# Division by zero

Emil Jeřábek

Institute of Mathematics of the Czech Academy of Sciences

Žitná 25, 115 67 Praha 1, Czech Republic, email: `jerabek@math.cas.cz`

September 24, 2016

*To Albert Visser*

## Abstract

For any sufficiently strong theory of arithmetic, the set of Diophantine equations provably unsolvable in the theory is algorithmically undecidable, as a consequence of the MRDP theorem. In contrast, we show decidability of Diophantine equations provably unsolvable in Robinson's arithmetic $Q$. The argument hinges on an analysis of a particular class of equations, hitherto unexplored in Diophantine literature. We also axiomatize the universal fragment of $Q$ in the process.

**Keywords:** Robinson arithmetic, Diophantine equation, decidability, universal theory

**MSC:** 03F30

## 1 Introduction

The standard Gödel–Church–Turing–Rosser undecidability theorem tells us that if $T$ is any consistent theory extending Robinson's arithmetic $Q$, the set of $\Pi_1$ consequences of $T$ is undecidable. Furthermore, the Matiyasevich–Robinson–Davis–Putnam theorem shows that every $\Pi_1$ formula is equivalent to unsolvability of a certain Diophantine equation. Since the MRDP theorem can be formalized in $I\Delta_0 + EXP$ due to Gaifman and Dimitracopoulos [5], we see that if $T$ extends $I\Delta_0 + EXP$, it is undecidable whether a given Diophantine equation is provably unsolvable in $T$, or dually, whether it has a solution in a model of $T$.

Surprisingly, Kaye [6, 7] proved that the same holds already for extensions of the weak theory $IU_1^-$ (induction for parameter-free bounded universal formulas), despite that it likely does *not* formalize the MRDP theorem as such. One can check that Kaye's methods also apply to extensions of Cook's theory $PV$ of polynomial-time functions (see e.g. Krajíček [8] for a definition).

Going further down, decidability of solvability of Diophantine equations in models of the theory $IOpen$ of quantifier-free induction has remained an intriguing open problem ever since it was posed by Shepherdson [12], see e.g. [13, 3, 10] for partial results.

The purpose of this note is to show that solvability of Diophantine equations in models of $Q$ is decidable, specifically NP-complete. Since $Q$ does not include ring identities that allow the usual manipulations of polynomials, it may be ambiguous what exactly is meant by Diophantine equations, so let us first state the problem precisely.

**Definition 1.1** A *Diophantine equation* is a formula of the form

$$t(\vec{x}) = u(\vec{x}),$$

where $t$ and $u$ are terms in the basic language of arithmetic $L_Q = \langle 0, S, +, \cdot \rangle$. If $T$ is a theory whose language contains $L_Q$, the *Diophantine satisfiability problem for $T$*, denoted $\mathrm{D}_T$, consists of all Diophantine equations $t = u$ satisfiable in a model of $T$ (shortly: $T$-satisfiable). That is,

$$\mathrm{D}_T = \{\langle t, u \rangle : T + \exists \vec{x}\, t(\vec{x}) = u(\vec{x}) \text{ is consistent}\}.$$

The decidability of $\mathrm{D}_Q$ is on the whole not so surprising, as $Q$ has models with "black holes" (to use Albert Visser's term) that can serve to equate nearly any pair of terms. It turns out however that while this argument yields a simple proof of decidability of Diophantine satisfiability for certain mild extensions of $Q$, it does not suffice for $Q$ itself: it only provides a reduction to systems of equations of a special form (see Eq. (1) below) that we have to investigate in detail.

We need to embed certain models in models of $Q$ as a part of our main construction, and to facilitate this goal, we will explicitly axiomatize the universal consequences of $Q$, which could be of independent interest.

See [1, 2] for related work.

## Acknowledgements

## 2 Robinson defeats Diophantus

We need a convenient way to refer to the individual axioms of $Q$, thus we can as well start by properly defining the theory, even though we trust it is familiar to the reader.

**Definition 2.1** $Q$ is the theory in language $L_Q$ with axioms

| | |
|---|---|
| (Q1) | $Sx \neq 0,$ |
| (Q2) | $Sx = Sy \rightarrow x = y,$ |
| (Q3) | $x = 0 \vee \exists y\, Sy = x,$ |
| (Q4) | $x + 0 = x,$ |
| (Q5) | $x + Sy = S(x + y),$ |
| (Q6) | $x \cdot 0 = 0,$ |
| (Q7) | $x \cdot Sy = x \cdot y + x.$ |

Let $t \simeq u$ denote that the terms $t$ and $u$ are syntactically identical. We define unary numerals $\overline{n} \simeq S^n 0$, and binary numerals

$$\underline{0} \simeq 0,$$
$$\underline{2n} \simeq \overline{2} \cdot \underline{n}, \qquad n > 0,$$
$$\underline{2n+1} \simeq S(\underline{2n})$$

for all natural numbers $n \in \mathbb{N}$. While unary numerals are easier to manipulate using axioms of the theory, we will need the much shorter binary numerals when discussing algorithmic complexity.

Of course, $Q$ proves $\underline{n} = \overline{n}$, and we will use both interchangeably in contexts where the distinction does not matter.

## 2.1 Black-hole models

As our starting point (already alluded to in the introduction), we can drastically reduce the complexity of the Diophantine satisfiability problem for $Q$ using black-hole models:

**Lemma 2.2** $\mathrm{D}_Q$ *is polynomial-time reducible to $Q$-solvability of equations of the form $t = \underline{n}$ with $n \in \mathbb{N}$.*

*Proof:* Consider the model $\mathbb{N}^\infty = \mathbb{N} \cup \{\infty\}$, where $S\infty = \infty + x = x + \infty = x \cdot \infty = \infty$ for all $x \in \mathbb{N}^\infty$, $\infty \cdot 0 = 0$, and $\infty \cdot x = \infty$ for $x \neq 0$. It is readily seen that $\mathbb{N}^\infty \vDash Q$.

When written in binary, the lengths of $n + m$ and $n \cdot m$ are bounded by the sum of lengths of $n$ and $m$. It follows by induction on the complexity of $t$ that given a term $t$ and $\vec{a} \in \mathbb{N}^\infty$, the length of the value of $t(\vec{a})$ in $\mathbb{N}^\infty$ is polynomial in the lengths of $t$ and $\vec{a}$, and we can compute $t(\vec{a})$ in polynomial time.

Crucially, the operations in $\mathbb{N}^\infty$ are defined so that they give a finite value only when forced so by the axioms of $Q$, hence we can show by induction on the complexity of $t$ that

$$t(\vec{\infty}) = n \in \mathbb{N} \implies Q \vdash t(\vec{x}) = \underline{n}.$$

For example, let $t \simeq u \cdot v$. Then $t(\vec{\infty}) \in \mathbb{N}$ only if both $u(\vec{\infty}), v(\vec{\infty}) \in \mathbb{N}$, or if $v(\vec{\infty}) = 0$. In the former case, the induction hypothesis gives

$$Q \vdash u(\vec{x}) = \underline{k}, \qquad Q \vdash v(\vec{x}) = \underline{l}$$

for some $k, l \in \mathbb{N}$, thus $Q \vdash t(\vec{x}) = \underline{n}$ with $n = kl$. In the latter case, $Q \vdash v(\vec{x}) = 0$ by the induction hypothesis, hence $Q \vdash t(\vec{x}) = 0$.

Thus, here is the promised reduction: given an equation $t_0 = t_1$, if $t_0(\vec{\infty}) = \infty = t_1(\vec{\infty})$, we have a witness that $t_0 = t_1$ is satisfiable, solving the problem outright; otherwise, at least one of the terms $t_i$ is provably equal to a numeral $\underline{n}$, which we can compute in polynomial time. The output of the reduction is (say) "$0 = 0$" in the former case, and "$t_{1-i} = \underline{n}$" in the latter case. $\square$

This is not yet the end of the story; we can further reduce the problem by unwinding the terms from top. For example, axioms Q1 and Q2 imply that an equation $St = \overline{n}$ is $Q$-satisfiable if and only if $n$ is nonzero, and $t = \overline{n-1}$ is satisfiable. Something to a similar effect also holds for the other function symbols, so let us see where it gets us.

**Definition 2.3** Let $Q_\forall$ denote the theory axiomatized by Q1, Q2, Q4–Q7, and

$$(Q8_n) \qquad\qquad x + y = \overline{n} \rightarrow \bigvee_{m \leq n} (y = \overline{m}),$$

$$(Q9_n) \qquad\qquad x \cdot y = \overline{n} \rightarrow x = 0 \vee \bigvee_{m \leq n} (y = \overline{m})$$

for $n \in \mathbb{N}$.

**Lemma 2.4**

   (i) $Q \vdash Q_\forall$.

(ii) *Let $n \in \mathbb{N}$. Then $Q_\forall$ proves*

$$x + y = \overline{n} \rightarrow \bigvee_{k+m=n} (x = \overline{k} \wedge y = \overline{m}),$$

$$x \cdot y = \overline{n} \rightarrow x = 0 \vee \bigvee_{km=n} (x = \overline{k} \wedge y = \overline{m}), \qquad n > 0.$$

*Proof:* (i): We prove $Q8_n$ by induction on $n$. Reason in $Q$, and assume $x + y = \overline{n}$. If $y = 0$, we are done. Otherwise $y = Sz$ for some $z$ by Q3, hence $S(x + z) = \overline{n}$ by Q5. This is only possible if $n > 0$ due to Q1, and we have $x + z = \overline{n-1}$ by Q2, thus $z = 0 \vee \cdots \vee z = \overline{n-1}$ by the induction hypothesis, and consequently $y = \overline{1} \vee \cdots \vee y = \overline{n}$. (Alternatively, notice that under the traditional definition of $u \leq v$ as $\exists w\, (v = w + u)$, $Q8_n$ may be read as the bounded sentence $\forall y \leq \overline{n} \bigvee_{m \leq n} y = \overline{m}$, hence its provability follows from the $\Sigma_1$-completeness of $Q$.)

The proof of $Q9_n$ is similar. Assuming $xy = \overline{n}$, we are done if $y = 0$, hence we can assume $y = Sz$. Then $xz + x = \overline{n}$ (Q7), thus $x = \overline{k}$ for some $k = 0, \ldots, n$ by $Q8_n$. If $k = 0$, we are done. Otherwise $xz = \overline{n-k}$ (Q4, Q5, Q2), where $n - k < n$, hence we can use the induction hypothesis to conclude $z = 0 \vee \cdots \vee z = \overline{n-k}$. This implies $y = \overline{1} \vee \cdots \vee y = \overline{n-k+1}$, where $n - k + 1 \leq n$.

(ii): If $x + y = \overline{n}$, we have $y = \overline{m}$ for some $m \leq n$ by $Q8_n$. Then $x + y = S^m x$ by Q4 and Q5, hence $x = \overline{n-m}$ by Q2.

Let $n \neq 0$, and reason in $Q_\forall$ again. Assume $xy = \overline{n}$. We have $\overline{n} \neq 0$ by Q1, hence $y \neq 0$ by Q6. Using $Q9_n$, either $x = 0$, or $y = \overline{m}$ for some $m = 1, \ldots, n$. In the latter case, $\overline{n} = x \cdot \overline{m-1} + x$ by Q7, hence $x = \overline{k}$ for some $k = 0, \ldots, n$ by $Q8_n$. Then $\overline{n} = xy = \overline{km}$ using Q4–7, hence $km = n$ by Q1, Q2. $\qquad \square$

**Proposition 2.5** *Let $Q^+$ denote $Q$ extended by the axiom $0 \cdot x = 0$. Then $D_{Q^+}$ is decidable.*

*Proof:* The proof of Lemma 2.2 works for $Q^+$, too, with $\mathbb{N}^\infty$ modified so that $0 \cdot \infty = 0$. We describe below a recursive procedure $\text{Sol}(E)$ that checks whether a finite set $E$ of equations of the form $t = \underline{n}$ is $Q^+$-satisfiable.

Let $t = \underline{n}$ be the first equation in $E$ such that $t$ is not a variable, and $E' = E \smallsetminus \{t = \underline{n}\}$:

(i) If $t \simeq t_0 \cdot t_1$ and $n \neq 0$, call $\text{Sol}(E' \cup \{t_0 = \underline{n_0}, t_1 = \underline{n_1}\})$ for every $n_0, n_1$ such that $n_0 n_1 = n$. Accept if any of the recursive calls accepted, otherwise reject.

(ii) If $t \simeq t_0 \cdot t_1$ and $n = 0$, call $\text{Sol}(E' \cup \{t_0 = 0\})$ and $\text{Sol}(E' \cup \{t_1 = 0\})$. Accept if any of the recursive calls accepted, otherwise reject.

(iii) If $t$ is of the form $t_0 + t_1$, $St_0$, or $0$, proceed similarly.

(iv) If the left-hand sides of all equations in $E$ are variables, reject if $E$ contains a pair of equations with the same left-hand sides and different right-hand sides, otherwise accept.

Each recursive call strictly decreases the total number of symbols on the left-hand sides, hence the algorithm terminates, and Lemma 2.4 and the extra axiom guarantee its correctness.

We note that $\text{Sol}(E)$ as presented is an exponential-time algorithm, but we can transform it into a nondeterministic polynomial-time algorithm by making only one, nondeterministically chosen, recursive call at each step. Thus, $D_{Q^+} \in \text{NP}$. $\qquad \square$

If we try to use $\text{Sol}(E)$ for $Q$, we run into trouble: while $Q$ proves $xy = 0 \rightarrow x = 0 \vee y = 0$ by Lemma 2.4, it does not prove the converse implication, hence the solvability of $E' \cup \{t_0 = 0\}$ does not imply the solvability of $E' \cup \{t_0 \cdot t_1 = 0\}$ in step (ii). Likewise, step (i) is incorrect, because $E' \cup \{t_0 \cdot t_1 = \overline{n}\}$ with $n \neq 0$ may be satisfied in such a way that $t_0 = 0$.

However, the other reductions remain valid, and this still proves useful: a variant of $\mathrm{Sol}(E)$ shows that $\mathrm{D}_Q$ reduces to $Q$-solvability of systems of equations of the form

$$
(1) \qquad
\begin{cases}
0 \cdot t_1(\vec{x}) = \overline{n_1}, \\
\quad \cdots \\
0 \cdot t_k(\vec{x}) = \overline{n_k}.
\end{cases}
$$

Diophantine systems of this type have not yet received the attention they deserve, so we are on our own. Their $Q$-satisfiability turns out to be an unexpectedly circuitous problem: on the one hand, we will see that nearly every such equation is satisfiable by itself in a suitable model, on the other hand there are subtle dependencies that make systems such as

$$
\begin{aligned}
0 \cdot (x + \overline{2}) &= \overline{5} \\
0 \cdot (y + 0 \cdot x) &= \overline{7} \\
0 \cdot Sy &= \overline{4}
\end{aligned}
$$

unsatisfiable[1]. One consequence is that we cannot make do with a one-size-fits-all model of $Q$ like in Lemma 2.2; we will need a variety of countermodels for different systems. This will be our task in the next two subsections.

## 2.2 Universal fragment of $Q$

We intend to use term models of a kind as our supply of models to satisfy various equations, but this approach is not very friendly to the predecessor axiom Q3, so to make our lives easier, we first determine what structures can be *extended* to models of $Q$ by adding predecessors (and other elements that are forced upon us). By general model theoretic considerations, these are exactly the models of the *universal fragment* of $Q$, hence we can reformulate the problem as a description of this universal fragment. Since we used very suggestive notation, the answer should come as no surprise:

**Proposition 2.6** $Q_\forall$ *is the universal fragment of $Q$. That is, every model of $Q_\forall$ embeds in a model of $Q$.*

*Proof:* Fix $M \vDash Q_\forall$. Identifying each $n \in \mathbb{N}$ with the corresponding numeral $\overline{n}^M \in M$, we may assume that $M$ includes the standard model $\mathbb{N}$; in particular, $M \smallsetminus \mathbb{N}$ is the set of nonstandard elements of $M$.

Let $A$ denote the set of nonzero elements of $M$ without a predecessor. We will embed $M$ in a structure with domain

$$
N = M \cup \{\infty\} \cup \{\langle a, k \rangle, \langle a, k, x \rangle : a \in A, k \in \mathbb{N}^{>0}, x \in M \smallsetminus \mathbb{N}\},
$$

where $\langle a, k \rangle$ should be thought of as $a - k$, and $\langle a, k, x \rangle$ as $x \cdot (a - k)$. We will define the interpretations of the $L_Q$-function symbols in $N$, and verify that $N \vDash Q$ along the way. All the function symbols are understood to retain their interpretations from $M$ on elements of $M$, thus we will not indicate such cases explicitly.

---

[1] By Q4–7, $0 \cdot (x + \overline{2}) = 0 \cdot SSx = (0 \cdot x + 0) + 0 = 0 \cdot x$, thus the first equation implies $0 \cdot x = \overline{5}$. Likewise, the third equation gives $0 \cdot y = \overline{4}$, while the second equation gives $\overline{7} = 0 \cdot (y + \overline{5}) = 0 \cdot y$. It is worth noting that the second equation does not imply anything about $0 \cdot y$ on its own—we need to know that $0 \cdot x$ is standard first.

*Successor:* we put

$$S^N \infty = \infty,$$
$$S^N \langle a, k, x \rangle = \langle a, k, x \rangle,$$
$$S^N \langle a, k \rangle = \begin{cases} \langle a, k-1 \rangle & k > 1, \\ a & k = 1. \end{cases}$$

We can see immediately that this makes $N$ a model of Q1–Q3. This also means we can unambiguously refer to $S^n x$ for $n \in \mathbb{Z}$ and $x \in N \smallsetminus \mathbb{N}$.

*Addition:* we put

$$\infty +^N y = \infty,$$
$$\langle a, k \rangle +^N y = \begin{cases} S^{n-k} a & y = n \in \mathbb{N}, \\ \infty & \text{otherwise.} \end{cases}$$

For $x \in M$, we define

$$x +^N \infty = \infty,$$
$$x +^N \langle a, k, y \rangle = \infty,$$
$$x +^N \langle a, k \rangle = S^{-k}(x +^M a).$$

Note that the last item makes sense: since $a \notin \mathbb{N}$, also $x +^M a \notin \mathbb{N}$ by Q8. Finally, we put

$$\langle a, k, x \rangle +^N n = \langle a, k, x \rangle, \qquad n \in \mathbb{N},$$
$$\langle a, k, x \rangle +^N S^n x = \begin{cases} \langle a, k-1, x \rangle & k > 1, \\ S^n(x \cdot^M a) & k = 1, \end{cases} \qquad n \in \mathbb{Z},$$
$$\langle a, k, x \rangle +^N y = \infty \qquad \text{for other } y.$$

Again, $x \cdot^M a \notin \mathbb{N}$ by Q9 as $x, a \notin \mathbb{N}$.

It is straightforward to check that $N$ validates Q4 and Q5.

*Multiplication:* for $x \in N \smallsetminus M$, we put

$$x \cdot^N n = (\cdots (0 +^N \underbrace{x) +^N \cdots +^N x) +^N x}_{n}, \qquad n \in \mathbb{N},$$
$$x \cdot^N y = \infty, \qquad y \notin \mathbb{N}.$$

For $x \in M$, we define

$$x \cdot^N \infty = \infty,$$
$$x \cdot^N \langle a, k, y \rangle = \infty,$$
$$x \cdot^N \langle a, k \rangle = \begin{cases} S^{-kn}(n \cdot^M a) & x = n \in \mathbb{N}, \\ \langle a, k, x \rangle & x \notin \mathbb{N}. \end{cases}$$

As above, $S^{-kn}(n \cdot^M a)$ exists: either $n = 0$ and the $S^{-kn}$ does nothing, or $n > 0$, in which case $n \cdot^M a \notin \mathbb{N}$ by Q9.

Again, it is straightforward to check Q6 and Q7. $\qquad \square$

## 2.3 Reduced terms

We now come to the crucial part of our construction: we establish that a system (1) is $Q$-satisfiable if the terms $t_i$ obey certain conditions that guarantee they do not interact with each other.

**Definition 2.7** An $L_Q$-term is *normal* if it contains no subterm of the form $t + 0$, $t + Su$, $t \cdot 0$, or $t \cdot Su$. A normal term is *irreducible* if it does not have the form $0$ or $St$. In other words, normal and irreducible terms are generated by the following grammar:

$$I ::= x_i \mid (N + I) \mid (N \cdot I)$$
$$N ::= I \mid 0 \mid SN$$

If $T$ is a set of terms, a normal term is $T$-*reduced* if it contains no subterm of the form $0 \cdot t$ for $t \in T$.

Each normal term can be uniquely written in the form $S^n 0$ or $S^n t$, where $n \in \mathbb{N}$, and $t$ is irreducible. A subterm of a normal ($T$-reduced) term is again normal ($T$-reduced, resp.).

**Lemma 2.8** *Let $T = \{t_i : i < k\}$ be a finite sequence of distinct irreducible terms such that $0 \cdot t_i$ is not a subterm of $t_j$ for any $i, j < k$ (i.e., the terms $t_i$ are $T$-reduced), and $\{n_i : i < k\} \subseteq \mathbb{N}$.*

(i) *The set of equations $\{0 \cdot t_i = \overline{n_i} : i < k\}$ is $Q$-satisfiable.*

(ii) *Given $T, \vec{n}$, and a term $t$, we can compute a $T$-reduced term $\widetilde{t}$ such that*

$$Q \vdash \bigwedge_{i<k} 0 \cdot t_i(\vec{x}) = \overline{n_i} \to t(\vec{x}) = \widetilde{t}(\vec{x}).$$

*Proof:* (i): We define a model $M$ whose domain consists of all $T$-reduced terms, and operations as follows. We put $0^M = 0$, and $S^M t = St$. If $t \in M$, $n \in \mathbb{N}$, and $u \in M$ is irreducible,

$$t +^M S^n 0 = S^n t,$$
$$t +^M S^n u = S^n (t + u).$$

If $t, u \in M$ are irreducible, and $n, m \in \mathbb{N}$, we put

$$S^n t \cdot^M S^m 0 = \underbrace{S^n(S^n(\ldots(S^n}_{m}(0 + \underbrace{t)\ldots) + t) + t}_{m}),$$

$$S^n t \cdot^M S^m u = \underbrace{S^n(S^n(\ldots(S^n}_{m}(S^n t \cdot u + \underbrace{t)\ldots) + t) + t}_{m}),$$

$$S^n 0 \cdot^M S^m 0 = S^{nm} 0,$$

$$S^n 0 \cdot^M S^m u = S^{nm}(S^n 0 \cdot u), \qquad n > 0,$$

$$0 \cdot^M S^m u = \begin{cases} S^{n_i} 0 & u = t_i, \\ 0 \cdot u & u \notin T. \end{cases}$$

It is readily checked that the operations are well-defined (i.e., the terms given above as their values are $T$-reduced), and that $M \vDash Q_\forall$. Let $v$ be the valuation in $M$ which assigns each variable $x_i$ to the corresponding element $x_i \in M$. Then $v(t) = t$ for every $T$-reduced term $t$, hence $v$ satisfies in $M$ the equations $0 \cdot t_i = \overline{n_i}$. By Proposition 2.6, we can embed $M$ into a model of $Q$.

(ii): Since the operations in $M$ are computable, we can compute the value $v(t) \in M$ by induction on the complexity of $t$. This value is a $T$-reduced term, so we can define $\widetilde{t} = v(t)$. Then we show that $Q$ proves the required implication

$$(2) \qquad \bigwedge_{i<k} 0 \cdot t_i(\vec{x}) = \overline{n_i} \to t(\vec{x}) = \widetilde{t}(\vec{x})$$

by induction on the complexity of $t$. For the induction steps, we observe the operations in $M$ are defined so that if $t +^M s = u$, then $Q \vdash t + s = u$, and likewise for $\cdot^M$ with the exception of the clause $0 \cdot^M S^m t_i = S^{n_i} 0$, which is handled by the premise of (2). $\square$

The reader might have realized that what just happened was term rewriting in thinly veiled disguise. Even though we will not need this point of view for our application, we make the digression to spell this connection out because of sheer curiosity.

**Definition 2.9** Let $R_Q$ denote the rewriting system for $L_Q$-terms generated by the rules

$$
(3) \qquad
\begin{cases}
\quad t + 0 \longrightarrow t, \\
\quad t + Su \longrightarrow S(t + u), \\
\quad t \cdot 0 \longrightarrow 0, \\
\quad t \cdot Su \longrightarrow t \cdot u + t.
\end{cases}
$$

More generally, if $\{t_i : i < k\}$ is a sequence of terms satisfying the conditions of Lemma 2.8, and $\{n_i : i < k\} \subseteq \mathbb{N}$, let $R_{\vec{t},\vec{n}}$ denote the rewriting system extending $R_Q$ with the rules

$$
(4) \qquad\qquad\qquad 0 \cdot t_i \longrightarrow \overline{n_i}, \qquad i < k
$$

(these rules are not supposed to allow substitution for variables inside $t_i$).

Notice that a term is normal in the sense of Definition 2.7 iff it is a normal form with respect to $R_Q$, and it is $T$-reduced (with $T = \{t_i : i < k\}$) iff it is a normal form with respect to $R_{\vec{t},\vec{n}}$ for an arbitrary choice of $\vec{n}$.

**Proposition 2.10** *For any $\vec{t}$ and $\vec{n}$ as in the definition, the rewriting system $T_{\vec{t},\vec{n}}$ is strongly normalizing and confluent. (That is, every term has a unique normal form, and every sequence of reductions will eventually reach it.)*

*Proof:* Put $c = 2 + \max_{i<k} n_i$, and define a "norm" function on terms by

$$
\begin{aligned}
\|x_i\| = \|0\| &= c, \\
\|St\| &= \|t\| + 3, \\
\|t + u\| &= \|t\| + 2\|u\|, \\
\|t \cdot u\| &= \|t\| \cdot \|u\|.
\end{aligned}
$$

Notice that $\|t\| \geq c$ for any term $t$, and the norm is strictly monotone in the sense that $\|u\| < \|v\|$ implies $\|t(u)\| < \|t(v)\|$. Using this, we can check easily that all $R_{\vec{t},\vec{n}}$-reduction steps strictly decrease the norm, thus there is no infinite sequence of reductions: in particular, we have

$$
\|\overline{n_i}\| = 3n_i + c \leq 4c - 6 < c^2 \leq \|0 \cdot t_i\|.
$$

This shows strong normalization of $R_{\vec{t},\vec{n}}$.

By Newman's lemma, confluence is implied by local confluence: that is, it suffices to show that if $s \longrightarrow v_0$ and $s \longrightarrow v_1$, then $v_0 \overset{*}{\longrightarrow} w$ and $v_1 \overset{*}{\longrightarrow} w$ for some term $w$, where $\longrightarrow$ denotes one-step reduction, and $\overset{*}{\longrightarrow}$ its reflexive transitive closure.

The local confluence property obviously holds if the two reductions $s \longrightarrow v_i$ are identical, or if they operate on disjoint terms. It also holds if $s \longrightarrow v_i$ is one of the $R_Q$-reductions as given in (3), and $s \longrightarrow v_{1-i}$ operates inside one of the terms $t, u$ on the left-hand side of (3): we can instead perform the reduction on their copies on the right-hand side.

This in fact covers all possibilities: the only redexes properly included inside the left-hand side of any $R_Q$-rule in (3) are inside $t$ or $u$, as there are no rules reducing $0$ or $Su$; there are no redexes properly included inside $0 \cdot t_i$ by the assumption that $t_i$ is $T$-reduced; and each redex can be reduced only in one way—the only possible clashes could be between (4) and the $R_Q$-rules for multiplication, but these are prevented as $t_i$ is assumed not to be of the form $0$ or $Su$. $\qquad \square$

Proposition 2.10 provides an alternative proof for most of Lemma 2.8: first, the $T$-reduced term $\widetilde{t}$ in 2.8 (ii) is just the $R_{\vec{t},\vec{n}}$-normal form for $t$. Second, we can use confluence (Church–Rosser property) to construct the model $M$ for (i) as the model of $R_{\vec{t},\vec{n}}$-normal terms, or equivalently, as the quotient of the free term model by the equivalence relation induced by reduction. It is automatically a model of axioms Q4–7 embodied in the reduction rules, and it is easily seen to satisfy Q1 and Q2 because there are no rules with redex $Su$. It would still take a little work to establish the validity of Q8 and Q9.

## 2.4 Witnessing satisfiability

To complete our analysis of $\mathrm{D}_Q$, we will now show that a general equation $t = \overline{n}$ can only be $Q$-satisfied if it is implied by a (suitably bounded) system of the form (1) that respects the assumptions of Lemma 2.8.

**Definition 2.11** For any term $u$, let $\widetilde{u}$ denote its $\varnothing$-reduced form as given by Lemma 2.8.

A *labelling* of a term $t$ is a partial map $\ell$ from subterms of $t$ to $\mathbb{N}$. If $\ell$ is a labelling of $t$, and $u$ a subterm of $t$ (written henceforth as $u \subseteq t$), let $u_\ell$ be the term obtained from $u$ by replacing all maximal proper labelled subterms of $u$ by numerals for their labels.

A *witness* for $t = \overline{n}$ is a labelling $\ell$ of $t$ by numbers $k \le n$ such that:

(i) $\ell(t) = n$.

(ii) If $u, v \subseteq t$ are such that $\widetilde{u_\ell} \simeq \widetilde{v_\ell}$, then $\ell(u) = \ell(v)$ (meaning both are undefined, or both are defined and equal).

(iii) If $u \in \mathrm{dom}(\ell)$, and $\widetilde{u_\ell} \simeq \overline{k}$ for some $k \in \mathbb{N}$, then $\ell(u) = k$.

(iv) If $u \in \mathrm{dom}(\ell)$, then all immediate subterms of $u$ are labelled, unless $u \simeq v \cdot w$, and $v$ or $w$ is labelled 0.

Note that (ii) implies that occurrences of the same subterm either all have the same label, or are all unlabelled. We also remark that in (iii), $k \le n$ is not a premise, but part of the conclusion.

**Example 2.12** Table 1 shows a labelling $\ell$ of the term $t \simeq x \cdot y + x \cdot SSSy$ that is a witness for satisfiability of the equation $t = \overline{8}$. For convenience, the table also lists for each term $u \subseteq t$ its set of maximal proper labelled subterms, as well as $u_\ell$ and $\widetilde{u_\ell}$, which makes it easy to check that conditions (i)–(iv) hold. In particular, for (ii), the two terms $u$ with $\widetilde{u_\ell} \simeq 0 \cdot y$ have the same label $\ell(u) = 4$; for (iii), the only applicable case is $\widetilde{t_\ell} = \overline{8}$, which agrees with $\ell(t) = 8$. We invite the reader to verify that $\ell$ is in fact the *only* possible witness for $t = \overline{8}$.

For this example, the set $E$ considered below in the proof of Lemma 2.13 consists of the single equation $0 \cdot y = \overline{4}$.

**Lemma 2.13** *An equation $t = \overline{n}$ is $Q$-satisfiable if and only if it has a witness.*

| $u$ | $\ell(u)$ | m.p.l.s. | $u_\ell$ | $\widetilde{u_\ell}$ |
|---|---|---|---|---|
| $x \cdot y + x \cdot SSSy$ | 8 | $x \cdot y, x \cdot SSSy$ | $\overline{4} + \overline{4}$ | $\overline{8}$ |
| $x \cdot y$ | 4 | $x$ | $0 \cdot y$ | $0 \cdot y$ |
| $x \cdot SSSy$ | 4 | $x$ | $0 \cdot SSSy$ | $0 \cdot y$ |
| $x$ | 0 | $-$ | $x$ | $x$ |
| $S^i y\ (i = 0, \ldots, 3)$ | $-$ | $-$ | $S^i y$ | $S^i y$ |

Table 1: Witness for $x \cdot y + x \cdot SSSy = \overline{8}$

*Proof:* Left-to-right: let $M \vDash Q$ and $\vec{a} \in M$ be such that $t^M(\vec{a}) = n$. Define a labelling of $t$ by putting $\ell(u) = u^M(\vec{a})$ if $u^M(\vec{a}) \in \{0, \ldots, n\}$, and $\ell(u)$ is undefined otherwise. Since a term equals its $\varnothing$-reduction provably in $Q$, we have $\widetilde{u_\ell}^M(\vec{a}) = u^M(\vec{a})$ for any $u \subseteq t$. It follows easily that $\ell$ is a witness for $t = \overline{n}$, using Lemma 2.4 for condition (iv).

Right-to-left: let $E$ denote the set of equations

$$\widetilde{u_\ell} = \overline{k}$$

where $u \simeq v \cdot w \subseteq t$, $\ell(u) = k$, $\ell(v) = 0$, and $\widetilde{u_\ell} \neq 0$ (which implies $w \notin \mathrm{dom}(\ell)$). Note that $\widetilde{u_\ell}$ then must be of the form $0 \cdot u^-$, where $u^-$ is an irreducible term, and $\widetilde{w_\ell} \simeq S^m u^-$ for some $m$.

**Claim 1** *If $E$ is satisfiable, then $t = \overline{n}$ is satisfiable.*

*Proof:* Fix a model $M \vDash Q$ and $\vec{a} \in M$ that satisfies $E$. Note that $E$ only contains unlabelled variables[2]; if $x_i$ is labelled, we make sure that $a_i = \ell(x_i)$ (this is independent of the choice of an occurrence of $x_i$ in $t$ by condition (ii)). We claim that

$$u \in \mathrm{dom}(\ell) \implies u^M(\vec{a}) = \ell(u),$$

which gives $t^M(\vec{a}) = n$ by condition (i). We prove this by induction on the complexity of $u$.

The statement holds for variables, and condition (iii) implies it holds for $u \simeq 0$.

If $u \in \mathrm{dom}(\ell)$ is of the form $Sv$ or $v + w$, then $v, w \in \mathrm{dom}(\ell)$ by (iv), and $\ell(u)$ equals $\ell(v) + 1$ or $\ell(v) + \ell(w)$ (resp.) by (iii), thus $u^M(\vec{a}) = \ell(u)$ by the induction hypothesis for $v$ and $w$.

The same argument applies if $u \simeq v \cdot w$, and both $v, w \in \mathrm{dom}(\ell)$, or $\ell(w) = 0$. Assume $\ell(v) = 0$ and $w \notin \mathrm{dom}(\ell)$. Using the induction hypothesis for subterms of $u$, and the soundness of reduction, we have $u^M(\vec{a}) = u_\ell^M(\vec{a}) = \widetilde{u_\ell}^M(\vec{a})$. If $\widetilde{u_\ell} = 0$, this means $u^M(\vec{a}) = 0 = \ell(u)$ by (iii); otherwise the equation $\widetilde{u_\ell} = \ell(u)$ is in $E$, hence it is satisfied by $\vec{a}$. $\qquad\square$ (Claim 1)

Condition (ii) ensures that $E$ does not contain two equations with the same left-hand side. Moreover, for any

$$0 \cdot u_0^- = \overline{k_0}$$
$$0 \cdot u_1^- = \overline{k_1}$$

in $E$, $0 \cdot u_0^-$ is not a subterm of $u_1^-$: writing $u_1 \simeq v_1 \cdot w_1$, inspection of the definition of reduction shows that this could only happen if $(w_1)_\ell$ contained a nonconstant subterm $s$ such that $\widetilde{s} \simeq 0 \cdot u_0^-$. But then $s \simeq r_\ell$ for some $r \subseteq w_1$ such that $r \notin \mathrm{dom}(\ell)$, whereas we should have $\ell(r) = k_0$ by (ii), a contradiction. Thus, $E$ is satisfiable by Lemma 2.8. $\qquad\square$

---

[2] Variables in $\widetilde{u_\ell}$ come from variables in $u$. However, a labelled variable in $u$ is a maximal proper labelled subterm of $u$, or is included in such a maximal subterm; consequently, it disappears in $u_\ell$ (and $\widetilde{u_\ell}$) by virtue of being replaced with a constant term (a numeral).

**Theorem 2.14** $\mathrm{D}_Q$ *is decidable.*

*Proof:* By Lemma 2.2, $\mathrm{D}_Q$ reduces to $Q$-satisfiability of equations of the form $t = \overline{n}$. These can be checked by the criterion from Lemma 2.13: a witness for $t = \overline{n}$ has size bounded by a computable function of $t$ and $n$, and using the computability of $\widetilde{u}$, we can algorithmically recognize a witness when we see it. $\hfill\square$

# 3 Computational complexity

Our arguments thus far give an exponential-time algorithm for checking if a given Diophantine equation is $Q$-satisfiable. We can in fact determine the complexity of $\mathrm{D}_Q$ precisely. First, a general lower bound follows from a beautiful result of Manders and Adleman [9] that there are very simple NP-complete Diophantine problems.

**Theorem 3.1 (Manders and Adleman)** *The following problem is* NP*-complete: given $a, b \in \mathbb{N}$ in binary, determine whether $x^2 + ay - b = 0$ has a solution in $\mathbb{N}$.* $\hfill\square$

(They state it with $ax^2 + by - c$, but it is easy to show that the version here is equivalent.)

**Corollary 3.2** *If $T$ is a consistent extension of $Q_\forall$, then $\mathrm{D}_T$ is* NP*-hard.*

*Proof:* If $a > 0$ (which we can assume without loss of generality), $x^2 + ay - b = 0$ is solvable iff the equation

$$(5) \qquad\qquad x \cdot x + \underline{a} \cdot y = \underline{b}$$

is in $\mathrm{D}_T$: on the one hand, a solution in $\mathbb{N}$ yields a solution in any model of $T$. On the other hand, (5) implies in $Q_\forall$ that $x \cdot x$ and $\underline{a} \cdot y$ are standard and bounded by $b$ using Q8, hence $y$ is standard by Q9. Also by Q9, $x = 0$, or $x = 0, \ldots, \underline{b}$; either way, $x$ is standard. Thus, if (5) is solvable in any model of $T \supseteq Q_\forall$, it is solvable in $\mathbb{N}$. $\hfill\square$

We will show that $\mathrm{D}_Q$ is as easy as possible, i.e., NP-complete. Now, the witnesses for satisfiability from Definition 2.11 are polynomial-size objects (if we write all numbers in binary), but it is not immediately clear they can be recognized in polynomial time. In particular, the conditions demand us to test $\widetilde{u}_\ell \simeq \widetilde{v}_\ell$ for subterms $u, v \subseteq t$, which naïvely takes exponential time as the $\widetilde{t}$ reduction from Lemma 2.8 can exponentially blow up sizes of terms (e.g., it unwinds a binary numeral term to the corresponding unary numeral). Fortunately, the offending overlarge pieces have a very boring, repetitive structure, hence we can overcome this obstacle by devising a succinct representation of terms such that on the one hand, the reduction of a given term has a polynomial-size representation, and on the one hand, we can efficiently test whether two representations describe the same term.

The representations we use below (called *descriptors*) have the syntactic form of terms over the language $L_Q$ augmented with extra function symbols $S_n(x)$, $A_{n,m}(x)$, and $B_{n,m}(x, y)$, where $n, m$ are integer indices written in binary. Their exact meaning is explained below, however, the intention is that they facilitate implementation of the operations (especially multiplication) introduced in the proof of Lemma 2.8.

**Definition 3.3** We define a set of expressions called *(term) descriptors*, and for each descriptor $t$ a term $\mathrm{d}(t)$ which it denotes, as follows.

- The constant 0 and variables $x_i$ are descriptors denoting themselves.

- If $t, u$ are descriptors, then $t + u$ and $t \cdot u$ are descriptors, and $\mathrm{d}(t + u) = \mathrm{d}(t) + \mathrm{d}(u)$, $\mathrm{d}(t \cdot u) = \mathrm{d}(t) \cdot \mathrm{d}(u)$.

- If $t$ is a descriptor, and $n \geq 1$ is written in binary, then $S_n(t)$ is a descriptor, and $\mathrm{d}(S_n(t)) = S^n(\mathrm{d}(t))$.

- If $u$ is a descriptor, and $n \geq 0$, $m \geq 2$ are written in binary, then $A_{n,m}(t)$ is a descriptor, and
$$\mathrm{d}(A_{n,m}(u)) = \underbrace{S^n(\ldots(S^n}_{m-1}(0 + \underbrace{\mathrm{d}(u))\ldots) + \mathrm{d}(u)) + \mathrm{d}(u)}_{m}.$$

- If $t, u$ are descriptors, and $n \geq 0$, $m \geq 1$ are written in binary, then $B_{n,m}(t, u)$ is a descriptor, and
$$\mathrm{d}(B_{n,m}(t, u)) = \underbrace{S^n(\ldots(S^n}_{m-1}(S^n(\mathrm{d}(u)) \cdot \mathrm{d}(t) + \underbrace{\mathrm{d}(u))\ldots) + \mathrm{d}(u)) + \mathrm{d}(u)}_{m}.$$

A descriptor is *minimal* if it contains no subdescriptors of the form

$$
\begin{array}{ll}
S_n(S_m(t)), & n, m \geq 1, \\
S_n(0 + u) + u, & n \geq 0, \\
S_n(u) \cdot t + u, & n \geq 0, \\
S_n(A_{n,m}(u)) + u, & n \geq 0,\ m \geq 2, \\
S_n(B_{n,m}(t, u)) + u, & n \geq 0,\ m \geq 1,
\end{array}
$$

where $S_0(t)$ is understood as $t$.

Notice that the definitions of $\mathrm{d}(A_{n,m}(u))$ and $\mathrm{d}(B_{n,m}(t, u))$ are short of an outer $S^n$ as compared to the relevant clauses in Lemma 2.8. The reason for this choice is that in the inductive construction of $\tilde{t}$, we need to be able to peel off easily the outer stack of $S$'s from the terms we got from the inductive hypothesis in order to proceed.

**Lemma 3.4**

(i) *Given a descriptor $t$, we can compute in polynomial time a minimal descriptor $t'$ such that $\mathrm{d}(t) \simeq \mathrm{d}(t')$.*

(ii) *Given a term $t$, we can compute in polynomial time a descriptor $t'$ such that $\mathrm{d}(t') \simeq \tilde{t}$.*

(iii) *If $t_0, t_1$ are minimal descriptors such that $\mathrm{d}(t_0) \simeq \mathrm{d}(t_1)$, then $t_0 \simeq t_1$.*

(iv) *Given descriptors $t$ and $u$, we can test in polynomial time whether $\mathrm{d}(t) \simeq \mathrm{d}(u)$.*

*Proof:* (i): We minimize the descriptor by applying the following rules to its subdescriptors in arbitrary order:

$$
\begin{aligned}
S_n(S_m(t)) &\longrightarrow S_{n+m}(t), \\
S_n(0 + u) + u &\longrightarrow A_{n,2}(u), \\
S_n(u) \cdot t + u &\longrightarrow B_{n,1}(t, u), \\
S_n(A_{n,m}(u)) + u &\longrightarrow A_{n,m+1}(u), \\
S_n(B_{n,m}(t, u)) + u &\longrightarrow B_{n,m+1}(t, u),
\end{aligned}
$$

where $n, m$ are as appropriate for each case according to Definition 3.3. Each rule strictly decreases the number of function symbols in the descriptor, hence the procedure stops after polynomially

many steps, and it clearly produces a minimal descriptor. Also, the maximal length (in binary) of numerical indices increases by at most 1 in each step, hence all descriptors produced during the process have polynomial size, and the algorithm runs in polynomial time.

(ii): By a straightforward bottom-up approach mimicking the definition in Lemma 2.8, we compute for each subterm $u \subseteq t$ a descriptor $u'$ such that $d(u') \simeq \widetilde{u}$, and $u'$ has the form $S_n(u'')$ where $u''$ is 0 or denotes an irreducible term. As in Lemma 2.2, all numerical indices appearing during the computation have bit-length bounded by the size of $t$. If $u \simeq u_0 + u_1$ or $u \simeq u_0 \cdot u_1$, then $u'$ can be expressed by at most one occurrence of each of $u_0''$, $u_1''$, and a bounded number of other symbols (by employing the $A_{n,m}$ and $B_{n,m}$ functions). It follows easily that all descriptors constructed during the computation have polynomial size, and the computation works in polynomial time.

(iii): By induction on the complexity of $d(t_0), d(t_1)$. If $t_0 \simeq u_0 \cdot v_0$, then $t_1$ must be of the form $u_1 \cdot v_1$, as other descriptors denote terms whose topmost symbols are different from $\cdot$. Then $d(u_0) \simeq d(u_1)$ and $d(v_0) \simeq d(v_1)$, hence $u_0 \simeq u_1$ and $v_0 \simeq v_1$ by the induction hypothesis, hence $t_0 \simeq t_1$. A similar argument applies when the topmost symbol of $t_0$ or $t_1$ is a variable, 0, or $S_n$ (in the last case, we use the fact that if $t_i \simeq S_n(u_i)$, then $u_i$ cannot have topmost symbol $S_m$ by minimality).

The remaining cases are when both $t_i$ are of the forms $u_i + v_i$, $A_{n_i,m_i}(v_i)$, or $B_{n_i,m_i}(u_i, v_i)$, so that the topmost symbol of $d(t_i)$ is $+$. We have $d(t_i) \simeq d(w_i) + d(v_i)$, where

$$
w_i \simeq \begin{cases}
u_i & t_i \simeq u_i + v_i, \\
S_{n_i}(A_{n_i,m_i-1}(v_i)) & t_i \simeq A_{n_i,m_i}(v_i), m_i \geq 3, \\
S_{n_i}(0 + v_i) & t_i \simeq A_{n_i,2}(v_i), \\
S_{n_i}(B_{n_i,m_i-1}(u_i, v_i)) & t_i \simeq B_{n_i,m_i}(u_i, v_i), m_i \geq 2, \\
S_{n_i}(v_i) \cdot u_i & t_i \simeq B_{n_i,1}(u_i, v_i).
\end{cases}
$$

Here $n_i$ may be 0, in which case $S_{n_i}$ is void. The descriptor $w_i$ as given here is minimal, except that in the last case, it might happen that $v_i$ itself starts with $S_{k_i}$ for some $k_i$; in that case, we modify $w_i$ in the obvious way. Since $d(w_0) \simeq d(w_1)$ and $d(v_0) \simeq d(v_1)$ are proper subterms of $d(t_0) \simeq d(t_1)$, we may now apply the induction hypothesis, yielding $v_0 \simeq v_1$, and $w_0 \simeq w_1$. By inspection, we see that for each of the five clauses of the definition of $w_i$, we can read off the original parameters ($n_i$, $m_i$, $u_i$; we already know $v_i$) from $w_i$. Moreover, two distinct clauses cannot result in the same $w_i$: the only problematic case is the first clause, where we need to use the minimality of $t_i$. Thus, all in all, $w_i$ and $v_i$ uniquely determine $t_i$, hence we obtain $t_0 \simeq t_1$.

(iv) follows from (i) and (iii). $\qquad \square$

**Theorem 3.5** $D_Q$ *is* NP-*complete.*

*Proof:* NP-hardness is Corollary 3.2, hence in view of Lemmas 2.2 and 2.13, it suffices to show that we can check the existence of a witness $\ell$ for $t = \underline{n}$ in NP. It is immediate from the definition that $\ell$ has size polynomial in $\log n$ and in the length of $t$ if we write labels in binary, so it remains to verify conditions (i)–(iv) in polynomial time.

Conditions (i) and (iv) are clearly polynomial-time. As for (ii), notice first that it makes no difference whether we use unary or binary numerals in the construction of $u_\ell$, as both end up the same after applying $\widetilde{\phantom{u}}$. Thus, in order to test $\widetilde{u}_\ell \simeq \widetilde{v}_\ell$ in polynomial time, we can compute $u_\ell, v_\ell$ using binary numerals, compute descriptors denoting $\widetilde{u}_\ell, \widetilde{v}_\ell$ using Lemma 3.4 (ii), and compare them using Lemma 3.4 (iv).

Condition (iii) is similar: given a term $u$, we can compute a minimal descriptor for $\widetilde{u}_\ell$ in polynomial time, and then check easily whether it has the form $S_k(0)$, and if so, extract $k$. $\qquad \square$

# 4 Conclusion

Unlike stronger theories of arithmetic, we have seen that satisfiability of Diophantine equations in models of $Q$ can be tested in NP, hence undecidability only sets in for more complicated $\Sigma_1$ sentences. The proof also revealed that Robinson's arithmetic can divide standard numbers by zero with ruthless efficiency (albeit in a lopsided way).

Some related questions suggest themselves, such as how far can we push the argument? On the one hand, the criterion in Lemma 2.13 does not use in any way that we are dealing with a single equation. Considering also that the models constructed in Lemma 2.8 only equate terms with the same reduced form, we obtain easily the following generalization:

**Proposition 4.1** *Q-satisfiability of existential sentences, all of whose positively occurring atomic subformulas are of the form $t = \underline{n}$, is decidable, and NP-complete.* □

On the other hand, the reduction in Lemma 2.2 breaks down already for conjunctions of two equations, hence we are led to

**Problem 4.2** *Is Q-satisfiability of existential sentences decidable?*

A question in another vein is how much stronger can we make the theory while maintaining decidability. Observe that the simple argument in Proposition 2.5 applies not just to $Q^+$ itself, but also to all its extensions valid in the variant $\mathbb{N}^\infty$ model used in the proof. This model is actually quite nice: a totally ordered commutative semiring, one pesky axiom short of the theory $PA^-$!

**Problem 4.3** *Is $\mathrm{D}_{PA^-}$ decidable?*

This problem appears to be essentially as hard as the decidability of $\mathrm{D}_{IOpen}$ mentioned in the introduction, cf. [13, 3].

# References

[1] Peter G. Doyle and John H. Conway, *Division by three*, arXiv:math/0605779 [math.LO], 1994, http://arxiv.org/abs/math/0605779.

[2] Peter G. Doyle and Cecil Qiu, *Division by four*, arXiv:1504.01402 [math.LO], 2015, http://arxiv.org/abs/1504.01402.

[3] Lou van den Dries, *Which curves over* **Z** *have points with coordinates in a discrete ordered ring?*, Transactions of the American Mathematical Society 264 (1981), no. 1, pp. 181–189.

[4] Jan van Eijck, Rosalie Iemhoff, and Joost J. Joosten (eds.), *Liber amicorum Alberti: A tribute to Albert Visser*, Tributes vol. 30, College Publications, London, 2016.

[5] Haim Gaifman and Constantinos Dimitracopoulos, *Fragments of Peano's arithmetic and the MRDP theorem*, in: Logic and algorithmic, Monographie de L'Enseignement Mathématique no. 30, Université de Genève, 1982, pp. 187–206.

[6] Richard Kaye, *Diophantine induction*, Annals of Pure and Applied Logic 46 (1990), no. 1, pp. 1–40.

[7] ———, *Hilbert's tenth problem for weak theories of arithmetic*, Annals of Pure and Applied Logic 61 (1993), no. 1–2, pp. 63–73.

[8] Jan Krajíček, *Bounded arithmetic, propositional logic, and complexity theory*, Encyclopedia of Mathematics and Its Applications vol. 60, Cambridge University Press, 1995.

[9] Kenneth L. Manders and Leonard M. Adleman, *NP-complete decision problems for binary quadratics*, Journal of Computer and System Sciences 16 (1978), no. 2, pp. 168–184.

[10] Margarita Otero, *On Diophantine equations solvable in models of open induction*, Journal of Symbolic Logic 55 (1990), no. 2, pp. 779–786.

[11] rainmaker, *Decidability of diophantine equation in a theory*, MathOverflow, 2015, `http://mathoverflow.net/q/194491`.

[12] John C. Shepherdson, *A nonstandard model for a free variable fragment of number theory*, Bulletin de l'Académie Polonaise des Sciences, Série des Sciences Mathématiques, Astronomiques et Physiques 12 (1964), no. 2, pp. 79–86.

[13] Alex J. Wilkie, *Some results and problems on weak systems of arithmetic*, in: Logic Colloquium '77 (A. Macintyre, ed.), North-Holland, 1978, pp. 285–296.