

The fact that Peano Arithmetic is expressive enough to represent all computable functions means that some of the unsolvable problems from the theory of computation carry over to first-order logic. We know that the theory of computation is undecidable, that is we can formulate problems about computable functions that cannot be decided by an algorithm in a uniform way. Since we have shown how to represent computable functions in logic, the same problems can be formulated in first-order logic as well, which means that first-order logic must be undecidable. We are now in a position to give a unified account of this and two other negative results of logic:

- *Church's Theorem*: the undecidability of logic
- *Tarski's Theorem*: the undefinability of truth
- *Gödel's Theorem*: the incompleteness of systems of arithmetic.

23.1 Gödel Numberings and Diagonalization

The key to all these results is an ingenious discovery made by Gödel in the 1930's: it is possible to effectively enumerate all computable functions in a uniform way (via so-called *Gödel numberings*, and it is possible to define computable functions by *diagonalization* through the table of computable functions. Assumptions about decidability and computability then enable us to define a new computable function that is different from all other computable functions and thus leads to a contradiction. We will briefly sketch both methods.

Definition: A *Gödel numbering* is a mapping from a set of expressions to \mathbb{N} that satisfies the following conditions

- (1) Different expressions receive different *Gödel numbers*. (*injective*)
- (2) The Gödel number of an expression can be effectively calculated. (*computable*)
- (3) It is effectively decidable whether a given number is a Gödel number or not.

It is easy to see that *the language of Peano-Arithmetic has a Gödel-numbering*. As we only use *finitely many different symbols* to write down arithmetic formulas we simply assign a number to every symbol. For instance, we may assign 0 to (, 1 to), 2 to \forall , 3 to \exists , etc. Thus every formula corresponds to a unique finite sequence of numbers that represents the string of symbols needed to write down the formula.

There are multiple ways to *encode sequences of numbers* by natural numbers. One possibility is using polynomials as described in the representation of primitive recursion. Another one is extending the bijective encoding $\langle \rangle$ of pairs of numbers to an encoding $\langle \rangle^3$ of triples, quadruples, etc, and then to a bijective encoding $\langle \langle \rangle \rangle$ of arbitrary lists of numbers as follows:

$$\langle x \rangle^1 \equiv x, \quad \langle x_1, \dots, x_{k+1} \rangle^{k+1} \equiv \langle \langle x_1, \dots, x_k \rangle^k, x_{k+1} \rangle, \quad \langle \langle x_1, \dots, x_n \rangle \rangle \equiv \langle n, \langle x_1, \dots, x_n \rangle^n \rangle$$

All these functions are computable and have computable inverses, which means that given the assignments of numbers to symbols (and the rules for parsing formulas) we can determine whether a number corresponds to a formula and we are able to construct the formula if this is the case.

Diagonalization is a method for deriving contradictions from certain assumptions, which is best illustrated when considering an infinite set of functions. We construct a new function by going diagonally through the table of all function values and modifying the entry at each diagonal point. As a result we get a function that cannot be represented by a row of the table. The assumptions are usually crucial for proving that this diagonal function exists at all.

	0	1	2	3	4	5	6	...
f_0	$f_0(0)$	$f_0(1)$	$f_0(2)$	$f_0(3)$	$f_0(4)$	$f_0(5)$	$f_0(6)$...
f_1	$f_1(0)$	$f_1(1)$	$f_1(2)$	$f_1(3)$	$f_1(4)$	$f_1(5)$	$f_1(6)$...
f_2	$f_2(0)$	$f_2(1)$	$f_2(2)$	$f_2(3)$	$f_2(4)$	$f_2(5)$	$f_2(6)$...
f_3	$f_3(0)$	$f_3(1)$	$f_3(2)$	$f_3(3)$	$f_3(4)$	$f_3(5)$	$f_3(6)$...
f_3	$f_3(0)$	$f_3(1)$	$f_3(2)$	$f_3(3)$	$f_3(4)$	$f_3(5)$	$f_4(6)$...
f_5	$f_5(0)$	$f_5(1)$	$f_5(2)$	$f_5(3)$	$f_5(4)$	$f_5(5)$	$f_5(6)$...
f_6	$f_6(0)$	$f_6(1)$	$f_6(2)$	$f_6(3)$	$f_6(4)$	$f_6(5)$	$f_6(6)$...
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

Cantor's famous theorem about the nondenumerability of the set of total functions from \mathbb{N} to \mathbb{N} (which is equivalent to the nondenumerability of the set of real numbers), for instance, uses this argument. If we assume that all functions from \mathbb{N} to \mathbb{N} are denumerable, then we can define a new function $f: \mathbb{N} \rightarrow \mathbb{N}$ as $f(i) := f_i(i) + 1$. Because of the assumption, f must be some f_j in the table and we get $f(j) = f_j(j) + 1 = f(j) + 1$, a clear contradiction. As a result, the set of functions from \mathbb{N} to \mathbb{N} must be nondenumerable.

Gödel's important modification to that argument was the insight that *diagonalization on computable functions is computable*, provided we use a Gödel-numbering of computable functions. Gödel originally expressed his construction without an explicit reference to computable functions (there was not yet a developed theory of computation). Instead he formulated it entirely within logic as an encoding of *self-reference* within formulas, which is the key issue in a diagonalization.

Definition: Let $\lceil X \rceil$ denote the representation of the Gödel number of the formula X (in which x may occur free). The *diagonalization of X* is the formula $(\exists x)(x = \lceil X \rceil \wedge X)$.

Lemma 1: Diagonalization is computable: there is a computable function *diag* such that $n = \lceil X \rceil$ implies $\text{diag}(n) = \lceil (\exists x)(x = \lceil X \rceil \wedge X) \rceil$, that is *diag*(n) is the Gödel number of the diagonalization of X whenever n is the Gödel number of the formula X .

Proof sketch: Given a number n we can effectively determine whether it is a Gödel number of some formula X . Once we have X we can construct the formula $(\exists x)(x = \underline{n} \wedge X)$ and compute its Gödel number.

The difficult part of the actual proof is recasting the argument so that it deals with natural numbers only. One needs a specific Gödel-numbering for this purpose.

Diagonal Lemma: If \mathcal{T} is a theory in which *diag* is representable, then for any formula $B(x)$ with exactly one free variable x there is a formula G such that $\models_{\mathcal{T}} G \Leftrightarrow B(\lceil G \rceil)$.

Proof: Assume that D represents *diag* in \mathcal{T} and let F be the formula $(\exists y)(D(x, y) \wedge B(y))$. Choose $G \equiv (\exists x)(x = \ulcorner F \urcorner \wedge F)$ as the diagonalization of F and let n and g be the Gödel numbers for F and G , respectively. By definition of *diag* we know $\text{diag}(\ulcorner F \urcorner) = \ulcorner G \urcorner$ and thus $D(\underline{n}, \underline{g})$ must be valid in \mathcal{T} .

Furthermore $G \equiv (\exists x)(x = n \wedge (\exists y)(D(x, y) \wedge B(y)))$ is logically equivalent to the formula $(\exists y)(D(\underline{n}, y) \wedge B(y))$. Because of the functionality of D and the validity of $D(\underline{n}, \underline{g})$ this formula is equivalent to $D(\underline{n}, \underline{g}) \wedge B(\underline{g})$, which in turn is equivalent to $B(\underline{g})$.

Thus G is logically equivalent to $B(\ulcorner G \urcorner)$ in \mathcal{T} and hence $\models_{\mathcal{T}} G \Leftrightarrow B(\ulcorner G \urcorner)$. \square

Note that the **diagonal lemma holds for Peano Arithmetic**, as *diag* is representable in any theory that can represent the computable functions.

23.2 Incompleteness Results

The diagonal lemma shows that in theories that can represent computability all formulas have a fixed point. Fixed point constructors, on the other hand, lead to inconsistencies, as they make it possible to define formulas that are equivalent to their own negation. Before we prove this, let us introduce some useful notation.

Definition: Let \mathcal{T} be a theory

- (1) \mathcal{T} is called *consistent*, if there is no theorem in \mathcal{T} whose negation is also in \mathcal{T} .
- (2) \mathcal{T} is called *complete*, if for every formula X in the language of \mathcal{T} either X or $\sim X$ is a theorem in \mathcal{T} .
- (3) A set S of formulas is called *decidable* if the set of Gödel numbers of S is decidable, i.e. if the characteristic function of that set is computable.
- (4) \mathcal{T} is called *axiomatizable*, if there is a decidable subset of \mathcal{T} whose logical consequences are exactly the theorems of \mathcal{T} . \mathcal{T} is *finitely axiomatizable* if it is axiomatizable with a finite set of axioms.
- (5) A set $S \subseteq \mathbb{N}$ is called *definable in \mathcal{T}* if there is unary predicate R_S in the formal language of \mathcal{T} , such that for all $y \in \mathbb{N}$: $y \in S$ implies $\models_{\mathcal{T}} R_S(\underline{y})$ and $y \notin S$ implies $\models_{\mathcal{T}} \sim R_S(\underline{y})$.
Definability extends naturally from sets (i.e. unary relations) to n -ary relations.
- (6) A theory \mathcal{T}' is an *extension of \mathcal{T}* if $\mathcal{T} \subseteq \mathcal{T}'$, i.e. if every theorem in \mathcal{T} is also one in \mathcal{T}' .
- (7) The theory of *arithmetic* is the theory over the language $\mathcal{L}(=, +, *, 0, 1)$ whose theorems are the formulas that are true in the *standard interpretation* $\langle \mathbb{N}, =, +, *, 0, 1 \rangle$.

Note that a set S is *definable in \mathcal{T} if its characteristic function f_S is representable in \mathcal{T}* . For if $R_f(x, y)$ represents f_S , then $n \in S$ implies $f_S(n) = 1$ and thus $\models_{\mathcal{T}} R_f(\underline{n}, \underline{1})$, and $n \notin S$ implies $0 = f_S(n) \neq 1$ and thus $\models_{\mathcal{T}} \sim R_f(\underline{n}, \underline{1})$. So the predicate $R_S(x) \equiv R_f(x, \underline{1})$ defines S .

An immediate consequence of this insight is that *every decidable set of numbers is definable in Peano arithmetic* and in any other theory that can represent the computable functions, as the characteristic function of every decidable set is computable.

Representability and definability is preserved under extensions. In fact, if f is representable in \mathcal{T} and \mathcal{T}' extends \mathcal{T} then f is representable in \mathcal{T}' by the same predicate. Arithmetic is an extension

of Peano Arithmetic since the Peano axioms are true in both standard and non-standard interpretations, so representability and definability carries over from Peano arithmetic to arithmetic. However, there are certain kinds of functions and sets that cannot be represented anymore once a theory is expressive enough to represent the computable functions. A key observation, for instance, is that such *a theory cannot define its own theorems* by a formula.

Lemma 2: If \mathcal{T} is a consistent theory in which *diag* is representable, then the set of Gödel numbers of \mathcal{T} -theorems is not definable in \mathcal{T} .

Proof: Assume that GN defines the set of Gödel numbers of \mathcal{T} -theorems in \mathcal{T} . By the diagonal lemma, there must be a formula G such that $\models_{\mathcal{T}} G \Leftrightarrow \sim GN(\ulcorner G \urcorner)$. We show that both G and $\sim G$ are \mathcal{T} -theorems, which contradicts the consistency of \mathcal{T} .

Assume G is not a \mathcal{T} -theorem. Then $\models_{\mathcal{T}} \sim GN(\ulcorner G \urcorner)$ by definition of GN and thus $\models_{\mathcal{T}} G$ because of the above equivalence. This, in turn, means that G is a \mathcal{T} -theorem. Thus by definition of GN we know $\models_{\mathcal{T}} GN(\ulcorner G \urcorner)$. Because of the above equivalence $\models_{\mathcal{T}} \sim G$ must hold. \square

The proof of Lemma 2 mimics in logic what Cantor's argument did to functions on natural numbers. The assumption that the predicate GN is definable corresponds to the assumption that we can construct a function f via $f(y)=f_y(y)$. Modifying $GN(y)$ into $\sim GN(y)$ captures the key idea behind defining $f(y)=f_y(y) + 1$, which makes f different from every f_y . The contradiction $G \Leftrightarrow \sim G$ then follows for the fixed point G of $\sim GN$ – while in Cantor's theorem it is stated as $f(j)=f(j) + 1$ for the index j of f .

An immediate consequence of Lemma 2 is that it is impossible to represent the arithmetically true statements by a formula in arithmetic.

Theorem: [Tarski's Indefinability Theorem] *Arithmetical truth is not arithmetically definable*, i.e. the set of Gödel numbers of formulas that are true in natural number arithmetic are not definable in the theory of arithmetic.

If a theory \mathcal{T} is decidable, then there is an effective method for deciding whether a given formula is a theorem in \mathcal{T} . One only has to calculate the Gödel number of that formula and test whether the value of the corresponding characteristic function is 1. Conversely, if a theory is undecidable, then there cannot be an effective method for checking whether a given formula is a theorem or not.

Theorem 1: *No theory that can represent the computable functions is decidable.*

Proof: Assume that \mathcal{T} is decidable and can represent the computable functions. Then the characteristic function of the set of Gödel numbers of \mathcal{T} is representable, which means that the set of Gödel numbers of \mathcal{T} is definable. Because of Lemma 2 this cannot be the case. \square

Corollary: *Arithmetic is undecidable.*

An immediate consequence of Theorem 1 is *Church's undecidability theorem*, which states that *first-order logic is undecidable*. We can easily reduce the decidability problem for first-order logic to the decidability problem for theories that can represent the computable functions, provided we find a finitely axiomatizable theory that has this property. Since the theories we considered so far are not finitely axiomatizable, we postpone the proof of Church's undecidability theorem until we have introduced the theory \mathcal{Q} .

The following lemma provides the key argument for Gödel’s first Incompleteness Result.

Lemma 3: Every axiomatizable and complete theory is decidable.

Proof: Let \mathcal{T} be an axiomatizable complete theory and S be its set of axioms.

If \mathcal{T} is inconsistent, then for some formula X both X and $\sim X$ are theorems of \mathcal{T} , thus consequences of the axioms. As a result, every formula in the language of \mathcal{T} is a logical consequence of the axioms, which means that \mathcal{T} is the set of all formulas, hence decidable.

Assume that \mathcal{T} is consistent. Since S is decidable, we can determine whether a given derivation of a formula from the axioms constitutes a valid logical inference. Therefore, it is possible to *enumerate* all valid derivations in \mathcal{T} and all logical consequences of the axioms.

Now let X be an arbitrary formula in the language of \mathcal{T} . Since \mathcal{T} is complete, the enumeration of the logical consequences of the axioms in S will yield either X or $\sim X$ after finitely many steps. Thus we can decide whether X is a theorem in \mathcal{T} or not. \square

Theorem: [Gödel’s First Incompleteness Theorem] There is no consistent, complete, and axiomatizable theory that can represent the computable functions.

Proof: This is an immediate consequence of Lemma 3 and the fact that no theory that can represent the computable functions is decidable (Theorem 1). \square

Corollary: Arithmetic is not axiomatizable.

Gödel’s incompleteness theorem is often described as “*any consistent and sufficiently strong formal theory of arithmetic is incomplete*”, where a formal theory is viewed as one whose theorems are derivable from an axiom system. For such theories there will always be formulas that are true (for instance, in the standard interpretation of arithmetic) but not theorems of the theories. When it comes to arithmetic, *truth and provability are in no sense the same*.

The above formulation is a sharpening of that statement, as it specifies “sufficiently strong” as “capable of representing computable functions”. As we will see later, this includes theories that are much weaker than Peano Arithmetic.

Gödel’s incompleteness result has a strong impact on the development of formal reasoning environments for mathematics. It is impossible to develop a consistent and complete proof system for any fragment of mathematics that includes arithmetic. As formal systems necessarily need to have decidable sets of axioms (even if we use axiom schemes), we have to sacrifice completeness, which means that there will always be theorems that are true but cannot be proven in the proof system.

On the other hand, one has to be aware that the incompleteness result has little effect on *practical theorem proving*. Most, if not all theorems in practical mathematics do not involve the issue of self-application and will therefore remain provable in a formal system.

Undecidability, however, is a bigger problem, as it shows the limitations of fully automated proof systems. While it is possible to find a proof for almost all (valid) practical problems in a formal proof system, the undecidability of first-order logic makes it impossible to check whether a proof attempt fails or just needs more time to complete. For this reason, all proof systems that aim at formalizing large fragments of mathematics and computer science (like Cornell’s Nuprl system) proceed interactively with some controlled automated support for decidable sub-problems.

The undecidability of logic, the undefinability of arithmetical truth, and the incompleteness of systems of arithmetic are not the only unsolvable problems for expressive logics. A fourth issue involves the *undefinability of provability*: it is not possible to describe a predicate $Prov$ that represents provability in a theory \mathcal{T} such that $\sim Prov(\ulcorner \text{false} \urcorner)$ is a theorem in \mathcal{T} . We call a predicate $Prov$ a *provability predicate for \mathcal{T}* if it satisfies the following conditions for all formulas X and Y .

- If $\models_{\mathcal{T}} X$ then $\models_{\mathcal{T}} Prov(\ulcorner X \urcorner)$
- $\models_{\mathcal{T}} Prov(\ulcorner X \supset Y \urcorner) \supset (Prov(\ulcorner X \urcorner) \supset Prov(\ulcorner Y \urcorner))$
- $\models_{\mathcal{T}} Prov(\ulcorner X \urcorner) \supset Prov(\ulcorner Prov(\ulcorner X \urcorner) \urcorner)$

The first condition states that every theorem should be provable, the second that the modus ponens holds for provability, and the third that provability is provable. Note that the second and third conditions are stronger than the first in the sense that the implication itself must be a theorem in \mathcal{T} .

Note that a condition like $\models_{\mathcal{T}} Prov(\ulcorner X \urcorner) \supset X$ is not included in the definition, since this requirement cannot be satisfied unless \mathcal{T} is inconsistent. In fact, **Löb's Theorem** shows that this condition implies $\models_{\mathcal{T}} X$.

Theorem: [**Löb's Theorem**] If $Prov$ is a provability predicate for a theory \mathcal{T} that can represent the computable functions then $\models_{\mathcal{T}} Prov(\ulcorner X \urcorner) \supset X$ implies $\models_{\mathcal{T}} X$ for any X .

Proof: Assume $\models_{\mathcal{T}} Prov(\ulcorner X \urcorner) \supset X$. By the diagonal lemma, the formula $Prov(y) \supset X$ must have a fixed point G , i.e. there is a formula G such that $\models_{\mathcal{T}} G \Leftrightarrow (Prov(\ulcorner G \urcorner) \supset X)$. Since $Prov$ is a provability predicate, we know $\models_{\mathcal{T}} Prov(\ulcorner G \supset (Prov(\ulcorner G \urcorner) \supset X) \urcorner)$ and also $\models_{\mathcal{T}} Prov(\ulcorner G \urcorner) \supset Prov(\ulcorner Prov(\ulcorner G \urcorner) \supset X \urcorner)$. and $\models_{\mathcal{T}} Prov(\ulcorner G \urcorner) \supset Prov(\ulcorner Prov(\ulcorner G \urcorner) \urcorner) \supset Prov(\ulcorner X \urcorner)$.

Since $\models_{\mathcal{T}} Prov(\ulcorner G \urcorner) \supset Prov(Prov(\ulcorner G \urcorner))$ we get $\models_{\mathcal{T}} Prov(\ulcorner G \urcorner) \supset Prov(\ulcorner X \urcorner)$.

Because of our assumption we conclude $\models_{\mathcal{T}} Prov(\ulcorner G \urcorner) \supset X$ and $\models_{\mathcal{T}} G$ follows from the definition of G . From this we conclude $\models_{\mathcal{T}} Prov(\ulcorner G \urcorner)$, which finally gives us $\models_{\mathcal{T}} X$. \square

An immediate consequence of Löb's Theorem is Gödel's Second Incompleteness Theorem, which states that provability predicates for consistent predicates cannot be complete.

Theorem: [**Gödel's Second Incompleteness Theorem**] If $Prov$ is a provability predicate for a consistent theory \mathcal{T} that can represent the computable functions then $\not\models_{\mathcal{T}} \sim Prov(\ulcorner \mathbf{0}=\mathbf{1} \urcorner)$.

Proof: Suppose $\models_{\mathcal{T}} \sim Prov(\ulcorner \mathbf{0}=\mathbf{1} \urcorner)$. By definition of negation $\models_{\mathcal{T}} Prov(\ulcorner \mathbf{0}=\mathbf{1} \urcorner) \supset \mathbf{0}=\mathbf{1}$. By Löb's theorem, we then conclude $\models_{\mathcal{T}} \mathbf{0}=\mathbf{1}$, which cannot be because \mathcal{T} is consistent. \square

In other words, no sufficiently expressive consistent theory can prove its own consistency.