A Machine-Assisted Proof of Gödel's Incompleteness Theorems

Lawrence C. Paulson, Computer Laboratory, University of Cambridge

The most misunderstood theorems in mathematics

- Gödel's theorems have highly technical, syntactic proofs.
 - 1. Every "strong enough" formal system is *incomplete*, in that at least one formula can neither be proved nor disproved.
 - 2. And if such a formal system admits a proof of its own consistency, then it is actually *inconsistent*.

- For the first time, *both* of Gödel's proofs have been mechanised, following a paper by Swierczkowski (2003)
- The machine proof, in the structured Isar language, is complete, almost readable, and can be perused interactively.

Hereditarily finite set theory

- * A *hereditarily finite set* is a finite set of HF sets.
- Many mathematical constructions, including natural numbers and sequences, can be defined as in standard set theory.
- * HF set theory is equivalent to Peano Arithmetic via the mapping

$$f(\mathbf{x}) = \sum \left\{ 2^{f(\mathbf{y})} \mid \mathbf{y} \in \mathbf{x} \right\}$$

Benefits of Using HF Set Theory

- Can use standard definitions of pairing and sequences.
- The first incompleteness theorem requires an HF development of the natural numbers, induction, etc., but not addition.

- The second incompleteness theorem requires operations on sequences and addition, but not multiplication.
- No need for least common multiples, prime numbers or the Chinese remainder theorem.

The Axioms of HF Set Theory

 $z = 0 \leftrightarrow \forall x [x \notin z]$ $z = x \triangleleft y \leftrightarrow \forall u [u \in z \leftrightarrow u \in x \lor u = y]$ $\phi(0) \land \forall xy [\phi(x) \land \phi(y) \rightarrow \phi(x \triangleleft y)] \rightarrow \forall x [\phi(x)]$

- 0 denotes the empty set
- * $x \triangleleft y$ denotes the set *x* extended with the element *y*.
- There are no other function symbols.
- * Union, intersection, etc can be shown to exist by induction.

Stages of the Proofs

- The *syntax* of a first-order theory is formalised: terms, formulas, substitution...
- A *deductive calculus* for sequents of the form Γ ⊢ α (typically for Peano arithmetic, but here HF)
- Meta-theory to relate truth and provability. E.g. "all true Σ formulas are theorems". Σ formulas are built using ∨ ∧ ∃ and bounded ∀.

- A system of coding to formalise the calculus within itself. The code of α is a term, written ¬α¬.
- Syntactic predicates to recognise codes of terms, substitution, axioms, etc.
- * Finally the predicate Pf, such that $\vdash \alpha \iff \vdash Pf \vdash \alpha \neg$.

First Incompleteness Theorem

- To prove Gödel's *first* incompleteness theorem, construct δ that expresses that δ is not provable.
- It follows (*provided* the calculus is consistent) that neither δ nor its negation can be proved.
- Need to show that substitution behaves like a function.
 - Requires a detailed proof in the calculus,
 - * ... alternatively, other detailed calculations.

Second Incompleteness Theorem

If α is a Σ sentence, then $\vdash \alpha \to \operatorname{Pf} \ulcorner \alpha \urcorner$.

- * This crucial lemma for Gödel's *second* incompleteness theorem is proved by induction over the construction of α as a Σ formula.
- * It requires generalising the statement above to allow the formula α to contain free variables.
 - complex technicalities
 - lengthy deductions in the calculus

Proving Theorems in the Calculus

- Gödel knew that formal proofs were difficult. They could be eliminated, but at what cost?
- By coding all predicates as executable functions, and proving a meta-theorem, Gödel reduced provability to truth.
- But then only bounded quantifiers can be used, with tricky arithmetical proofs that the bounds are adequate.

- With Σ formulas, provability is reduced to truth for most desired properties, with no tricky proofs about bounds.
- Instead, some straightforward inductions need to be formalised in the calculus.
- The second theorem requires working in the calculus anyway.

Isabelle/HOL and Nominal

- a proof assistant for higherorder logic
- much automation to hide the underlying proof calculus
- support for recursive functions and inductive sets
- the *nominal package*, for working with named variables

- Free names are significant, but not once they are bound.
- Syntax involving variable binding can be defined using recursion, provided variables are used "sensibly".
- During proof by induction, bound variable names can be guaranteed not to clash with specified other terms.

De Bruijn Indexes

- * This approach to variable binding replaces names by numbers.
- * 0 denotes the innermost bound variable, 1 for the next, etc.
- This approach destroys readability, but substitution and abstraction are very easy to define.
- During coding, formulas are translated into the de Bruijn format.
- * And so there is no need to formalise the nominal theory within HF.

Defining Terms and Formulas

nominal_datatype tm = Zero | Var name | Eats tm tm

nominal_datatype fm =

- Mem tm tm (infixr "IN" 150) | Eq tm tm (infixr "EQ" 150) | Disj fm fm (infixr "OR" 130) | Neg fm
- | Ex x::name f::fm binds x in f

Variable binding formalised using nominal

Defining Substitution

nominal_primrec subst :: "name ⇒ tm ⇒ tm ⇒ tm"
where
 "subst i x Zero = Zero"
 | "subst i x (Var k) = (if i=k then x else Var k)"
 | "subst i x (Eats t u) = Eats (subst i x t) (subst i x u)"

```
nominal_primrec subst_fm :: "fm \Rightarrow name \Rightarrow tm \Rightarrow fm"
where
Mem: "(Mem t u)(i::=x) = Mem (subst i x t) (subst i x u)"
| Eq: "(Eq t u)(i::=x) = Eq (subst i x t) (subst i x u)"
| Disj: "(Disj A B)(i::=x) = Disj (A(i::=x)) (B(i::=x))"
| Neg: "(Neg A)(i::=x) = Neg (A(i::=x))"
| Ex: "atom j \ddagger (i, x) \Longrightarrow (Ex j A)(i::=x) = Ex j (A(i::=x))"
```

The variable *j* must be fresh for *i* and *x*

Properties of substitution have simple proofs.

Defining the HF Calculus

```
inductive hfthm :: "fm set \Rightarrow fm \Rightarrow bool" (infixl "\vdash" 55)
  where
     Hyp: "A \in H \implies H \vdash A"
   | Extra: "H \vdash extra_axiom"
   | Bool: "A \in boolean_axioms \implies H \vdash A"
    \  \  \, I \  \, Eq: \qquad "A \ \in \  \, equality\_axioms \ \Longrightarrow \  \, H \  \, \vdash \  \, A"
   | Spec: "A \in special_axioms \implies H \vdash A"
   | HF: "A \in HF_axioms \implies H \vdash A"
   | Ind: "A \in induction_axioms \implies H \vdash A"
   / MP: \qquad "H \vdash A IMP B \implies H' \vdash A \implies H \cup H' \vdash B"
   | Exists: "H \vdash A IMP B \Longrightarrow
               atom i \ddagger B \implies \forall C \in H. atom i \ddagger C \implies H \vdash (Ex i A) IMP B"
       The variable i must be
           fresh for B and H
```

Early Steps in the HF Calculus

- the *deduction theorem* (yielding a sequent calculus)
- *derived rules* to support explicit formal proofs
 - * for defined connectives, including $\land \rightarrow \forall$
 - * for equality, set induction, ...
- definitions and proofs for *subsets*, *extensionality*, *foundation* and natural number *induction*

S Formulas

Strict Σ formulas only contain variables and are the basis for the main induction of the second incompleteness theorem. We can still derive the general case of Σ formulas.

inductive ss_fm :: "fm \Rightarrow bool" where MemI: "ss_fm (Var i IN Var j)" DisjI: "ss_fm A \implies ss_fm B \implies ss_fm (A OR B)" ConjI: "ss_fm A \implies ss_fm B \implies ss_fm (A AND B)" ExI: "ss_fm A \implies ss_fm (Ex i A)" All2I: "ss_fm A \implies atom j \ddagger (i,A) \implies ss_fm (All2 i (Var j) A)"

"Sigma_fm A \longleftrightarrow (\exists B. ss_fm B & supp B \subseteq supp A & {} \vdash A IFF B)"

theorem "[Sigma_fm A; ground_fm A; eval_fm e0 A]] \implies {} \vdash A"

True Σ formulas are theorems!

Coding Terms and Formulas

must first translate from nominal to de Bruijn format

* the actual coding is a simple recursive map:

*
$$r 0 = 0$$
, $r x_k = k$, $r x \triangleleft y = \langle r \triangleleft , r x \neg, r y \neg \rangle$, ...

also define (in HF) predicates to recognise codes

- *abstraction* over a variable needed to define Form(*x*), the predicate for formulas
- *substitution* for a variable

Example: Making a Formula

definition MakeForm :: "hf \Rightarrow hf \Rightarrow hf \Rightarrow bool" where "MakeForm y u w \equiv y = q_Disj u w \lor y = q_Neg u \lor (\exists v u'. AbstForm v 0 u u' \land y = q_Ex u')"

> $y = u \lor w$, or $y = \neg u$, or $y = (\exists v) u$ with an *explicit* abstraction step on u

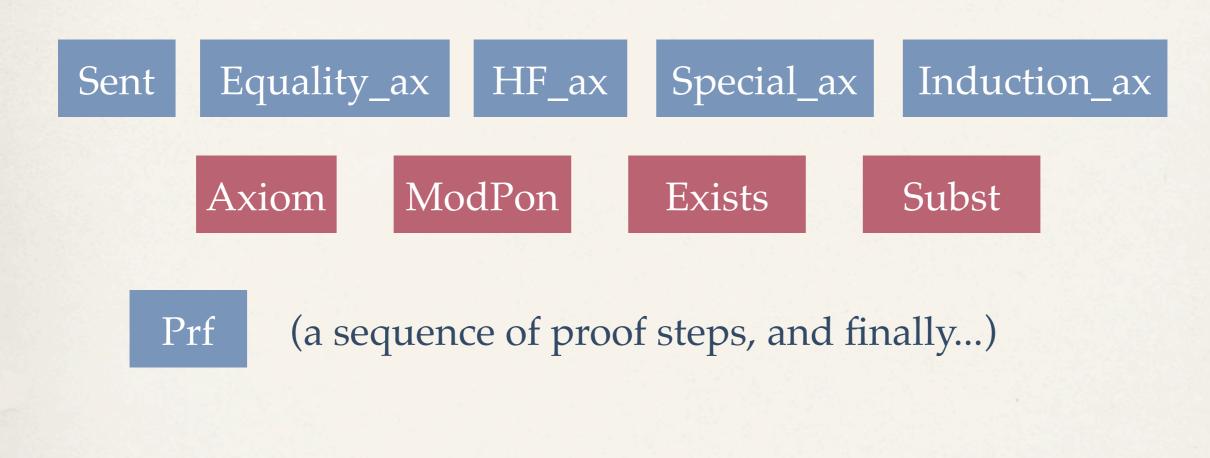
nominal_primrec MakeFormP :: "tm \Rightarrow tm \Rightarrow tm \Rightarrow fm"
where "[atom v \$\\$ (y,u,w,au); atom au \$\\$ (y,u,w)] \Longrightarrow MakeFormP y u w =
 y EQ Q_Disj u w OR y EQ Q_Neg u OR
 Ex v (Ex au (AbstFormP (Var v) Zero u (Var au) AND y EQ Q_Ex (Var au)))"

The "official" version as an HF formula, not a boolean

Those Coding Predicates



... And Proof Predicates



Steps to the First Theorem

- * We need a function *K* such that $\vdash K(\ulcorner \phi \urcorner) = \ulcorner \phi(\ulcorner \phi \urcorner) \urcorner$
- * ... but we have no function symbols. Instead, define a relation, KRP:
 lemma prove_KRP: "{} ⊢ KRP 「Var i┐ 「A┐ 「A(i::=「A┐)¬"
- Proving its functional behaviour takes 600 HF proof steps.
 lemma KRP_unique: "{KRP v x y, KRP v x y'} ⊢ y' EQ y"
- Finally, the diagonal lemma:

lemma diagonal: obtains δ where "{} $\vdash \delta$ IFF $\alpha(i::=\lceil \delta \rceil)$ " "supp δ = supp α - {atom i}"

```
theorem Goedel_I:
  assumes Con: "\neg {} \vdash Fls"
  obtains \delta where "{} \vdash \delta IFF Neg (PfP \lceil \delta \rceil)"
                       "¬ {} \vdash \delta" "¬ {} \vdash Neg \delta"
                       "eval_fm e \delta " "ground_fm \delta "
proof -
  obtain \delta where "{} \vdash \delta IFF Neg ((PfP (Var i))(i::=\lceil \delta \rceil))"
              and [simp]: "supp \delta = supp (Neg (PfP (Var i))) - {atom i}"
    by (metis SyntaxN.Neg diagonal)
  hence diag: "{} \vdash \delta IFF Neg (PfP \lceil \delta \rceil)"
    by simp
  hence np: "\neg \{\} \vdash \delta"
     by (metis Con Iff_MP_same Neg_D proved_iff_proved_Pf)
  hence npn: "\neg {} \vdash Neg \delta" using diag
     by (metis Iff_MP_same NegNeg_D Neg_cong proved_iff_proved_Pf)
  moreover have "eval_fm e \delta" using hfthm_sound [where e=e, OF diag]
     by simp (metis Pf_quot_imp_is_proved np)
  moreover have "ground_fm \delta "
     by (auto simp: ground_fm_aux_def)
  ultimately show ?thesis
     by (metis diag np npn that)
qed
```

Steps to the Second Theorem

* Coding must be generalised to allow *variables* in codes.

*
$$[x \triangleleft y]_V = \langle [\triangleleft \neg, [x \neg, [y \neg] \rangle] \rangle$$
 codes of variables
* $[x \triangleleft y]_V = \langle [\triangleleft \neg, x, y \rangle$ are integers

- * Variables must be renamed, with the intent of creating "pseudo-terms" like ⟨¬ ⊲ ¬, Q x, Q y⟩.
- * Q is a magic function: Q x = r t r where *t* is some canonical term denoting the set *x*.

Complications

- * Q must be a *relation*.
 - Function symbols cannot be added...
 - * Sets do not have an easily defined canonical ordering.
- * QR(0,0)
- * QR(x,x'), QR(y,y') \Longrightarrow QR($x \triangleleft y, \langle \neg \neg, x', y' \rangle$)

One of the Final Lemmas

 $QR(x, x'), QR(y, y') \vdash x \in y \to Pf [x' \in y']_{\{x', y'\}}$ $QR(x, x'), QR(y, y') \vdash x \subseteq y \to Pf [x' \subseteq y']_{\{x', y'\}}$ $QR(x, x'), QR(y, y') \vdash x = y \to Pf [x' = y']_{\{x', y'\}}$

- * The first two require simultaneous induction, yielding the third.
- * Similar proofs for the symbols $\vee \wedge \exists$ and bounded \forall .
- * The proof in the HF calculus needs under 450 lines.
- Fills a major gap in various presentations, including Swierczkowski's.

theorem Goedel_II:

```
assumes Con: "\neg {} \vdash Fls"
```

```
shows "\neg {} \vdash Neg (PfP \lceilFls\rceil)"
```

proof -

from Con Goedel_I obtain δ

```
where diag: "{} \vdash \delta IFF Neg (PfP \lceil \delta \rceil)" "\neg {} \vdash \delta"
```

and gnd: "ground_fm δ "

by metis

have "{PfP $\lceil \delta \rceil$ } \vdash PfP $\lceil PfP \lceil \delta \rceil$ "

by (auto simp: Provability ground_fm_aux_def supp_conv_fresh) moreover have $"{PfP \lceil \delta \rceil} \vdash PfP \lceil Neg (PfP \lceil \delta \rceil) \rceil"$

apply (rule MonPon_PfP_implies_PfP [OF _ gnd])

apply (auto simp: ground_fm_aux_def supp_conv_fresh) using diag

by (metis Assume ContraProve Iff_MP_left Iff_MP_left' Neg_Neg_iff) moreover have "ground_fm (PfP $\lceil \delta \rceil$)"

by (auto simp: ground_fm_aux_def supp_conv_fresh) ultimately have "{PfP $\lceil \delta \rceil$ } \vdash PfP $\lceil Fls \rceil$ " using PfP_quot_contra

by (metis (no_types) anti_deduction cut2)

thus " \neg {} \vdash Neg (PfP \lceil Fls \rceil)"

by (metis Iff_MP2_same Neg_mono cut1 diag) qed

What Did We Learn?

- * Some highly compressed proofs were finally made explicit.
- * The entire proof development can be examined interactively.
- * The nominal package can cope with very large developments...

(BUT: performance issues, some repetitive notation, complications in accepting function definitions)

- * <9 months for the first theorem, a further 4 for the second</p>
- * Under 16 000 lines of proof script in all.

Conclusions

- the first-ever machine formalisation of Gödel's second incompleteness theorem
- * using both nominal and de Bruijn syntax for bound variables
- within an axiom system for hereditarily finite set theory
- * conducted using Isabelle/HOL.

Many thanks to Christian Urban for help at critical points!

Also Jesse Alama, Dana Scott.