

**Demonstrating the Undecidability and
Incompleteness of Peano Arithmetic with a
Higher Order Logic Theorem Prover**

**Robert T. Bauer
Dr. Raymond J. Toal**

12/3/1997

**Department of Electrical Engineering
and Computer Science
College of Science and Engineering
Loyola Marymount University**

Acknowledgments

The work presented in this paper could not have been accomplished without Dr. Ray Toal's help and encouragement; however, it is important to recognize that Loyola Marymount University has been very kind to allow me to pursue a course of study dominated by theoretical topics in Computer Science, Metamathematics, Logic, and Distributed Computing.

Dr. Ray Toal, Dr. Phil Dorin, and Dr. Tony Karrer have each provided me with an example of excellence, dedication, and scholarship; to each of them, I offer my deepest and most humble respect.

In no particular order, I would like to acknowledge the following persons: Lisa and Mike Malgeri for their friendship and Ray Bala II for showing that excellence needs no loudspeaker. I especially thank Albert James Ingram (age twelve) who stated "A number is a thing that you use to count other things." I hope that one day he reads this work and comes to understand both the limitations and significance of his statement.

Finally, a personal statement: In the course of this work, I have walked in the footsteps of great intellectual achievement and have been protected from the vastness of confusion by the shadow of genius. I have chased the diagonal, touched the edge of infinity, and experienced the awe of understanding. I cannot claim the clarity of vision nor the depth of understanding of Galileo, Cantor, Dedekind, Frege, Russell, Church, Turing, Godel, or Milner; however, I am satisfied that this work will help others on their journey to understanding.

Table of Contents

1. INTRODUCTION.....	2
2. TYPE THEORY	4
2.1. HOL	4
2.1.1. <i>Defining New Logical Types</i>	5
2.1.2. <i>Tactics, Tacticals and Backward Proofs</i>	7
2.2. ANDREW'S SYSTEM Q	9
2.2.1. <i>Primitive Basis of Q_0</i>	9
2.2.2. <i>Axioms of Q_0</i>	9
2.2.3. <i>Elementary Logic in Q_0</i>	10
2.2.4. <i>Formalized Number Theory in Q_0</i>	11
2.3. COMPARING FORMALIZED NUMBERS IN HOL AND Q_0	14
3. GODEL'S FIRST INCOMPLETENESS THEOREM.....	16
3.1. GODEL'S SENTENCE G	16
3.2. REPRESENTING PROOF	18
3.3. REPRESENTING AXIOMS	20
3.3.1. <i>Axiom 1: $g_{\alpha\alpha}T_\alpha \wedge g_{\alpha\alpha}F_\alpha = !x_\alpha . g_{\alpha\alpha}x_\alpha$</i>	20
3.3.2. <i>Axiom 2: $[X_\alpha = Y_\alpha] \supset . h_{\alpha\alpha}X_\alpha = h_{\alpha\alpha}Y_\alpha$</i>	20
3.3.3. <i>Axiom 3: $f_{\alpha\beta} = g_{\alpha\beta} = !X_\beta . f_{\alpha\beta}X_\beta = g_{\alpha\beta}X_\beta$</i>	22
3.3.4. <i>Axiom 4: $[\lambda x_\alpha . B_\beta]A_\alpha = B_\beta$</i>	22
3.3.5. <i>Axiom 5: $[\lambda x_\alpha . x_\alpha]A_\alpha = A_\beta$</i>	24
3.3.6. <i>Axiom 6: $[\lambda x_\alpha . B_\beta C_\gamma]A_\alpha = [([\lambda x_\alpha . B_\beta]A_\alpha)]([\lambda x_\alpha . C_\gamma]A_\alpha)$</i>	24
3.3.7. <i>Axiom 7: $[\lambda x_\alpha . \lambda y_\gamma . B_\beta]A_\alpha = [\lambda y_\gamma . ([\lambda x_\alpha . B_\beta]A_\alpha)]$</i>	24
3.3.8. <i>Axiom 8: $[\lambda x_\alpha . \lambda x_\alpha . B_\beta]A_\alpha = [\lambda x_\alpha . B_\beta]$</i>	24
3.3.9. <i>Axiom 9: $[!_{\alpha(\alpha)} [Q_{\alpha\alpha} y] = y_\alpha]$</i>	25
3.4. REPRESENTING RULES OF INFERENCE.....	25
4. DEMONSTRATING GIT IN HOL	26
5. CONCLUSION.....	29
5.1. SUMMARY	29
5.2. PROBLEMS ENCOUNTERED	29
5.3. SUGGESTIONS FOR FURTHER WORK	30
6. REFERENCES.....	31

Abstract

Godel's Incompleteness Theorem (GIT) expresses the undecidability of any formal system that is a recursively axiomatizable extension of First Order Peano Arithmetic. The goal of this work is to demonstrate a proof of GIT using the Higher Order Logic (HOL) Theorem Prover. Since HOL is based on a variant of Church's Type Theory, we provide a comprehensive background in type theory and HOL. We then present a backwards demonstration and proof of GIT in type theory. Finally we take initial steps towards the demonstration of GIT in the HOL Theorem Prover.

1. Introduction

Godel's First Incompleteness Theorem stands as the most important 20th Century landmark in Mathematical Logic. In 1931, Kurt Godel [9] demonstrated that any reasonably powerful formal system of axioms and inference rules contained true sentences that could not be proved – i.e., the formal system was incomplete. Godel constructed a true sentence that could not be proved and whose complement (by definition, false) could not be refuted. Within the formal system, this sentence is undecidable. In a relatively short paper, Kurt Godel crushed all hopes of using *algorithms* to correctly identify whether a given formal statement was a theorem in a given formal theory of mathematics.

Following Smullyan [18], consider a computer program that examines sentences. If the sentence is true, the computer program prints the sentence. Given a finite alphabet, we could generate all of the sentences of one symbol, then two symbols, etc., – and our program would only print those that were true. Given enough time, our program enumerates all truths within our language. To test whether a sentence represents a true statement, we simply count the number of symbols appearing in the sentence, index to the appropriate part of the print-out (or wait for the sentence to be printed, confirming the old adage that “truth comes to those who wait”) and then find the sentence being examined or determine that it is not (will never be) printed. If the sentence appears on the print-out, it represents truth, otherwise it represents falsehood.

Let P mean “printable” in the following sense. Suppose that some symbol, say Z has the value “false,” then our program will not print Z ; however, $P Z$ is also not printable because it is not true that Z is printable ($P Z = \text{false}$). Let \sim represent negation and note that $\sim P Z$ will be printed (because $P Z$ is not printable) and also that $\sim Z$ will be printed. Let “(“ and “)” be used to group symbols, e.g., $\sim P (P Z)$ – a true statement, meaning that $P Z$ is not printable. Other (true) groupings include $P \sim (P Z)$ and $P (\sim P Z)$. Let $N(X)$ stand for $X(X)$. So, for example, $P N (P N)$ is printable if $N(P N)$ is true and $N(P N)$ is $P N(P N)$; therefore, if we interpret $N(P N)$ as true, then $P N(P N)$ is printed, otherwise it is not. Now consider the sentence $\sim P N(\sim P N)$. Suppose we interpret $N(\sim P N)$ as true, then $\sim P N(\sim P N)$ is not printed; however, $N(\sim P N)$ is $\sim P N(\sim P N)$ – resulting in a true sentence that is not printed. Suppose we interpret $N(\sim P N)$ as false. Then $\sim P N(\sim P N)$ is printed; however, $N(\sim P N)$ is $\sim P N(\sim P N)$ – resulting in a false sentence that is printed. If our program is consistent, it cannot print a false sentence; therefore, we must admit that while our machine only prints “true” sentences, there are true sentences that it does not print.

A formal system is said to be complete if and only if every “true” sentence is derivable from the axioms. It is said to be sound if and only if every sentence derivable from the axioms is “true.” The formal system is said to be consistent if only if every sentence derivable from the axioms implies that its complement is not derivable from the axioms. The sentence above, shows that our program is either incomplete (a true sentence not printed) or inconsistent (a false sentence that is printed). Now consider the sentence $P N(\sim P N)$.

Suppose $N(\sim P N)$ is true, then $P N(\sim P N)$ is printed; however $N(\sim P N)$ is $\sim P N(\sim P N)$ – which if true, means that the sentence we just printed is false, and if false, means we should not have printed the sentence. Suppose that $N(\sim P N)$ is false, then $P N(\sim P N)$ is not printed; but $N(\sim P N)$ is $\sim P N(\sim P N)$ and if false, means that we did not print a true sentence. $P N(\sim P N)$ is an example of an undecidable sentence: A formal system is said to be decidable if and only if, for any sentence we can determine whether it or its complement is derivable from the axioms.

The goal of this project is to formulate Godel's First Incompleteness Theorem within the framework of the Higher Order Logic (HOL) theorem prover. The significance of this effort

is entirely technical; yet, it lays the foundation for further metamathematical work within the HOL system. In preparation of this effort, we present a Type-theoretic version of Godel's First Incompleteness Theorem. Whereas the traditional presentation generally proceeds as a forward proof, our presentation follows the spirit of HOL and presents a backward proof. As we work backwards through the proof, we identify the requirements for new tacticals and/or the selection of existing tacticals.

Although Andrews [1] does a very fine job of describing type theory including the development Godel's First and Second Incompleteness Theorems, we relied on Smullyan [18], Nagel [16], Shankar [17], Davis [7], and Godel [9] to fully understand and appreciate Godel's work. Throughout the backwards proof, we fill in many details that will not be found in Andrews; thus, this work may be considered a definitive, type-theoretic exposition of Godel's First Incompleteness Theorem.

The remainder of this work is organized as follows: Section 2 provides an overview of Higher Order Logic and especially the definition of new types and a detailed description of tactics. Section 2 also provides an overview of Andrew's type theory (Q) including the development of natural numbers (Peano Axioms) and Axioms of Infinity. We also compare and contrast Peano Arithmetic as expressed in HOL. In Section 3, we illustrate a type-theoretic proof of Godel's Incompleteness Theorem. Section 4 outlines Godel's Incompleteness in terms of HOL. Lastly, Section 5 presents our conclusions and provides ideas and directions for further research.

2. Type Theory

Church's λ -calculus is a type free theory about functions as rules rather than as graphs. Dirichlet is usually credited with the important idea that functions could be considered as sets of pairs of argument and value (extending Descartes' notion of pairs), but the λ -calculus retains the notion of functions as rules to stress the computational aspects of functions [2]. λ -calculus was developed to provide a general theory about functions and their manipulation. It was hoped that it would also help provide a foundation for (parts of) mathematics.

While λ -calculus was successful in describing functions, all attempts to provide a foundation in mathematics failed. Typed λ -calculus was introduced in order to prevent the various paradoxes (e.g., Russell's) from demonstrating the inconsistencies of these mathematical systems. All "type" systems, including typed λ -calculus, allow new types to be specified:

- The set of types are inductively defined as follows.
 1. There exists a base type (in Church's original paper, the base types consisted of propositional and individual corresponding to Andrew's individuals and truth (boolean); in HOL, only individuals are in the base type, all other types are derived).
 2. If σ, τ are types, then $\sigma \rightarrow \tau$ (meaning a function mapping type σ to type τ is also a type).
- Thus, all functions, variables, and constants (e.g., $+$) are strongly typed. For example, the "plus" constant for adding two objects of type number is specified as $+_{\sigma \rightarrow \sigma}$ where σ means number.

$(\alpha\beta)$ means $\beta \rightarrow \alpha$. Also, $(\alpha\beta\delta)$ stands for $\delta \rightarrow \beta \rightarrow \alpha$. Note that types are left associative. Thus, $\delta \rightarrow \beta \rightarrow \alpha$ and $\delta \rightarrow (\beta \rightarrow \alpha)$ are equivalent; however, $(\delta \rightarrow \beta) \rightarrow \alpha$ denotes a function that takes (as an argument) functions of the type $\delta \rightarrow \beta$ and returns an α . $\delta \rightarrow \beta \rightarrow \alpha$, on the other hand, denotes a function that takes a δ and returns a function that takes a β and returns an α .

In the remainder of this section, we discuss the type systems of HOL and Andrew's system Q. We conclude this section with a comparison of formalized numbers in HOL and Q.

2.1. HOL

Higher Order Logic (HOL) is called higher order because:

- Variables are allowed to range over functions and predicates;
- Functions can take functions as arguments and yield functions as results; and
- The notation of the λ -calculus can be used to write terms that denote functions.

In HOL, each variable has an associated logical type that specifies the kind of values it ranges over. Sets and even recursive structures (lists and trees) can be represented in HOL by extending the syntax of types. This is done by first defining these new types in terms of already existing types and then deriving properties about these new types by formal proof. The idea is that formal proof will guarantee that adding a new type will not introduce inconsistency.

The syntax of HOL includes terms that correspond to conventional predicate calculus – $P\ x$ expresses the proposition that x has property P . $R(x,y)$ means that relation R holds between

x and y. Examples of P include *Man Socrates* to express the notion that *Socrates* is human and $<(3,5)$ to denote $3 < 5$.

In HOL, T indicates truth, F falsity, ~ negation, \ / disjunction, /\ conjunction, ==> implication, = equality, ! universality, ? existential, @ hilbert's e operator, => conditional, and ?! unique existence. For example, a theorem of primitive recursion states that for any equation written inductively (base and induction statement), there exists a unique recursive function computing the same value:

$$\vdash \text{?F} . !x \text{ f} . ((F \text{ f}) 0 = x) /\ !n . (F \text{ f}) (n+1) = f ((F \text{ f}) n)$$

For any variable x, for any function f: function F, which takes a function as an argument and yields a function as a result, has the following properties. Given the base of the induction ($f\ 0 = x$), F takes f and 0 as arguments and yields x; thus, F is equivalent in the base case. And (\wedge) for all n (!n), F takes f and n+1 as arguments and computes the same value as F taking f and n as arguments, the result being applied to f (recursively).

HOL is a mechanized, interactive, proof-assistant that has been primarily used to reason about the correctness of digital hardware; however, much of what has been developed in HOL for hardware verification – the theory of arithmetic, is also fundamental for other applications...

HOL supports secure theorem proving by representing its logic in a ML, a strongly-typed functional programming language. Propositions and theorems of the logic are represented by ML abstract data types and interaction with the theorem prover takes place by executing ML procedures that operate on values of these data types. Users can write arbitrarily complex programs to implement proof strategies and because the logic is represented using ML abstract data types, the user-defined proof strategies are guaranteed to perform only valid logical inferences.

2.1.1. Defining New Logical Types

The use of types in higher order logic eliminates the potential for the inconsistency that comes with allowing higher order variables. For example, Russell's Paradox: $P\ x = \sim\ x\ x$, then substitute P for x, leading to $P\ P = \sim\ P\ P$, is avoided using typed variables.

In HOL, three steps are needed to define a new type:

1. find an appropriate subset of an existing type to represent the new type;
2. extend the syntax of logical types to include a new type symbol, and using a *type definition axiom* to relate this new type to its representation; and
3. deriving from the type definition axiom and the properties of the representing type a set of theorems that serves as an 'axiomatization' of the new type.

In the first step, a model for the new type is given by specifying a set of values – this is done by defining a predicate P on an existing type such that the set of values satisfying P denotes the properties that the new type is expected to have. In the second step, a new type constant or type operator is added to the logic – this *type definition axiom* serves to relate values of the new type to the corresponding values of the existing type that represents them.

In the last step, a collection of theorems is proved that abstractly characterizes the new type. Here "abstractly" means that essential properties of the new type are stated without reference to the way the values are represented; we use the term 'axiomatized' because the collection of theorems are derived from the definition of the subset predicate defining the set of values for the new type and from the *type definition axiom*. This last step gives a

consistency proof of the axioms (theorems?) for the new type by showing that there is a model for them.

Suppose that α is an existing type of the logic and $P:\alpha \rightarrow \text{bool}$ is a predicate on values of type α that defines the subset for the new type. Since the semantics of the hilbert e operator require all types to denote non-empty sets, it must be the case that " $\vdash ?x:\alpha. P x$ " and we must prove it before the type definition axiom can be added to the system.

The type definition axiom asserts the existence of a representation function that maps a value of the new type to the value of the "base" type that represents it:

$$\vdash ?f:\beta \rightarrow \alpha . (!a_1 a_2 . f a_1 = f a_2 \supset a_1 = a_2) \wedge (!r . P r = (?a . r = f a))$$

This axiom states that function f which maps from the new type (β) to the "base" type (α) is one-to-one and onto the subset defined by P . Note that new operators can be defined in a similar fashion. If t is of type (α, β) , then $P:t(\alpha, \beta) \rightarrow \text{bool}$ is the desired predicate and if $?x.P x$, then a new binary type operator $(\alpha, \beta)op$ is asserted:

$$\vdash ?f:(\alpha, \beta)op \rightarrow t(\alpha, \beta) . (!a_1 a_2 . f a_1 = f a_2 \supset a_1 = a_2) \wedge (!r . P r = (?a . r = f a))$$

To formulate the abstract axioms required in step three, it is convenient to have logical constants which in fact denote the isomorphism (and its inverse). A representation function that maps values in the new type to the value in the "base" type is defined:

$$\vdash \text{REP} = @f . (!a_1 a_2 . f a_1 = f a_2 \supset a_1 = a_2) \wedge (!r . P r = (?a . r = f a))$$

Once REP is defined, the e -operator can be used to defined the inverse abstraction function $\text{ABS}:\text{base_type} \rightarrow \text{new_type}$:

$$\vdash !r . \text{ABS } r = (@a . r = \text{REP } a)$$

From these definitions, it follows that ABS is the left inverse of REP and for REP is the left inverse of ABS :

$$\vdash !a . \text{ABS}(\text{REP } a) = a$$

$$\vdash !r . P r = (\text{REP}(\text{ABS } r) = r)$$

2.1.1.1.A Simple Type Definition

The simplest and smallest type possible in HOL is the type constant *one* denoting a set containing exactly one element. To represent one, any singleton subset of an existing type will do. Since T (truth within the type bool) is one of the two elements that are represented in the type bool , we can use the predicate $\lambda b:\text{bool}. b$ denoting the identity function on bool . This set is non-empty, $\vdash ?x. (\lambda b. b) x$ follows from $\vdash (\lambda b. b) T$ which is equivalent to $\vdash T$.

Now that we have established that $\lambda b. b$ specifies a non-empty set of booleans, the type constant *one* can be defined:

$$\vdash ?f:\text{one} \rightarrow \text{bool} . (!a_1 a_2 . f a_1 = f a_2 \supset a_1 = a_2) \wedge (!r . (\lambda b. b) r = (?a . r = f a))$$

and we specify f as $\text{REP_one}:\text{one} \rightarrow \text{bool}$. The axiomatization of the type *one* consists of a single theorem:

$$\vdash !f:\alpha \rightarrow \text{one} . !g:\alpha \rightarrow \text{one} . (f = g)$$

This theorem says that all functions mapping values of type α to values of type *one* are equal. Thus, there is only one value of type *one* (otherwise $f \neq g$ at least once!). Note that we have specified an essential property of the type without reference to how the type is represented.

Of course we need to prove the theorem we have just stated:

$$\vdash !a_1 a_2 . \text{REP_one } a_1 = \text{REP_one } a_2 \supset a_1 = a_2$$

$$\vdash !r . r = (\text{?}a . r = \text{REP_one } a)$$

Specializing r in the second equation (if it is true for all r , then it is certainly true for any r) to $\text{REP_one}(f \ x)$ results in:

$$\vdash \text{REP_one}(f \ x) = (\text{?}a . \text{REP_one}(f \ x) = \text{REP_one } a)$$

Note that the right-hand side is \top , since our new type is non-empty, so that $\text{REP_one } a$ maps to \top any a which is of type *one*. The theorem simplifies to $\vdash \text{REP_one}(f \ x)$ and similar reasoning yields $\vdash \text{REP_one}(g \ x)$. Then since equal things are equal, we have $\vdash \text{REP_one}(f \ x) = \text{REP_one}(g \ x)$. From this and the first equation which states that REP_one is one-to-one, we have $\vdash f \ x = g \ x$ and therefore $\vdash !f \ g . (f = g)$.

Finally we note that once we prove the axiom about *one*, we can easily demonstrate that there is only one value of type *one*. Here's the definition of a constant called **one**:

$$\vdash \text{one} = @x:\text{one} . \top$$

From the previously proved axiom about *one*, we have " $\vdash \lambda x:\alpha . v = \lambda x:\alpha . \text{one}$ " – resulting in " $\vdash v = \text{one}$ ". Since v is not free in **one**, we generalize:

$$\vdash !v:\text{one} . v = \text{one}$$

stating that every value v of type *one* is equal to the constant **one**.

2.1.2. Tactics, Tacticals and Backward Proofs

HOL uses sequents to express logical constructs. A sequent is a pair (Γ, t) where Γ is a finite set of formulae called the set of assumptions and t is a single formula called the conclusion. A deductive system D is a set of pairs $(L, (\Gamma, t))$ where L is a possibly empty list of sequents. A sequent follows from a set of sequents, Δ , if and only if there exist sequents $(\Gamma_1, t_1), \dots, (\Gamma_n, t_n)$ such that:

1. $(\Gamma, t) = (\Gamma_n, t_n)$; and
2. for all i such that $1 \leq i \leq n$:
 - $(\Gamma_i, t_i) \in \Delta$; or
 - $(L_i (\Gamma_i, t_i)) \in D$ for some list of L_i of members of $\Delta \cup \{(\Gamma_1, t_1), \dots, (\Gamma_n, t_n)\}$.

The sequence $(\Gamma_1, t_1), \dots, (\Gamma_n, t_n)$ is called a proof of (Γ, t) from Δ with respect to (model) D . As is usual, HOL specifies its deductive system as a number of (schematic) rules of inference:

- Assumption Introduction – Assume t results in the sequent " $t \vdash t$ ".
- Reflexivity – " $\vdash t = t$ ".
- Beta-conversion – " $\vdash t (\lambda x . t) u = t[u/x]$ "; substitute u for x .
- Substitution – " $\Gamma_1 \vdash t_1 = u_1, \dots, \Gamma_n \vdash t_n = u_n, \Gamma \vdash t[t_1, \dots, t_n] \rightarrow \Gamma_1 \cup \dots \cup \Gamma_n \cup \Gamma \vdash t[t_1, \dots, t_n]$ ".
- Abstraction – " $\Gamma \vdash t = u \rightarrow \Gamma \vdash (\lambda x . t) = (\lambda x . u)$ " provided x is not free in Γ .
- Type Instantiation – " $\Gamma \vdash t \rightarrow \Gamma \vdash t[\alpha_1, \dots, \alpha_n / \beta_1, \dots, \beta_n]$ " where $t[\alpha_1, \dots, \alpha_n / \beta_1, \dots, \beta_n]$ is the result of substituting, in parallel, the types α for the type variables β with two conditions:
 1. None of the type variables, β , occur in Γ ; and

2. No distinct variables in t become identified after the instantiation. For example, let $F(v)$ stand for $\exists t.(t \neq v)$. Clearly, by $F(t)$ we do not mean the incorrect formula $\exists t.(t \neq t)$, but rather $\exists z.(z \neq t)$.

- Discharging an Assumption – “ $\Gamma \vdash t \rightarrow \Gamma - \{u\} \vdash u \rightarrow t$ ”.
- Modus Ponens – “ $(\Gamma_1 \vdash t_1 \rightarrow t_2, \Gamma_2 \vdash t_1) \rightarrow \Gamma_1 \cup \Gamma_2 \vdash t_2$ ”.

All provable statements of HOL may be derived from the rules of inference given above and the five axioms listed below:

- $\vdash \forall b.(b = T) \vee (b = F)$
- $\vdash \forall b_1 b_2.(b_1 \rightarrow b_2) \rightarrow (b_2 \rightarrow b_1) \rightarrow (b_1 = b_2)$
- $\vdash \forall f_{\alpha \rightarrow \beta} . (\lambda x. f x) = f$
- $\vdash \forall P_{\alpha \rightarrow \text{bool}} x. P x \rightarrow P (@ P)$
- $\vdash \exists f_{\text{ind_ind}} . \text{One_One } f \wedge \neg(\text{Onto } f)$

Note, these axioms are extended through the development of new theories (c.f., Section 3.5 where we detail the theory of numbers in HOL).

A theorem is a sequent that follows from the empty set of sequents; thus, a theorem is the last element of a proof from the empty set of sequents. A forward proof proceeds from the “axioms” to the theorem being proved. A backwards proof proceeds from the theorem back to the axioms.

Although one can perform forward proofs in HOL, the HOL theorem prover was designed to support backward (goal directed) proofs. A *tactic* in HOL is an ML function that when applied to a goal reduces it to a list of sub-goals and a justification function that maps a list of theorems to a theorem. For example, suppose the “goal” was to prove a simple goal such as “ $T \wedge \neg F$ ” is straight-forward.

In HOL there is a tactic called CONJ_TAC. It takes a conjunction, such as the goal specified above, and breaks it into two sub-goals and a justification function. The justification function takes a list of “values” and generates a theorem. Let Truth stand for “ $\vdash T$ ”. Note that “Truth” has no assumptions. Let Not_False stand for “ $\vdash \neg F$ ”. Again note that there are no assumptions. The use of the justification function for CONJ_TAC is shown below:

CONJ_TAC_JUSTIFICATION [Truth, Not_False] $\Rightarrow \vdash T \wedge \neg F$

Tactics are ML functions, so HOL users are free to create their own tactics. The ML type inference system guarantees soundness and the HOL system actually carries out the functions specified to ensure that proofs are evaluated.

A tactic solves a goal if it reduces the goal to the empty set of sub-goals. To do this requires at least one tactic that maps a goal to the empty sub-goal list. ACCEPT_TAC is a function that maps a goal to the empty list of sub-goals.

Tacticals differ from tactics. A tactical is an ML function that returns a tactic or list of tactics as a result. The most common example is the tactical ORELSE. If S and T are tactics, $S \text{ ORELSE } T$ evaluates to a tactic that applies S unless that fails. If S fails, ORELSE applies T . The tactical FIRST is similar, except that it applies the first tactic in a list of tactics that succeed.

2.2. Andrew's System Q

Andrew's type theory, like all type theoretic systems, is based on Church's typed λ -calculus where every individual, variable, function, and predicate has a type. For example, the function $+$ that operates on pairs of integers to yield their sum is denoted by $+_{\sigma\sigma}$ where σ is the type of integers. In Andrew's type theory, the type o denotes the type of truth values (i.e., boolean) so any function of type $(o\alpha)$ identifies a set of elements of type α or not type α .

2.2.1. Primitive Basis of Q_0

$\alpha, \beta, \dots, \zeta$, etc., are variables that range over type symbols. Type symbols are defined inductively:

o is a type symbol denoting the type of truth values¹.

i is a type symbol denoting the type of individuals.

If α and β are type symbol, then $(\alpha\beta)$ is a type symbol denoting the type of functions from elements of type β to elements of type α (e.g., $f:\beta \rightarrow \alpha$).

$Q_{\alpha\alpha}$ denotes the binary relation of equality. If A_o is a wff of type o in which the variable y_i occurs free, and if there is a unique element y_i of which A_o is true, we denote that element by $[y_i A_o]$ – meaning the y_i such that A_o . In type theory, $[A_\alpha = B_\alpha]$ abbreviates $[Q_{\alpha\alpha} A_\alpha B_\alpha]$ and the truth value True abbreviates $[Q_{ooo} = Q_{ooo}]$. Falsity, the truth value of False abbreviates $[\lambda x_o.T] = [\lambda x_o.x_o]$ (we define false as the statement “every individual of type truth is true,” since we can easily demonstrate the false truth value).

The logical constants of Q_0 are equality and selection:

- $Q_{(o\alpha)\alpha}$; and
- $I_{(i(o\alpha))}$

By wff $_\alpha$, we mean a wff of type α . We write A_α, B_α , etc., as variables ranging over wff's of type α . Wff's are defined inductively:

1. Any primitive variable or constant of type α is a wff $_\alpha$;
2. $[A_{\alpha\beta} B_\beta]$ is a wff $_\alpha$; and
3. $[\lambda x_\beta A_\alpha]$ is a wff $_{\alpha\beta}$.

2.2.2. Axioms of Q_0

1. $g_{ooo}.T_o \wedge g_{ooo}.F_o = !\chi_o . g_{ooo}\chi_o$ – meaning that truth and falsehood are the only truth values.
2. $[X_\alpha = Y_\alpha] \supset . h_{\alpha\alpha}X_\alpha = h_{\alpha\alpha}Y_\alpha$ – the results of function application are equal when the arguments are equal.
3. $(f_{\alpha\beta} = g_{\alpha\beta}) = !X_\beta . f_{\alpha\beta}X_\beta = g_{\alpha\beta}X_\beta$ – so that equal functions produce equal results.
4. $[\lambda x_\alpha B_\beta]A_\alpha = B_\beta$ where B_β is a primitive constant or variable distinct from x_α – showing that the entire range is mapped by the entire domain.
5. $[\lambda x_\alpha x_\alpha]A_\alpha = A_\alpha$ – λ identity function abstraction.

¹ Unfortunately Andrew's chose not to refer to this set as the set of boolean values. In previous sections we refer to "truth" values as boolean; however, since this is Andrew's theory, we keep his terminology.

6. $[\lambda x_\alpha . B_\beta C_\gamma] A_\alpha = [[\lambda x_\alpha B_\beta] A_\alpha] [[\lambda x_\alpha C_\gamma] A_\alpha]$ – substitution is simultaneous, even in the presence of composition.
7. $[\lambda x_\alpha . \lambda y_\gamma B_\beta] A_\alpha = [\lambda y_\gamma . [\lambda x_\alpha B_\beta] A_\alpha]$ where y_γ is distinct from x_α and from all variables in A_α .
8. $[\lambda x_\alpha . \lambda x_\alpha B_\beta] A_\alpha = [\lambda x_\alpha B_\beta]$
9. $\iota_{(0)} [Q_{ou} y_\gamma] = y_\gamma$

Q_0 has only one rule of inference: From C and $A_\alpha = B_\alpha$ to infer the result of replacing one occurrence of A_α in C by an occurrence of B_α , provided that the occurrence of A_α in C is not an occurrence of a variable immediately preceded by λ .

2.2.3. Elementary Logic in Q_0

The class of propositional wffs is the smallest class of wffs containing T_0 , F_0 , every variable of type 0 and if it contains A and B also contains:

- $[\sim A]$
- $[A \wedge B]$
- $[A \vee B]$
- $[A \rightarrow B]$
- $[A = B]$

Recall that the operators in the above sentential forms are abbreviations for the primitive constructions given in the previous section.

Variables of type 0 may be assigned truth values; thus if A represents a wff, then $\mathcal{V}_\phi A$ refers to the value of the wff A given some assignment of truth values to the variables of A . We now give the denotation of the logical basis of Q :

1. $\mathcal{V}_\phi T_0 = T$ (truth)
2. $\mathcal{V}_\phi F_0 = F$ (falsehood)
3. $\mathcal{V}_\phi p_0 = \phi p_0$ (the “value” of a propositional variable is the value it is assigned)
4. $\mathcal{V}_\phi (\sim A) = T$ if $\mathcal{V}_\phi (A) = F$ and $\mathcal{V}_\phi (\sim A) = F$ if $\mathcal{V}_\phi (A) = T$
5. $\mathcal{V}_\phi [A \wedge B] = T$ if $\mathcal{V}_\phi (A) = T$ and $\mathcal{V}_\phi (B) = T$ otherwise $\mathcal{V}_\phi [A \wedge B] = F$
6. $\mathcal{V}_\phi [A \vee B] = T$ if $\mathcal{V}_\phi (A) = T$ or $\mathcal{V}_\phi (B) = T$ otherwise $\mathcal{V}_\phi [A \vee B] = F$
7. $\mathcal{V}_\phi [A \rightarrow B] = T$ if $\mathcal{V}_\phi (A) = F$ or $\mathcal{V}_\phi (B) = T$ otherwise $\mathcal{V}_\phi [A \rightarrow B] = F$
8. $\mathcal{V}_\phi [A = B] = T$ if $(\mathcal{V}_\phi (A) = \mathcal{V}_\phi (B))$ otherwise $\mathcal{V}_\phi [A = B] = F$

Although we do not prove the soundness and completeness of Q with respect to elementary logic², we state the main conclusion: if A is a tautology ($\mathcal{V}_\phi (A) = T$ for all truth-value assignments ϕ), then A is a theorem of Q : $\text{Tautology}(A) \rightarrow \vdash A$.

² Andrew [1] is the obvious source for information regarding Q ; however, [12] is a very good source for metatheoretic results concerning propositional and first order logic.

2.2.4. Formalized Number Theory in Q_0

2.2.4.1. Definition of natural numbers

Two sets are said to be equivalent³ if the members of each set can be mapped one-to-one. The cardinal number of a set is "that which is common" to all equivalent sets. Thus, if x is equivalent to y , the cardinality of x must equal the cardinality of y : $x \approx y \rightarrow \bar{x} = \bar{y}$. It is interesting to note that set theory has no convenient way to express the notion of equivalence. In Q_0 equivalence is expressed:

$$[\lambda p_{\alpha\beta} \lambda q_{\alpha\alpha} \exists s_{\alpha\beta} \cdot \forall x_{\beta} [p_{\alpha\beta} x_{\beta} \rightarrow q_{\alpha\alpha} \cdot s_{\alpha\beta} x_{\beta}] \wedge \forall y_{\alpha} \cdot q_{\alpha\alpha} y_{\alpha} \rightarrow \exists ! x_{\beta} \cdot p_{\alpha\beta} x_{\beta} \wedge y_{\alpha} = s_{\alpha\beta} x_{\beta}]$$

which has the type $((\beta \rightarrow \alpha) \rightarrow (\alpha \rightarrow \alpha)) \rightarrow \alpha$ and which "asserts" the existence of a function s that establishes a one-to-one and onto relationship between the sets p and q . If the sets are not equivalent, s does not exist so that the lambda-expression is false.

Letting $E_{((\beta \rightarrow \alpha) \rightarrow (\alpha \rightarrow \alpha)) \rightarrow \alpha}$ stand for the equivalence relation, cardinal numbers are expressed:

$$[\lambda u_{\alpha(\alpha\beta)} \exists p_{\alpha\beta} \cdot u_{\alpha(\alpha\beta)} = E_{\alpha(\alpha\beta)(\alpha\beta)} p_{\alpha\beta}]$$

with the type $((\beta \rightarrow \alpha) \rightarrow \alpha) \rightarrow \alpha$. As defined earlier, a cardinal number is that which is common to all equivalent sets. If u is a cardinal number, then it represents such a commonality as follows: u is the cardinal number for some set (of β individuals) if there exists some other set of β individuals equivalent to the one referenced by u .

Having defined cardinal numbers, we wish to define ordinal numbers (so that we can say, for example, that one number is less than another). In traditional set theory, we would define zero as a symbol that associates with the empty set (i.e., $0 \equiv \{\}$). The symbol one would represent the set containing the empty set ($1 \equiv \{\{\}\}$). Clearly, only sets with a single element are equivalent to 1. Inductively, we define the set n to be the set containing the set zero, 1, ..., $n-1$ ($n \equiv \{\{\}, "1", "2", \dots "n-1"\}$). Where "1" means equivalence class of sets containing exactly one element.

That 2 is greater than 1, is argued as follows: The ordinal number 2 is a symbol associated with a set of sets equivalent to the set containing the empty set and the sets equivalent to 1. Since the set represented by 2 contains the set represented by 1, the set represented by 2 is larger than 1.

In type theory, the concept of ordinality is much simpler. The successor function in type theory is constructed around the notion that the cardinal number associated with some set A is a successor to the cardinal number associated with a set B , if by removing one element from set A we obtain a set, say A' with the same cardinality as that of B . With this notion, type theory easily expresses the notion of order.

The symbol σ is used to represent the equivalence classes of sets of individuals $(\iota \rightarrow \alpha) \rightarrow \alpha$ and may be interpreted as "number" throughout the remainder of this paper. The equivalence class of sets with zero members is used to represent zero:

$$0_{\sigma} \text{ stands for } Q_{\sigma(\sigma\iota)} [\lambda x_{\iota} F_{\sigma}]$$

The successor function rewrites a lambda expression to generate the successor:

$$S_{\sigma\sigma} \text{ stands for } [\lambda n_{\sigma(\sigma\iota)} \lambda p_{\sigma\iota} \exists x_{\iota} \cdot p_{\sigma\iota} x_{\iota} \wedge n_{\sigma(\sigma\iota)} [\lambda t_{\iota} \cdot t_{\iota} \neq x_{\iota} \wedge p_{\sigma\iota} t_{\iota}]]$$

Consider $S_{\sigma\sigma} S_{\sigma\sigma} 0_{\sigma}$ (representing the number 2):

³ The set-theoretic notion of equivalence is very different from ordinary usage of the term. Although one would like to use equinumerous instead of equivalence neither numbers nor equality has been defined!

$$[\lambda p_o. \exists x_i. p_o. x_i \wedge [\lambda q_o. \exists z_i. q_o. z_i \wedge Q_{o(o)(o)} [\lambda w_i. F_o] [\lambda u_i. u_i \neq z_i \wedge q_o. u_i]] [\lambda t_i. t_i \neq x_i \wedge p_o. t_i]]$$

$$[\lambda p_o. \exists x_i. p_o. x_i \wedge [\exists z_i. [\lambda t_i. t_i \neq x_i \wedge p_o. t_i] z_i \wedge Q_{o(o)(o)} [\lambda w_i. F_o] [\lambda u_i. u_i \neq z_i \wedge [\lambda t_i. t_i \neq x_i \wedge p_o. t_i] u_i]]]$$

$$[\lambda p_o. \exists x_i. p_o. x_i \wedge [\exists z_i. z_i \neq x_i \wedge p_o. z_i] \wedge Q_{o(o)(o)} [\lambda w_i. F_o] [\lambda u_i. u_i \neq z_i \wedge u_i \neq x_i \wedge p_o. u_i]]$$

The set p such that x is in p , z is in p , u is in p and the set u is empty and there is an x in p which is not the same as any z in p and neither z nor x are equivalent to the empty set. Clearly this represents the set with exactly two (distinct) elements. More precisely, the lambda expression represents the characteristic function of the equivalence relation: It is true precisely when applied to argument that is a set with two distinct elements.

The type theoretic representation of the natural numbers:

$$N_{oo} \text{ stands for } [\lambda n_o. \forall p_{oo}. [p_{oo} 0_o \wedge \forall x_o. p_{oo} x_o \rightarrow p_{oo} S_{oo} x_o] \rightarrow p_{oo} n_o]$$

If p is true for zero and if being true for some number it is also true of that number's successor, we say p is a property of numbers. If we consider the set of all such properties of numbers, then the set of numbers may be defined as the set of things that have all of those properties.

As discussed previously, cardinal numbers represent that which is common to equivalent sets. In this sense, when we write the symbol 2, we are using a sign to represent that notion of "in-common." Indeed, it would very inconvenient to write $Q[\lambda x. F]$ for zero or $SSQ[\lambda x. F]$ for 2. Since natural numbers as defined above are simply equivalence classes and because every equivalence class has a cardinality, then every natural number is a cardinal number:

$$\text{Finite}_{o(o)} \text{ stands for } [\lambda p_o. \exists n_o. N_{oo} n_o \wedge n_o p_o]$$

For any set, p , there exists a number n that represents "how many" things are in the set.

2.2.4.2. Peano Axioms

The axiomatization of the natural numbers is usually accomplished through the Peano Axioms:

- (a) Zero is a natural number.
- (b) If n is a natural number so is S_n (S_n is the successor of n , often written n').
- (c) If $S_m = S_n$, then $m = n$.
- (d) There is no number n for which $S_n = 0$.
- (e) Axiom of complete induction:

If a property P of the natural numbers satisfies the following two conditions, then P holds for every natural number:

- (1) P holds for 0.
- (2) For every natural number n , if P holds for n , then P holds for S_n .

Expressed in type theory, the Peano Axioms become:

$$(a) \vdash N_{oo} 0_o$$

Proof. By definition of N .

$$(b) \vdash \forall x_o. N_{oo} x_o \rightarrow N_{oo} S_{oo} x_o$$

Proof. Suppose the converse. Then, $\exists y . \neg N y \rightarrow N S y$, asserting the existence of an equivalence relation upon the successor but not the predecessor. By definition of the successor function the equivalence of the predecessor is established.

$$(c) \vdash \forall n_o \forall m_o . S_{oo} n_o = S_{oo} m_o \rightarrow n_o = m_o$$

Proof. Suppose the converse – that $S n = S m$ and $n \neq m$. Since “=” stands for Q_{ooo} and because the successor function S uses Q , we have the curious phenomena that the Q on the lhs is somehow different from the Q on the rhs.

$$(d) \vdash \forall n_o . S_{oo} n_o \neq 0_o$$

Proof. Consider the definition of the successor expression and the definition of zero. Zero represents all sets that are equivalent to the empty set. The successor function specifies that the successor set can be made equivalent to the predecessor set by removing an individual from the successor set. If zero is the successor of some number, then it is possible to remove an individual from the empty set.

$$(e) \vdash \forall p_{oo} . [p_{oo} 0_o \wedge \forall x_o . p_{oo} x_o \rightarrow p_{oo} . S_{oo} x_o] \rightarrow \forall x_o p_{oo} x_o$$

Proof. The definition of N says that if a property is valid for the equivalence relation with the empty set and if it is true for some equivalence relation and the successor of the equivalence relation, then it is true for a specific equivalence relation. The definition of complete induction says that given any property of natural numbers, if that property is valid for zero and given any natural number, that property is valid on the successor, then the property is valid for all natural numbers.

2.2.4.3. Axioms of Infinity

Infinity has played a very fundamental role in the modern development of mathematics and computer science. Intuitionism attacked the entire notion of infinity by claiming that exhaustive searches fail with infinite domains and that proof-by-contradiction cannot yield a constructive demonstration of existence in an infinite domain. Formalism in mathematics was a reaction to intuitionism – the assumption was that infinities were convenient idealizations and that conclusions based on infinities could be shown to be correct. Out of this formalism grew metamathematics where the whole scope of mathematics became the object of algorithms (symbols and syntactic rules for forming sentences and proofs).

Thus, from Dedekind, Cantor, Peano, Ferge, Hilbert, Whitehead & Russell, and finally to Kurt Godel, we travel from abstract notions of infinity to a concrete demonstration of undecidability which in its own way expresses an entirely different kind of infinity. Yet, it is almost startling to learn that Galileo was aware of the one-one correspondence between the set of natural numbers and the set of squares of the natural numbers.

Following Dedekind, it is still convenient to use the characteristic property that any infinite set is equivalent to at least one of its proper subsets. Finite sets can be defined as any set that cannot be placed in one-one correspondence with any of its proper subsets. The axiom used in type theory, on the other hand, is one that is a bit more natural: there is an irreflexive partial ordering of individuals with respect to which there is no maximal element:

$$\exists r_{ou} \forall x_i \forall y_i \forall z_i [\exists w_i . r_{ou} x_i w_i \wedge \neg r_{ou} x_i x_i \wedge . r_{ou} x_i y_i \rightarrow . r_{ou} y_i z_i \rightarrow r_{ou} x_i z_i]$$

When formalizing natural numbers, the axiom of infinity is not needed except to show:

$$\vdash \forall n_o \forall m_o . S_{oo} n_o = S_{oo} m_o \rightarrow n_o = m_o$$

2.3. Comparing Formalized Numbers in HOL and Q_0

HOL formalizes numbers as a theory that is derived from the theory of individuals. The theory of individuals is derived from the theory of booleans. In this section we quickly recap the theory of individuals (there are an infinite number of them) and the theory of booleans. We then compare and contrast how numbers are expressed in these two type-theoretic systems.

The theory of booleans has only implication (\rightarrow), equality ($=$), choice (which uses the symbol $@$ for Hilbert's ϵ -operator). Hilbert's ϵ -operator has already been introduced; however, we augment that introduction by noting that an HOL metatheorem states:

$$[\exists x . t \ x] = T \rightarrow [t \ @x . t \ x] = T$$

Our final remark concerning Hilbert's ϵ -operator is that if $!x . t \ x = F$, then $@x . t \ x$ denotes an unspecified value (of the proper type). According to the HOL documentation, $@x$ looks suspicious, but is very useful for naming things one knows to exist but for which one has no name!

The Hilbert ϵ -operator is equivalent to the Axiom of Choice (select the value of x such that x has the desired property) which when reformulated in HOL becomes $[@x . P \ x]$.

With only three logical constants, here is how HOL defines boolean expressions and abbreviates the "standard" boolean connectives:

- T stands for $((\lambda x . x) = (\lambda x . x))$ – Q defines T and F as primitives in the sense that T and F are the (only) values that may be assigned to boolean variables. So, in HOL, we have a precise notion of the value `True`.
- The binder $!$ (the HOL representation of \forall) stands for $[\lambda P . P = [\lambda x \ T]]$.
- The binder $?$ (the HOL representation of \exists) stands for $[@x . \lambda P . P \ x = T]$
- F stands for $[!b . b]$
- Negation (\sim) stands for $[\lambda b . b \rightarrow F]$
- And (\wedge) stands for $[\lambda x \ \lambda y . !z . (x \rightarrow (y \rightarrow z)) \rightarrow z]$
- Or (\vee) stands for $[\lambda x \ \lambda y . !z . (x \rightarrow z) \rightarrow ((y \rightarrow z) \rightarrow z)]$
- Equal ($=$) stands for $[\lambda x \ \lambda y . (x \rightarrow y) \wedge (y \rightarrow x)]$
- Unique existence ($?!$) $[\lambda P . ? P \wedge !x . !y . P \ x \wedge P \ y . \rightarrow . x = y]$

The axioms related to theory of booleans:

- $!b . (b = T) \vee (b = F)$
- $!x . !y . (x \rightarrow y) \rightarrow (y \rightarrow x) \rightarrow (x = y)$
- $!f . (\lambda x . f \ x) = f$
- $!P \ x . P \ x \rightarrow P (@ P)$

Other constants of the theory `bool` include `One_One`, `Onto`, and `Onto_Subset`, etc., some of these have already been discussed in regard to how new types are defined within HOL.

The theory of individuals extends the two element type `boolean`. If there are x distinct elements of type α and y distinct elements of type β , then there are y^x distinct elements of type $\alpha \rightarrow \beta$. Since there are infinitely many numbers, there is no hope of constructing a representing type for numbers from `bool` and implication (\rightarrow).

Individuals are a primitive type in HOL that contain infinitely many distinct elements. The axiom of infinity used in HOL is of a mapping from the class of individuals to the class of individuals such that the mapping is a function (One_One) but not Onto (does not map the entire range):

$$\text{?F:ind} \rightarrow \text{ind} . \text{One_One } F \wedge \neg(\text{Onto } F)$$

This axiom is derives the natural numbers as follows:

The representing type for zero is

$$\text{Zero_Rep} = @y:\text{ind} . !x . \neg(y = \text{Suc_Rep } x)$$

meaning zero is the “value” of the representing type (individuals) such that this value is not the value of the representing type for the successor of any other individual.

Then,

$$!x . \neg (\text{Zero_Rep} = \text{Suc_Rep } x)$$

So that, $1 = \text{Suc_Rep } (\text{Zero_Rep})$, $2 = (\text{Suc_Rep } (\text{Suc_Rep } (\text{Zero_Rep})))$, etc.

The proceeding demonstrates that in HOL numbers are represented by the subset of individuals (ind) consisting of Zero_Rep and all elements of the form $\text{Suc_Rep}^n \text{Zero_Rep}$. If we could define the function Is_Num as $[\lambda x . (x = \text{Zero_Rep}) \vee ?n . x = \text{Suc_Rep}^n \text{Zero_Rep}]$ it would be very convenient; however, the definition assumes that n is a number. In HOL

$$\text{Let Is_Num stand for } [\lambda x . !P . (P \text{ Zero_Rep}) \wedge (!y . P y \rightarrow P(\text{Suc_Rep } y)) \rightarrow P x]$$

Thus, we have

$$\text{Is_Num Zero_Rep; and}$$

$$!x . \text{Is_Num } x \rightarrow \text{Is_Num}(\text{Suc_Rep } x)$$

Therefore, the representing type for num is the class of individuals. The subset predicate is Is_Num. The subset is non-empty since Is_Num Zero_Rep . The representation function is $\text{Rep_Num: num} \rightarrow \text{ind}$.

Finally, in a manner reminiscent of our derivation of the Peano postulates in type theory, we develop them (Peano postulates) for HOL.

Define the abstract for numbers as $\text{Abs_Num} = \text{Inverse Rep_Num}$. Then,

- $0 = \text{Abs_Num Zero_Rep}$
- $\text{Suc} = \text{Abs_Num}(\text{Suc_Rep } (\text{Rep_Num}))$

The first axiom, zero is a natural number follows by definition: Is_Num Zero_Rep ;

The second axiom, if n is a number, then the successor of n is a number, follows from the definition of Suc.

The third axiom, if the successor of m equals the successor of n, then m and n are equal, can be proved based on the definition of One_One: $[\lambda f:\alpha \rightarrow b . !x . !y . (f x = f y) \rightarrow (x = y)]$.

The fourth axiom, zero is not the successor of any number, is the definition of zero in HOL.

The fifth axiom, complete induction, also follows from the definition of Is_Num:

$$[!P:\text{num} \rightarrow \text{bool} . P 0 \wedge (!m . P m \supset P (\text{Suc } m)) \supset !m . P m]$$

3. Godel's First Incompleteness Theorem

3.1. Godel's Sentence G

Previous sections have introduced the reader to basic concepts, notation, etc., of type theory and higher order logic. In this section, we start with an explanation of Godel's incompleteness theorem and work backwards through the proof. Our approach here is a bit unusual in that we begin with a type-theoretic version of Godel's sentence G and use it to introduce syntax and semantics:

Let G stand for $[\lambda x_{\sigma} \forall y_{\sigma} \neg \text{Proof}_{\sigma\sigma}^S [\text{"} * x_{\sigma} * [\text{Num}_{\sigma\sigma} x_{\sigma}] * \text{"}] y_{\sigma}]_{\sigma \rightarrow o}$.

Numbers have type σ and booleans have type o . Proof is a function that takes two numbers and returns a boolean value. The superscript S represents that *Proof* is specific to some formal system called S. The function Proof works as follows:

- Every sentence in the language is comprised of a sequence of one or more symbols from the alphabet of the language.
- Each sentence can be assigned a unique number called the Godel number of the sentence. We can determine whether an arbitrary number is a Godel number and we can recover the sentence encoded by the Godel number.
- Each unique series of sentences has a unique Godel number.
- A "proof" is a series of sentences with the following properties:
 1. Each sentence in the series either begins or ends a subseries in the sentence.
 2. Whether an arbitrary sentence is in the series is decidable.
 3. Each sentence in the series is either an axiom of the formal system or is the result of applying a rule of inference to one or more earlier sentences in the series.
 4. The last sentence in the series represents that which is being proved.
- Proof is true if and only if the second argument represents a "proof" whose last sentence is equivalent to the first argument.

The function proof contains the very powerful notion that an algorithm exists that can determine whether an arbitrary number represents the Godel encoding of a valid proof within the formal system being considered (in this case S). Indeed, Godel demonstrated just such an algorithm for Whitehead & Russell's Principia Mathematica.

The sequence of symbols $\text{"} * x_{\sigma} * [\text{Num}_{\sigma\sigma} x_{\sigma}] * \text{"}$ is quite complex: " represents the Godel number associated with the symbol $\{$. Note, we use this notation quite often: $\forall x . \text{Suc } x \neq 0$ means the Godel number associated with the sentence enclosed in quotes. Concatenation is an operation of juxtaposition denoted by $*$. The semantics of concatenation are straightforward: $\text{"} A * B \text{"}$ must be equivalent to $\text{"} AB \text{"}$.

Num is a tricky function to understand. In a previous section we described the formalization of numbers within type theory (and Higher Order Logic). Generally, a number is represented as so many successors to zero; thus, our common notion of the number 5 is expressed as $\text{Suc Suc Suc Suc Suc } 0$ or equivalently $\text{SSSSS}0$ and sometimes (set-theoretic) $0''''''$. The function Num transforms the series of symbols into a single number-symbol. In the literature, one finds a difference between 5 and Num 5. 5 is some

symbol (whose Godel encoding does not have to be the number 5) whereas Num 5 is a symbol representing the concept of the number 5.

Consider the number 1003. It is comprised of a series of 4 symbols. We write 1003 because it is more convenient than writing SSSS...S0. However, Num 1003 is not 1003 (there are 1004 symbols in the representation!). It is possible, though, that some sentence has the Godel number 1003, e.g., " $\forall x = x$ " \equiv 1003 where 1003 is a single number-symbol. To avoid confusion, " $\forall x = x$ " represents the Godelization of the sentence; Num " $\forall x = x$ " represents the number associated with the Godelization. Note: Num "1003" \neq 1003 and that Num 1003 = 1003, since 1003 is a single number-symbol whereas "1003" is the Godel number associated with the series '1' '0' '0' '3' of four symbols.

Getting back to sentence G, consider application to "G"; "G" represents the Godelization of G:

$$G \text{ "G"} \equiv [\lambda x_o \forall y_o \neg \text{Proof}_{ooo}^{\circ} [\text{"["} * x_o * [\text{Num}_{oo} x_o] * \text{"}]"] y_o]_{o \rightarrow o} \text{"G"}$$

$$G \text{ "G"} \equiv [\forall y_o \neg \text{Proof}_{ooo}^{\circ} [\text{"["} * \text{"G"} * [\text{Num}_{oo} \text{"G"}] * \text{"}]"] y_o]_{o \rightarrow o}$$

If G "G" is a theorem, then "G" * [Num_{oo} "G"] must be the last sentence of some valid proof. So if we assume $\vdash G \text{ "G"}$ meaning that G "G" is a theorem, then we have

$$\vdash [\forall y_o \neg \text{Proof}_{ooo}^{\circ} [\text{"["} * \text{"G"} * [\text{Num}_{oo} \text{"G"}] * \text{"}]"] y_o]_{o \rightarrow o}$$

Since $\vdash G \text{ "G"}$, then $\exists n . \text{Proof} [\text{"["} * \text{"G"} * [\text{Num}_{oo} \text{"G"}] * \text{"}]"] n$

Specialize y to n:

$$\vdash [\neg \text{Proof}_{ooo}^{\circ} [\text{"["} * \text{"G"} * [\text{Num}_{oo} \text{"G"}] * \text{"}]"] n]_{o \rightarrow o}$$

Simplify:

$$\vdash [\neg T_o]$$

$$\vdash F_o$$

If we derive a falsehood, as we have done, then our system is inconsistent; rejecting inconsistency, we must accept $\neg \vdash G \text{ "G"}$; so that G "G" is not a theorem (i.e., cannot be proved).

If G "G" is not a theorem, it is possible that $\neg G \text{ "G"}$ is a theorem. Assume $\vdash \neg G \text{ "G"}$, then

$$\vdash \neg [\forall y_o \neg \text{Proof}_{ooo}^{\circ} [\text{"["} * \text{"G"} * [\text{Num}_{oo} \text{"G"}] * \text{"}]"] y_o]_{o \rightarrow o}$$

Then, noting the equivalence, $\sim \forall \sim \equiv \exists$, we have

$$\vdash [\exists y_o \text{Proof}_{ooo}^{\circ} [\text{"["} * \text{"G"} * [\text{Num}_{oo} \text{"G"}] * \text{"}]"] y_o]_{o \rightarrow o}$$

If there existed a "proof", then G "G" would be a theorem, therefore we reject $\vdash \neg G \text{ "G"}$ and accept $\neg \vdash \neg G \text{ "G"}$.

Finally, we inquire as to the "truth" of G "G". Suppose G "G" were false, then there would exist some proof (as G "G" asserts that for all y, there is no proof); however, we have established that the existence of such a proof leads to inconsistency. Hence, G "G" is true (but not provable). Thus, with sentence G "G" we have demonstrated a true sentence that cannot be proved and with $\neg G \text{ "G"}$ a false sentence that cannot be refuted.

Having worked through the top level of Godel's First Incompleteness Theorem, we now introduce some relevant definitions and concepts that have not been introduced:

Num, *, and Proof are representable⁴ in Q^- . Then, if the formal system \mathcal{A} is any recursively axiomatized⁵ pure extension of Q^- , Num, *, and Proof are representable in \mathcal{A} . Let \mathcal{A} be an extension of Q^- :

1. \mathcal{A} is *consistent* iff $\neg \vdash_{\mathcal{A}} F_0$.
2. \mathcal{A} is ω -*consistent* iff there is no closed wff B_{00} of \mathcal{A} such that $\vdash_{\mathcal{A}} \exists y_0 B_{00}y_0$ and for each natural number k , $\vdash_{\mathcal{A}} \neg B_{00}k$. Another way of stating this is that the formal system \mathcal{A} is ω -*inconsistent* if there is a formula $F(\omega)$ in one free variable ω , such that the sentence $\exists \omega F(\omega)$ is provable, yet all the sentences $F(0), F(1), \dots, F(n), \dots$ are refutable. A ω -*inconsistent* cannot be correct (only true sentences are provable), but it may be consistent. A system that is not ω -*inconsistent* is ω -*consistent* – meaning that whenever the sentence $\exists \omega F(\omega)$ is provable, then for at least one number n , $F(n)$ is not refutable.
3. \mathcal{A} is *complete* iff for each sentence A_0 of \mathcal{A} , $\vdash_{\mathcal{A}} A_0$ or $\vdash_{\mathcal{A}} \neg A_0$.
4. \mathcal{A} is ω -*complete* iff $\vdash_{\mathcal{A}} \forall y_0 B_{00}y_0$ for each closed wff B_{00} of \mathcal{A} such that $\vdash_{\mathcal{A}} B_{00}k$ for each natural number k .
5. \mathcal{A} is *essentially incomplete* iff \mathcal{A} is consistent and recursively axiomatized and every consistent recursively axiomatized pure extension of \mathcal{A} is incomplete.

3.2. Representing Proof

Consider the following demonstration of Modus Ponens:

- (1) Assume $t1$
 - (1.1) $t1 \vdash t1$
- (2) Assume $t1 \rightarrow t2$
 - (2.1) $t1 \rightarrow t2 \vdash t1 \rightarrow t2$
- (3) Conclude (Modus Ponens) $\vdash t2$
 - (3.1) $t1, t1 \rightarrow t2 \vdash t2$

Now consider the same demonstration written in a functional form:

$$\text{Modus_Ponens}(\text{Assume}(t1), \text{Assume}(t1 \rightarrow t2)) = \vdash t2$$

Godel's idea was that checking a proof could be done by an algorithm; and even more efficient, the proof itself could be written as algorithm (e.g., program). But perhaps, the most important result is that Godel precisely defined the concept of proof and then showed that it was recursive. Godel's definition of a recursive number-theoretic function extended the concept of general recursion to include the notion that a proof is also recursive: If a proof is a finite series of sentences such that the last sentence in the series is the conclusion (i.e., the sentence proved by the proof) and where each sentence is either an axiom or the result of applying a rule of inference to one or more earlier sentences, then clearly proof is recursive (indeed, it is obvious that is recursively enumerable).

⁴ A function, F , is said to represent a number set A , $n \in A \leftrightarrow \vdash F(n)$. In type theory: $\vdash n. A\ n = \vdash F\ n$.

⁵ Godel's definition of recursively axiomatized is "Function ψ is called recursive, if there exists a finite series of number-theoretic functions $\psi_1, \psi_2, \psi_3, \dots, \psi_n$, which ends in ψ and has the property that every function ψ_k of the series is either recursively defined by two of the earlier ones, or is derived from any of the earlier ones by substitution or is a constant or the successor function."

With this background, let us delve into the details of the function proof. Informally, we may define proof as $[\lambda m, n. \exists \text{ a proof } P^0, \dots, P^k \text{ such that } m = "P^k" \text{ and } n = "P^0, \dots, P^k"]_{\text{agg}}$ – however, we need a more precise definition:

$$[\lambda m, n. (m = E(n) \mathcal{G}L n)]_0 \wedge$$

$$(\forall k \leq E(n)) [Axiom(k \mathcal{G}L n) \vee (\exists i < k)(\exists j < k) Rule(i \mathcal{G}L n, j \mathcal{G}L n, k \mathcal{G}L n)]_{\text{agg}}]_{\text{agg}}$$

Where:

- $E(x)$ stands for $(\text{Max } n \leq x)[\text{Pr}(n) \text{ Divides } x]$ – Pr enumerates the prime numbers (e.g., $\text{Pr}(0)$ stands for 2, $\text{Pr}(5)$ stands for 13, etc.).
 $x \text{ Divides } y$ stands for the sentence $[\lambda x, y. \exists z. (z \leq y) \wedge (y = x \text{ Multiply } z)]$.
 $x \text{ Multiply } y$ stands for the function:
 $[\lambda x, y. @z. (y = 0 \wedge z = 0) \vee (z = (x \text{ Multiply } y-1) \text{ Plus } x)]$.⁶
 $x \text{ Plus } y$ stands for $[\lambda x, y. @z. (y = 0 \wedge z = x) \vee (z = (\text{Suc } x) \text{ Plus } y-1)]$.
 Pr stands for the sentence:
 $[\lambda x. @z. (y = 0 \wedge z = 2) \vee (\text{Pr}(\text{Suc } z) = (\text{Min } m \leq \text{Suc}(\text{Pr}(z)!)) [\text{Prime } m \wedge (\text{Pr}(z) < m)]]$.
Finally, Prime stands for $[\lambda x. (2 \leq x) \wedge (\forall y \leq x)[y \text{ Divides } x \rightarrow y = x \vee y = 1]]$.
- $n \mathcal{G}L x$ stands for $(\text{Max } y \leq x)[\text{Pr}(n)^y \text{ Divides } x]$ – Here we use exponentiation.
 x^y stands for $[\lambda x, y. @z. (y = 0 \wedge z = 1) \vee (z = x \text{ Multiply } (x^{y-1}))]$.⁷
- $Axiom(n)$ stands for $[\lambda n. \exists z. axiom(z, n)]$ – We devote an entire subsection to this function. For the present purpose, it is sufficient to note that this boolean term is true if and only if n represents a valid encoding of one the axioms.
- $Rule(m, n, p)$ stands for $[\exists \text{ exist wff's } A, B, C, \text{ and } D \text{ such that } m = "[A = B]" \text{ and } n = "C" \text{ and } p = "D", \text{ and } D \text{ is obtained from } [A = B] \text{ and } C \text{ by rule } R]$. Note rule R is the only rule of inference of Q_0 and Q^* . As with $Axiom$, we devote an entire subsection to this function. For the present, $Rule$ identifies whether the sentence m is derived from other sentences via application of a (the) rule of inference.

With this background, let us revisit the expression for proof. The first part of the expression contains $(m = E(n) \mathcal{G}L n)_0$ to ensure that the sentence encoded by m is indeed the last sentence in the series encoded by n .

The next part contains $(\forall k \leq E(n)) [Axiom(k \mathcal{G}L n)]$. To understand this term, recall that Godelization of a series of sentences requires that the Godel encoding of each sentence be concatenated by using them as exponents to successively greater prime numbers. For example consider the sentence A . Its Godel encoding is " A ". Now consider the sentence B . Its Godel encoding is " B ". For the concatenation of these sentences we would like AB to be encoded as " AB ". We also want this encoding to be composable so that " A " * " B " = " AB " = AB . Here's how: Passing from left to right, we associate each sentence with successive prime numbers. Then we use the Godel number of each sentence as the exponent to the prime number associated with each sentence. Thus, if " A " = 10000 and " B " = 100, then " AB " is $2^{10000} \times 3^{100}$.

⁶ The definitions here for Multiply , Plus , and Pr differ from Andrew's Theory Q and from the HOL definition. We derived these definitions for their constructive nature and to demonstrate that using the Hilbert " e " operator can simplify the specification of a primitive recursive relationship.

3.3. Representing Axioms

In the following subsections we define the expressions “Axiom x” in great detail.

The following shows the transformation of Axiom 1 from the form given to its more natural functional form:

The definitions of $*$ and $'''$ are explained in the next subsection.

⁷ Exponentiation is somewhat complicated by the definition that $0^0 = 1$ but $0^{\text{Suc } y} = 0$.

If we properly bracket Axiom 2 based on function application and partial function application:

$$[[\supset_{ooo} [[Q_{o\alpha\alpha} X_\alpha] Y_\alpha]] [[Q_{ooo} [h_{o\alpha} X_\alpha]] [h_{o\alpha} Y_\alpha]]]$$

The Godel encoding must consider the polymorphism of the type α and we also need to verify that the type substitution represents a valid type. Thus, we require a function to determine whether some non-zero number represents a valid type and this is not so simple as might be first thought. Types are Godelized as follows:

- Type o is Godelized, “ o ”, as 3
- Type ι is Godelized, “ ι ”, as 5
- Type $\alpha\beta$ is Godelized, “ $\alpha\beta$ ”, as $3^{“\alpha”} \times 5^{“\beta”}$

Thus, to Godelize a multi-symbol type, we traverse the series from left to right and compute $\prod \text{Pr}(i)^{\text{symbol}(i)}$ where i refers to the index position of each symbol and $\text{symbol}(i)$ refers to the symbol at position i in the series. This Godelization process quickly demonstrates an important issue related to Godelization: The numbers quickly become very, very large.

Consider the type σ that we have used in previous sections to represent the type number as a shorthand for the type oot . It should not make a difference which we choose to use: We use σ to remind us of “number” and oot to remind us that a number is just “that which is common” to all equivalent sets. The Godel number of the (simple) type $\sigma\sigma$ is difficult to use because it inherently means $((\sigma)(\sigma))$ which maps to $((oot)(oot))$. This means that “ $\sigma\sigma$ ” is $3^{“\sigma”} \times 5^{“\sigma”}$ where “ σ ” is “ oot ” $\rightarrow 3^3 \times 5^3 \times 7^3$; now, this means that “ oot ” * “ oot ” = “ $ootoot$ ” which is $3^3 \times 5^3 \times 7^3 \times 9^3 \times 11^3 \times 13^3$; so that “ oot ” = 56,723,625 and “ $ootoot$ ” = 2.0436×10^{19} but, “ σ ” is $3^{56723625}$ and “ $\sigma\sigma$ ” is $3^{56723625} \times 5^{56723625}$ = a really big number! Note also that “ $((oot)(oot))$ ” is $3^{56723625} \times 5^{56723625}$.

Thus, determining whether a given number represents a type (i.e., some sequence of o and/or ι) or a type symbol (i.e., a symbol representing a type) is without too much complication.

If the number represents a type, the lambda expression

$$[\lambda n . \forall i . ((i \leq E(n) \wedge (\text{Pr}(i)^3 \text{ Divides } n \vee \text{Pr}(i)^5 \text{ Divides } n))$$

will be true. If the number represents a type symbol, the lambda expression

$$[\lambda m . \forall i . ((i \leq E(m)) \wedge \exists n . ((\forall j . ((j \leq E(n) \wedge (\text{Pr}(j)^3 \text{ Divides } n \vee \text{Pr}(j)^5 \text{ Divides } n)) \wedge \text{Pr}(i)^n \text{ Divides } m)))$$

will be true. Now, letting Type stand for

$$[\lambda n . \forall i (i \leq E(n) \wedge (\text{Pr}(i)^3 \text{ Divides } n \vee \text{Pr}(i)^5 \text{ Divides } n))]$$

and TypeSymbol stand for

$$[\lambda m . \forall i . ((i \leq E(m)) \wedge \exists n . ((\forall j . ((j \leq E(n) \wedge (\text{Pr}(j)^3 \text{ Divides } n \vee \text{Pr}(j)^5 \text{ Divides } n)) \wedge \text{Pr}(i)^n \text{ Divides } m)))]$$

We let ValidType⁸ stand for

$$[\lambda n . \text{Type}(n) \vee \text{TypeSymbol}(n)]$$

⁸ Our ValidType is a big improvement over Andrews [1]. In [1], his motivation is to show that these expressions were recursive. By building expressions that represent how a decision algorithm would operate, we demonstrate recursiveness in a manner more natural for the computer scientist.

and note that $\text{ValidType}(n) = T$, if and only if, n represents a type or type symbol. With all of this, we are almost in a position to demonstrate the expression used to determine whether a given number represents a valid encoding of an instance of Axiom 2. We must first identify the Godel encoding of the “constant” symbols. The “constant” symbols of Axiom 2 (shown below) and their Godel encodings are computed below:

$$\text{Axiom 2: } [[\supset_{ooo} [[Q_{o\alpha\alpha} X_\alpha] Y_\alpha]] [[Q_{ooo} [h_{o\alpha} X_\alpha] [h_{o\alpha} Y_\alpha]]]]$$

Where,

$$“h_{o\alpha}” = 31^{“o\alpha”} = 31 \exp(3^3 \times 5^{“o\alpha”})$$

$$“X_\alpha” = 103^{“\alpha”} = 103 \exp(3^{“\alpha”})$$

$$“Y_\alpha” = 107^{“\alpha”} = 107 \exp(3^{“\alpha”})$$

$$“Q_{o\alpha\alpha}” = 19^{“o\alpha\alpha”} = 19 \exp(3^3 \times 5^{“o\alpha”} \times 7^{“\alpha”})$$

Finally⁹,

$$\text{Let axiom}(2,n) \text{ stand for } [\lambda n . \exists a . a \leq n . \text{ValidType}_o(a) \wedge n = “[[\supset_{ooo} [[” * 19^{(3^3 \times 5^{“a”} \times 7^{“a”})} * 103^{“a”} * “]” * 107^{“a”} * “]]” * 19^{(3^3 \times 5^3 \times 7^3)} * 31^{(3^3 \times 5^{“a”})} * 103^{“a”} * “]]” * 31^{(3^3 \times 5^{“a”})} * 107^{“a”} * “]]”]]$$

3.3.3. Axiom 3: $f_{\alpha\beta} = g_{\alpha\beta} = !X_\beta . f_{\alpha\beta} X_\beta = g_{\alpha\beta} X_\beta$

Axiom 3, $f_{\alpha\beta} = g_{\alpha\beta} = !X_\beta . f_{\alpha\beta} X_\beta = g_{\alpha\beta} X_\beta$, when rewritten in functional form becomes:

$$[[Q_{ooo} [[Q_{o(\alpha\beta)(\alpha\beta)} F_{(\alpha\beta)}] G_{(\alpha\beta)}]] [[Q_{o(o\beta)(o\beta)} [\lambda x_\beta T_o]] [\lambda x_\beta [[Q_{o\alpha\alpha} [F_{(\alpha\beta)} x_\beta]] [G_{(\alpha\beta)} x_\beta]]]]]$$

and its Godelization:

$$\text{Let axiom}(3,n) \text{ stand for } [\lambda n . \exists \alpha, \beta . \alpha \leq n \wedge \beta \leq n . \text{ValidType}_o(\alpha) \wedge \text{ValidType}_o(\beta) \wedge n = “[[” * “Q”^{“o\alpha” * “o\alpha” * “o\alpha”} * “[[” * “Q”^{“o\alpha” * “(\alpha\beta)” * (3^{“o\alpha”} \times 5^{“\beta”})} * “F”^{“(\alpha\beta)”} * “]” * “G”^{“(\alpha\beta)”} * “]]” * “Q”^{“o\alpha” * “(o\beta)” * (3^{“o\alpha”} \times 5^{“\beta”})} * “[\lambda” * “x”^{“\beta”} * “T”^{“o\alpha”} * “]]” * “[\lambda” * “x”^{“\beta”} * “[[Q”^{“o\alpha” * “o\alpha” * “o\alpha”} * “[” * “F”^{“(\alpha\beta)”} * “x”^{“\beta”} * “]]” * “G”^{“(\alpha\beta)”} * “x”^{“\beta”} * “]]]]”]$$

Note that in axiom(3,n) we have made the following simplifications:

- We use “Q”, “x”, etc., to represent the Godel number/encoding for the object being quoted;
- λx introduces an abstraction – x represents any valid variable name and axiom(3,n) should have the conjunctive term $(\exists x \leq n . \text{ValidVariable}(x))$ and everywhere we have used “x” we should have used x (the Godel number associated with the valid variable name); and,
- Except for one instance, we have favored the simpler form “ $(\alpha\beta)$ ” for types over the more complex (and less obvious) $(3^{“\alpha”} \times 5^{“\beta”})$.

For completeness, let $\text{ValidVariable}(n)$ stand for $[\lambda n . \exists a, i . a \leq n \wedge i > \text{'number of primitive symbols'} \wedge \text{ValidType}(a) \wedge n = \text{Pr}(i)^a]$.

3.3.4. Axiom 4: $[\lambda x_\alpha B_\beta] A_\alpha = B_\beta$

Axiom 4, $[\lambda x_\alpha B_\beta] A_\alpha = B_\beta$, rewritten in functional form: $[[Q_{o\beta\beta} [[\lambda x_\alpha B_\beta] A_\alpha]] B_\beta]$

⁹ Our axiom 2 corrects a number of serious errors in Andrews. His encoding of Axiom 2 is not composable, as is ours.

Let axiom(4,n) stand for $[\lambda n . \exists x, B, A, \alpha, \beta . \alpha \leq n \wedge \beta \leq n \wedge A \leq n \wedge B \leq n \wedge x \leq n \wedge \text{ValidVariable}(x) \wedge \text{WFF}(A) \wedge \text{TypeOfWFF}(A) = \text{TypeOfVariable}(x) \wedge \text{ValidSymbol}(B) \wedge B \neq x \wedge \text{ValidType}(\alpha) \wedge \text{ValidType}(\beta) \wedge n = "[[" * "Q"^{a_n} * "a"^{a_n} * "[[" * "x"^{a_n} * "B"^{a_n} * "]" * "A"^{a_n} * "]" * "B"^{a_n} * "]"]]$

Where,

- $\text{ValidSymbol}(n)$ stands for $[\lambda n . \text{ValidVariable}(n) \vee n = \text{Pr}(\text{index of primitive symbol } i)^{a_n} \vee (\exists a . a \leq n \wedge \text{ValidType}(a) \wedge n = "Q"^{a_n} * "a"^{a_n})]$
- $\text{TypeOfWFF}(n)$ stands for $[\lambda n . @y . \min y \leq n . \exists z \leq n . n = z^y \wedge \text{WFF}(z)]$
- $\text{TypeOfVariable}(n)$ stands for $[\lambda n . @y . \min y \leq n . \exists z \leq n . n = z^y \wedge \text{Prime}(z)]$

3.3.4.1. Parsing as a Type-Theoretic Expression

The function $\text{WFF}_{\rightarrow o}$ is a very complicated expression. In what follows we demonstrate that WFF is a computable function of numbers to boolean. We then show how one would express this function as a type-theoretic expression.

Consider the following type-theoretic expression¹⁰:

$$[\exists n_o, m_o . \forall x_i . [m_o x_i \rightarrow n_o x_i] \wedge \exists y_i . [n_o y_i \wedge \sim m_o y_i] \wedge E_{o(o \cup o)} n_o m_o]$$

It asserts the existence of two sets (m and n) such that for all individuals, x, if x is a member of m, then x is a member of n (m is a subset of n). The expression continues by asserting the existence of an individual y that is in n, but not in m; thus, m is a proper subset of n. The last term says that the two sets can be brought into one-to-one correspondence. Now that we are convinced that the sentence makes sense, how do we “prove” that it is a wff?

In order for a sentence to be well-formed, it must be syntactically correct and for a type-theoretic system, types are part of the syntax. Thus, in the above expression, we can state that if all of the sub-terms are wff, then the sentence is wff. For example, consider the sub-term involving the existential y as it more properly written: $[\wedge_{o \cup o} [[n_o y_i] [\sim m_o y_i]]]$. This sub-term is wff, since each of the factors resolve to the proper type¹¹.

For type theory, we can build a parser to verify that a sentence is a wff. Here's how:

- Simplest approach: Do not try to parse the number theoretic sentence. As part of the Godelization process, do the parse; thus, if and only if the sentence is wff will our Godelizer generate a number.
- More difficult: Although complex, we can extract the unique sentence represented by the Godel number being considered. This sentence is then passed to the parser which returns success if the sentence is wff and fails otherwise.

If well-formedness is essentially the parsing problem, then we can imagine a Turing machine “programmed” with our parser. Cook's Theorem tells us that we can translate the parser to an array of machine descriptions, such that the first row is the instantaneous description of the initial machine configuration with the input sentence. Each subsequent row is a legal successor to the previous row. The last row contains either accept or reject.

¹⁰ This expression is another representation of infinity. It says that a set is infinite if and only if it can be put into one-to-one correspondence with a proper subset of itself.

¹¹ HOL allows us to write these functions, including the parser for wff in ML. If we were not able to do this, we would be in the tedious position of having to actually implement a propositional logic formulation of the parser – we discuss this application of Cook's Theorem in subsequent paragraphs.

Finally, we define propositional variables (using only propositional logic) to “simulate” this array. Thus, variable $q(t, i)$ is true if our Turing machine simulated tape contains variable q on square i at time t . Since each square (we know in advance that the upper bound of the number of squares need to parse a sentence is related to the size of the sentence) contains at least symbol at all times, then we have $q(t, i) \vee a(t, i) \vee \dots \vee Q(t, i)$. Since each square can only have one symbol at each time, we also have $\sim(q(t, i) \vee a(t, i))$ for each pair of symbols.

At this point the reader should be convinced that we can build a propositional logic expression that describes the operation of the parser in such a manner that we can determine whether the parser accepted or rejected the sentence. Actually, the situation is much more profound: We can build such an expression for any recursive set including our metatheory. Gödel’s Incompleteness Theorem teaches us, though, that when presented with sentence G , that the system will not be able to decide on its status as a theorem.

3.3.5. Axiom 5: $[\lambda x_\alpha x_\alpha]A_\alpha = A_\beta$

Axiom 5, $[\lambda x_\alpha x_\alpha]A_\alpha = A_\beta$, in functional form: $[[Q_{\alpha\alpha}][[\lambda x_\alpha x_\alpha]A_\alpha]] A_\alpha$

Let axiom(5,n) stand for $[\lambda n . \exists A, x, \alpha . \alpha \leq n \wedge x \leq n \wedge A \leq n \wedge \text{ValidVariable}(x) \wedge \text{ValidType}(\alpha) \wedge \text{WFF}(A) \wedge \text{TypeOf}(A) = \text{PrimeLog}(x) \wedge n = \text{“}[[\text{“} * \text{“}Q\text{”}^{\alpha\alpha} * \text{“}[[\lambda * x^{\alpha\alpha} * x^{\alpha\alpha} * \text{“}]] * A^{\alpha\alpha} * \text{“}]] * A^{\alpha\alpha} * \text{“}]]\text{”}]$

3.3.6. Axiom 6: $[\lambda x_\alpha . B_\beta C_\gamma]A_\alpha = [[\lambda x_\alpha B_\beta]A_\alpha][[\lambda x_\alpha C_\gamma]A_\alpha]$

Axiom 6, $[\lambda x_\alpha . B_\beta C_\gamma]A_\alpha = [[\lambda x_\alpha B_\beta]A_\alpha][[\lambda x_\alpha C_\gamma]A_\alpha]$, in functional form is quite daunting:

$$[[Q_{\alpha\beta\gamma}][[\lambda x_\alpha [B_\beta C_\gamma]]A_\alpha]][[\lambda x_\alpha B_\beta]A_\alpha][[\lambda x_\alpha C_\gamma]A_\alpha]]$$

Let axiom(6,n) stand for $[\lambda n . \exists A, B, C, x, \alpha, \beta, \gamma . \alpha \leq n \wedge \beta \leq n \wedge \gamma \leq n \wedge x \leq n \wedge A \leq n \wedge B \leq n \wedge C \leq n \wedge \text{ValidVariable}(x) \wedge \text{ValidType}(\alpha) \wedge \text{ValidType}(\beta) \wedge \text{ValidType}(\gamma) \wedge \text{WFF}(A) \wedge \text{WFF}(B) \wedge \text{WFF}(C) \wedge \text{TypeOf}(A) = \text{PrimeLog}(x) \wedge \text{TypeCompatible}(\text{TypeOf}(B), \text{TypeOf}(C)) \wedge n = \text{“}[[\text{“} * \text{“}Q\text{”}^{\alpha\alpha} * \text{“}[[\lambda * x^{\alpha\alpha} * \text{“}[[\text{“} * B^{\beta\beta} * C^{\gamma\gamma} * \text{“}]] * A^{\alpha\alpha} * \text{“}]] * [[\lambda * x^{\alpha\alpha} * B^{\beta\beta} * \text{“}]] * A^{\alpha\alpha} * \text{“}]] * [[\lambda * x^{\alpha\alpha} * C^{\gamma\gamma} * \text{“}]] * A^{\alpha\alpha} * \text{“}]]\text{”}]$

TypeCompatibility anticipates extensions of well-formedness to include parsing the entire sentence so that we may claim that if a sentence is well-formed then so are its sub-terms.

3.3.7. Axiom 7: $[\lambda x_\alpha . \lambda y_\gamma B_\delta]A_\alpha = [\lambda y_\gamma [\lambda x_\alpha B_\delta]A_\alpha]$

The functional form for Axiom 7: $[[Q_{\alpha(\delta\gamma)\delta\eta}][[\lambda x_\alpha [\lambda y_\gamma B_\delta]]A_\alpha]][\lambda y_\gamma [[\lambda x_\alpha B_\delta]A_\alpha]]$

and let axiom(7,n) stand for $[\lambda n . \exists A, B, x, y, \alpha, \delta, \gamma . \alpha \leq n \wedge \delta \leq n \wedge \gamma \leq n \wedge x \leq n \wedge y \leq n \wedge A \leq n \wedge B \leq n \wedge \text{ValidVariable}(x) \wedge \text{ValidVariable}(y) \wedge y \neq x \wedge \text{ValidType}(\alpha) \wedge \text{ValidType}(\delta) \wedge \text{ValidType}(\gamma) \wedge \text{WFF}(A) \wedge \text{WFF}(B) \wedge \text{TypeOf}(A) = \text{PrimeLog}(x) \wedge n = \text{“}[[\text{“} * \text{“}Q\text{”}^{\alpha\alpha} * \text{“}[[\lambda * x^{\alpha\alpha} * \text{“}[[\text{“} * y^{\gamma\gamma} * B^{\delta\delta} * \text{“}]] * A^{\alpha\alpha} * \text{“}]] * [[\lambda * y^{\gamma\gamma} * \text{“}[[\text{“} * x^{\alpha\alpha} * B^{\delta\delta} * \text{“}]] * A^{\alpha\alpha} * \text{“}]]\text{”}]$

3.3.8. Axiom 8: $[\lambda x_\alpha . \lambda x_\alpha B_\delta]A_\alpha = [\lambda x_\alpha B_\delta]$

In functional form, Axiom 8: $[[Q_{\alpha(\delta\alpha)(\delta\alpha)}][[\lambda x_\alpha B_\delta]A_\alpha]][\lambda x_\alpha B_\delta]$

Let axiom(8,n) stand for $[\lambda n . \exists A, B, x, \alpha, \delta . \alpha \leq n \wedge \delta \leq n \wedge x \leq n \wedge A \leq n \wedge B \leq n \wedge \text{ValidVariable}(x) \wedge \text{ValidType}(\alpha) \wedge \text{ValidType}(\delta) \wedge \text{WFF}(A) \wedge \text{WFF}(B) \wedge \text{TypeOf}(A) = \text{PrimeLog}(x) \wedge n = \text{“}[[\text{“} * \text{“}Q\text{”}^{\alpha\alpha} * \text{“}[[\lambda * x^{\alpha\alpha} * B^{\delta\delta} * \text{“}]] * A^{\alpha\alpha} * \text{“}]] * [[\lambda * x^{\alpha\alpha} * B^{\delta\delta} * \text{“}]]\text{”}]$

3.3.9. Axiom 9: $[l_{(o)}[Q_{oi}y]=y_l]$

The functional representation of Axiom 9 is similar to axiom 2: $[[Q_{oi}[y_l]][l_{(o)}[Q_{oi}[y_l]]]$.

Let axiom(9,n) stand for $[\lambda n . \exists y . y \leq n \wedge \text{ValidVariable}(y) \wedge n = "[[" * "Q"^{o_n} * "o_l" * "o_l" * "[" * y^{o_l} * "]" * "]" * "l"^{o_l} * "o_l" * "[" * "]" * "Q"^{o_n} * "o_l" * "o_l" * "[" * y^{o_l} * "]" * "]" * "l"^{o_l} * "o_l" * "]"]]$

3.4. Representing Rules of Inference

Q_0 has only one rule of inference: From C and $A_\alpha = B_\alpha$ to infer the result of replacing one occurrence of A_α in C by an occurrence of B_α , provided that the occurrence of A_α in C is not an occurrence of a variable immediately preceded by λ .

This means that there is some m such that $m = "[[Q_{oi} A_\alpha] B_\alpha]"$ and formulas (not necessarily well-formed) J and K (possibly empty), such that $C = "J" * A_\alpha * "K"$ and $D = "J" * B_\alpha * "K"$.

Let Rule(m,n,p) stand for $[\lambda m . \lambda n . \lambda p . \exists A, B, \alpha, z . A \leq m \wedge B \leq m \wedge \alpha \leq m \wedge \text{wff}(A) \wedge \text{wff}(B) \wedge \text{TypeOf}(A) = \text{TypeOf}(B) \wedge \text{TypeOf}(A) = z \wedge m = "[[" * "Q"^{o_n} * "z" * "z" * A^z * "]" * B^z * "]" \wedge \text{TypeOf}(n) = "o" \wedge \text{TypeOf}(p) = "o" \wedge \exists j, k . j \leq n \wedge k \leq n \wedge n = j * A^z * k \wedge p = j * B^z * k]$

4. Demonstrating GIT in HOL

Section 3, "Godel's Incompleteness Theorem" expressed a backwards proof of GIT in HOL. In this section, we are concerned with the implementation details of using the HOL theorem prover to demonstrate that we have a correct proof.

Godel's sentence $G, [\lambda x. \forall y. \sim \text{Proof}_{\omega\omega}^S [\text{"["} * x * [\text{Num}_{\omega\omega} x] * \text{"}"]] y]_{\omega \rightarrow \omega}$, demands that there be no integer y that encodes a Proof of "G" G . Since all HOL proofs are constructive, we must show that any such Proof of "G" G leads to contradiction.

An interesting approach to the problem would be to show that Godel's Incompleteness Theorem is reducible to another problem. For example, consider the following program.

Program P:

```

  For i = 0 to  $\infty$ 
  do
    if i encodes a proof that P does not halt, then halt
  od;
End P.

```

By construction: P halts iff there is a proof that P does not halt.

Let $\psi \equiv \text{"P does not halt"}$

Then,

$$\psi \rightarrow \sim \text{Proof}(\sim \text{Halt P}) \equiv \psi \rightarrow \sim \text{Proof}(\psi)$$

and

$$\sim \psi \rightarrow \text{Proof}(\sim \text{Halt P}) \equiv \sim \psi \rightarrow \text{Proof}(\psi)$$

This program fragment shows a reducibility of Godel's sentence G to the halting problem. The proof relies on a violation of ω -consistency, $\vdash \exists \omega F(\omega) \rightarrow \exists \omega F(\omega)$, as in Godel's original proof.

Whether we choose Godel's original sentence G or a reduction of it, we shall always be faced with the impossibility of showing constructive proof. The problem, of course, is that we need to show that each integer does not encode a valid proof. Even though the cardinality of this set is "countably" infinite, it is still infinite none-the-less. Unfortunately, incompleteness means that within the formal system, Godel's sentence G is undecidable – so that neither $\text{Proof "G" } G$ nor $\sim \text{Proof "G" } G$ can be determined within the formal system.

Thus, given a correctly implemented proof, the result cannot be decided constructively. To implement proof-by-contradiction in HOL, requires a different strategy. One first expresses an implication (e.g., $p \rightarrow q$) whose negation is provable: $\sim(p \rightarrow q) \equiv \sim(T \rightarrow F) \equiv \sim F$.

The following example shows how, at a high-level, Godel's Incompleteness Theorem would be demonstrated using the HOL theorem prover.

Every backward proof begins with a goal. In HOL, a goal is specified using the `set_goal` function. This function has two arguments: 1) A list of terms corresponding to the assumptions of the goal; and 2) a term specifying the goal.

```
set_goal([], string_to_term "?G.  $\sim(\text{lg:num} . G \text{ g } x) \implies (!y . \sim G \text{ g } y)$ ");
```

HOL places the goal on the goal stack:

```

val it =
  Status: 1 proof.
  1. Incomplete:

```

```

Initial goal:
  (--'?G. ~(!g. (?x. G g x) ==> (!y. ~(G g y)))'--)
: proofs

```

For illustrative purposes we instantiate G using the "EXISTS_TAC". Note that function e is an abbreviation for expand, a function that applies tactics (and tacticals) to the goal stack.

```
- e(EXISTS_TAC(string_to_term "\\(x:num) (y:num). T"));
```

The "EXISTS_TAC" replaces every occurrence of G with the specified lambda-expression. HOL responds with:

```

OK..
1 subgoal:
val it =
  (--'?~(!g. (?x. (\x y. T) g x) ==> (!y. ~((\x y. T) g y)))'--)
: goalstack

```

We then facilitate a β -conversion reduction using "BETA_TAC":

```
- e(BETA_TAC);
```

HOL responds by "doing" the conversion and displaying the result:

```

OK..
1 subgoal:
val it =
  (--'?~(!g. (?x. T) ==> (!y. ~T))'--)
: goalstack

```

Next we convert the subgoal to prenex format using the PRENEX_CONV conversion. In HOL conversions rewrite a term and then assert a theorem stating that the original term implies the rewritten term. The CONV_TAC tactic puts the original term in the assumption (undischarge) so that the rewritten term is the new subgoal.

```
- e(CONV_TAC PRENEX_CONV);
```

The result shows that all quantifiers have been moved to the extreme left-hand side of the equation:

```

OK..
1 subgoal:
val it =
  (--'?g y x. ~(T ==> ~T)'--)
: goalstack

```

Next the existentials are removed. REPEAT is used to cause EXISTS_TAC to be applied as many times as possible.

```
- e(REPEAT (EXISTS_TAC (string_to_term "1")));
```

Since all quantifiers were moved to the extreme left (prenex form) and all of them were existential, they are all removed by repeatedly applying "EXIST_TAC":

```

OK..
1 subgoal:
val it =
  (--'?~(T ==> ~T)'--)
: goalstack

```

Finally, the tautology is eliminated:

- e(CONV_TAC Taut.TAUT_CONV);

The HOL theorem prover recognizes that the goal is proved, informs the user, and then displays the proof:

OK..

Goal proved. $\vdash \neg(T \implies \neg T)$

Goal proved. $\vdash \neg \exists y x. \neg(T \implies \neg T)$

Goal proved. $\vdash \neg(!g. (\exists x. T) \implies (!y. \neg T))$

Goal proved. $\vdash \neg(!g. (\exists x. (\lambda x y. T) g x) \implies (!y. \neg((\lambda x y. T) g y)))$

val it =

Initial goal proved.

$\vdash \neg G. \neg(!g. (\exists x. G g x) \implies (!y. \neg(G g y))) : \text{goalstack}$

To actually verify GIT, one needs to instantiate G with Godel's sentence and then continue with definitions for Proof, etc., as we did in section 3.

5. Conclusion

5.1. Summary

Even with the passage of nearly four decades, Godel's Incompleteness Theorem remains difficult. Other researchers have reworked and reinterpreted the theorem, but the details are still formidable and quite challenging. Godel has shown that any recursively axiomatizable formal system (i.e., any formal system whose axioms and rules of inference are computable) is incomplete or inconsistent. The ramifications run deep and wide: Can a rational system capture its own metatheory? If not, is the irrationality of man simply a reflection of Godel's Incompleteness Theorem? Does the mimicry of intelligence stand for intelligence or do we remain forever undecided?

In this work we have taken the novel and original approach of presenting Godel's First Incompleteness Theorem in reverse. Even more fundamental is that we show how to compute the sub-goals. This is a radical departure from every other presentation of Godel's Incompleteness Theorem – traditionally, the terms, predicates, etc., comprising Godel's Theorem are shown to be general recursive.

The Church-Turing thesis states that our notion of computability is expressed equally and completely by general recursion or Turing machines; therefore, Godel's Incompleteness Theorem may be expressed equally and completely within the context of a program instead of a set of general recursive functions.

In this regard, HOL is an ideal theorem prover. The HOL theorem prover relies on the programming language, ML (Meta-Language) as the meta-language for HOL. ML was developed as part of the LCF (Logic for Computable Functions) system for interactive automated reasoning about higher-order, recursively defined functions. Thus, HOL and ML are ideal for expressing metamathematics and related metatheory. ML is especially useful as a bridge between general recursion and Von-Neumann style imperative programming paradigms.

5.2. Problems Encountered

Godel's Incompleteness Theorem is very complex. A number of authors [16] have attempted to make Godel's work accessible to the non-mathematician. Unfortunately, these attempts generally fail to capture important notions such as proof, validity, consistency, soundness, and completeness. Other authors [1,3,7,12,17,18], writing for mathematicians (and/or computer scientists), are quite liberal in their approach and notation; thus, making each of them incomparable. In the end, one is forced to go through the original proof and through each of the various authors; eventually, of course, one understands.

Type theory is relatively simple, but has a moderate learning curve. HOL, on the other hand, is quite complex and has a steep learning curve. The HOL documentation is well-written, but for a system called HOL-88, not HOL-90 which was used in this work. Unhappily, the differences between HOL-90 and HOL-88 syntax are not always easily determined. Overtime, as one works through the dual learning curves of HOL and the syntax differences, the HOL system becomes less difficult.

One area of difficulty that was not actually encountered (but still very important) involves the use of numbers. In Godel's Theorem, we "Godelize" expressions, equations, etc., in order to demonstrate the concept of proof, etc., numbers. As we have discussed, our approach is not number theoretic but program-theoretic; thus, we can avoid the problems of big numbers by handling "Godelization" symbolically.

5.3. *Suggestions for Further Work*

Godel's First Incompleteness Theorem when expressed solely in terms of recursive function theory is extremely complex. A type-theoretic presentation of the theorem is simpler to understand; however, the relationship between type theory and general recursion is not obvious. Due to time limitations, etc., it was not possible to completely encode Godel's First Incompleteness Theorem with the HOL theorem prover.

A useful and worthy project would be to complete the encoding and to build a parser for symbolic Godelization.

6. References

- [1] P.B. Andrews, *An Introduction to Mathematical Logic and Type Theory: To Truth Through Proof*. Academic Press, 1986.
- [2] H.P. Barendregt, *The Lambda Calculus Its Syntax and Semantics*, North-Holland, 1984.
- [3] G. Boolos, *The Logic of Provability*, Cambridge University Press, 1993.
- [4] R. Carnap, *Introduction to Symbolic Logic*, Dover, 1958.
- [5] A. Church, "A formulation of the simple theory of types," *J. Symbolic Logic*, Vol. 5, pp. 56-68, 1940.
- [6] A. Church, *Introduction to Mathematical Logic*, Princeton University Press, 1996.
- [7] M. Davis, *Computability and Unsolvability*, Dover, 1973.
- [8] A. DiCanzio, *Galileo*, ADASI, 1996
- [9] K. Godel, *On Formally Undecidable Propositions of Principia Mathematica and Related Systems*, Dover, 1992.
- [10] M. Gordon, "HOL - A Machine Oriented Formulation of Higher Order Logic," University of Cambridge Computer Laboratory Technical Report No. 68, July 1985.
- [11] M.J.C. Gordon and T.F. Melham, *Introduction to HOL*. Cambridge University Press, 1993.
- [12] G. Hunter, *Metalogic An Introduction to the Metatheory of Standard First Order Logic*, University of California Press, 1996.
- [13] R.M. Karp, "Lecture Notes on Reducibility and NP-Completeness," UC Berkeley.
- [14] L.C. Paulson, *ML for the Working Programmer*. Cambridge University Press, 1991.
- [15] S. MacLane, *Categories for the Working Mathematician*, Springer-Verlag, 1971.
- [16] E. Nagel and J.R. Newman, *Godel's Proof*, New York University Press, 1986.
- [17] N. Shankar, *Metamathematics, Machines, and Godel's Proof*. Cambridge University Press, 1994.
- [18] R.M. Smullyan, *Godel's Incompleteness Theorems*. Oxford University Press, 1992.
- [19] A. Tarski, *Introduction to Logic and to the Methodology of Deductive Sciences*, Dover, 1995.
- [20] R. Toal, "Toward Automated Compiler Verification," Ph.D. Dissertation, UCLA, 1993.
- [21] J. von Wright, "Representing higher-order logic proofs in HOL," *The Computer Journal*, Vol. 38, No. 2, pp. 171-179, July 1995.