Compilers and Formal Languages

Email: christian.urban at kcl.ac.uk

Office: N7.07 (North Wing, Bush House)

Slides: KEATS (also home work is there)

Compilers & Boeings

777

First flight in 1994. They want to achieve triple redundancy in hardware faults.

They compile 1 Ada program to

- Intel 80486
- Motorola 68040 (old Macintosh's)
- AMD 29050 (RISC chips used often in laser printers)
 using 3 independent compilers.

Compilers & Boeings

777

First flight in 1994. They want to achieve triple redundancy in hardware faults.

They compile 1 Ada program to

- Intel 80486
- Motorola 68040 (old Macintosh's)
- AMD 29050 (RISC chips used often in laser printers)
 using 3 independent compilers.

Airbus uses C and static analysers. Recently started using CompCert.

seL4 / Isabelle

- verified a microkernel operating system (\approx 8000 lines of C code)
- US DoD has competitions to hack into drones; they found that the isolation guarantees of seL4 hold up
- CompCert and seL4 sell their code

POSIX Matchers

 Longest match rule ("maximal munch rule"): The longest initial substring matched by any regular expression is taken as the next token.

 Rule priority: For a particular longest initial substring, the first regular expression that can match determines the token.

POSIX Matchers

 Longest match rule ("maximal munch rule"): The longest initial substring matched by any regular expression is taken as the next token.

 Rule priority: For a particular longest initial substring, the first regular expression that can match determines the token.

Kuklewicz: most POSIX matchers are buggy

http://www.haskell.org/haskellwiki/

$$\begin{array}{lll} \operatorname{der} c \ (\mathbf{0}) & \stackrel{\operatorname{def}}{=} & \mathbf{0} \\ \operatorname{der} c \ (\mathbf{1}) & \stackrel{\operatorname{def}}{=} & \mathbf{0} \\ \operatorname{der} c \ (d) & \stackrel{\operatorname{def}}{=} & \operatorname{if} c = d \ \operatorname{then} \ \mathbf{1} \ \operatorname{else} \ \mathbf{0} \\ \operatorname{der} c \ (r_1 + r_2) & \stackrel{\operatorname{def}}{=} & (\operatorname{der} c \ r_1) + (\operatorname{der} c \ r_2) \\ \operatorname{der} c \ (r_1 \cdot r_2) & \stackrel{\operatorname{def}}{=} & \operatorname{if} \ \operatorname{nullable}(r_1) \\ \operatorname{then} \ ((\operatorname{der} c \ r_1) \cdot r_2) + (\operatorname{der} c \ r_2) \\ \operatorname{else} \ (\operatorname{der} c \ r_1) \cdot r_2 \\ \operatorname{der} c \ (r^*) & \stackrel{\operatorname{def}}{=} & (\operatorname{der} c \ r_1) \cdot r_2 \\ \operatorname{der} c \ (r^{\{n\}}) & \stackrel{\operatorname{def}}{=} & \operatorname{if} \ n = 0 \ \operatorname{then} \ \mathbf{0} \\ \operatorname{else} \ (\operatorname{der} c \ r) \cdot (r^{\{n-1\}}) \\ \operatorname{der} c \ (r^{\{\uparrow n\}}) & \stackrel{\operatorname{def}}{=} & \operatorname{if} \ n = 0 \ \operatorname{then} \ \mathbf{0} \\ \operatorname{else} \ (\operatorname{der} c \ r) \cdot (r^{\{\uparrow n-1\}}) \\ \end{array}$$

Proofs about Rexps

Remember their inductive definition:

If we want to prove something, say a property P(r), for all regular expressions r then ...

Proofs about Rexp (2)

- P holds for 0, 1 and c
- P holds for $r_1 + r_2$ under the assumption that P already holds for r_1 and r_2 .
- P holds for $r_1 \cdot r_2$ under the assumption that P already holds for r_1 and r_2 .
- P holds for r* under the assumption that P already holds for r.

...

Proofs about Strings

If we want to prove something, say a property P(s), for all strings s then ...

- P holds for the empty string, and
- P holds for the string c::s under the assumption that P already holds for s

Correctness of the Matcher

We want to prove

```
matches r s if and only if s \in L(r)
```

```
where matches\ r\ s \stackrel{\text{def}}{=} nullable(ders\ s\ r)
```

Correctness of the Matcher

We want to prove

matches
$$r$$
 s if and only if $s \in L(r)$

where $matches\ r\ s \stackrel{\text{def}}{=} nullable(ders\ s\ r)$

We can do this, if we know

$$L(der\ c\ r) = Der\ c\ (L(r))$$

Some Lemmas

- Der $c (A \cup B) =$ (Der c A) \cup (Der c B)
- If [] ∈ A then

$$Der c (A @ B) = (Der c A) @ B \cup (Der c B)$$

- If $[] \notin A$ then

 Der c (A @ B) = (Der c A) @ B
- Der $c(A^*) = (Der c A) @ A^*$ (interesting case)

Why?

Why does

Der
$$c(A^*) = (Der c A) @ A^*$$
 hold?

$$Der c (A^*) = Der c (A^* - \{[]\})$$

$$= Der c ((A - \{[]\}) @ A^*)$$

$$= (Der c (A - \{[]\})) @ A^*$$

$$= (Der c A) @ A^*$$

using the facts

Der c A = Der c
$$(A - \{[]\})$$
 and $(A - \{[]\}) @A^* = A^* - \{[]\}$

POSIX Spec

$$\frac{[] \in \mathbf{1} \to \mathsf{Empty}}{\mathsf{c} \in \mathsf{c} \to \mathsf{Char}(\mathsf{c})}$$

$$\frac{s \in r_1 \to \nu}{s \in r_1 + r_2 \to Left(\nu)}$$
$$\frac{s \in r_2 \to \nu \quad s \notin L(r_1)}{s \in r_1 + r_2 \to Right(\nu)}$$

$$s_{1} \in r_{1} \to \nu_{1}$$

$$s_{2} \in r_{2} \to \nu_{2}$$

$$\neg (\exists s_{3} s_{4}. s_{3} \neq [] \land s_{3} @ s_{4} = s_{2} \land s_{1} @ s_{3} \in L(r_{1}) \land s_{4} \in L(r_{2}))$$

$$s_{1} @ s_{2} \in r_{1} \cdot r_{2} \to Seq(\nu_{1}, \nu_{2})$$

Sulzmann & Lu Paper

 I have no doubt the algorithm is correct — the problem is I do not believe their proof.

"How could I miss this? Well, I was rather careless when stating this Lemma:)

Great example how formal machine checked proofs (and proof assistants) can help to spot flawed reasoning steps."

Sulzmann & Lu Paper

 I have no doubt the algorithm is correct — the problem is I do not believe their proof.

"How could I miss this? Well, I was rather careless when stating this Lemma:)

Great example how formal machine checked proofs (and proof assistants) can help to spot flawed reasoning steps."

Sulzmann & Lu Paper

 I have no doubt the algorithm is correct — the problem is I do not

Lemma 3 (Projection and Injection). Let r be a regular expression, l a letter and v a parse tree.

```
1. If \vdash v : r and |v| = lw for some word w, then \vdash proj_{(r,l)} v : r \setminus l.
```

2. If
$$\vdash v : r \setminus l$$
 then $(proj_{(r,l)} \circ inj_{r \setminus l}) \ v = v$.

3. If
$$|v| = v$$
 and $|v| = v$ for some word v , then $(inj_{r \setminus l} \circ proj_{(r,l)}) v = v$.

MS:BUG[Come accross this issue when going back to our constructive reg-ex work] Consider $\vdash [Right\ (), Left\ a]: (a+\epsilon)^*$. However, $proj_{((a+\epsilon)^*,a)}\ [Right\ (), Left\ a]$ fails! The point is that proj only works correctly if applied on POSIX parse trees.

MS:Possible fixes We only ever apply *proj* on Posix parse trees.

For convenience, we write " $\vdash v : r$ is POSIX" where we mean that $\vdash v : r$ holds and v is the POSIX parse tree of r for word |v|.

Lemma 2 follows from the following statement.