Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

# On the hierarchy of natural theories

James Walsh

Sage School of Philosophy
Cornell University

*jameswalsh@cornell.edu*

4/6/22

**Introduction**
Set theory as a case study
The consistency operator
Second-order arithmetic

The starting points are Gödel's incompleteness theorems.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

The starting points are Gödel's incompleteness theorems.

### Theorem (Gödel)

*No reasonable axiomatic theory is complete.*

**Introduction**
Set theory as a case study
The consistency operator
Second-order arithmetic

The starting points are Gödel's incompleteness theorems.

### Theorem (Gödel)

*No reasonable axiomatic theory is complete.*

### Theorem (Gödel)

*No reasonable axiomatic theory proves its own consistency.*

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

The starting points are Gödel's incompleteness theorems.

### Theorem (Gödel)

*No reasonable axiomatic theory is complete.*

### Theorem (Gödel)

*No reasonable axiomatic theory proves its own consistency.*

No axiom system suffices for the development of all of mathematics;
how should we navigate the vast array of axiomatic theories?

**Introduction**
Set theory as a case study
The consistency operator
Second-order arithmetic

The so-called **consistency strength hierarchy** maps out the reasonable axiomatic theories and their relations.

**Introduction**
Set theory as a case study
The consistency operator
Second-order arithmetic

The so-called **consistency strength hierarchy** maps out the reasonable axiomatic theories and their relations.

### Definition

For a base theory $B$, we say that $T \leq_{\text{Con}}^{B} U$ if $B$ proves that the consistency of $U$ implies the consistency of $T$.

**Introduction**
Set theory as a case study
The consistency operator
Second-order arithmetic

The so-called **consistency strength hierarchy** maps out the reasonable axiomatic theories and their relations.

### Definition

For a base theory $B$, we say that $T \leq_{Con}^{B} U$ if $B$ proves that the consistency of $U$ implies the consistency of $T$.

### Definition

$T <_{Con}^{B} U$ if $T \leq_{Con}^{B} U$ and $U \nleq_{Con}^{B} T$.

**Introduction**
Set theory as a case study
The consistency operator
Second-order arithmetic

The so-called **consistency strength hierarchy** maps out the reasonable axiomatic theories and their relations.

### Definition

For a base theory $B$, we say that $T \leq_{\text{Con}}^B U$ if $B$ proves that the consistency of $U$ implies the consistency of $T$.

### Definition

$T <_{\text{Con}}^B U$ if $T \leq_{\text{Con}}^B U$ and $U \not\leq_{\text{Con}}^B T$.

### Definition

$T$ and $U$ are **equiconsistent** over $B$ if $T \leq_{\text{Con}}^B U$ and $U \leq_{\text{Con}}^B T$.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

### Theorem (Folklore)

$<_{\mathsf{Con}}$ is not **pre-linear**, i.e., there are non-equiconsistent $T$ and $U$ such that $T \not<_{\mathsf{Con}} U$ and $U \not<_{\mathsf{Con}} T$.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

### Theorem (Folklore)

$<_{\mathsf{Con}}$ is not **pre-linear**, i.e., there are non-equiconsistent $T$ and $U$ such that $T \not<_{\mathsf{Con}} U$ and $U \not<_{\mathsf{Con}} T$.

### Theorem (Folklore)

The ordering $<_{\mathsf{Con}}$ is **ill-founded**, i.e., there is a sequence $T_0 >_{\mathsf{Con}} T_1 >_{\mathsf{Con}} T_2 >_{\mathsf{Con}} ...$ where each $T_i$ is consistent.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

All known instances of non-linearity and ill-foundedness are ad hoc; they were discovered by applying logical techniques.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

All known instances of non-linearity and ill-foundedness are ad hoc; they were discovered by applying logical techniques.

**Empirical Observation:** The restriction of $<_{\mathsf{Con}}$ to the theories that arise in practice is a *well-ordering*.

$$\mathsf{EA}, \mathsf{EA}^+, \mathsf{PRA}, I\Sigma_n, \mathsf{PA}, \mathsf{ATR}_0, \Pi^1_n\mathsf{CA}_0, \mathsf{PA}_n, \mathsf{ZF}, \mathsf{AD}^{L(\mathbb{R})}$$

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

All known instances of non-linearity and ill-foundedness are ad hoc; they were discovered by applying logical techniques.

**Empirical Observation:** The restriction of $<_{Con}$ to the theories that arise in practice is a *well-ordering*.

$$EA, EA^+, PRA, I\Sigma_n, PA, ATR_0, \Pi_n^1 CA_0, PA_n, ZF, AD^{L(\mathbb{R})}$$

Explaining this contrast is widely regarded as a major outstanding conceptual problem in mathematical logic.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

*The fact that "natural" theories, i.e. theories which have something like an "idea" to them, are almost always linearly ordered with regard to logical strength has been called one of the great mysteries of the foundations of mathematics.*

S. Friedman, Rathjen, Weiermann

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

1. Introduction

2. Set theory as a case study

3. The consistency operator

4. Second-order arithmetic

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Three reasons for discussing set theory.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Three reasons for discussing set theory.

1. Set theory has proceeded in an explicitly axiomatic way since the isolation of ZFC.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Three reasons for discussing set theory.

1. Set theory has proceeded in an explicitly axiomatic way since the isolation of ZFC.
2. ZFC is highly general.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Three reasons for discussing set theory.

1. Set theory has proceeded in an explicitly axiomatic way since the isolation of ZFC.

2. ZFC is highly general.

3. ZFC is insufficient for answering many of the problems that motivated the early development of set theory:
   - The Continuum Hypothesis
   - Projective Measure
   - Suslin's Hypothesis

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Set theorists have investigated a wide array of extensions of ZFC.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Set theorists have investigated a wide array of extensions of ZFC.

- large cardinal axioms

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Set theorists have investigated a wide array of extensions of ZFC.

- large cardinal axioms
- determinacy axioms

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Set theorists have investigated a wide array of extensions of ZFC.

- large cardinal axioms
- determinacy axioms
- forcing axioms

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Set theorists have investigated a wide array of extensions of ZFC.

- large cardinal axioms
- determinacy axioms
- forcing axioms

Can we make rational judgments about the correctness of these principles or their consequences?

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Set theorists have investigated a wide array of extensions of ZFC.

- large cardinal axioms
- determinacy axioms
- forcing axioms

Can we make rational judgments about the correctness of these principles or their consequences?

Steel has promoted the following maxim:

<div align="center">

MAXIMIZE STRENGTH.

</div>

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Steel's Maxim echoes Cantor's dictum of mathematical freedom.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Steel's Maxim echoes Cantor's dictum of mathematical freedom.

The $<_{Con}$ tells us what mathematics can be developed on the basis of one theory rather than another; (more or less) if $Con(T)$ implies $Con(U)$ then $T$ can **interpret** $U$ and **not** vice-versa.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Steel's Maxim echoes Cantor's dictum of mathematical freedom.

The $<_{\mathsf{Con}}$ tells us what mathematics can be developed on the basis of one theory rather than another; (more or less) if $\mathrm{Con}(T)$ implies $\mathrm{Con}(U)$ then $T$ can **interpret** $U$ and **not** vice-versa.

- Poincaré interpreted two dimensional hyperbolic geometry in the Euclidean geometry of the unit circle.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Steel's Maxim echoes Cantor's dictum of mathematical freedom.

The $<_{\mathrm{Con}}$ tells us what mathematics can be developed on the basis of one theory rather than another; (more or less) if $\mathrm{Con}(T)$ implies $\mathrm{Con}(U)$ then $T$ can **interpret** $U$ and **not** vice-versa.

- Poincaré interpreted two dimensional hyperbolic geometry in the Euclidean geometry of the unit circle.
- Dedekind interpreted analysis in set theory.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Steel's Maxim echoes Cantor's dictum of mathematical freedom.

The $<_{\mathsf{Con}}$ tells us what mathematics can be developed on the basis of one theory rather than another; (more or less) if $\mathsf{Con}(T)$ implies $\mathsf{Con}(U)$ then $T$ can **interpret** $U$ and **not** vice-versa.

- Poincaré interpreted two dimensional hyperbolic geometry in the Euclidean geometry of the unit circle.
- Dedekind interpreted analysis in set theory.
- Gödel interpreted proof theory in arithmetic.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Is Steel's Maxim coherent?

Let's consider some sentence $\varphi$ that is independent of ZFC.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Is Steel's Maxim coherent?

Let's consider some sentence $\varphi$ that is independent of ZFC.

1. $\varphi$ increases strength but $\neg\varphi$ does not.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Is Steel's Maxim coherent?

Let's consider some sentence $\varphi$ that is independent of ZFC.

1. $\varphi$ increases strength but $\neg\varphi$ does not.
2. $\neg\varphi$ increases strength but $\varphi$ does not.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Is Steel's Maxim coherent?

Let's consider some sentence $\varphi$ that is independent of ZFC.

1. $\varphi$ increases strength but $\neg\varphi$ does not.
2. $\neg\varphi$ increases strength but $\varphi$ does not.
3. Neither $\varphi$ nor $\neg\varphi$ increases strength.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Is Steel's Maxim coherent?

Let's consider some sentence $\varphi$ that is independent of ZFC.

1. $\varphi$ increases strength but $\neg\varphi$ does not.
2. $\neg\varphi$ increases strength but $\varphi$ does not.
3. Neither $\varphi$ nor $\neg\varphi$ increases strength.
4. Both $\varphi$ and $\neg\varphi$ increases strength.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Is Steel's Maxim coherent?

Let's consider some sentence $\varphi$ that is independent of ZFC.

1. $\varphi$ increases strength but $\neg\varphi$ does not.
2. $\neg\varphi$ increases strength but $\varphi$ does not.
3. Neither $\varphi$ nor $\neg\varphi$ increases strength.
4. Both $\varphi$ and $\neg\varphi$ increases strength.

It turns out that **all four** possibilities are realized; in the fourth case we **cannot** follow Steel's Maxim.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

1. $\varphi$ increases strength but $\neg\varphi$ does not.

2. $\neg\varphi$ increases strength but $\varphi$ does not.

3. neither $\varphi$ nor $\neg\varphi$ increases strength.

4. both $\varphi$ and $\neg\varphi$ increases strength.

When we restrict our attention to natural theories, only the first three possibilities are realized.

This is just to say that natural theories are linearly ordered by consistency strength.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Consider again the axiom systems extending ZFC:

- large cardinal axioms
- axioms of definable determinacy
- forcing axioms

These systems have different motivations, but they are well-ordered by consistency strength.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Consider again the axiom systems extending ZFC:

- large cardinal axioms
- axioms of definable determinacy
- forcing axioms

These systems have different motivations, but they are well-ordered by consistency strength.

They converge on statements about $\mathbb{N}$; in fact, they converge on statements about $\mathbb{R}$.

Introduction
**Set theory as a case study**
The consistency operator
Second-order arithmetic

Consider again the axiom systems extending ZFC:

- large cardinal axioms
- axioms of definable determinacy
- forcing axioms

These systems have different motivations, but they are well-ordered by consistency strength.

They converge on statements about $\mathbb{N}$; in fact, they converge on statements about $\mathbb{R}$.

> *At the level of sentences about $\mathbb{R}$, we know of only one road upward. We are led to it many different ways.*

Steel

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Fix a sufficiently strong, sound, effectively axiomatized theory $T$,
e.g., *elementary arithmetic*, *Peano Arithmetic*, ...

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Fix a sufficiently strong, sound, effectively axiomatized theory $T$, e.g., *elementary arithmetic*, *Peano Arithmetic*, ...

$T$ is incomplete by Gödel's first theorem; $T$ does not prove $\text{Con}_T$ by Gödel's second theorem.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Fix a sufficiently strong, sound, effectively axiomatized theory $T$, e.g., *elementary arithmetic*, *Peano Arithmetic*, ...

$T$ is incomplete by Gödel's first theorem; $T$ does not prove $\mathrm{Con}_T$ by Gödel's second theorem.

Are there any proper extensions of $T$ that are strictly weaker than $T + \mathrm{Con}_T$?

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Rosser introduced a trick whereby we can find sentences strictly weaker than $\mathrm{Con}_{\mathcal{T}}$.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Rosser introduced a trick whereby we can find sentences strictly weaker than $\mathrm{Con}_T$.

$$T \vdash \left( R_T \leftrightarrow \forall x \big( \mathrm{Pf}_T(x, \ulcorner R_T \urcorner) \to \exists y < x \, \mathrm{Pf}_T(y, \ulcorner \neg R_T \urcorner) \big) \right)$$

$R_T$ "says": If there are any proofs of $R_T$, then they are preceded by proofs of $\neg R_T$.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Rosser introduced a trick whereby we can find sentences strictly weaker than $\mathsf{Con}_T$.

$$T \vdash \left( R_T \leftrightarrow \forall x \big( \mathsf{Pf}_T(x, \ulcorner R_T \urcorner) \rightarrow \exists y < x \mathsf{Pf}_T(y, \ulcorner \neg R_T \urcorner) \big) \right)$$

$R_T$ "says": If there are any proofs of $R_T$, then they are preceded by proofs of $\neg R_T$.

We can use Rosser's trick to produce independent sentences strictly weaker than $\mathsf{Con}_T$.

$$\mathsf{Con}_T \vee R_{T + \neg \mathsf{Con}_T}$$

Yet these sentences are highly unnatural.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Rosser introduced a trick whereby we can find sentences strictly weaker than $\mathrm{Con}_T$.

$$T \vdash \Big( R_T \leftrightarrow \forall x \big( \mathrm{Pf}_T(x, \ulcorner R_T \urcorner) \rightarrow \exists y < x \mathrm{Pf}_T(y, \ulcorner \neg R_T \urcorner) \big) \Big)$$

$R_T$ "says": If there are any proofs of $R_T$, then they are preceded by proofs of $\neg R_T$.

We can use Rosser's trick to produce independent sentences strictly weaker than $\mathrm{Con}_T$.

$$\mathrm{Con}_T \vee R_{T + \neg \mathrm{Con}_T}$$

Yet these sentences are highly unnatural.

Self-reference, dependence on a seemingly arbitrary numeration of proofs,...

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Instead of focusing on specific theories, we focus on algorithms for **uniformly** extending theories.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Instead of focusing on specific theories, we focus on algorithms for **uniformly** extending theories.

We are particularly interested in $\mathfrak{g}$ that are *monotone*:

If $T$ proves $\varphi \to \psi$, then $T$ proves $\mathfrak{g}(\varphi) \to \mathfrak{g}(\psi)$.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

There are many monotone algorithms for uniformly extending
theories.

Introduction
Set theory as a case study
**The consistency operator**
Second-order arithmetic

There are many monotone algorithms for uniformly extending
theories.

$\mathfrak{id} : \varphi \mapsto \varphi$

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

There are many monotone algorithms for uniformly extending theories.

$\mathfrak{id} : \varphi \mapsto \varphi$

$\mathfrak{Con} : \varphi \mapsto \mathsf{Con}_T(\varphi)$

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

There are many monotone algorithms for uniformly extending theories.

$\mathfrak{id} : \varphi \mapsto \varphi$

$\mathfrak{Con} : \varphi \mapsto \mathsf{Con}_T(\varphi)$

Rosser's trick engenders an algorithm for extending theories, but it is **not** monotone.

Introduction
Set theory as a case study
**The consistency operator**
Second-order arithmetic

There are many monotone algorithms for uniformly extending theories.

$\mathfrak{id} : \varphi \mapsto \varphi$

$\mathfrak{Con} : \varphi \mapsto \mathsf{Con}_T(\varphi)$

Rosser's trick engenders an algorithm for extending theories, but it is **not** monotone.

Indeed, the Rosser algorithm is not monotone in virtue of the pathological properties flagged earlier.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

The goal is to prove that the consistency operator is the unique
weakest monotone algorithm for uniformly extending theories.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

The goal is to prove that the consistency operator is the unique weakest monotone algorithm for uniformly extending theories.

What does "the unique weakest" mean?

We can make sense of this claim **only** modulo a suitable equivalence relation.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

The goal is to prove that the consistency operator is the unique
weakest monotone algorithm for uniformly extending theories.

What does "the unique weakest" mean?

We can make sense of this claim **only** modulo a suitable
equivalence relation.

Let $\varphi$ be a true sentence. Then the set of sentences that implies $\varphi$
is a **cone**.

$$\{\psi : T + \psi \text{ proves } \varphi\}$$

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

Let's call a function $\mathfrak{g}$ bounded if there exists a $k \in \mathbb{N}$ such that, for every $\varphi$, $\mathfrak{g}(\varphi) \in \Pi^0_k$.

For technical reasons, we restrict our attention to bounded functions.

Introduction
Set theory as a case study
**The consistency operator**
Second-order arithmetic

### Theorem (W.)

*Let $\mathfrak{g}$ be a bounded, computable, and monotone. Then one of the following holds:*

1. *There is a cone $\mathfrak{C}$ such that for all $\varphi \in \mathfrak{C}$, $T + \varphi \vdash \mathfrak{g}(\varphi)$.*

2. *There is a cone $\mathfrak{C}$ such that for all $\varphi \in \mathfrak{C}$,*
   $T + \varphi + \mathfrak{g}(\varphi) \vdash \mathrm{Con}_T(\varphi)$.

Introduction
Set theory as a case study
**The consistency operator**
Second-order arithmetic

### Theorem (W.)

*Let $\mathfrak{g}$ be a bounded, computable, and monotone. Then one of the following holds:*

1. *There is a cone $\mathfrak{C}$ such that for all $\varphi \in \mathfrak{C}$, $T + \varphi \vdash \mathfrak{g}(\varphi)$.*

2. *There is a cone $\mathfrak{C}$ such that for all $\varphi \in \mathfrak{C}$, $T + \varphi + \mathfrak{g}(\varphi) \vdash \mathrm{Con}_T(\varphi)$.*

That is, either $\mathfrak{g}$ is as weak as the identity on a cone or as strong as the consistency operator on a cone.

Introduction
Set theory as a case study
**The consistency operator**
Second-order arithmetic

### Theorem (W.)

*Let $\mathfrak{g}$ be a bounded, computable, and monotone. Then one of the following holds:*

1. *There is a cone $\mathfrak{C}$ such that for all $\varphi \in \mathfrak{C}$, $T + \varphi \vdash \mathfrak{g}(\varphi)$.*

2. *There is a cone $\mathfrak{C}$ such that for all $\varphi \in \mathfrak{C}$,*
   $T + \varphi + \mathfrak{g}(\varphi) \vdash \mathsf{Con}_T(\varphi)$.

That is, either $\mathfrak{g}$ is as weak as the identity on a cone or as strong as the consistency operator on a cone.

The consistency operator is the unique weakest method for uniformly extending theories.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

This contributes to a partial explanation of the well-ordering
phenomenon.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

This contributes to a partial explanation of the well-ordering phenomenon.

It suggests that the iterates of the consistency operator form a spine of axiomatic theories that is, in some sense, canonical.

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

We now shift our attention to second-order arithmetic, the joint
theory of the natural numbers and the real numbers.

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

We now shift our attention to second-order arithmetic, the joint theory of the natural numbers and the real numbers.

$ACA_0$ is our base theory; it is a second-order pendant of PA.

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

The $\Pi_1^0$ formulas are the formulas $\forall x \in \mathbb{N} \ \varphi$ where $\varphi$ is computable.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

The $\Pi_1^0$ formulas are the formulas $\forall x \in \mathbb{N} \; \varphi$ where $\varphi$ is computable.

The $\Sigma_1^0$ formulas are the formulas $\exists x \in \mathbb{N} \; \varphi$ where $\varphi$ is computable.

Introduction
Set theory as a case study
The consistency operator
Second-order arithmetic

The $\Pi_1^0$ formulas are the formulas $\forall x \in \mathbb{N} \; \varphi$ where $\varphi$ is computable.

The $\Sigma_1^0$ formulas are the formulas $\exists x \in \mathbb{N} \; \varphi$ where $\varphi$ is computable.

The $\Pi_1^1$ formulas are the formulas $\forall x \in \mathbb{R} \; \varphi$ where $\varphi$ has no quantifiers over $\mathbb{R}$.

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

The $\Pi_1^0$ formulas are the formulas $\forall x \in \mathbb{N}\ \varphi$ where $\varphi$ is computable.

The $\Sigma_1^0$ formulas are the formulas $\exists x \in \mathbb{N}\ \varphi$ where $\varphi$ is computable.

The $\Pi_1^1$ formulas are the formulas $\forall x \in \mathbb{R}\ \varphi$ where $\varphi$ has no quantifiers over $\mathbb{R}$.

The $\Sigma_1^1$ formulas are the formulas $\exists x \in \mathbb{R}\ \varphi$ where $\varphi$ has no quantifiers over $\mathbb{R}$.

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

### Definition

A theory $T$ is $\Gamma$-sound if every $\Gamma$ sentence that $T$ proves is true.

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

### Definition

A theory $T$ is $\Gamma$-sound if every $\Gamma$ sentence that $T$ proves is true.

$$\mathsf{RFN}_\Gamma(T) := \forall \varphi \in \Gamma \big( \mathsf{Pr}_T(\varphi) \to \mathsf{True}_\Gamma(\varphi) \big)$$

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

### Definition

A theory $T$ is $\Gamma$-sound if every $\Gamma$ sentence that $T$ proves is true.

$$\mathsf{RFN}_\Gamma(T) := \forall \varphi \in \Gamma \big( \mathsf{Pr}_T(\varphi) \to \mathsf{True}_\Gamma(\varphi) \big)$$

**Fact:** A theory is consistent just in case it is $\Pi_1^0$-sound.

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

**Definition**

A theory $T$ is $\Gamma$-sound if every $\Gamma$ sentence that $T$ proves is true.

$$\mathrm{RFN}_\Gamma(T) := \forall \varphi \in \Gamma \big( \mathrm{Pr}_T(\varphi) \to \mathrm{True}_\Gamma(\varphi) \big)$$

**Fact:** A theory is consistent just in case it is $\Pi_1^0$-sound.

**Definition**

$T \vdash^\Gamma \varphi$ if there is a true $\psi \in \Gamma$ such that $T + \psi \vdash \varphi$.

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

### Definition

A theory $T$ is $\Gamma$-sound if every $\Gamma$ sentence that $T$ proves is true.

$$\mathsf{RFN}_\Gamma(T) := \forall \varphi \in \Gamma \big( \mathsf{Pr}_T(\varphi) \to \mathsf{True}_\Gamma(\varphi) \big)$$

**Fact:** A theory is consistent just in case it is $\Pi^0_1$-sound.

### Definition

$T \vdash^\Gamma \varphi$ if there is a true $\psi \in \Gamma$ such that $T + \psi \vdash \varphi$.

**Fact:** For any $T$ and $\varphi$, $T \vdash \varphi$ if and only if $T \vdash^{\Sigma^0_1} \varphi$.

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

### Definition

$T \leq_{\mathsf{Con}} U := \mathsf{ACA}_0 \vdash \mathsf{Con}(U) \to \mathsf{Con}(T)$.

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

### Definition

$T \leq_{\mathsf{Con}} U := \mathsf{ACA}_0 \vdash^{\Sigma_1^0} \mathsf{Con}(U) \to \mathsf{Con}(T).$

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

### Definition

$T \leq_{\mathsf{Con}} U := \mathsf{ACA}_0 \vdash^{\Sigma^0_1} \mathsf{RFN}_{\Pi^0_1}(U) \to \mathsf{RFN}_{\Pi^0_1}(T).$

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

### Definition

$T \leq_{\mathsf{Con}} U := \mathsf{ACA}_0 \vdash^{\Sigma_1^0} \mathsf{RFN}_{\Pi_1^0}(U) \to \mathsf{RFN}_{\Pi_1^0}(T).$

### Definition

$T \leq_{\Pi_1^1}^{\Sigma_1^1} U := \mathsf{ACA}_0 \vdash^{\Sigma_1^1} \mathsf{RFN}_{\Pi_1^1}(U) \to \mathsf{RFN}_{\Pi_1^1}(T).$

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

### Definition

$T \leq_{\mathsf{Con}} U := \mathsf{ACA}_0 \vdash^{\Sigma^0_1} \mathsf{RFN}_{\Pi^0_1}(U) \to \mathsf{RFN}_{\Pi^0_1}(T).$

### Definition

$T \leq^{\Sigma^1_1}_{\Pi^1_1} U := \mathsf{ACA}_0 \vdash^{\Sigma^1_1} \mathsf{RFN}_{\Pi^1_1}(U) \to \mathsf{RFN}_{\Pi^1_1}(T).$

### Theorem (W.)

*The relation $\leq^{\Sigma^1_1}_{\Pi^1_1}$ pre-well-orders the $\Pi^1_1$-sound extensions of $\mathsf{ACA}_0$.*

Introduction
Set theory as a case study
The consistency operator
**Second-order arithmetic**

## Thanks!

📄 J. Walsh (2020)

A note on the consistency operator.

*Proceedings of the American Mathematical Society.* 148(6):2645–2654

📄 J. Walsh (2022)

On the hierarchy of natural theories.

*arXiv.*

📄 J. Walsh (2022)

A robust proof-theoretic well-ordering.

*arXiv.*