

# Recording of Scientific Notes Notarized by Dan E. Willard on May10, 2014

Our proof of  $P \neq NP$  will certainly be non-trivial. It has, thus, taken us 21+ years to discover how [37]'s initial 1993 formalism for generating quite innocent-looking boundary-case forms of exceptions to the Second Incompleteness Theorem can be revised to provide a proof of  $P \neq NP$ . A nice aspect of this proof is that it is much easier to appreciate with good retrospective hindsight than in its original (where one begins with a full-blown start that will inevitably cause many reader turns to be made). This is because there are multiple misleading misgivings, that lie along the roadways towards proving  $P \neq NP$ .

Our formal proof of  $P \neq NP$  will be contained in Sections 777-777. The next three sections will explain how our core methodologies are related to some seminal observations by Pudlik, Salover, Nelson and Wilkie-Parić [27, 22, 23, 35], along with some related work by Dus˘-Igajatović and Švajdar [6, 33].

Their results establish that the logical axioms system  $\alpha$ , employing a Hilbert-style deductive methods, obey more powerful generalizations of the Second Incompleteness Theorem made by Pudlik, Salover, Nelson and Wilkie-Parić [27, 22, 23, 35]. The appreciation of the importance of the seminal observations by Pudlik, Salover, Nelson and Wilkie-Parić [27, 22, 23, 35] was accepted on May 2 (and rewriting it required my attention).

**NEW SECTION 7.10: THE BACKGROUND NOTES**

Throughout this article,  $\alpha$  will denote a recursively enumerable set of "proper axioms", and  $d$  will denote a deduction method, used for deriving theorems from  $\alpha$ . Such  $d$  will typically include some rules of inference (such as modus ponens) accompanied by some "logical axioms", denoted as  $L_d$ , that we consider to be part of  $d$ 's infrastructure rather than  $\alpha$ 's infrastructure. (This seemingly minor notational point, which views  $L_d$  as part of the deduction method  $d$ , meticulously separated from the set of proper axioms  $\alpha$ , shall greatly simplify several aspects of our proof. We therefore emphasize that by the axiom system " $\alpha$ ", our discourse will always be referring to a set of proper axioms that is, automatically, further extended by  $L_d$ 's logical axioms.)

**Example 0.1** The deduction methods  $d$ , that we will use, are quite conventional. They can, for instance, correspond to the paradigms that were used in the textbooks by say of Enderton, Mendelson or Papadimitriou [7, 7, 7]. Thus,  $d$  can represent:

1. Enderton's 6-part set of logical axioms, augmented by solely modus ponens as a rule of inference, when one uses [7] as a framework to define  $L_d$ . (This framework was also employed in Papadimitriou's textbook [7].)

2. Mendelson's framework [7], which has two inference rules of modus ponens and generalization, and which would have  $L_d$  include the five pure logical axioms of A1-A5 (from section 2.3 of [7]) plus the two equality axioms A6 & A7 (defined in its Section 2.8).

In general, we will use the term Hilbert Deduction to refer to a deductive methodology that includes a modus ponens rule and whose efficiency differs from the above two proof methodologies by no more than a polynomial increase in proof length. Our proof of  $P \neq NP$  will employ essentially any Hilbert-style deductive method. We will focus our attention mostly on Item 1's paradigm for the sake of simplicity because they are probably best known to a wide audience of readers. (The best source about the more advanced properties of Hilbert deduction and probably be found in the Hájek-Pudlik textbook [13], but its more advanced techniques will not be used during our proof of  $P \neq NP$ .)

**Remark 0.2** Some deduction methods  $d$  fall into a non-Hilbert category, where they possess the same logical power as Example 0.1's "Hilbert" methods, while they are super-exponentially inefficient in the worst case. This category includes the semantic tableau and resolution methods. If an automated theorem proving [7]. Despite their enormous disadvantage, these techniques produce, often, good heuristics in automated deduction. [37, 40, 48, 49] illustrate how these tableau-style methods can help approximate how humans motivate themselves to cogitate. These results will interest logicians, historians of logic, and especially self-justifying logics, but they will not be germane to our proof of  $P \neq NP$ . For simplicity, the latter proof will employ only the better known Hilbert-style methods, defined in Example 0.1.

**REMOVE** Thus happily, this second lengthy group of articles, published by us, can be fully ignored by a reader who wishes to focus mostly on our proof of  $P \neq NP$ .)

**Definition 0.3** Once again, let  $\alpha$  denote an axiom system, and  $d$  denote a deduction method. The ordered pair  $(\alpha, d)$  will be called Self-Justifying when:

- I. one of  $\alpha$ 's theorems (or at least one of its axioms) will state that the deduction method  $d$ , applied to the system  $\alpha$ , will produce a consistent set of theorems, and
- II. the axiom system  $\alpha$  is in fact consistent.

**Example 0.4** For say,  $(\alpha, d)$ , it is easy to construct a system  $\alpha^d \supseteq \alpha$  that satisfies the Part-I condition. For instance,  $\alpha^d$  could consist of all of  $\alpha$ 's axioms plus an added "SelfRef( $\alpha, d$ )" sentence, defined as stating:

- There is no proof (using  $d$ 's deduction method) of  $0 = 1$  from the union of the axiom system  $\alpha$  with this sentence "SelfRef( $\alpha, d$ )" (looking at itself).

Kleene [17] discussed how to encode approximate analogs of this "SelfRef( $\alpha, d$ )" statement. Each of Kleene, Rogers and Jeroslow [17, 20, 16] noted  $\alpha^d$  may, however, be inconsistent (despite SelfRef( $\alpha, d$ )'s assertion), thus causing it to violate Part-II of self-justification's definition. This is because if the ordered pair  $(\alpha, d)$  is too strong then a classic Gödel-style diagonalization argument can be applied to the axiom system  $\alpha^d = \alpha + \text{SelfRef}(\alpha, d)$ , where the added presence of the statement SelfRef( $\alpha, d$ ) causes this extended version of  $\alpha$  to logically become automatically inconsistent. Thus, the machinery of the sentence "SelfRef( $\alpha, d$ )", while relatively easy to encode via an application of the Fixed Point Theorem, is ironically most often useless for purposes of demonstrating  $P \neq NP$ .

For the sake of clarity, our proof of  $P \neq NP$  will ultimately be a theorem proven by a fully classical logic, whose correctness is beyond doubt. The reason for our interest in classical logic, that evades the Second Incompleteness Theorem, is that these strange-looking entities will be a useful intermediate step to examine during our proof of  $P \neq NP$ .

**NEW Section 7.2: THE CLASSIFIED WORK OF PUDLIK, SLOVER, NELSON AND WILKIE-PARIĆ**

Throughout this article, Add( $x, y, z$ ) and Mult( $x, y, z$ ) will denote two 3-way predicate symbols specifying  $x + y = z$  and  $x \cdot y = z$ . An axiom system  $\alpha$  will be said to recognize successor, addition and multiplication as Total Functions iff it includes 1-3 as axioms:

$$\begin{aligned} \text{Va } 3x \quad & \text{Add}(x, 1, z) \\ \text{Va } Vx \quad & \text{Add}(x, y, z) \\ \text{Va } Vx \quad & \text{Mult}(x, y, z) \end{aligned}$$

Also, an axiom system  $\alpha$  will be called Type-M iff it contains (1) – (3) as axioms, Type-N iff it contains only (1) and (2) as axioms, and Type-S iff it contains only (1) as an axiom. Moreover,  $\alpha$  will be called Type-NS iff it contains none of these axioms.

Our initial paper [37] about self-justification involved Type-A logics that can formalize their own consistency under semantic tableau deduction. Shortly after we published this result, Robert Salover telephoned us, in April of 1994, indicating his intent to strengthen a variant of the Second Incompleteness Theorem due to Pavel Pudlik [27], with additional techniques developed by Nelson and Wilkie-Parić [23, 29], to show that the reasonable Type-S axiom systems (that treat Addition and Multiplication as 3-way relations) can verify its own Hilbert consistency. This result, summarized by Theorem 3.1, improves upon Beeson-Sherman-Shepherdson's earlier work in [3].

**Theorem 3.1** (Salover's 1994 modification [23] of Pudlik's 1985 formalism [27] using the added methodologies of Nelson and Wilkie-Parić [23, 35]) Let  $\alpha$  denote any arithmetical Type-S system that can verify addition and multiplication satisfy their usual associativity, commutativity, distributive and identity axioms. Then  $\alpha$  cannot prove a theorem affirming its own consistency under any of Beeson-SJS's forms of Hilbert-style deduction.

We will never use Theorem 3.1 during our proof of  $P \neq NP$ . Its importance in achieving this result cannot, however, be understated. This is because Theorem 3.1 indicates that only Type-NS arithmetics are plausibly capable of verifying their own Hilbert consistency.

If this indicates any hope in using self-justifying logics to prove  $P \neq NP$  must center around the daunting prospect of astonishingly weak Type-NS systems. (I have no doubt that my 21+ year effort to prove  $P \neq NP$  would have taken more than roughly — giggle, giggle — say 210+ years — if I was not aware of Theorem 3.1's crucial negative result and the tiny class of Type-NS boundary-case exceptions to the Second Incompleteness Theorem that it implicitly permits.)

**SOME RELEVANT BACKGROUND INFORMATION** One certainly salient aspect of Theorem 3.1 is that there is currently no available published proof of its statement. Essentially what had transpired was I called Pudlik in 1993 a proof that relatively rich Type-A logics can verify their own consistency under semantic tableau and resolution forms of deduction [37]. Robert Salover subsequently telephoned me in April of 1994 [23] to tell me that he knew how to generalize a theorem of Pudlik to show that a variant of Theorem 3.1 would pertain to Hilbert deduction.

That day was reaching the age where Berkeley would grant him a very lucrative early retirement in 1994. He showed no interest in writing up either Theorem 3.1's proof or spending more time fine-tuning his dissertation on his subject matter (together with sending me a 10-page email, outlining the core aspects of his fascinating methodology) and suggesting I read the less-general results by Pudlik, Nelson and Wilkie-Parić [27, 23, 35].

In a 2009 blog, Richard Lipton [21] has indicated that Salover was always reluctant to write up the proofs of his results. For instance, this blog reports Salover was the only mathematician to submit manuscripts to journals before he published a single paper. Also, many logicians [6, 13, 18, 28, 27, 28, 32, 34, 35] have commented about Salover's unusual reluctance to publish many results he derived, analogous to Theorem 3.1, that logicians have greatly valued.

Thus, on the suggestion of an anonymous referee, first JSL article [36] included a deliberately abbreviated 4-page Appendix A that formalized a slightly weaker but-more-easily comprehensible version of Theorem 3.1. Other logicians (its editor, and not fully analogous to Theorem 3.1), were published by Dus˘-Igajatović, Švajdar and ourselves in [6, 33, 43]. Also, Pudlik's initial paper [27] comes very close to establishing Theorem 3.1's result. Furthermore, Pudlik specifies in [28] that some also-unpublished work by Harvey Friedman [6] runs in some partially but-not-quite analogous directions to Pudlik's work on proof lengths.

In any case, the reader will be happily relieved to learn our proof of  $P \neq NP$  will never use the formalism of the mysterious Theorem 3.1, lying in an analog of the tales "Neverland" from "Peter Pan" novel by J. M. Barrie, as an intermediate step. Instead, the sole function of this non-published theorem is to serve as a signal about what steps a proof of  $P \neq NP$  need absolutely avoid.

This is because there are no examples of viable Type-S logics that verify their Hilbert-style consistency. Our proof of  $P \neq NP$  will thus use Theorem 3.1 as a beacon, signaling that only the amazingly weak Type-NS formalisms can plausibly support a proof of  $P \neq NP$ , via self-justifying logics. (Beyond this useful guiding hint, our proof of  $P \neq NP$  will not use Theorem 3.1's special machinery.)

**NEW Section 7.3: THIS GIVING STARTED**

Our first paper on self-justifying axiom systems were the conference papers [37, 38] whose unconventional results were handily readily accepted by the respective teams of editors of Göttscche-Lötzsch & Henk and Henk & Henk. These papers [37] established that relatively strong axiom systems that satisfied Equation (2)'s Type-A requirement could verify their own semantic tableau consistency. The latter article [38] indicated some Type-NS systems could corroborate their own Hilbert consistency and speculated whether a proof of  $P \neq NP$  could be established via a stepped-up version of these results.

The Sections 7-10 of [38] briefly outlined the approximate reasons for [38]'s core conjecture. The 18 years that had separated the years 1994 from 2014 obviously indicates that much additional work was needed to complete this task.

Dan E. Willard  
Notary Public, State of New York  
Qualified in Albany County  
Commission Expires 01/22/2016

No. 01SH6181014

5-10-14

Signature of Notary Public

It turns out that there are many different variants of Type-NS arithmetic that can verify different forms of their own Hilbert consistency. It was not until the year 2000 [47] that we would introduce in [47] two versions of Type-NS self-verifying arithmetics that could recognise their own Hilbert consistency. The remainder of this section will re-weigh summarise our results. Our earlier 1993 and 1996 announcements in [37, 38] are of historical interest, but only [47]'s more mature year-2000 formalisms will be central for proving  $P \neq NP$ . During our discourse, the term "Introspective Semantics" will refer to the analog of our self-justifying axiom systems from Definition 0.3 and Example 0.4 that have the same names always begin with the acronym "IS", as useful reader.

In a position where  $\beta$  is an initial axiom system, that contains at least the logical power of Peano Arithmetic (PA), the two main variants of introspective systems from [47] carried the names of "ISCE( $\beta$ )" and "ISINF( $\beta$ )". Both were Type-NS self-justifying formalisms that could recognise their own Hilbert consistency and also prove all the  $\Pi_1^0$  theorems that  $\beta$  could prove under a slightly modified language, where  $\text{Add}(x, y, z)$  represents 3-way predicate symbol that formalises addition and multiplication. The remainder of this section will offer a encapsulated summary of ISCE( $\beta$ )'s and ISINF( $\beta$ )'s properties and discuss [47]'s related generalisation of the Second Incompleteness Theorem, which demonstrates that these two formalisms are near-maximal.

The formalism ISCE( $\beta$ ) shall avoid using Equation (1)'s Type-3 axiom sentence, declaring successor is a total function, by instead employing an infinite number of built-in constant symbols  $C_0, C_1, C_2, C_3 \dots$  for defining the set of positive integers. The constant symbols  $C_0$  and  $C_1$  will represent the integers of 0 and 1. Each other  $C_j$  will be defined to represent the quantity  $2^{j-1}$ . These integers will be defined via essentially an "Additive Naming Convention" indicating  $C_{j+1} = C_j + C_j$ . Since ISCE( $\beta$ ) technically does not contain an additive function symbol, its definition of  $C_{j+1}$  will formally rest upon using Equation (4)'s 3-way addition predicate symbol:

$$\text{Add}(C_j, C_j, C_{j+1}) \quad (4)$$

We will assume the "name" of the built-in constant symbol " $C_j$ " is encoded using  $O(\log(j+2))$  bits. The advantage of using an "additive naming convention", which assigns names to only integers which are powers of 2, is that its methodology will nicely ensure that the powers of 2 have integer names that are shorter than the lengths of their binary encodings.

Since the ISCE( $\beta$ ) system will contain a built-in function symbol for representing integer-subtraction as a total function (where  $x - y$  is defined to be equal to zero when  $x < y$ ), Equation (4)'s additive naming convention can clearly define any integer that is not a power of 2. For instance since  $10 = 16 - 4 - 2$ , the integer 10 can be represented as " $C_4 - C_2 - C_1$ ". A detailed definition of ISCE( $\beta$ ) is provided in [47]. It uses an analog of Example 0.4's "SelfRef" axiomatic sentence to corroborate its own consistency. The Theorem 3 of [47] indicated ISCE( $\beta$ ) is a self-justifying formalism that verifies its own Hilbert consistency. In essence, ISCE( $\beta$ ) avades the Padiák-Solovay variant of the Second Incompleteness Theorem because it is a Type-NS system. At the same time, the Theorem 4 from [47] demonstrated that if one used an alternate naming convention, where say  $C_j^*$  equals  $2^{j-2}$  when  $j \geq 2$ , then the force of the Second Incompleteness Theorem would return, even though the said formalism is a Type-NS arithmetic. In particular, [47]'s generalisation of the Second Incompleteness Theorem will apply to settings where one replaces Equation (4)'s additive naming convention with (6)'s "Multiplicative Naming Convention":

$$\text{Mult}(C_j^*, C_j^*, C_{j+1}^*) \quad (5)$$

It turns out that if a reader wishes to glance at our article [47], then he can afford to entirely omit its Section 5 and Theorems 4 & 5. This is because the latter, unlike [47]'s Theorem 3, are unrelated to the proof of  $P \neq NP$ . Indeed, I would recommend that readers, interested in mainly NP's characterisation, initially omit [47]'s Section 5 and its Theorems 4 & 5 because the latter involve a complicated proof that is ultimately unrelated to NP's fundamental properties.

A peculiar aspect of [47] is that it does contains one other result, that we recently discovered to be crucial for proving  $P \neq NP$ , although we previously presumed Theorem 6 was too specialised for it to have much significance. Its awkward-but-useful formalism involves the following definition:

**Definition 0.5** Let us assume that  $\alpha$  is an axiom system that contains a Predecessor function symbol, where  $\text{Pred}(a) = \text{Min}(a - 1, 0)$ , as well as contains a constant symbol  $C_1$  for representing the value "1". Also, let  $\text{Pred}^k(a)$  denote a functional operation that consists of  $j$  iterations of such a predecessor function (e.g. this notation implies that  $k$  is the unique integer that satisfies the identity  $\text{Pred}^{k-1}(k) = C_1$ ). Then  $\alpha$  will be said to possess Infinite Far Reach iff there exists some finite subset of  $\alpha$ 's set of proper axioms, called say  $\gamma$ , such that for every integer  $b$  the finite system  $\gamma$  is capable of proving (6)'s invariant (that intuitively states the integer quantity  $\gamma$  does exist).

$$3 = \text{Pred}^{b-1}(a) = C_1 \quad (6)$$

The Theorem 6 of [47] shows it is possible to construct awkward-but-viable self-justifying arithmetics with infinite far reach that are capable of corroborating their own Hilbert consistency and proving the validity of all of Peano Arithmetic's  $\Pi_1^0$  theorems, using again a slightly modified language that treats addition and multiplication as 3-way relations (rather than as function primitives). In particular, [47]'s ISINF formalism can achieve this property without violating the Padiák-Solovay version of the Second Incompleteness Theorem because it is a Type-NS formalism, specifically incapable of verifying any of the operations of SAT. Additionally, its identification are total functions; moreover, this ISINF class of formalisms has a fundamentally different axioms from [47]'s alternate ISCE variant of self-justifying logic because the latter uses an infinite number of different instances of the axiom schema (4) to construct the full infinite range of integers.

Thus while many aspects of [47]'s ISINF style formalism are awkward and super-impractical (due to the unwieldy long proofs it produces), this framework differs from the alternative ISCE mechanism by, at least, showing that self-justifying axiom systems with Infinite Far Reach are theoretically capable of confirming their own Hilbert consistency.

The intuition behind our proof of  $P \neq NP$  is that we noticed that if this invariant was false then was then a hybridisation of the ISCE and ISINF frameworks will produce a contradiction (showing that the invariant  $P = NP$  cannot possibly hold). In particular assuming that there is available an algorithm  $a$  that can solve length-n SAT problems on a Turing machine in say  $n^b$  time, one can apply the ordered pair  $(a, b)$  to develop two types of self-justifying axiom systems, called "Is.Bulky" and "Is.Super.Bulky", that are natural hybrids of the ISCE and ISINF frameworks, with the following pleasing combination of properties:

- A natural generalisation of the techniques, used to prove [47]'s Theorems 3 and 6, will imply that both Is.Bulky and Is.Super.Bulky are efficient and consistent, when one uses the ordered pair  $(a, b)$  to construct the axiomatic input  $\beta$  for Is.Bulky and Is.Super.Bulky frameworks.
- In contrast to Item I's positive result, a generalisation of Gödel-style diagonalisation argument will imply both of Is.Bulky and Is.Super.Bulky are inconsistent (essentially because they are too efficient to escape the reach of the Second Incompleteness Theorem).

Our proof of  $P \neq NP$  will essentially rest on the fact that Items I and II are incompatible with each other whenever an ordered pair  $(a, b)$  formally represents a polynomial solution for SAT problems. (We will, thus, derive the result  $P \neq NP$ , via a proof by contradiction, that shows an ordered pair  $(a, b)$  would otherwise produce two incompatible results.)

**Cheerful News about Misingled Mirages** As earlier section of this article indicated that there existed several types of mirages, which would tend to make most computer researchers and logicians overlook the essence of the proof of  $P \neq NP$ . Some cheerful news, that will be now revealed, is that our proof of  $P \neq NP$  becomes much easier to conceptualise, intuitively, when one is aware of two mirages that can be initially quite misleading.

The first mirage concerns the question about whether Gödel's Second Incompleteness Theorem is 100% rather than 99% ubiquitous, when it describes the inability of formalisms to corroborate their own consistency. A theme of our research in [37]-[40] is that the the Second Incompleteness Theorem is obviously 100% technically correct, as well as 100% germane to conventional mathematical paradigms. However, our research perspective has been that seemingly tiny 1% categories of boundary-case exceptions can be helpful in solving many open questions, including  $P \neq NP$ .

**REMOVE** However, we view the tiny 1% category of boundary-case exceptions as also revealing. In a sense the same passion to examine that final seemingly minuscule tiny remaining 1% category of possibilities, that led to the Friedman-Willard Fusion Treat at the 1990 STOC and FOCS conferences [7, 7], will now lead to a proof of  $P \neq NP$ .

The second surprising mirage also caught this author off guard. It concerns [47]'s "ISINF" axiomatic framework, and our APAL-2006 paper had called "awkward" and "artificial". The intuition behind this label, originally, was that ISINF was so fully inefficient and purely artificial in its defining structure that one can confidently predict THAT IT IS INCAPABLE of ever producing any useful result.

Such italicised and bold-faced words, which I certainly did not dare insert into [47]'s published text, are, technically, accurate. Yet in April of 2014, I discovered, that they were also a mirage! This is because the statement " $P \neq NP$ " is a negative assertion, discussing the inability of polynomial time algorithms to simulate a non-deterministic Turing machine. The latter class of negational statements have quite different semantics and logical properties than positive-styled theorems!

Thus although ISINF's anatomy is unlikely to produce any useful strictly positive results, it is a useful vehicle in a proof by contradiction, which shows that a hybridisation of the ISINF and ISCE frameworks make it impossible for  $P = NP$  to hold (e.g. see Items I & II on the prior page).

(via the 2-part argument that was roughly summarised by the Items I and II, earlier in this section).

**NEW Section's Title Added Notation and Further Intuition**

Imply the combination of Item I's efficiency results, together with prohibitions imposed by the Incompleteness Theorem will imply that either and both of the either  $P \neq NP$  or any other fundamental fact.

d4d2.3

For any  $(\alpha, d)$ , it is easy to construct a system  $\alpha^d \supseteq \alpha$  that satisfies the Part-I condition. For instance,  $\alpha^d$  could consist of all of  $\alpha$ 's axioms plus an added "SelfRef( $\alpha, d$ )" sentence, defined as stating:

This article will confirm the validity of Gödel's second  $P \neq NP$  conjecture, a topic whose foundational implications have been discussed also by Karp and Levin. We had conjectured at a 1996 DIMACS Workshop, co-chaired by Sam Buss and Paul Beame, that a proof of  $P \neq NP$  could be obtained by somehow hybridising the mechanics of Gödel's Diagonalisation Arguments with those of the unusual boundary-case exceptions pertaining to Gödel's Second Incompleteness Theorem, that we have called "Self-justifying arguments".

The above conjecture was, of course, highly counter-intuitive because Gödel's First Incompleteness Theorem (which asserts the inherent undecidability of conventional arithmetic) and his Second Incompleteness Theorem (which establishes the inability of conventional logic to corroborate their own consistency) are both known to be amazingly robust and resilient twin results. However after researching this topic during the last essentially 514 years (and publishing six papers about it in the Journal of Symbolic and Applied Logic), we have now, finally, developed a method to demonstrate  $P \neq NP$  by means of a proof-by-contradiction that exploits the inherent tension and divergent properties that separate the seemingly minuscule-but-valid boundary-case exceptions for Gödel's Second Incompleteness Theorem from sophisticated generalisations of Gödel-style diagonalisation arguments.

Our discourse will presume that the reader has a sufficient knowledge about Mathematical Logic to understand the main themes in at least one of the graduate-level textbooks by say Budzianowicz, Mendelson or Papadimitriou [7, 7, 7].

**NEW Section's Title NEED**

Throughout this article,  $\alpha$  will denote an axiom system, and  $d$  will denote a deduction method. An ordered pair  $(\alpha, d)$  will be called Self Justifying when:

- one of  $\alpha$ 's theorems will state that the deduction method  $d$ , applied to the system  $\alpha$ , will produce a consistent set of theorems, and
- the axiom system  $\alpha$  is in fact consistent.

For any  $(\alpha, d)$ , it is easy to construct a system  $\alpha^d \supseteq \alpha$  that satisfies the Part-I condition. For instance,  $\alpha^d$  could consist of all of  $\alpha$ 's axioms plus an added "SelfRef( $\alpha, d$ )" sentence, defined as stating:

• There is no proof (using  $\alpha$ 's deduction method) of  $0 = 1$  from the union of the axiom system  $\alpha$ , with this sentence " $\text{SelRef}(\alpha, d) \wedge \text{looking at itself}$ ".

Kleene [17] discussed how to encode approximate analogs of this "SelRef( $\alpha, d$ )" statement. Each of Kleene, Rogers and Jeroslow [17, 21, 16] noted  $\alpha^d$  may, however, be inconsistent (despite  $\text{SelRef}(\alpha, d)$ 's assertion).

This problem arises in settings more general than Gödel's paradigm, where  $\alpha$  was an extension of Peano Arithmetic. There are many settings where the Second Incompleteness Theorem does generalize [1, 2, 3, 4, 5, 6, 8, 10, 11, 12, 13, 14, 15, 16, 17] and [17, 21, 16] noted  $\alpha^d$  may, however, be inconsistent (despite  $\text{SelRef}(\alpha, d)$ 's assertion). Such result formalizes a paradigm where self-justification is feasible, due to a diagonalization issue. Many logicians have, thus, hesitated to employ a  $\text{SelRef}(\alpha, d)$  axiom because  $\alpha^d = \alpha + \text{SelRef}(\alpha, d)$  is typically inconsistent.

Our research explored special circumstances [10, 15, 16, 17] where it is feasible to construct self-justifying formalisms. These paradigms break down when properties specifying  $a + b = a$  and  $a \cdot b = a$ . Then a logic will be said to recognize multiplication (to avoid the preceding difficulties). To be more precise, let  $\text{Add}(a, b, c)$  and  $\text{Mult}(a, b, c)$  denote two 3-way predicates specifying  $a + b = c$  and  $a \cdot b = c$ . Then a logic will be said to recognize addition and multiplication or Total Functions (if it includes 1-3 as axioms).

$$\begin{aligned} \forall a \exists x \quad & \text{Add}(a, 1, x) \\ \forall a \forall y \exists z \quad & \text{Add}(a, y, z) \\ \forall a \forall y \exists z \quad & \text{Mult}(a, y, z) \end{aligned} \quad \begin{matrix} (7) \\ (8) \\ (9) \end{matrix}$$

A logic system  $\alpha$  will be called Type-M if it contains (1) – (3) as axioms, Type-A if it contains only (1) and (2) as axioms, and Type-S if it contains only (1) as an axiom. Also a system is called Type-NS if it contains none of these axioms. The significance of these constructs is explained by items (a) and (b):

- a. The existence of Type-A systems that can recognize their own consistency under semantic tableau deduction, while proving analogs of all Peano Arithmetic's  $\Pi_1$  theorems (in a slightly different language), was demonstrated in [16]. Also, [17] noted that some specialised forms of Type-NS systems can likewise recognize their own Hilbert consistency.
- b. The above versions of the Second Incompleteness Theorem are known to be near-maximal in a mathematical sense. This is because the combined work of Pudlák, Solovay, Nelson and Wilkie-Pairis [22, 27, 32, 36] implied no natural Type-S system can recognize its Hilbert consistency, and Willard subsequently [11, 18, 19] hybridized their formalisms with some techniques of Adamowics-Bierwirth [1, 2] to establish that most Type-M systems cannot recognize their own semantic tableau consistency.

An interesting aspect of (a) and (b) is that there is a tight match between their positive and negative results from a quantitative perspective, though suggested by Hilbert in  $\alpha + \phi$  and Gödel in  $\alpha$ , which has not been addressed. This remaining gap will be explored by our Theorems 77, 77, 77 and 77.

Gödel's Incompleteness Theorem is a 2-part result. Its first half indicates no decision procedure can identify all the true statements of arithmetic. Its Second Incompleteness Theorem specifies sufficiently strong systems cannot verify their own consistency. Gödel was careful to insert the following caveat into his historic paper [10], indicating a diluted form of Hilbert's Consistency Program could be successful:

"It must be expressly noted that Proposition XI (e.g., the Second Incompleteness Theorem) represents no contradiction of finite means, and there might conceivably be finite proofs which cannot be stated in  $P$  (or in  $M$  or in  $A$ )."

Years ago has summarized, in detail, Gödel's considerations about this subject. Thus [36] indicated that "for several years" after [10]'s publication, Gödel "was cautious not to pre-judge" whether some unusual formalism, different from Peano Arithmetic, might provide some type of proof of its own consistency.

#### OLDER NOTES 18 NOTES

On Tuesday March 19 while going to Dr. Salter's office, I corrected my old proof of the Pudlák conjecture to make it work. I also consider it possible, albeit not likely, that it could lead to a proof of Cook's conjecture. I can prove Pudlák's conjecture both for Atai-style pipedream growth functions and their generalization for sum functions that maps powers of constant  $c$  onto other powers of  $c$ . (Indeed, it should work for numbers of the form  $\sum P(i)$  where  $P(i)$  is a polynomial of  $i$ .) This finishes my old paper about Pudlák's conjecture about generalized Atai operation, which I will henceforth call "opaque" functions. A function that full defines all values of  $P(i)$  will be a "completely transparent" function.

It is best to prove Pudlák's conjecture by using a modified form of semantic tableau, called X-tableau, that allows for all instances of the Law of Excluded Middle to be used as axioms. All other Prolog-style proofs can reduce to this method except that the Cook-Li axioms are replaced by Mendelson's first three logical axiom schema. Also we will call a proof a Normalized Tableaux proof relative to an axiom  $\Phi$  if  $\Phi$  appears only

once in the proof tree and nowhere else.

Only for the sake of notational simplicity, will we assume our semantic trees are so normalized and that they have no growth function symbol(s.e.g. they represent a function  $F(x_1, x_2, \dots, x_j)$  as a  $j+1$  way relation  $R_f(x_1, x_2, \dots, x_j, y)$  where  $y$  is the function's output).

Then all the results of my old paper at the GWT Math conference will nicely generalize. I think my new method will called something like Judicialized Blind Proving. If such a X-tab "candidate tree" has  $n$  nodes then we will say it has an L-model "M" with  $L = n + 2$  when  $M$  is a model of the natural numbers that assumes that all integers  $\leq L + 1$  exists and our opaque Atai-like function  $P$  satisfies the following rules

$$V = L \quad F(n) \leq L + 1 \quad (11)$$

Moreover, we will say a node  $v$  in the candidate tree is L-consistent in a 0-blid if there exists some L-Model M where the sentence  $v$  and all its ancestors with the possible exception of the true root  $\top$  satisfy the model M under an interpretation where each of the constants C in this proof tree assume values  $\leq L$ . The magic bullet, used in the proof is the observation that if G is the global number of the minimal proof of absurdity in our self-verifying axiomatic system, whose "I am consistent" declaration is stored in its root node  $\top$  then at least one top down tree branch ending at node v will be L-consistent in a 0-blid model where  $L$  is a number less than any  $\frac{G}{2}$ . (Indeed with a little extra work I will not do now, this proof will work when  $L$  equals  $\text{Log}(G)$ .) The point is that I know how to make this methodology produce my needed proof by induction.

Moreover, I can generalize my method in many directions. One lesson is that it is probably best for my paper to concentrate on the simplest version of Atai function that are analogous to the incremental naming convention, and then to proceed to the additive naming convention at the end of the article (which obviously will also support my method when one walks through powers of 2 because I have enough wiggle room available). Other terms are listed below.

1. Method generalizes when a proof is defined not by one integer p but rather a k-tuple of such integers
2. If  $\Phi$  is 0 then this implies that my method will know that there exists no simultaneous proof of a sentence and its negation (both encoded in prenex normal form for simplicity).
3. Also, it will apply to many versions of the tameness reflection principle, including the Tarski-Rest principle and stronger

Somewhere I will want to use the words from 2 days ago of "blind specimen as starting with a permutation  $p$  of  $L$  integers, blind case, blind specimen, Starting Blind Set, and make Reduced Blind SET".

**NOTES ADDED MARCH 25** afternoon (starting at 1:30pm after exercise for 75 minutes at YMCA) I saw Nilsen on Friday at 11am, and she told me that her guess was that a co-author would make the writing of the paper harder (rather than easier). (Therefore, I will not ask AG for help in writing the paper, but I might ask him to read it and think it through in the Acknowledgement, similar to my prior papers.) Also, I went over the basic idea on Sunday morning and I convinced myself they were correct. (Interestingly, I could not reconstruct my prior idea without first reading my March 21 notes (below) and then filling in the missing pieces. This is partly because I took a 2-day rest from my new project and returned to the old Waller project.) Another new point is that I decided to call my core idea "Opaque Pseudo-Model Analysis (OPMA)". It rests on the idea of taking an X-tableau proof where only axioms you derive in  $\Omega_1$  sentence called  $\Phi$  and making  $\Phi$  the root of the proof tree (appearing nowhere else. Then you entertain the possibility that  $\Phi$  is false in the Standard Model. In that case you can create a model consisting of all  $\alpha$ 's axioms except that  $\Phi$  is replaced by  $\neg\Phi$ . The call of OPMA analysis is to show that the existence of the resulting model is contradictory, which is why I "judicialize" call it a pseudo-model. The trick is to use the combination of the false presumption of the existence of the pseudo-model M to construct a path  $p$  through the proof tree where all nodes, except "its root"  $\Phi$  hold in this model M. That is impossible because the proof tree has each node end with a contradiction. (The contradiction essentially arises because the bottom nodes of the proof tree, being  $\Delta_0$  sentences stemming from  $\Phi$  must be correct. This is because if  $\Phi$  is untrue of the form  $V = \Phi(a)$  then all constants symbols  $C$  constructed by the proof are small enough to have  $\Phi(a)$  be true when we assume our proof tree T is MINIMUM CONTRADICTORY PROOF.) I wonder if this technique is related to Percing in Cohen's discussion of the Axiom of Choice?

**NOTES ADDED APRIL 11** (11:45 am, running tomorrow to Washington DC to join my son Robert and celebrate my wife's birthday. I am certain my new method will show that if  $P \neq NP$  then there will exist a transfinite sequence of proofs of infinite length, and they likely can also prove  $P \neq NP$  itself. The diagonalization sentence needed to be this is the statement: "No short proof of these create a short proof of this sentence" (This sentence will be called  $\psi$ ).

By short proof  $p$  we mean a proof  $p$  having the property that a temporally-hypothesized algorithm  $\alpha$  for solving SAT problems it time  $b^k$  cannot find in  $p$  units of time existence of a proof of  $0=1$  from our self-referencing axiomatic system and ALSO having the minor additional constraint (from temporally-hypothesized LogPhi) is greater than some trivially specified constant K.

I will use the term ISWISH(a,b,A) to refer temporally to my analog of the old ISIMP(A) system in this short note. Its name will be later presumably changed in the future to perhaps "Is\_Problematic" or better yet "Is\_Troubling" or "Is\_Vexing". It will need to prove  $P \neq NP$  by showing that it displays contradictory behaviour when  $(a,b)$  is a generic method for solving SAT problems in polynomial time. It will differ from ISIMP in the following respects:

A) It is far more complex and much matters simpler for ISWISH(a,b,A) to include the axioms of the multiplicative naming convention.

B) The Group-3 axioms of IsWish and include the self-referencing statement that on each occasion when it is proven that our (a,b) methodology in time t is able to demonstrate the non-existence of a proof of  $0 = 1$  of length less than  $L$ , the statement  $\beta^L$  holds.

C) It is possible that IsWISH's group-3 will also create an analog of ISIMP's statement "There is no proof of  $0 = 1$  from me, but I am not sure if that statement is necessary, but I am 99 % certain that it is accurate."

Interestingly I have been going in circles during the last week until I realized the above triple diagonalisation sentence in  $\Omega$  seems to work, while the alternate version below (that has one extra use of the word "no") is false:

"No short proof of there exists NO short proof of this sentence" (e.g. second use of the word "no" should be avoided, and it has been causing me difficulties for many years).

The intention of the word "no" is to be able to establish it by contradiction. Thus if "no" was false then there would be a short proof  $P$  of the inside part of  $\Omega$  (which I will henceforth call  $\Psi$ ). This short proof of  $\Psi$  will be correct. This is because if  $\Psi$  is untrue of the form  $V = \Psi(a)$  then all constants symbols  $C$  constructed by the proof are small enough to have  $\Psi(a)$  be true when we assume our proof tree T is MINIMUM CONTRADICTORY PROOF.) I wonder if this technique is related to Percing in Cohen's discussion of the Axiom of Choice?

Its definition will require using a SAT(a,b) formalism for the sake of contradiction, (defined on April 11). Its Group-3 will use this system to prove larger numbers, and it will use a variant of the additive naming convention that defines a series of fixed constants  $c_1, c_2, \dots, c_{32}$  where the VERY HELPFUL upper bound  $b$  depends on  $(a,b)$ . (This bound was not in my April 11 notes.) It will also contain the new group-3 clause that states if a proof of  $0=1$  of length  $L$  has been found then no integer exists that is larger than say  $\text{Max}(C_m, 2\sqrt{L})$ . (This clause makes it unnecessary for group-3 to include the sentence "There is no proof of  $0=1$  from me" as an axiom because it can make such into a theorem. However, the retaining this as an axiom is probably helpful in simplifying the proof.)

<sup>1</sup> Typically,  $\alpha^d = \alpha + \text{SelRef}(\alpha, d)$  will be inconsistent because a standard Gödel-like self-referencing construction will produce a proof of  $0 = 1$  from  $\alpha^d$ , even when  $\alpha$  is consistent.

