and are a bit
awkward to use in a purely functional language

# Sources

The inventors of splay trees:
Daniel Sleator and Robert Tarjan.
Self-adjusting Binary Search Trees. *J. ACM*, 1985.

The formalization is based on
Berry Schoenmakers. A Systematic Analysis of Splaying.
*Information Processing Letters*, 1993.

# Archive of Formal Proofs

```
https://www.isa-afp.org/entries/Pairing_
Heap.shtml
```

# Implementation type

**datatype** $'a\ heap\ =\ Empty\ |\ Hp\ 'a\ ('a\ heap\ list)$

Heap invariant:

$pheap\ Empty\ =\ True$
$pheap\ (Hp\ x\ hs)\ =$
$(\forall\ h\in set\ hs.\ (\forall\ y\in\#mset\_heap\ h.\ x \leq y)\ \wedge\ pheap\ h)$

Also: $Empty$ must only occur at the root

# *insert*

*insert x h = merge (Hp x [])* h

*merge h Empty = h*
*merge Empty h = h*
*merge (Hp x hsx =: hx) (Hp y hsy =: hy) =*
(if $x < y$ then *Hp x (hy # hsx)* else *Hp y (hx # hsy)*))

Like function *link* for binomial heaps

# del_min

$del\_min\ Empty\ =\ Empty$
$del\_min\ (Hp\ x\ hs)\ =\ pass_2\ (pass_1\ hs)$

$pass_1\ (h_1\ \#\ h_2\ \#\ hs)\ =\ merge\ h_1\ h_2\ \#\ pass_1\ hs$
$pass_1\ hs\ =\ hs$

$pass_2\ []\ =\ Empty$
$pass_2\ (h\ \#\ hs)\ =\ merge\ h\ (pass_2\ hs)$

# Fusing $pass_2 \circ pass_1$

$merge\_pairs\ [] = Empty$
$merge\_pairs\ [h] = h$
$merge\_pairs\ (h_1\ \#\ h_2\ \#\ hs) =$
$merge\ (merge\ h_1\ h_2)\ (merge\_pairs\ hs)$

## Lemma
$pass_2\ (pass_1\ hs) = merge\_pairs\ hs$

# Functional correctness proofs

Straightforward

# 24 Pairing Heap
Amortized Analysis

# Analysis

Analysis easier (more uniform) if a pairing heap is viewed as a binary tree:

$homs :: {}'a\ heap\ list \Rightarrow {}'a\ tree$
$homs\ [] = \langle\rangle$
$homs\ (Hp\ x\ hs_1\ \#\ hs_2) = \langle homs\ hs_1,\ x,\ homs\ hs_2\rangle$

$hom :: {}'a\ heap \Rightarrow {}'a\ tree$
$hom\ Empty = \langle\rangle$
$hom\ (Hp\ x\ hs) = \langle homs\ hs,\ x,\ \langle\rangle\rangle$

Potential function same as for splay trees

Verified:

> The functions $insert$, $del\_min$ and $merge$ all have $O(\log_2 n)$ amortized complexity.

These bounds are not tight.
Better amortized bounds in the literature:
$insert \in O(1)$, $del\_min \in O(\log_2 n)$, $merge \in O(1)$

The exact complexity is still open.

# Archive of Formal Proofs

```
https://www.isa-afp.org/entries/Amortized_
            Complexity.shtml
```

# Sources

The inventors of the pairing heap:

M. Fredman, R. Sedgewick, D. Sleator and R. Tarjan.
The Pairing Heap: A New Form of Self-Adjusting Heap.
*Algorithmica*, 1986.

The functional version:

Chris Okasaki. *Purely Functional Data Structures.*
Cambridge University Press, 1998.

# More trees

Huffman Trees

Finger Trees

B Trees

$k$-d Trees

Optimal BSTs

Priority Search Trees

Treaps

…

# Graph algorithms

Floyd-Warshall

Dijkstra Dijkstra

Maximum Network Flow

Strongly Connected Components

Kruskal Kruskal

Prim Prim

# Algorithms

Knuth-Morris-Pratt

Median of Medians

Approximation Algorithms

FFT

Gauss-Jordan

Simplex

QR-Decomposition

Smith Normal Form

Probabilistic Primality Testing

…

# Dynamic programming

- Start with recursive function
- Automatic translation to memoized version incl. correctness theorem
- Applications
  - Optimal binary search tree
  - Minimum edit distance
  - Bellman-Ford (SSSP)
  - CYK
  - …

# Infrastructure

Refinement Frameworks by Lammich:

Abstract specification
⤳ functional program
⤳ imperative program

using a library of collection types

# Model Checkers

- SPIN-like LTL Model Checker:
  Esparza, Lammich, Neumann, Nipkow, Schimpf,
  Smaus 2013
- SAT Certificate Checker:
  Lammich 2017; beats unverified standard tool

Mostly in the Archive of Formal Proofs