

Obfuscating Set Representations

Stephen Drape

Oxford University Computing Laboratory

with thanks to Jeff Sanders

What is obfuscation?

Obfuscation is a program transformation:

- Used to make a program "harder to understand"
- Try to make reverse engineering harder
- Must preserve functionality
- Concerns about efficiency

Why is obfuscation needed?

Obfuscation is usually applied to **object-oriented languages** such as Java and C#.

When compiling these languages, an **intermediate representation** is produced.

It is possible to recover the original code from this representation – **obfuscation can make this process harder.**

Fresh Approach

Instead of obfuscating an imperative program, we consider obfuscating operations of a **data-type** – we can then exploit properties of that data-type (see later!).

We model the data-types and the operations in **Haskell**.

Deriving and Proving

We want to obfuscate some set operations.

Using functional programs, we can:

- **Derive** obfuscations
- Easily establish **proofs of correctness**

Both of these are **difficult** to do in imperative languages.

List splitting

Adapt "array splitting" – consider a particular example "alternating split"

Write $xs \sim \langle l, r \rangle_a$ to denote xs is split into two lists l and r – xs is *data refined* by $\langle l, r \rangle_a$

$$[5, 7, 5, 4, 3, 1, 1] \sim \langle [5, 5, 3, 1], [7, 4, 1] \rangle_a$$

Invariant: $|r| \leq |l| \leq |r| + 1$

Splitting Function

$$\text{split}([]) = \langle [], [] \rangle_a$$

$$\text{split}([p]) = \langle [p], [] \rangle_a$$

$$\text{split}(p:q:xs) = \langle p:l, q:r \rangle_a$$

$$\text{where } \langle l, r \rangle_a = \text{split}(xs)$$

$$\text{unsplit } \langle [], [] \rangle_a = []$$

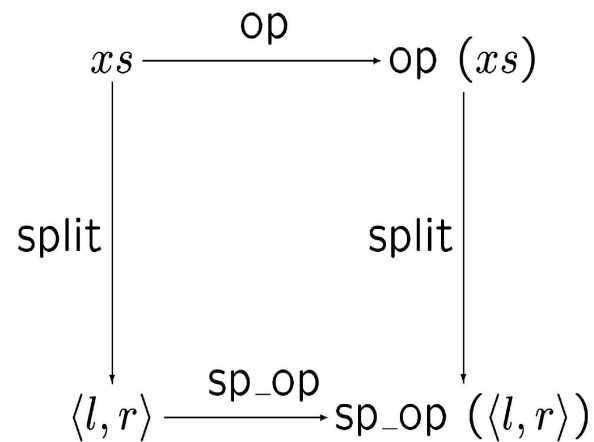
$$\text{unsplit } \langle [p], [] \rangle_a = [p]$$

$$\text{unsplit } \langle p:l, q:r \rangle_a = \\ p:q:\text{unsplit}(\langle l, r \rangle_a)$$

Derivations

$$\text{split} \cdot \text{op} = \text{sp_op} \cdot \text{split}$$

Data
Refinement



$$\text{sp_op} = \text{split} \cdot \text{op} \cdot \text{unsplit}$$

$$\text{op} = \text{unsplit} \cdot \text{sp_op} \cdot \text{split}$$

List operations

$$p : \langle l, r \rangle_a = \langle p : r, l \rangle_a$$

$$\langle l_0, r_0 \rangle_a ++ \langle l_1, r_1 \rangle_a$$

$$| |l_0| == |r_0| \quad = \langle l_0 ++ l_1, r_0 ++ r_1 \rangle_a$$

$$| \text{otherwise} \quad = \langle l_0 ++ r_1, r_0 ++ l_1 \rangle_a$$

: and ++ distribute over split

Unordered Lists without Duplicates

```
member p xs = or (map (==p) xs)
```

```
insert p xs = if member p xs then xs  
            else p:xs
```

```
delete p xs = ys ++ (if null zs  
                    then zs  
                    else tail zs)  
            where (ys, zs) = span (/=p) xs
```

Deriving Delete

Deriving the delete operation for split lists

`deletea p = split.delete p.unsplit`

Let $l = [l_0, l_1, \dots, l_j, l_{j+1}, \dots, l_n]$ $r = [r_0, r_1, \dots, r_{n'}]$

and $xs \sim \langle l, r \rangle_a$

We have three cases

(a) $p \in l$

(b) $p \in r$

(c) $p \notin \langle l, r \rangle_a$

Case (a) $(p \in l)$

1

Suppose that $l_j = p$

Let $(ly, lz) = \text{span } (/=p) l$
 $= ([l_0, l_1, \dots, l_{j-1}], [l_j, l_{j+1}, \dots, l_n])$

$(ry, rz) = \text{splitAt } |ly| r$
 $= ([r_0, r_1, \dots, r_{j-1}], [r_j, \dots, r_n])$

Case (a) $(p \in l)$

2

$\text{delete}_a p \langle l, r \rangle_a$

= {derivation equation}

$\text{split}(\text{delete } p(\text{unsplit} \langle l, r \rangle_a))$

= {definition of unsplit }

$\text{split}(\text{delete } p [l_0, r_0, l_1, \dots])$

= {definition of delete , $l_j = p$ }

$\text{split}([l_0, r_0, \dots, r_{j-1}] ++ [r_j, \dots])$

Case (a) $(p \in l)$

3

$\text{split}([l_0, r_0, \dots, r_{j-1}] ++ [r_j, \dots])$

= { split distributes over $++$ }

$\text{split}([l_0, r_0, \dots]) ++ \text{split}([r_j, \dots])$

= {definition of split }

$\langle [l_0, \dots], [r_0, \dots] \rangle_a ++ \langle [r_j, \dots, r_{n'}], [l_{j+1}, \dots, l_n] \rangle_a$

= {earlier definitions}

$\langle ly, ry \rangle_a ++ \langle rz, \text{tail } lz \rangle_a$

Case (a) $(p \in l)$

4

$$\begin{aligned} & \langle ly, ry \rangle_a ++ \langle rz, \text{tail } lz \rangle_a \\ & = \{\text{definition of } ++, |ly| = |ry|\} \\ & \langle ly ++ rz, ry ++ \text{tail } lz \rangle_a \end{aligned}$$

We cannot simplify this further,
but as lists are unordered:

$$\langle ly, ry \rangle_a \text{ is equivalent to } \langle ry, ly \rangle_a$$

Case (a) $(p \in l)$

5

$$\begin{aligned} & \langle ry, ly \rangle_a ++ \langle rz, \text{tail } lz \rangle_a \\ & = \{\text{definition of } ++\} \\ & \langle ry ++ rz, ly ++ \text{tail } lz \rangle_a \\ & = \{\text{definitions}\} \\ & \langle r, \text{delete } p \ l \rangle_a \end{aligned}$$

□

Case (b) $(p \in r)$

1

$\text{delete}_a p \langle l, r \rangle_a$

$= \{ l = (\text{head } l) : (\text{tail } l), l \neq [] \}$

$\text{delete}_a p \langle (\text{head } l) : (\text{tail } l), r \rangle_a$

$= \{\text{definition of } : \}$

$\text{delete}_a p ((\text{head } l) : \langle r, \text{tail } l \rangle_a)$

$= \{ \text{head } l \neq p \}$

$(\text{head } l) : (\text{delete}_a p \langle r, \text{tail } l \rangle_a)$

Case (b) $(p \in r)$

2

$(\text{head } l) : (\text{delete}_a p \langle r, \text{tail } l \rangle_a)$

= {previous definition of delete_a }

$(\text{head } l) : (\langle \text{tail } l, \text{delete } p r \rangle_a)$

= {definition of $:$ }

$\langle (\text{head } l) : (\text{delete } p r), \text{tail } l \rangle_a$

□

Finally

(c) $p \notin l$ and $p \notin r$

$$\text{delete}_a p \langle l, r \rangle_a = \langle l, r \rangle_a \quad \square$$

Final definition

$\text{delete}_a p \langle l, r \rangle_a$

| $\text{member } p l = \langle r, \text{delete } p l \rangle_a$

| $\text{member } p r =$

$\langle (\text{head } l) : (\text{delete } p r), \text{tail } l \rangle_a$

| otherwise = $\langle l, r \rangle_a$

Insert operation

$\text{insert}_a p \langle l, r \rangle_a =$
 if $\text{member}_a p \langle l, r \rangle_a$
 then $\langle l, r \rangle_a$
 else $\langle p:r, l \rangle_a$

We will now prove that

$\text{insert } p = \text{unsplit}. (\text{insert}_a p). \text{split}$

Case for $p \in xs$ is trivial.

Otherwise, suppose that:

$$xs \sim \langle l, r \rangle_a \text{ and } p \notin xs$$

`unsplit (inserta p split (xs))`

$$= \{xs \sim \langle l, r \rangle_a \}$$

`unsplit (inserta p (⟨l, r⟩a))`

Proof

2

$\text{unsplit}(\text{insert}_a p (\langle l, r \rangle_a))$

= {definition of insert_a }

$\text{unsplit}(\langle p:r, l \rangle_a)$

= {definition of $:$ }

$\text{unsplit}(p:\langle l, r \rangle_a)$

= {property of $:$ }

$p:\text{unsplit}(\langle l, r \rangle_a)$

$p : \text{unsplit} (\langle l, r \rangle_a)$

$= \{xs \sim \langle l, r \rangle_a\}$

$p : xs$

$= \{\text{definition of insert}\}$

$\text{insert } p \ xs$

□

Complexity

1

```
delete p xs = ys ++ (if null zs
                      then zs
                      else tail zs)
              where (ys,zs) = span (≠p) xs
```

```
delete_a p ⟨l,r⟩_a
| member p l = ⟨r, delete p l⟩_a
| member p r = ⟨(head l):(delete p r), tail l⟩_a
| otherwise = ⟨l,r⟩_a
```

Both functions have linear complexity.

Complexity

2

```
insert  $p$   $xs$  = if member  $p$   $xs$   
                then  $xs$   
                else  $p:xs$ 
```

```
insert $a$   $p$   $\langle l, r \rangle_a$  = if member $a$   $p$   $\langle l, r \rangle_a$   
                          then  $\langle l, r \rangle_a$   
                          else  $\langle p:r, l \rangle_a$ 
```

Again, these functions have linear complexity.

"Obfuscating Set Representations"

The paper looks at three representations:

- Unordered with duplicates
- Unordered without duplicates
- Strictly-increasing

Proofs and derivations of `delete` and `insert` are given for the other representations.

Also, another split is considered.

At the beginning, it was stated we obfuscate data-types directly so that we can exploit properties of the data-type.

Suppose that we want to split a matrix and we want to develop a transpose operation for the split matrix.

Suppose we flatten the matrix to an array and then split this array.

Using arrays means that we lose the "shape" of the matrix and so we have difficulty in constructing a transpose operation.

Using matrices directly:

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix}^T = \begin{pmatrix} A^T & C^T \\ B^T & D^T \end{pmatrix}$$

Conclusions

We have seen that using data-types and functional programming, we can

- derive obfuscations
- prove correctness

Our operations make little change to the complexity

Have to keep `split` secret

Future Work

Possible areas for future work

- Other obfuscations
- Other data-types (matrices, trees)
- Automation
- Obfuscation definition

