

Konrad Slind

School of Computing, University of Utah

November 20, 2006



▲ロト▲聞と▲臣と▲臣と 臣 のへの

- Logic built on top of typed lambda calculus
- Originally due to Church (1940s)
- First implemented by Gordon (early 1980s), by adapting LCF implementation of Milner and colleagues
- Now we have HOL-4, HOL-Light, ProofPower, and Isabelle/HOL, all vital systems
- Basically a kind of typed set theory that builds in functions
- Not clear that HOL has anything to do with FP.

Basic FP in HOL

- 1980s
- Gordon's initial work developed some basic types (numbers, pairs, lists) sufficient to do hardware verification examples.
- Melham (thesis) implemented a package for definition of inductive datatypes
- Each such definition provided induction and a so-called *primitive recursion* principle

Theorem (Primitive Recursion Theorem for lists)

```
|- !e f. ?!fn. (fn [] = e) /\
(!n. fn (h::t) = f h t (fn t))
```

◆□▶ ◆□▶ ▲□▶ ▲□▶ ▲□ ◆ ◇ ◇ ◇

Basic FP in HOL

- Illustrates standard methodology for HOL developments : do not extend logic with new axioms, but instead derive tools on top that use inference to mechanize general theorems
- For example,

is introduced by constructing appropriate e and f instances for the P.R. theorem and then deriving the specified equations

- Time consuming and possibly slow,
- BUT cuts down on soundness bugs, and gives a nice assurance story
- Reasoning about prim. rec. functions over inductive datatypes supported by custom induction principles.

1990s

- Good start
- BUT, from the point of view of FP'ers this is an impoverished setting in which to write programs:
 - Only very simple patterns
 - Only very restricted kinds of recursion
 - Numerous other irritations
- In 1990s systems emerged that dealt with some of these problems
- Fourman's LAMBDA system (defunct)
- TFL
 - Complex patterns
 - Arbitrary recursion (termination proofs required)
 - Per-function induction principles (following Boyer and Moore)

・ロット (雪) (日) (日) (日)

Based on

Theorem (Wellfounded Recursion theorem)

 $\vdash \textbf{WF} \ R \land (f = \textbf{WFREC} \ R \ M) \Rightarrow \exists f. \forall x. \ f(x) = M \ (f \mid_{Rx}) \ x$

- Works by instantiating and manipulating WF Rec. thm (proved in OL)
- A parameterized implementation, instantiated to HOL-4 and Isabelle/HOL
- Handles deep patterns, *e.g.*Okasaki-style Red-Black trees
- Deals well with mutually recursive functions
- Deals with nested recursive functions, but not well (since improved by Matthews and Krstic, and recently by Krauss)

(日)

FLAT again

The following version of **FLAT** has more complex patterns and also needs a termination proof in order to be admitted.

Definition (FLAT)

```
|- (FLAT [] = []) /\
  (FLAT([]::rst) = FLAT rst) /\
  (FLAT ((h::t)::rst) = h :: FLAT (t::rst))
```

Theorem (FLAT induction)

```
|- !P. P [] /\
    (!rst. P rst ==> P ([]::rst)) /\
    (!h t rst. P (t::rst) ==> P((h::t)::rst))
    ==>
    !list. P list
```

◆□▶ ◆□▶ ▲□▶ ▲□▶ ▲□ ◆ ○ ◆

Depth first fold with graph represented as a function of type

$\alpha \to \alpha \text{ list}$

which takes a node and delivers the children of the node.

$$\mathsf{DFSp}: (\alpha \to \alpha \mathsf{ list}) \to (\alpha \to \beta \to \beta) \to \alpha \mathsf{ list} \to \alpha \mathsf{ list} \to \beta \to \beta$$

- Folds function f over directed, possibly cyclic graph
- Applies *f* to each node and accumulating parameter *acc*
- By instantiating *f* can get map, search, max, filter, *etc* functions for such graphs
- Perfectly acceptable functional program

- The functional representation of the graph allows infinite graphs
- Makes the function partial (if graph has an infinite number of reachable nodes)
- For example $\lambda x.[x + 1]$
- HOL only supports total functions so DFSp wouldn't be admitted
- How to repair (totalize)?

< ロ > < 同 > < 回 > < 回 > 、 回

The fold will always terminate given a finite set of reachable nodes. How to define reachability?

```
Definition (Reachability)

\mathbf{R}_G \ x \ y \equiv \mathbf{mem} \ y \ (G \ x)

\mathbf{reach}_G \equiv \mathbf{RTC} \ R_G

\mathbf{reachlist}_G \ nodes \ y \equiv \exists x. \ \mathbf{mem} \ x \ nodes \land \mathbf{reach}_G \ x \ y
```

Thus, we want to constrain **DFSp** by finiteness of nodes reachable from root nodes of graph.

```
\mathsf{DFS}: (\alpha \to \alpha \mathsf{ list}) \to (\alpha \to \beta \to \beta) \to \alpha \mathsf{ list} \to \alpha \mathsf{ list} \to \beta \to \beta
```

```
DFS G f seen to visit acc =
  if Finite (reachlist<sub>G</sub> to visit)
      then case to visit
             of ] \Rightarrow acc
              |(h::t) \Rightarrow
                 if mem h seen
                     then DFS G f seen t acc
                     else DFS G f (h :: seen)
                                      (G h ++ t)
                                      (f h acc)
```

else **ARB**

▲ロト▲聞と▲臣と▲臣と 臣 のへの

```
\vdash DFS G f seen [] acc = acc
```

くロ とく得 とくき とくき とうき

- In first recursive call, list of seen nodes doesn't change, but nodes still to visit shrinks
- In second recursive call, seen nodes gets bigger and nodes to visit can increase in size
- So simple measures don't work.
- Recall that the set of reachable nodes is finite.
- Idea. Let ≺ be the lexicographic combination of the number of reachable nodes not yet seen and the number of nodes in *to_visit*.

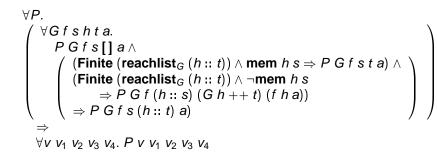
Termination relation

Definition

$$(G, f, seen', to_visit', acc') \prec (G, f, seen, to_visit, acc)$$

iff
 $(||reachlist_G to_visit' \setminus ListToSet seen'||, length to_visit') <_{lex}$
 $(||reachlist_G to_visit \setminus ListToSet seen||, length to_visit)$

- In first recursive call h has been previously seen, so the set of unseen reachable nodes does not change, and t is smaller than h :: t.
- In the second call, h is added to the seen list, and all of the nodes reachable from the children of h are also reachable from h itself. Thus the set of reachable nodes gets no additions, and since the addition to the seen list was previously reachable the size of the calculated set decreases.



< ロ > < 同 > < 回 > < 回 > .

3

What does it mean for this fold on graphs to be correct?

- All reachable nodes are visited
- No unreachable nodes are visited
- No reachable node is visited twice
- How, though, do we capture the notion of visits?
 - We capture this notion by using **cons** as the folding function given to **DFS**, so that the returned list is just the visited nodes.

DFS with folding function *f* is equal to gathering all the visited nodes and then folding *f* over the resulting list.

Theorem (DFS Fold)

Finite (reachlist_G to_visit) \Rightarrow DFS G f seen to_visit acc =

foldr f acc (DFS G cons seen to_visit [])

・ロット (雪) (日) (日) 日

Correctness

With this understanding, it suffices to prove that the invocation **DFS** *G* **cons** *seen to_visit* [] contains no duplicate entries, contains each node reachable from *to_visit*, and contains no nodes not so reachable. The first property is

Theorem (DFS Distinct)

Finite (reachlist_G to_visit) \Rightarrow all_distinct (DFS G cons seen to_visit [])

and the other two are phrased as

Theorem (DFS Reach)

Finite (reachlist_G to_visit) \Rightarrow $\forall x. \text{ reachlist}_G \text{ to}_visit x$ \Leftrightarrow mem x (DFS G cons [] to_visit [])

Konrad Slind FP in HOL

- Possible to go on and instantiate the various parameters of DFS to get various simplifications
 - Simpler constraint assuring (but not characterizing) termination : finite number of parent nodes in graph
 - DFS with adjacency lists

∃ >

An adjacency list

 $A: (\alpha \times \alpha \text{ list})$ list

gives a listing of nodes alongside their children.

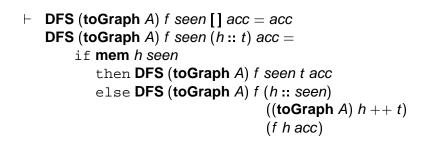
• toGraph converts an adjacency list into a graph.

```
Definition (toGraph)
```

```
toGraph al n =
case filter (\lambda(x, \_). (x = n)) al
of [] \rightarrow []
| (__, x) :: t \rightarrow x
```

 Can then prove that DFS terminates when called on any graph derived from an adjacency list

・ロン・(部)・・ヨン・ヨン 三連



・ロト・(部・・注)・・注)

- Fun tutorial study
- Formalization challenges (partiality, termination, *visits*, ...)
- Need to do math in order to work with such programs (*e.g.*, reachability)
- Possibly of future use; could be added to a library

∃ >

Current state of affairs

- Programming total functions over inductive datatypes is pretty well handled in most proof systems
- Could always be improved, of course
- Isabelle/HOL has a nice development of domain theory for applications that need it.
- But domains make life more complicated (lifting)
- Support for lazy datatypes and functions over them, not using domains, was pioneered by John Matthews, but is still not mechanized well in any HOL implementation.
- HOL systems have only simple types, and there doesn't seem to be much momentum for supporting more expressive type systems.

< ロ > < 同 > < 回 > < 回 > < □ > <

Critique and a Response

- Functions essentially trapped inside the formal system ("Case studies are boring")
- Sterile environment?
- But ACL2 community has shown that breaking free of the proof system is possible
- Emitting formal programs into outside world has many benefits
- Other systems have followed: PVS, Isabelle/HOL, HOL-4 all provide export for formal programs
- Coq has a variety of solutions, both internal and external
- Other systems (e.g., Matlab) also export programs and/or hardware

・ロット (雪) (日) (日) (日)

What then?

- Observation: programs are exported to the metalanguage
- Lisp for ACL2, ML for Isabelle/HOL and HOL
- Nice research project: export to mainstream languages like Java or even C
- However, current program export facilities exploit the fact that the conceptual gap between the formal program and the host PL is small
- But what if we want to export to Java, C, or even hardware?
- End up re-capitulating phases of compilation
- But then the small gap gets ever wider ...

3) 3

Compilation-by-proof

- Our current research investigates ways to
 - specify functional programs as mathematics
 - prove correctness properties at the mathematics level
 - translate to assembly or hardware
 - translation done by proof, so result is guaranteed to return the correct answers.
- Amounts to compilation of logic functions, inside the logic
- Two approaches to providing this:
 - Verified compiler. This is what is done traditionally
 - Translation validation. Recent alternative proposed by Pnueli

Compilation-by-proof

- Have built two prototype TV compilers for a very simple functional language
 - hardware (with Mike Gordon)
 - ARM assembly (with Owens, Li, Tuerk)
- Work is still very much in progress
- Target example: Elliptic Curve Crypto (relatively efficient replacement for RSA)
 - Formal theory of elliptic curves (on top of finite field theory)
 - Define recursive functions that implement, *e.g.*, addition of points on elliptic curves
 - Compile these to ARM assembly
 - formal ARM model in HOL-4

- Try to do as much as possible by source-to-source translations.
- Start by translating to combinator form, then to ANF (administrative normal form)
- These end up being *semantic* versions of the standard syntax bashing done in CPS translation
- Register allocation done by standard graph-colouring algorithm. Used to deliver an α-convertible version (nice trick from Jason Hickey)
- Maintenance of equality, by proof, from starting program
- That's the front end

・ロット (雪) (日) (日) 日

```
 \vdash \text{Rounds} (n, (y, z), (k_0, k_1, k_2, k_3), s) = \\ \text{if } n = 0w \text{ then } ((y, z), (k_0, k_1, k_2, k_3), s) \\ \text{else Rounds } (n - 1w, \\ \text{let } s' = s + 2654435769w \text{ in} \\ \text{let } y' = y + \text{ShiftXor}(z, s', k_0, k_1) \\ \text{in } ((y', z + \text{ShiftXor}(y', s', k_2, k_3)), (k_0, k_1, k_2, k_3), s')
```

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ()

After Front-end processing

```
\vdash Rounds(r_0, (r_8, r_5), (r_4, r_3, r_2, r_6), r_7) =
      let v_9 = (\mathbf{op} =) (r_0, 0w)
      in if V_9 then ((r_8, r_5), (r_4, r_3, r_2, r_6), r_7)
           else let m_2 = (\mathbf{op} -) (r_0, \mathbf{1}w) in
                     let m_4 = (op +) (r_7, 2654435769w) in
                     let r_1 = ShiftXor (r_5, m_4, r_4, r_3) in
                     let r_9 = (\mathbf{op} +) (r_8, r_1) in
                     let r_1 = ShiftXor (r_9, m_4, r_2, r_6) in
                     let r_1 = (\mathbf{op} +) (r_5, r_1) in
                     let ((m_5, m_3), (m_1, m_0, m_6, r_1), r_0) =
                                 Rounds (m_2, (r_9, r_1), (r_4, r_3, r_2, r_6), m_4)
                     in ((m_5, m_3), (m_1, m_0, m_6, r_1), r_0)
```

▲ロト ▲御 ト ▲ 国 ト ▲ 国 ト 二 連

- Back end poof *synthesizes* a counterpart function to the front end function. Then a tactic is executed to show the two are equal.
- In general, the backend is pretty conventional compiler verification technology
- We synthesize the IL semantics from the ARM semantics, rather than relating two operational semantics
- Main difficulty is dealing with memory
- In particular, function call is hard.

- To be done in near future: defunctionalization to support higher order
- Not mentioned: work on translating FP by proof to hardware
- Early days in this area, but I think it's quite exciting

・ロット (雪) (日) (日) 日

• **TFL** supports *recursion schemes*, by allowing free variables in rhs, for example

Definition (While-loops)

While s = if B s then While (C s) else s.

- While can also be defined directly
- Connection with FP: Lewis, Shields, Meijer, Launchbury, Implicit parameters: Dynamic scoping with static types

Polytypism

- Polytypism (type-indexed functions) is becoming a basic FP tool
- Applications in logic:
 - Termination proofs
 - Normalization by Evaluation
 - Translation between representations (mapping to binary format, to SAT, to LISP, *etc*)
- In HOL systems, two ways to support it:
 - A polytypic function *f* is represented by a meta-level function parameterized by a P.R. theorem (HOL-4).
 - Explicit definitions over type structure (Isabelle/HOL).

3

- Actually, not so recent ...
- Used extensively in Isabelle/HOL, but not the other HOL systems
- Supports some abstract algebra and number theory hierarchies
- Recent work from CMU translates type classes to ML functors, offering a way to map formal developments from Isabelle/HOL to HOL-Light

3

THE END



◆□▶▲圖▶▲≣▶▲≣▶ ▲国▼