

Combining decision procedures for the reals

Jeremy Avigad and Harvey Friedman*

September 15, 2006

Abstract

We address the general problem of determining the validity of boolean combinations of equalities and inequalities between real-valued expressions. In particular, we consider methods of establishing such assertions using only restricted forms of distributivity. At the same time, we explore ways in which “local” decision or heuristic procedures for fragments of the theory of the reals can be amalgamated into global ones.

Let $T_{add}[\mathbb{Q}]$ be the first-order theory of the real numbers in the language with symbols $0, 1, +, -, <, \dots, f_a, \dots$ where for each $a \in \mathbb{Q}$, f_a denotes the function $f_a(x) = ax$. Let $T_{mult}[\mathbb{Q}]$ be the analogous theory for the language with symbols $0, 1, \times, \div, <, \dots, f_a, \dots$. We show that although $T[\mathbb{Q}] = T_{add}[\mathbb{Q}] \cup T_{mult}[\mathbb{Q}]$ is undecidable, the universal fragment of $T[\mathbb{Q}]$ is decidable. We also show that terms of $T[\mathbb{Q}]$ can fruitfully be put in a normal form. We prove analogous results for theories in which \mathbb{Q} is replaced, more generally, by suitable subfields F of the reals. Finally, we consider practical methods of establishing quantifier-free validities that approximate our (impractical) decidability results.

1 Introduction

This paper is generally concerned with the problem of determining the validity of boolean combinations of equalities and inequalities between real-valued expressions. Such computational support is important not only for the formal verification of mathematical proofs, but, more generally, for any application which depends on such reasoning about the real numbers.

Alfred Tarski’s proof [23] that the theory of the real numbers as an ordered field admits quantifier-elimination is a striking and powerful response to the problem. The result implies decidability of the full first-order theory, not just the quantifier-free fragment. George Collins’s [10] method of cylindrical algebraic decomposition made this procedure feasible in practice, and ongoing research in computational real geometry has resulted in various optimizations and alternatives (see e.g. [14, 6, 5]). Recently, a proof-producing version of an

*Work by both authors has been supported by NSF grant DMS-0401042. We are grateful to three anonymous referees for numerous comments, suggestions, and corrections.

elimination procedure due to Paul Cohen has even been implemented in the framework of a theorem prover for higher-order logic [20].

There are two reasons, however, that one might be interested in alternatives to q.e. procedures for real closed fields. The first is that their generality means that they can be inefficient in restricted settings. For example, one might encounter an inference like

$$0 < x < y \rightarrow (1 + x^2)/(2 + y)^{17} < (1 + y^2)/(2 + x)^{10},$$

in an ordinary mathematical proof. Such an inference is easily verified, by noticing that all the subterms are positive and then chaining through the obvious inferences. Computing sequences of partial derivatives, which is necessary for the full decision procedure, seems misguided in this instance. A second, more compelling reason to explore alternatives is that decision procedures for real closed fields are not extensible. For example, adding trigonometric functions or an uninterpreted unary function symbol renders the full first-theory undecidable. Nonetheless, an inference like

$$0 < x < y \rightarrow (1 + x^2)/(2 + e^y) < (2 + y^2)/(1 + e^x)$$

is also straightforward, and it is reasonable to seek procedures that capture such inferences.

The unfortunate state of affairs is that provability in most interesting mathematical contexts is undecidable, and even when decision procedures are available in restricted settings, they are often infeasible or impractical. This suggests, instead, focusing on heuristic procedures that traverse the search space by applying a battery of natural inferences in a systematic way (for some examples in the case of real arithmetic, see [7, 17, 25]). There has been, nonetheless, a resistance to the use of such procedures in the automated reasoning community. For one thing, they do not come with a clean theoretical characterization of the algorithm's behavior, or the class of problems on which one is guaranteed success. This is closely linked to the fact that the algorithms based on heuristics are brittle: small changes and additions as the system evolves can have unpredictable effects.

The strategy we pursue here is to develop a theoretical understanding that can support the design of such heuristic procedures, by clarifying the possibilities and limitations that are inherent in a method, and providing a general framework within which to situate heuristic approaches. One observation we exploit here is that often distributivity is used only in restricted ways in the types of verifications described above. Arguably, any inference that requires factoring a complex expression does not count as "obvious." Conversely, multiplying through a sum can result in the loss of valuable information, as well as lead to increases in the lengths of terms. As a result, steps like these are usually spelled out explicitly in textbook reasoning when they are needed. It is therefore natural to ask whether one can design procedures that reasonably handle those inferences that do not make use of distributivity, relying on the user or other methods to then handle the latter.

The “distributivity-free” fragment of the theory of the reals as an ordered field can naturally be viewed as a combination of the additive and multiplicative fragments, each of which is easily seen to be decidable. This points to another motivation for our approach. A powerful paradigm for designing useful search procedures involves starting with procedures that work locally, for restricted theories, and then amalgamating them into a global procedure in some principled way. For example, Nelson-Oppen methods are currently used to combine decision procedures for theories that are disjoint except for the equality symbol, yielding decision procedures for the universal fragment of their union. Shostak methods perform a similar task more efficiently by placing additional requirements on the theories to be amalgamated. (See [18, 22] and the introduction to [3] for overviews of the various approaches.) Such methods are appealing, in that they allow one to unify such decision procedures in a uniform and modular way. This comes closer to what ordinary mathematicians do: in simple, domain-specific situations, we know exactly how to proceed, whereas in more complex situations, we pick out the fragments of a problem that we know how to cope with and then try to piece them together. One would therefore expect the notion of amalgamating decision procedures, or even heuristic procedures, to be useful when there is more significant overlap between the theories to be amalgamated. For example, the Nelson-Oppen procedure has been generalized in various ways, such as to theories whose overlap is “locally finite” [15]. Our results here show what can happen when one tries to amalgamate decision procedures for theories where the situation is not so simple.

Sections 2 and 3, below, provide general background. In Section 2, we discuss the theoretical results that underly Nelson-Oppen methods for combining decision procedures for theories that share only the equality symbol, or for theories with otherwise restricted overlap. In Section 3, we describe some particular decision procedures for fragments of the reals, which are candidates for such a combination.

In Section 4, we define the theories $T[F]$, which combine the additive and multiplicative fragments of the theory of the reals, allowing multiplicative constants from a field F . The theory $T[F]$, in particular, can, alternatively, be thought of as the theory of real closed fields minus distributivity, except for constants in F . Because of the nontrivial overlap, Nelson-Oppen methods no longer apply. In Section 5, we provide two examples that clarify what these theories can do. On the positive side, we show that when a multivariate polynomial has no roots on a compact cube, $T[\mathbb{Q}]$ is strong enough to prove that fact. On the negative side, we show that the theories $T[F]$ cannot prove $x^2 - 2x + 1 \geq 0$, a fact which is easily proved using distributivity.

In Sections 6–8 we establish our decidability results. Using a characterization of the universal fragment of $T[F]$ developed in Section 6, we show, in Section 7, that whenever F is an appropriately computable subfield of \mathbb{R} , the universal fragment of $T[F]$ is decidable. So, in particular, the universal fragment of $T[\mathbb{Q}]$ is decidable. In Section 8, we describe normal forms for terms of $T[F]$, which make it easy to determine whether two terms are provably equal. We also show that these provable equalities are independent of the parts of the theory that

have to do with the ordering.

In Sections 9–11, we establish our undecidability results. In Section 9, we present a flexible technique that will allow us to build suitable models of the theories $T[F]$. In Section 10, we use this technique to reduce the problem of determining the truth of an existential sentence over the field F to that of the provability of a related formula in $T[F]$. As a result, if Diophantine equations in the rationals are unsolvable (which is generally believed to be the case), then so is the set of existential consequences of $T[\mathbb{Q}]$. In Section 11, we reduce the problem of determining the solvability of a Diophantine equation in the integers to the provability of a related $\forall\forall\forall\exists^*$ -sentence in any $T[F]$. As a result, we have an unconditional undecidability result for that fragment.

The procedure implicit in our decidability results is not useful in practice: it works by reducing the question as to whether a universal sentence is provable in $T[F]$ to the question as to whether a more complex sentence is provable in the theory of real closed fields, and then appeals to the decidability of the latter. In Sections 12–14, we consider the problem of designing pragmatic procedures that approximate our decidability results, are more flexible than decision procedures for real closed fields, and work reasonably well on ordinary textbook inferences. In Section 12, we suggest a restriction of the theories $T[F]$ which avoids disjunctive case splits, which are a key source of infeasibility. In Section 13, we describe a search procedure that works along these lines, making use of the normal forms introduced in Section 8. In Section 14, we indicate a number of directions in which one might extend and improve our crude algorithm.

Finally, in Section 15, we offer some final thoughts and conclusions.

2 Combining decision procedures

In this section, we briefly review the mathematical foundation for the Nelson-Oppen combination procedure [21]. For more detail, see [4, 11, 16, 24]; an important program verification system based on these methods is described in [13].

Let Δ be the set of first-order formulas in the language of equality asserting that the universe is infinite. A theory T is said to be *stably infinite* if whenever $T \cup \Delta$ proves a universal sentence φ , then T proves it as well. Equivalently, T is stably infinite if whenever a quantifier-free formula is satisfied in any model of T , it is satisfied in some infinite model of T . In particular, if T only has infinite models, then T is stably infinite.

The Nelson-Oppen procedure for combining decidable theories of equality is based on the following:

Theorem 2.1. *Suppose T_1 is a theory in a language L_1 , T_2 is a theory in a language L_2 , T_1 and T_2 are stably infinite, and the languages L_1 and L_2 are disjoint except for the equality symbol. Suppose the universal fragments of T_1 and T_2 are decidable. Then the universal fragment of $T_1 \cup T_2$ is decidable.*

The proof of Theorem 2.1 is not difficult. The question as to whether $T_1 \cup T_2$ proves a universal formula is equivalent to the question as to whether it proves the quantifier-free matrix. (One can treat the free variables as new constants, if one prefers, but here and below we will speak in terms of proving or refuting sets of formulas with free variables.) Since any quantifier-free formula can be put in conjunctive normal form, the problem reduces to that of determining provability of disjunctions of literals, or, equivalently, that of determining whether $T_1 \cup T_2$ refutes a conjunction of literals.

Let Γ be a set of literals. The first step in the procedure is to introduce new variables to “separate terms.” For example, the universal closure of a formula of the form $\varphi(f(s_1, \dots, s_k))$ is equivalent to the universal closure of $x = f(s_1, \dots, s_k) \rightarrow \varphi(x)$, where x is a new variable. This is, in turn, equivalent to the universal closure of $y_1 = s_1 \wedge \dots \wedge y_k = s_k \wedge x = f(y_1, \dots, y_k) \rightarrow \varphi(x)$. By introducing new variables in this way, we can obtain sets of equalities Π_1 and Π_2 in L_1 and L_2 respectively, and a set of literals, Π_3 , in which no function symbols occur, such that $T_1 \cup T_2$ refutes Γ if and only if it refutes $\Pi_1 \cup \Pi_2 \cup \Pi_3$. Let Γ_1 be Π_1 together with the literals in Π_3 that are in L_1 , and let Γ_2 be Π_2 together with the literals in Π_3 that are in L_2 . Then each Γ_i is in the language of T_i , and $T_1 \cup T_2$ refutes Γ if and only if $T_1 \cup T_2$ refutes $\Gamma_1 \cup \Gamma_2$.

By the Craig interpolation theorem, $T_1 \cup T_2$ refutes $\Gamma_1 \cup \Gamma_2$ if and only if there is a quantifier-free interpolant θ in the common language (i.e. involving only the equality symbol and variables common to both Γ_1 and Γ_2) such that

$$T_1 \cup \Gamma_1 \vdash \theta$$

and

$$T_2 \cup \Gamma_2 \cup \{\theta\} \vdash \perp.$$

By the assumption that T_1 and T_2 are stably infinite, we can assume without loss of generality that each includes Δ . Since the theory of equality in an infinite structure has quantifier-elimination, θ is equivalent to a quantifier-free formula. In fact, we can assume without loss of generality that θ is in disjunctive normal form. So we are looking for a sequence $\theta_1, \dots, \theta_n$ of finite conjunctions of literals such that for each i ,

$$T_1 \cup \Gamma_1 \vdash \theta_1 \vee \dots \vee \theta_n$$

and

$$T_2 \cup \Gamma_2 \cup \{\theta_i\} \vdash \perp$$

for each i .

Each disjunct θ_i describes relationships between the variables \vec{x} of $\Gamma_1 \cup \Gamma_2$, in the language $L_1 \cap L_2$, which has only the equality symbol. The key point is this: over Δ , every “complete type” (that is, complete, consistent set of formulas with free variables \vec{x}) is determined by an exhaustive description of which of the variables are equal to one another and which are not. Furthermore, there are only finitely many such descriptions. Without loss of generality, we can assume that each θ_i is of this form, because otherwise it can be rewritten as a disjunction of such. Thus we simply need to use the decision procedure for T_2 to determine

all the complete types θ_i that can be refuted by $T_2 \cup \Gamma_2$, and then use the decision procedure for T_1 to determine whether $T_1 \cup \Gamma_1$ proves their disjunction. Equivalently, we can use the decision procedures to determine all the complete types that are consistent with either side; Γ can be refuted if and only if there is no complete type that is consistent with both.

This naive procedure is not very efficient. In fact, the Nelson-Oppen procedure iteratively searches for a disjunction of equalities derivable from either $T_1 \cup \Gamma_1$ or $T_2 \cup \Gamma_2$, adds this disjunction to the hypotheses, and then splits across the cases. It is not hard to show that this variant is complete; one can view it in terms using both $T_1 \cup \Gamma_1$ and $T_2 \cup \Gamma_2$ to derive a sequence of increasingly strong disjunctions of conjunctions of positive literals, until either a contradiction is reached or no further strengthening can be found. In the latter case, one can read off a complete type consistent with both $T_1 \cup \Gamma_1$ and $T_2 \cup \Gamma_2$. The procedure is much more efficient if either of the theories T_i is *convex*, that is, whenever φ is a conjunction of literals and $T_i \cup \varphi \vdash x_1 = y_1 \vee \dots \vee x_k = y_k$ then $T_i \cup \varphi \vdash x_i = y_i$ for some i . The linear theory of the reals has this property, though the multiplicative theory does not. Shostak's procedure provides further optimization under the assumptions that terms in the theory are "canonizable" and "solvable," again, features that are commonly satisfied.

For future use, we record the effects of "separating terms," as described above. We no longer assume L_1 and L_2 are disjoint languages.

Proposition 2.2. *Let φ be any universal sentence in the language $L_1 \cup L_2$. Then φ is equivalent to a sentence of the form*

$$\forall \vec{x} (\theta_1(\vec{x}) \wedge \theta_2(\vec{x}) \rightarrow \theta_3(\vec{x})),$$

where θ_1 is a conjunction of equalities in L_1 , θ_2 is a conjunction of equalities in L_2 , and θ_3 is a quantifier-free formula in $L_1 \cup L_2$ with no function symbols. As a result, φ can be written as a conjunction of formulas of the form

$$\forall \vec{x} (\varphi_1(\vec{x}) \vee \varphi_2(\vec{x})), \tag{1}$$

where each φ_i is a quantifier-free formula in L_i . If all the relation symbols in $L_1 \cup L_2$ are common to both L_1 and L_2 , or if the matrix of φ is equivalent to a disjunction of literals, one conjunct of the form (1) suffices.

3 Decision procedures for fragments of the reals

The method described in Section 2 requires only that the universal fragments of the theories T_1 and T_2 are decidable, and that for any sequence of variables, there are only finitely many complete types in the common language, each of which can be described by a single quantifier-free formula. In particular, we have the following:

Theorem 3.1. *Let T_1 and T_2 be theories extending the theory of dense linear orders without endpoints, with only $<$ and $=$ in the common language. If the*

universal fragments of T_1 and T_2 are decidable, then the universal fragment of $T_1 \cup T_2$ is also decidable.

As was the case when equality was the only common symbol, this theorem can be stated even more generally: we only need assume that T_1 and T_2 satisfy the property obtained by replacing Δ by the theory of dense linear orders without endpoints in the definition of “stably infinite” above. Of course, Theorem 3.1 can be iterated to combine theories T_1, T_2, T_3, \dots with the requisite properties.

Let us consider some examples of fragments of the reals that admit quantifier-elimination, and are hence decidable. Note that to eliminate quantifiers from any formula it suffices to be able to eliminate a single existential quantifier, i.e. transform a formula $\exists x \varphi$, where φ is quantifier-free, to an equivalent quantifier-free formula. Since $\exists x (\varphi \vee \psi)$ is equivalent to $\exists x \varphi \vee \exists x \psi$, we can always factor existential quantifiers through a disjunction. In particular, since any quantifier-free formula can be put in disjunctive normal form, it suffices to eliminate existential quantifiers from conjunctions of atomic formulas and their negations. Also, since $\exists x (\varphi \wedge \psi)$ is equivalent to $\exists x \varphi \wedge \psi$ when x is not free in ψ , we can factor out any formulas that do not involve x . Furthermore, whenever we can prove $\forall x (\theta \vee \eta)$, $\exists x \varphi$ is equivalent to $\exists x (\varphi \wedge \theta) \vee \exists x (\varphi \wedge \eta)$; so we can “split across cases” as necessary. We will use all of these facts freely below.

Proposition 3.2. *The theory of $\langle \mathbb{R}, 0, 1, +, -, < \rangle$ admits elimination of quantifiers, and hence is decidable.*

This is theory commonly known as linear arithmetic, and is the same as the theory of divisible ordered abelian groups. The universal fragment coincides with that of the theory of ordered abelian groups. The method of eliminating an existentially quantified variable implicit in the proof is known as the *Fourier-Motzkin procedure*.

Proof (sketch). It is helpful to extend the language to include multiplication by rational coefficients, though we can view this as nothing more than a notational convenience: for example, if n is a natural number, we can take nx to abbreviate $x + x + \dots + x$, and when n, m, k, l are natural numbers with m and l nonzero we can take $(n/m)s = (k/l)t$ to abbreviate $nls = kmt$.

Consider a sentence $\exists x \varphi$, where φ is quantifier-free. Writing $s \neq t$ as $s < t \vee t < s$ and $s \not< t$ as $t < s \vee t = s$, we can assume without loss of generality that φ is a positive boolean combination of atomic formulas of the form $s = t$ and $s < t$. Putting φ in disjunctive normal form and factoring the existential quantifier through the disjunction we can assume φ is a conjunction of atomic formulas. Solving for x , we can express each of these in the form $x = s$, $x < s$, or $s < x$, where s does not involve x (atomic formulas that do not involve x can be brought outside of the existential quantifier).

If any of the conjuncts is of the form $x = s$, then $\exists x \varphi(x)$ is equivalent to $\varphi(s)$, which is quantifier-free. So we are reduced to the case where φ is of the form $(\bigwedge_i s_i < x) \wedge (\bigwedge_j x < t_j)$. In that case, it is not hard to verify that $\exists x \varphi$ is equivalent to $\bigwedge_{i,j} s_i < t_j$. \square

For more on the Fourier-Motzkin procedure, see [1]. In fact, more efficient elimination procedures are available, and are not much more complicated; see [19, 27].

Proposition 3.3. *The theory of $\langle \mathbb{R}, 0, 1, -1, \times, \div, < \rangle$ with the convention $x \div 0 = 0$ admits elimination of quantifiers, and hence is decidable.*

Proof (sketch). Since $\langle \mathbb{R}^{>0}, 1, \times, \div, < \rangle$ is isomorphic to $\langle \mathbb{R}, 0, +, -, < \rangle$, the previous argument shows that the theory of this structure has quantifier-elimination. For the larger structure, consider $\exists x \varphi$, where φ is quantifier-free. As above, we can assume φ is a conjunction of equalities and strict inequalities. Introducing case splits we can assume that φ determines which variables are positive, negative, or 0. Temporarily replacing negative variables by their negations, we can further assume that φ implies that all the variables are positive. Bringing negation symbols to the front of each term, we are left with a conjunction of atomic formulas of the form $\pm s < \pm t$, where s and t are products of variables assumed to be positive. But then $-s < t$ is equivalent to \top ; $s < -t$ is equivalent to \perp ; and $-s < -t$ is equivalent to $t < s$. Similarly, $-s = -t$ is equivalent to $s = t$, and both $s = -t$ and $-s = t$ are equivalent to \perp . So, we are reduced to the case where all the variables are positive. \square

Proposition 3.4. *The theory of $\langle \mathbb{R}, \exp, \ln, 0, 1, < \rangle$, where $\exp(x) = e^x$ and $\ln(x) = 0$ for non-positive x , admits quantifier-elimination, and hence is decidable.*

Proof (sketch). Once again, we are reduced to the case of eliminating a quantifier of the form $\exists x \varphi$ where φ is a conjunction of equalities and strict inequalities. Expressions of the form $\ln(\exp(s))$ simplify to s , and across a case split of the form $s > 0 \vee s \leq 0$ an expression of the form $\exp(\ln(s))$ simplifies to s or 0 . Using the equivalences $s < t \leftrightarrow \exp(s) < \exp(t)$ and introducing case splits as necessary, we are reduced to the case where φ is a conjunction of terms of the form $u < \exp^n(v)$, $u > \exp^n(v)$, and $u = \exp^n(v)$, where u and v are variables and $\exp^n(u)$ denotes n applications of \exp to u . If there is an equality using x , we can use that to eliminate the existential quantifier. Otherwise, for suitable k we can arrange that φ is a conjunction of formulas of the form $s_i < \exp^k(x)$ and $\exp^k(x) < t_j$, in which case $\exists x \varphi$ is equivalent to $(\bigwedge_{i,j} s_i < t_j) \wedge (\bigwedge_j 0 < t_j)$. \square

From Theorem 3.1 we have:

Corollary 3.5. *The universal fragment of the union of the three theories above is decidable.*

The decision procedure implicit in the proof of Corollary 3.5 is, unfortunately, not very useful. There is a sense in which it does too little, and another sense in which it does too much.

A sense in which the procedure does too little is that the union of the three theories is too weak. For example, it is not hard to show (either using the interpolation theorem or a model-theoretic argument) that the theory does not

prove $\bar{2} \times \bar{2} = \bar{4}$, where $\bar{2}$ abbreviates the term $1 + 1$, and $\bar{4}$ abbreviates $1 + 1 + 1 + 1$. Similarly, it fails to prove $x + x = \bar{2}x$. In the next section, we will focus on the additive and multiplicative fragments of the reals, and respond to this problem by augmenting the structures to allow multiplication by arbitrary rational constants, or, more generally, constants from a suitably computable subfield F of the reals. Unfortunately, this means that the two structures share a language with infinitely many function symbols, and so the methods described in the last section can no longer be used. We will have to do a good deal of additional work to establish decidability in this case.

A sense in which the algorithm implicit in the proof of Corollary 3.5 does too much is that even in the absence of the new multiplicative constants, it is inefficient: the combination procedure relies on the fact that one can enumerate all possible descriptions of equalities and inequalities between variables, and, in general, the number of possibilities grows exponentially. Our proof of decidability for the augmented theories involves a reduction to the theory of real-closed fields, and so it does not represent a practical advance either. In Sections 12–14, we will address the issue of developing practical procedures that approximate the theories we describe here.

4 The theories $T[F]$

Let F denote any subfield of the reals. Let $T_{add}[F]$ be the theory of the real numbers for the language with symbols

$$0, 1, +, -, <, \dots, f_a, \dots$$

where for $a \in F$, f_a denotes the function $f_a(x) = ax$. Let $T_{mult}[F]$ be the analogous theory for the language with symbols

$$0, 1, \times, \div, <, \dots, f_a, \dots$$

where $x \div y$ is interpreted as 0 when $y = 0$. Our central concern in this paper is the union of these two theories, $T[F] = T_{add}[F] \cup T_{mult}[F]$. It will also be useful to denote their intersection, $T_{add}[F] \cap T_{mult}[F]$, by $T_{comm}[F]$. It often makes sense to restrict one's attention to computable subfields F of the real numbers; in particular, \mathbb{Q} , the minimal such subfield, is a natural choice. We will see below that, in a sense, the field of real algebraic numbers \mathbb{A} represents a maximal choice. Intermediate choices are also possible; for example, one might consider the smallest field containing \mathbb{Q} and closed under taking roots of positive numbers. It should be clear that each $T[F]$ proves, for example, $\bar{2} \times \bar{2} = \bar{4}$ and $x + x = \bar{2}x$.

We claim that the theories $T[F]$ are natural, and are sufficient to justify many of the inferences that come up in ordinary mathematical texts. The latter claim is an empirical one, however, and we will not try to justify it here.

Each of $T_{comm}[F]$, $T_{add}[F]$, and $T_{mult}[F]$ has quantifier-elimination, and hence is complete. The elimination procedures sketched in Section 3 can easily be extended to $T_{add}[F]$ and $T_{mult}[F]$, assuming the operations on F are

computable, in which case these theories are decidable as well. Similarly, a quantifier-elimination procedure for $T_{comm}[F]$ is easily obtained by extending the usual procedure for dense linear orders without endpoints, so this theory is also complete, and decidable when F is computable.

Reflecting these elimination procedures yields complete axiomatizations of the relevant theories. The theory $T_{comm}[F]$ is axiomatized by the following:

1. $<$ is a dense linear order
2. $0 < 1$
3. $f_a(f_b(x)) = f_{ab}(x)$, for every $a, b \in F$
4. $f_0(x) = 0$, $f_1(x) = x$
5. $x < y \leftrightarrow f_a(x) < f_a(y)$ for $0 < a \in F$
6. $x < y \leftrightarrow f_a(x) > f_a(y)$ for $0 > a \in F$
7. $0 < x \rightarrow x < f_a(x)$ for $1 < a \in F$

One obtains an axiomatization for $T_{add}[F]$ by adding the following:

1. $0, +, <$ is an ordered abelian group
2. $x - y = z \leftrightarrow x = y + z$
3. $f_a(x + y) = f_a(x) + f_a(y)$
4. $f_{a+b}(x) = f_a(x) + f_b(x)$

Similarly, one obtains an axiomatization of $T_{mult}[F]$ by adding the following to $T_{comm}[F]$:

1. $1, \times, <$ is a divisible ordered abelian group on the positive elements
2. $x/y = z \leftrightarrow (y = 0 \wedge z = 0) \vee x = yz$
3. $f_a(xy) = f_a(x)y$

In Sections 9–11, we will prove undecidability results for fragments of $T[F]$. We will find it useful to work with the following alternative system, $T[F]^*$, based on the symbols $0, 1, +, \times, <$ together with constant symbols c_a for $a \in F$. The axioms of $T[F]^*$ fall naturally into four groups:

1. $0, +, <$ is an ordered abelian group
2. $1, \times, <$ is a divisible ordered abelian group on the positive elements
3. (a) $c_{a+b} = c_a + c_b$, for $a, b \in F$
(b) $c_{ab} = c_a \times c_b$, for $a, b \in F$
(c) $0 < c_a$ for $0 < a$, $a \in F$

4. (a) $c_{a+b} \times x = (c_a \times x) + (c_b \times x)$, for $a, b \in F$
 (b) $c_a \times (x + y) = (c_a \times x) + (c_a \times y)$, for $a \in F$

Note that the extra symbols in the language of $T[F]$ are easily definable in $T[F]^*$. It is straightforward to verify the following.

Lemma 4.1. *Let φ be a formula in the language of $T[F]$ without $-$, \div . Let φ^* be the result of replacing each occurrence of $f_a(t)$ with $c_a \times t$, inductively, from innermost to outermost. Then φ is provable in $T[F]$ if and only if φ^* is provable in $T[F]^*$.*

Lemma 4.2. *Let φ be a formula in the language of $T[F]^*$. Let φ' be the result of replacing each occurrence of c_a with $f_a(1)$. Then φ is provable in $T[F]^*$ if and only if φ' is provable in $T[F]$.*

Theorem 4.3. *$T[F]$ and $T[F]^*$ prove the same sentences involving only the symbols $0, 1, +, \times, <$.*

Below we will call the symbols f_a the *auxiliary function symbols* and the symbols c_a the *auxiliary constant symbols*. For readability, we will write ax instead of $f_a(x)$ or c_ax when the context makes the meaning clear.

The following shows that as far as provability of formulas in the language of real closed fields is concerned, there is never a need to go beyond the real algebraic numbers in choosing F .

Theorem 4.4. *$T[\mathbb{R}]$ is a conservative extension of $T[\mathbb{A}]$.*

Proof. Since \div and $-$ are definable in terms of the other symbols of $T[F]$, we can focus on sentences in which these symbols do not occur, and use Theorem 4.3.

Let d be a proof of a sentence φ in $T[\mathbb{R}]^*$, where φ is in the language of $T[\mathbb{A}]^*$. Assign variables \vec{y} to the auxiliary constant symbols occurring in φ , and let $\psi(\vec{y})$ define the corresponding real algebraic numbers in the language of real closed fields. Assign variables \vec{z} to all the additional auxiliary constant symbols occurring in d , and let $\theta(\vec{y}, \vec{z})$ be the conjunction of all the axioms of $T[\mathbb{R}]$ used in d , with the constants replaced by the corresponding variables. The assertion $\exists \vec{y}, \vec{z} (\theta(\vec{y}, \vec{z}) \wedge \psi(\vec{y}))$ is true of the real numbers, and so, by transfer (i.e. the completeness of the theory of real closed fields, of which both the reals and the real algebraic numbers are a model), it is true of \mathbb{A} as well. Let \vec{a}, \vec{b} be real algebraic numbers witnessing the existential quantifiers. Because $\psi(\vec{a})$ determines \vec{a} uniquely, \vec{a} corresponds to the original auxiliary constant symbols in φ . Thus we have the even stronger result that d can be interpreted as a proof in $T[\mathbb{A}]^*$, taking the constant symbols to denote \vec{a}, \vec{b} . \square

This argument shows, more generally, that to prove a sentence with auxiliary function symbols f_{a_1}, \dots, f_{a_n} , there is no need to go beyond the real algebraic closure of $\{a_1, \dots, a_n\}$.

5 Examples

To provide a better feel for the theories $T[F]$, in this section we consider some theorems that clarify their strength. The first theorem provides a lower bound by showing that a decision procedure for the universal fragment of any $T[F]$ implies a decision procedure for the existence of roots of a multivariate polynomial on the unit cube.

Theorem 5.1. *Let F be any subfield of the real numbers, and let $f(x_1, \dots, x_k)$ be a multivariate polynomial with coefficients in F . Let $I = [0, 1]^k$ be the compact k -dimensional unit cube. Then f is nonzero on I if and only if $T[F]$ proves that fact.*

Proof. The “if” direction follows from the fact that the axioms of $T[F]$ are true of the real numbers. On the other hand, by the intermediate value theorem, if a polynomial function f is nonzero on I , then it is either strictly positive or strictly negative on I . So it suffices to show that if f is strictly positive on I , then $T[F]$ proves that this is the case.

Suppose $f(\vec{x}) = \sum_{i < n} t_i(\vec{x})$, where each t_i is a monomial in x_1, \dots, x_l with a coefficient in F , and suppose f is strictly positive on I . Given a point $\langle a_1, \dots, a_k \rangle$ in I , let $r_{\vec{a}} = f(\vec{a}) > 0$, and for each i , let $r_{\vec{a},i} = t_i(\vec{a})$. By continuity, we can find an open neighborhood $U_{\vec{a}}$ of \vec{a} , such that for each $\vec{b} \in U_{\vec{a}}$, $t_i(\vec{b}) > r_{i,\vec{a}} - r_{\vec{a}}/3n$. Shrinking $U_{\vec{a}}$ if necessary, we can assume that $U_{\vec{a}}$ is a product of open intervals with rational endpoints.

By compactness, I is covered by a finite set of these open neighborhoods, say $U_{\vec{a}_1}, \dots, U_{\vec{a}_m}$. Then:

1. $T[F]$ proves $\forall \vec{x} (\vec{x} \in I \rightarrow \vec{x} \in U_{\vec{a}_1} \vee \dots \vee U_{\vec{a}_m})$. In fact, this can be proved by $T_{comm}[F]$, since it is purely a property of the ordering on the rational numbers.
2. For each $j < m$ and $i < n$, $T_{mult}[F]$ proves $\vec{x} \in U_{\vec{a}_j} \rightarrow t_i(\vec{x}) > q_{i,j}$, where $q_{i,j}$ is any rational number less than $r_{\vec{a}_j,i} - r_{\vec{a}_j}/3n$ and greater than $r_{\vec{a}_j,i} - r_{\vec{a}_j}/2n$.
3. Using these lower bounds, for each $j < m$, $T_{add}[F]$ can prove $\vec{x} \in U_{\vec{a}_j} \rightarrow f(\vec{x}) > \sum_{i < n} q_{i,j}$.

The result follows from the fact that in the last claim,

$$\sum_{i < n} q_{i,j} > \sum_{i < n} (r_{\vec{a}_j,i} - r_{\vec{a}_j}/2n) = r_{\vec{a}_j} - r_{\vec{a}_j}/2 = r_{\vec{a}_j}/2 > 0.$$

This completes the proof. □

As an example of something $T[F]$ cannot do, consider the inequality $x^2 - 2x + 1 \geq 0$. That this is generally valid is clear from writing $x^2 - 2x + 1 = (x - 1)^2$, but this equality is a consequence of distributivity, which is not available in $T[F]$. In fact, we have:

Theorem 5.2. *For any F , $T[F]$ proves $\forall x (x^2 - 2x + 1 \geq \varepsilon)$ if and only if $\varepsilon < 0$. In particular, $T[F]$ does not prove $\forall x (x^2 - 2x + 1 \geq 0)$.*

Moreover, proofs of $\forall x (x^2 - 2x + 1 \geq \varepsilon)$ in $T[F]$ necessarily get longer as ε approaches 0, and the results that follow provide explicit lower bounds. Focusing on the domain of the function $x^2 - 2x + 1$ instead of the range, we also have:

Theorem 5.3. *For any F ,*

1. $T[F]$ proves $\forall x (x \leq r \rightarrow x^2 - 2x + 1 \geq 0)$ if and only if $r < 1$.
2. $T[F]$ proves $\forall x (x \geq r \rightarrow x^2 - 2x + 1 \geq 0)$ if and only if $r > 1$.

Theorem 5.3 implies Theorem 5.2. Assuming $x \in [1 - \delta, 1 + \delta]$ for a small rational constant δ , $T[F]$ can easily show $x^2 \geq 1 - 2\delta + \delta^2$ and $2x \leq 2 + 2\delta$, and hence $x^2 - 2x + 1 \geq -4\delta + \delta^2 \geq -4\delta$. So, taking r to be $1 - \delta$ and $1 + \delta$, respectively, in the two clauses Theorem 5.3, we have the “if” direction of Theorem 5.2. But the “only if” direction is a consequence of the fact that $T[F]$ does not prove $\forall x (x^2 - 2x + 1 \geq 0)$, which is immediate from Theorem 5.3.

The two clauses of Theorem 5.3 are proved in a similar way, and so we will only prove the first. Since $T[F]$ easily proves $x < 0 \rightarrow x^2 - 2x + 1 \geq 0$, we can replace the first statement in Theorem 5.3 by $\forall x (0 \leq x \leq r \rightarrow x^2 \geq 2x - 1)$. $T[F]$ proves this if and only if it refutes the set of formulas

$$\{0 \leq x, x \leq r, u = x^2, u < 2x - 1\}.$$

Recall that this happens if and only if there is an interpolant, θ , in disjunctive normal form, such that

$$T_{mult}[F] \cup \{0 \leq x, x \leq r, u = x^2\} \vdash \theta \tag{2}$$

and

$$T_{add}[F] \cup \{u < 2x - 1\} \cup \theta \vdash \perp. \tag{3}$$

So it suffices to show:

Theorem 5.4. *There is a DNF formula θ with at most n disjuncts satisfying (2) and (3) if and only if $r \leq n/(n+1)$.*

Proof. We will first show that if θ has n disjuncts and satisfies (2) and (3) then $r \leq n/(n+1)$. We will then show that, in fact, for $r = n/(n+1)$ such a θ exists.

Write $\theta = \theta_1 \vee \dots \vee \theta_n$, where each θ_i is a conjunction of literals involving only x and u . It is not hard to see that each θ_i is equivalent to a conjunction of literals of the form

$$a \triangleleft x \triangleleft b \wedge c \triangleleft u \triangleleft d \wedge ex \triangleleft u \triangleleft fx$$

where each \triangleleft is either $<$ or \leq (and some of the conjuncts may be absent). $T_{mult}[F] \cup \{0 \leq x, x \leq r, u = x^2\}$ proves this equivalent to a conjunction of the form

$$a \triangleleft x \triangleleft b \wedge a^2 \triangleleft u \triangleleft b^2 \wedge ax \triangleleft u \triangleleft bx \tag{4}$$

for some a, b in $[0, 1]$, and from the point of view of $T_{add}[F] \cup \{2x - 1 < u\}$, each of these disjuncts is no weaker than the original. Thus it suffices to prove the claim for interpolants that are of the form (4).

Now, $T_{add}[F] \cup \{u < 2x - 1\}$ refutes θ if and only if it refutes each disjunct. Thus the following lemma is crucial to our analysis.

Lemma 5.5. *For a, b in $[0, 1]$, $T_{add}[F] \cup \{2x - 1 < u\}$ refutes (4), for any versions of the relation \triangleleft , if and only if $b \leq 1/(2 - a)$.*

Proof. If $b < a$, $T_{add}[F] \cup \{2x - 1 < u\}$ easily refutes (4), and $b \leq 1/(2 - a)$ holds. So it suffices to consider the case $a \leq b$.

We need only work through the Fourier-Motzkin procedure by hand. Eliminating u , we obtain the inequalities $a^2 < 2x - 1$ and $ax < 2x - 1$. (Note that we get strict inequality, whether the initial \triangleleft 's are strict inequalities or not.) Solving for x , we obtain $(a^2 + 1)/2 < x$ and $1/(2 - a) < x$. Eliminating x , we get $(a^2 + 1)/2 < b$ and $1/(2 - a) < b$. This yields a contradiction if and only if b is less than or equal to the minimum of $(a^2 + 1)/2$ and $1/(2 - a)$. A calculation shows that the latter is always smaller for $a \in [0, 1]$, so we have the desired conclusion. \square

We can now finish off the proof of Theorem 5.4. Suppose $T_{mult}[F] \cup \{0 \leq x, x \leq r, u = x^2\}$ proves a disjunction $\theta_1 \vee \dots \vee \theta_n$ with each θ_i of the form (4) for some a_i and b_i . If any of the intervals (a_i, b_i) overlap, we can strengthen some disjuncts (and eliminate redundant ones) and obtain an equivalent interpolant where the intervals (a_i, b_i) are disjoint and are listed so that for each i , $a_i < a_{i+1}$. On the other hand, $T_{add}[F] \cup \{2x - 1 < u\}$ refutes θ if and only if it refutes each θ_i , and if this is the case, it is certainly true for any θ'_i such that $T_{add}[F] \cup \{2x - 1 < u\}$ proves $\theta'_i \rightarrow \theta_i$. Thus, from the point of view of proving the “only if” direction of the theorem, we may assume, without loss of generality, that θ is a disjunction of formulas of the form (4), and the intervals (a_i, b_i) corresponding to the a and b in each θ_i are increasing and disjoint.

But then it is clear that $T_{mult}[F] \cup \{0 \leq x, x \leq r, u = x^2\}$ proves $\theta_1 \vee \dots \vee \theta_n$ if and only if

1. $a_0 = 0$,
2. $b_i = a_{i+1}$, for each $i < n$,
3. $a_n = r$,

and the \triangleleft 's are chosen suitably. Lemma 5.5 guarantees that for each i , $a_{i+1} \leq 1/(2 - a_i)$. The largest possible value of r occurs when the inequality is replaced by an equality $a_{i+1} = 1/(2 - a_i)$, and a calculation shows that in that case, $a_i = i/(i + 1)$ for each $i \leq n$.

This proves the “only if” direction of the theorem, establishing an upper bound on the possible values of r . But the proof in fact yields an interpolant that shows that the upper bound can be obtained: if each θ_i is the formula

$$a_i \leq x \leq a_{i+1} \wedge a_i^2 \leq u \leq a_{i+1}^2 \wedge a_i x \leq u \leq a_{i+1} x$$

with $a_i = i/(i + 1)$, then $T_{mult}[F] \cup \{0 \leq x, x \leq r, u = x^2\}$ proves $\theta_1 \vee \dots \vee \theta_n$, and $T_{add}[F] \cup \{2x - 1 < u\}$ refutes each θ_i . \square

6 Provability of a universal sentence in $T[F]$

In this section, we will provide various characterizations of provability of a universal sentence in $T[F]$. These will be used in Section 7 to establish our decidability results.

By Proposition 2.2, if φ is a universal sentence in the language of some $T[F]$, φ is equivalent to a formula of the form $\forall \vec{x} (\varphi_{add}(\vec{x}) \vee \varphi_{mult}(\vec{x}))$, where φ_{add} and φ_{mult} are in the language of $T_{add}[F]$ and $T_{mult}[F]$, respectively.

Proposition 6.1. *Let $\varphi \equiv \forall \vec{x} (\varphi_{add}(\vec{x}) \vee \varphi_{mult}(\vec{x}))$ be as above. Then the following are equivalent:*

1. $T[F]$ proves φ .
2. There is a quantifier-free formula $\theta(\vec{x})$ in the language $T_{comm}[F]$ such that $T_{add}[F] \cup \{\theta(\vec{x})\} \vdash \varphi_{add}(\vec{x})$ and $T_{mult}[F] \cup \{\neg\theta(\vec{x})\} \vdash \varphi_{mult}(\vec{x})$.
3. There is a quantifier-free formula $\theta(\vec{x})$ in the language $T_{comm}[F]$ such that

$$\forall \vec{x} (\theta(\vec{x}) \rightarrow \varphi_{add}(\vec{x})) \quad \text{and} \quad \forall \vec{x} (\neg\theta(\vec{x}) \rightarrow \varphi_{mult}(\vec{x}))$$

hold of the real numbers, with the intended interpretation of the auxiliary function symbols.

Proof. If 2 holds, then clearly $T[F]$ proves $\varphi_{add}(\vec{x}) \vee \varphi_{mult}(\vec{x})$. Thus 2 implies 1. Conversely, if $T[F]$ proves φ , it proves $\neg\varphi_{mult}(\vec{x}) \rightarrow \varphi_{add}(\vec{x})$. Treating \vec{x} as new constants and applying the Craig interpolation lemma, we get an interpolant $\theta(\vec{x})$ in the language of $T_{comm}[F]$ satisfying the conclusion of 2. Since $T_{comm}[F]$ has quantifier-elimination, we can assume without loss of generality that $\theta(\vec{x})$ is quantifier-free.

The equivalence of 2 and 3 follows easily from the fact that each of $T_{add}[F]$ and $T_{mult}[F]$ is a complete theory that holds of the reals numbers with the intended interpretation of the auxiliary function symbols. \square

From a model-theoretic perspective, it is useful to replace provability by nonexistence of a countermodel. When we say $\Gamma(\vec{x})$ is a *type* over a theory T , we mean that Γ is a set of formulas in the language of T , involving only the free variables \vec{x} , such that Γ is consistent with T . Saying $\Gamma(\vec{x})$ is a *complete type* means that for every formula $\psi(\vec{x})$, either $\psi(\vec{x})$ or $\neg\psi(\vec{x})$ is in $\Gamma(\vec{x})$.

Proposition 6.2. *Let $\varphi \equiv \forall \vec{x} (\varphi_{add}(\vec{x}) \vee \varphi_{mult}(\vec{x}))$ be as above. Then the following are equivalent:*

1. $T[F]$ does not prove φ .
2. $T[F] \cup \{\neg\varphi\}$ is consistent.

3. The union of $T_{add}[F] \cup \{\neg\varphi_{add}(\vec{x})\}$ and $T_{mult}[F] \cup \{\neg\varphi_{mult}(\vec{x})\}$ is consistent.

4. There is a complete type $\Gamma(\vec{x})$ over $T_{comm}[F]$ such that

$$T_{add}[F] \cup \Gamma(\vec{x}) \cup \{\neg\varphi_{add}(\vec{x})\} \quad \text{and} \quad T_{mult}[F] \cup \Gamma(\vec{x}) \cup \{\neg\varphi_{mult}(\vec{x})\}$$

are both consistent.

5. There is a complete type $\Gamma(\vec{x})$ over $T_{comm}[F]$ such that for every finite $\Gamma'(\vec{x}) \subseteq \Gamma(\vec{x})$,

$$T_{add}[F] \vdash \exists \vec{x} (\bigwedge \Gamma'(\vec{x}) \wedge \neg\varphi_{add}(\vec{x}))$$

and

$$T_{mult}[F] \vdash \exists \vec{x} (\bigwedge \Gamma'(\vec{x}) \wedge \neg\varphi_{mult}(\vec{x})).$$

6. There is a complete type $\Gamma(\vec{x})$ over $T_{comm}[F]$ such that for every finite $\Gamma'(\vec{x}) \subseteq \Gamma(\vec{x})$, there are real numbers \vec{x} and \vec{y} satisfying

$$\Gamma'(\vec{x}) \wedge \neg\varphi_{add}(\vec{x}) \wedge \Gamma'(\vec{y}) \wedge \neg\varphi_{mult}(\vec{y}).$$

Proof. In light of the soundness and completeness of first-order logic, 1 is just a restatement of 2, and the equivalence with 3 follows from the definition of φ in terms of φ_{add} and φ_{mult} . The equivalence of 3 with 4 follows by the Robinson joint consistency theorem, or, equivalently, from the Craig interpolation theorem, using compactness.

That statement 4 implies statement 5 follows from the fact that $T_{add}[F]$ and $T_{mult}[F]$ are both complete theories; for example, $T_{add}[F] \cup \Gamma'(\vec{x}) \cup \{\neg\varphi_{add}(\vec{x})\}$ is consistent if and only if $T_{add}[F]$ proves $\exists \vec{x} (\bigwedge \Gamma'(\vec{x}) \wedge \neg\varphi_{add}(\vec{x}))$. The converse is immediate.

The equivalence of 5 and 6 follows from the fact that each of $T_{add}[F]$ and $T_{mult}[F]$ is the theory of the real numbers in the respective languages. \square

Note that the equivalence of 1–4 holds, in general, for any two theories. The equivalence with 5 relies only on the fact that $T_{add}[F]$ and $T_{mult}[F]$ are complete, and the equivalence with 6 relies only on the additional fact that they are satisfied by the reals.

Statement 6 provides a nice characterization of provability in $T[F]$. A universal sentence φ is true of the reals if and only if every sequence \vec{x} of reals satisfies either $\varphi_{add}(\vec{x})$ or $\varphi_{mult}(\vec{x})$. But a universal sentence φ is provable in $T[F]$ if and only if for every complete type $\Gamma(\vec{x})$ in the language of $T_{comm}[F]$, there is a finite subset $\Gamma'(\vec{x})$ such that either

$$\forall \vec{x} (\bigwedge \Gamma'(\vec{x}) \rightarrow \varphi_{add}(\vec{x})) \quad \text{or} \quad \forall \vec{x} (\bigwedge \Gamma'(\vec{x}) \rightarrow \varphi_{mult}(\vec{x}))$$

holds in the reals. In particular, this has to hold whenever $\Gamma(\vec{x})$ is the type corresponding to a sequence of real numbers; but we will see below that there

are types in the language of $T_{comm}[F]$ that are not of this form. Thus, provability in $T[F]$ imposes a stronger requirement.

In the remainder of this section, we consider various representations of the quantifier-free formulas $\varphi_{add}(\vec{x})$, $\varphi_{mult}(\vec{x})$, and the possible interpolants $\theta(\vec{x})$. We also consider representations of the types $\Gamma(\vec{x})$. The former will be relevant to the discussion of heuristic algorithms in Sections 12–14, whereas the latter will be used in our decidability proofs in Section 7.

Let $\varphi \equiv \forall \vec{x} (\varphi_{add}(\vec{x}) \vee \varphi_{mult}(\vec{x}))$ be as above. Since $\forall y \psi(y)$ is equivalent to $\forall y > 0 \psi(y) \wedge \psi(0) \wedge \forall y > 0 \psi(-y)$, as in the proof of Proposition 3.3, any universal sentence φ is equivalent to a conjunction of formulas of the form $\forall \vec{x} > 0 (\varphi_{add}(\vec{x}) \vee \varphi_{mult}(\vec{x}))$. We can absorb the condition $\vec{x} > 0$ into both $\varphi_{add}(\vec{x})$ and $\varphi_{mult}(\vec{x})$. By adding a new variable if necessary, we can also assume that each includes the condition $x_1 = 1$, and it will be notationally convenient to do so. Thus, for the rest of this section, we will assume that φ is a universal formula of the form $\forall \vec{x} (\varphi_{add}(\vec{x}) \vee \varphi_{mult}(\vec{x}))$ where $\varphi_{add}(\vec{x})$ and $\varphi_{mult}(\vec{x})$ are quantifier-free in the language of $T_{add}[F]$ and $T_{mult}[F]$, respectively, and $\neg\varphi_{add}(\vec{x})$ and $\neg\varphi_{mult}(\vec{x})$ each implies $\vec{x} > 0$ and $x_1 = 1$. The question as to the decidability of the universal fragment of $T[F]$ reduces to the question as to whether one can determine whether $T[F]$ proves a sentence of this form. Let $\Delta(\vec{x})$ be the set $\{\vec{x} > 0, x_1 = 1\}$.

Proposition 6.3. *Under hypotheses $\Delta(\vec{x})$, a quantifier-free formula in the language of $T_{comm}[F]$ can be put in any of the following forms:*

1. *a conjunction of disjunctions of atomic formulas of the form $x_i < ax_j$ or $x_i \leq ax_j$, with $a > 0$.*
2. *a conjunction of disjunctions of atomic formulas of the form $x_i < ax_j$, with $a > 0$, or of the form $x_i = ax_j$ with $a > 0$ and $i < j$.*
3. *either 1 or 2, with “conjunction” and “disjunction” switched.*

Proof. Let θ be quantifier-free. First, put θ in negation-normal form, so that it is built up from atomic formulas and negations of atomic formulas using \wedge and \vee . Replace $s \not< t$ by $t \leq s$, replace $s \not\leq t$ by $t < s$, and replace $s \neq t$ by $s < t \vee t < s$. As a result, all atomic literals occur positively. One can further eliminate either $s \leq t$ in favor of $s < t \vee s = t$, or one can eliminate $s = t$ in favor of $s \leq t \wedge t \leq s$. The resulting formula can then be put in either disjunctive or conjunctive normal form, without introducing negations.

In the end, all the atomic formulas are of the form $ax_i < bx_j$, $ax_i \leq bx_j$, or $ax_i = bx_j$. Dividing through by b (and reversing an inequality when b is negative), we can assume that in each case $b = 1$. With the assumptions in Δ , each atomic formula in which a is negative can be replaced by either \top or \perp . Then inequalities $ax_i < x_j$ (resp. $ax_i \leq x_j$) can be expressed as $x_i < (1/a)x_j$ (resp. $x_i \leq (1/a)x_j$), as necessary, and equalities $x_j = ax_i$ can be rewritten $x_i = (1/a)x_j$ when $i < j$. \square

Such normal forms can be useful in reducing the problem of proof search to restricted cases. From an implementation point of view, not all these reductions

are wise, however; for example, using case splits to ensure that the x 's are all positive or to eliminate $s \leq t$ in favor of $s < t$ or $s = t$ can result in an exponential blowup. In the absence of sign information, the normal forms are more complicated. For example, although $x_2 > 2x_3$ can be expressed as $x_3 < (1/2)x_2$, $x_2 > -x_3$ cannot be expressed in the form $x_i < ax_j$. Also, in the absence of sign information, neither of $x_2 < x_3$ and $x_2 < 2x_3$ implies the other. In that case, one has to consider normal forms with atomic formulas from among $x_i < ax_j$, $x_i \leq ax_j$, $x_i > ax_j$, and $x_i \geq ax_j$. A little thought shows that in a single conjunction or disjunction, for each pair i, j , no more than two such formulas are needed; see also the proof of Proposition 12.2.

We can similarly classify the complete types over $T_{\text{comm}}[F]$. Let $\Gamma(\vec{x}) \supseteq \Delta(\vec{x})$ be such a type. Since $T_{\text{comm}}[F]$ has quantifier elimination, Γ is determined by the atomic formulas that it contains. Hence it is also determined by its subsets $\Gamma_{i,j}(x_i, x_j)$, with $i < j$, where $\Gamma_{i,j}$ consists of the atomic formulas involving both x_i and x_j . If $\Gamma_{i,j}$ contains a formula of the form $x_i = ax_j$, that determines the set $\Gamma_{i,j}$ uniquely. We denote this type by $\Gamma_{x_i/x_j=a}$. Otherwise, $\Gamma_{i,j}$ contains the formula $x_i \neq ax_j$ for every a in F , and so $\Gamma_{i,j}$ is determined by the set of elements a such that $\Gamma_{i,j}$ contains the formula $x_i < ax_j$. This set is a downwards-closed subset of the positive part of F ; think of it as the set of a such that $x_i/x_j < a$. If this set is empty, that determines $\Gamma_{i,j}$ uniquely, and we denote the corresponding type $\Gamma_{x_i/x_j \approx \infty}$. Otherwise, the set has a greatest lower bound in the real numbers, say, r . If r is not an element of F , then $\Gamma_{i,j}$ contains $x_i < ax_j$ exactly when $r < a$, and this determines $\Gamma_{i,j}$ exactly; we denote the resulting type by $\Gamma_{x_i/x_j \approx r}$. If, on the other hand, r is an element a of F , there are two possibilities: $\Gamma_{i,j}$ contains the formula $x_i < ax_j$, or it does not (in which case it contains the formula $x_j < (1/a)x_i$). Denote the first type by $\Gamma_{x_i/x_j \approx a^-}$, and denote the second by $\Gamma_{x_i/x_j \approx a^+}$.

In short, we have shown the following:

Proposition 6.4. *Let $\Gamma(\vec{x})$ be any complete type over $T_{\text{comm}}[F]$ that includes $\Delta(\vec{x})$. Then for each $i < j$, Γ includes exactly one of the following:*

1. $\Gamma_{x_i/x_j=a}$, for some a in F
2. $\Gamma_{x_i/x_j \approx r}$, for some r in $\mathbb{R} \setminus F$
3. $\Gamma_{x_i/x_j \approx \infty}$
4. $\Gamma_{x_i/x_j \approx a^-}$, for some a in F
5. $\Gamma_{x_i/x_j \approx a^+}$, for some a in F

These data determine Γ uniquely.

Note that not every collection of sets Γ_{x_i/x_j} determines a consistent type over $T_{\text{comm}}[F]$; for example, the sets $\Gamma_{x_1/x_2=2}$, $\Gamma_{x_2/x_3=2}$, and $\Gamma_{x_1/x_3=2}$ are jointly inconsistent.

In the next section, we will combine the analysis given by Proposition 6.4, together with equivalence 6 of Proposition 6.2, to show that, with general conditions on F , the universal fragment of $T[F]$ is decidable.

7 Decidability

Let $\varphi \equiv \forall \vec{x} (\varphi_{add}(\vec{x}) \vee \varphi_{mult}(\vec{x}))$ be as in the previous section, so that φ_{add} and φ_{mult} are quantifier-free formulas in the language of $T_{add}[F]$ and $T_{mult}[F]$ respectively, and each of $\neg\varphi_{add}(\vec{x})$ and $\neg\varphi_{mult}(\vec{x})$ implies $\vec{x} > 0$ and $x_1 = 1$. We have seen that the decidability of the universal fragment of $T[F]$ reduces to the problem of determining whether $T[F]$ proves a formula φ of this sort; and that $T[F]$ does *not* prove such a φ if and only if

there is a complete type $\Gamma(\vec{x})$ over $T_{comm}[F]$ such that for every finite $\Gamma'(\vec{x}) \subseteq \Gamma(\vec{x})$, the sentence

$$\exists \vec{x} \left(\bigwedge \Gamma'(\vec{x}) \wedge \neg\varphi_{add}(\vec{x}) \right) \wedge \exists \vec{x} \left(\bigwedge \Gamma'(\vec{x}) \wedge \neg\varphi_{mult}(\vec{x}) \right)$$

is true of the real numbers.

Call this the “consistency criterion for $\neg\varphi$.” We also have a complete classification of the relevant types $\Gamma(\vec{x})$. In this section, we will use the latter to show that when F is a computable subfield of \mathbb{R} and membership of a real algebraic number in F is decidable, the consistency criterion for $\neg\varphi$ is decidable.

Fix φ and F , and hence $\varphi_{add}(\vec{x})$ and $\varphi_{mult}(\vec{x})$. If $\Gamma(\vec{x})$ is any set of atomic formulas in the language of $T_{comm}[F]$ involving the variables \vec{x} and $i < j$, let $\Gamma_{i,j}$ denote the set of formulas in Γ involving x_i and x_j . Let \mathcal{S} be the collection of sets Γ such that for each $i < j$, $\Gamma_{i,j}$ is one of the types described in Proposition 6.4. Since each such Γ consistent with $T_{comm}[F]$ uniquely determines the complete type that extends it, we can replace “complete type $\Gamma(\vec{x})$ over $T_{comm}[F]$ ” by “ $\Gamma \in \mathcal{S}$ ” in the consistency criterion for $\neg\varphi$.

We now show that we can modify the collection of sets \mathcal{S} to avoid the restrictions “ $a \in F$ ” in the clauses of Proposition 6.4. To do so, we consider types in the larger language, $T_{comm}[\mathbb{R}]$. Let the types $\hat{\Gamma}_{x_i/x_j=a}$, $\hat{\Gamma}_{x_i/x_j \approx r}$, $\hat{\Gamma}_{x_i/x_j \approx \infty}$, $\hat{\Gamma}_{x_i/x_j \approx a^-}$, and $\hat{\Gamma}_{x_i/x_j \approx a^+}$ be defined as in the paragraph before Proposition 6.4, except with respect to the language of $T_{comm}[\mathbb{R}]$. Let $\hat{\mathcal{S}}$ be the sets $\hat{\Gamma}$ of atomic formulas in $T_{comm}[\mathbb{R}]$ such that for each $i < j$, $\hat{\Gamma}_{i,j}$ is one of the following:

1. $\hat{\Gamma}_{x_i/x_j=a}$, for some a in \mathbb{R}
2. $\hat{\Gamma}_{x_i/x_j \approx r}$, for some r in $\mathbb{R} \setminus F$
3. $\hat{\Gamma}_{x_i/x_j \approx \infty}$
4. $\hat{\Gamma}_{x_i/x_j \approx a^-}$, for some a in \mathbb{R}
5. $\hat{\Gamma}_{x_i/x_j \approx a^+}$, for some a in \mathbb{R}

Note that we have replaced “ $a \in F$ ” by “ $a \in \mathbb{R}$ ” in the first item and in the last two items, but we have left $\mathbb{R} \setminus F$ alone in the second item.

Lemma 7.1. *The consistency criterion for $\neg\varphi$ is satisfied by a set $\Gamma \in \mathcal{S}$ if and only if it is satisfied by a set $\hat{\Gamma} \in \hat{\mathcal{S}}$.*

Proof. Suppose the consistency criterion is satisfied by some $\Gamma \in \mathcal{S}$. It is easy to check that it is then satisfied by the corresponding set $\hat{\Gamma} \in \hat{\mathcal{S}}$.

In the other direction, note that if a is in $\mathbb{R} \setminus F$, then each of $\hat{\Gamma}_{x_i/x_j=a}$, $\hat{\Gamma}_{x_i/x_j \approx a^-}$, and $\hat{\Gamma}_{x_i/x_j \approx a^+}$ includes $\Gamma_{x_i/x_j \approx a}$. Thus every set $\hat{\Gamma} \in \hat{\mathcal{S}}$ includes a set $\Gamma \in \mathcal{S}$. So, if the consistency criterion is satisfied by some $\hat{\Gamma} \in \hat{\mathcal{S}}$, it is satisfied by some $\Gamma \in \mathcal{S}$. \square

We now parameterize the finite subsets of each $\hat{\Gamma} \in \hat{\mathcal{S}}$. For each $\varepsilon > 0$, we define a formula

$$\hat{\Gamma}[\varepsilon] = \bigwedge_{i < j} \hat{\Gamma}_{i,j}[\varepsilon],$$

where

1. $\hat{\Gamma}_{x_i/x_j=a}[\varepsilon]$ is the formula $x_i = ax_j$
2. $\hat{\Gamma}_{x_i/x_j \approx r}[\varepsilon]$ is $(r - \varepsilon)x_j < x_i < (r + \varepsilon)x_j$
3. $\hat{\Gamma}_{x_i/x_j \approx \infty}$ is $x_i > (1/\varepsilon)x_j$
4. $\hat{\Gamma}_{x_i/x_j \approx a^-}$ is $(a - \varepsilon)x_j < x_i < ax_j$
5. $\hat{\Gamma}_{x_i/x_j \approx a^+}$ is $ax_j < x_i < (a + \varepsilon)x_j$

For every ε , $\hat{\Gamma}[\varepsilon]$ is implied by some finite subset of $\hat{\Gamma}$. Conversely, every finite subset of $\hat{\Gamma}$ is implied by $\hat{\Gamma}[\varepsilon]$ for some $\varepsilon > 0$, and, in fact, for an ε of the form $1/n$ for some $n \in \mathbb{N}$. Thus the consistency criterion for $\neg\varphi$ is equivalent to the following:

there is a set $\hat{\Gamma} \in \hat{\mathcal{S}}$ such that for every $\varepsilon > 0$, the sentence

$$\exists \vec{x} (\hat{\Gamma}[\varepsilon] \wedge \neg\varphi_{add}(\vec{x})) \wedge \exists \vec{x} (\hat{\Gamma}[\varepsilon] \wedge \neg\varphi_{mult}(\vec{x}))$$

is true of the real numbers.

The sets $\hat{\Gamma} \in \hat{\mathcal{S}}$, and the corresponding formulas $\hat{\Gamma}[\varepsilon]$, are parameterized by tuples of symbols from the set

$$\{ '=a' \mid a \in \mathbb{R} \} \cup \{ '\approx r' \mid r \in \mathbb{R} \setminus F \} \cup \{ '\infty' \} \cup \{ '\approx a^-', \approx a^+' \mid a \in \mathbb{R} \}.$$

When $F = \mathbb{R}$, there are no sets with parameters of the second kind, and so the consistency criterion can be expressed in the language of real closed fields. By Theorem 4.4, $T[\mathbb{R}]$ is a conservative extension of $T[\mathbb{A}]$. Thus we have:

Theorem 7.2. *The universal fragment of $T[\mathbb{A}]$ is decidable.*

When F is a proper subfield of \mathbb{R} , the revised consistency criterion for $\neg\varphi$ can be expressed as a sentence of the form

$$\exists \vec{r} \in \mathbb{R} \setminus F \exists \vec{a} \in \mathbb{R} \forall \varepsilon > 0 \exists \vec{x}, \vec{x}' \theta$$

where θ is a quantifier-free formula in the language of real closed fields. By quantifier-elimination for real closed fields, this is equivalent to a sentence of the form $\exists \vec{r} \in \mathbb{R} \setminus F \ \eta$, where η is a quantifier-free formula in the language of real closed fields. Say F is a *sufficiently computable* subfield of \mathbb{R} if F is a computable subfield of \mathbb{R} and there is an algorithm to determine whether a real algebraic number a (described in terms of a definition, say, in the language of real closed fields) is in F .

Theorem 7.3. *For any sufficiently computable $F \subseteq \mathbb{R}$, the universal fragment of $T[F]$ is decidable.*

By our analysis of the consistency criterion, Theorem 7.3 is a consequence of the following:

Theorem 7.4. *For any sufficiently computable $F \subseteq \mathbb{R}$, there is an algorithm to decide whether a sentence of the form $\exists \vec{x} \in \mathbb{R} \setminus F \ \varphi(\vec{x})$ holds of the reals, where φ is a formula in the language of real closed fields.*

We will prove something more general. Let R be any real closed field. A function $h(\vec{x})$ or a predicate $E(\vec{x})$ on R is said to be *semialgebraic* if it is definable in the language of real-closed fields without parameters.

Theorem 7.5. *Let R be any real closed field, and let F be any proper subfield of R . If E, h_1, \dots, h_m are semialgebraic, then*

$$\exists x_1 \notin F \ \dots \exists x_n \notin F \ (E(\vec{x}, \vec{y}) \wedge h_1(\vec{x}, \vec{y}) \notin F \wedge \dots \wedge h_m(\vec{x}, \vec{y}) \notin F)$$

is equivalent to a positive boolean combination of assertions of the form $D(\vec{y})$ and $g(\vec{y}) \notin F$, where D and g are semialgebraic. Furthermore, there is an algorithm for determining an expression of this form from (presentations of) E, h_1, \dots, h_m . This algorithm does not depend on R or F .

In particular, when there are no variables \vec{y} , Theorem 7.5 asserts that any assertion of the form $\exists \vec{x} \in \mathbb{R} \setminus F \ E(\vec{x})$ is effectively equivalent to a boolean combination of sentences in the language of real-closed fields and assertions of the form $g \notin F$, where g is a real algebraic constant. Thus Theorem 7.5 implies Theorem 7.4.

Proof of Theorem 7.5. We use induction on n . When $n = 0$ there is nothing to do. Suppose the theorem is true for n . Then

$$\exists x_1 \notin F \ \dots \exists x_{n+1} \notin F \ (E(\vec{x}, \vec{y}) \wedge h_1(\vec{y}) \notin F \wedge \dots \wedge h_m(\vec{y}) \notin F)$$

is equivalent to $\exists x_1 \notin F \ \psi(x_1, \vec{y})$, where ψ has the requisite form. We can then write ψ as a disjunction of formulas of the form

$$D(x_1, \vec{y}) \wedge g_1(x_1, \vec{y}) \notin F \wedge \dots \wedge g_l(x_1, \vec{y}) \notin F$$

where D, g_1, \dots, g_l are semialgebraic. Since we can factor the existential quantifier $\exists x_1$ across the disjunction, it suffices to prove Theorem 7.5 for the special case $n = 1$.

So, resorting to the original notation, let $E(x, \vec{y}), h_1(x, \vec{y}), \dots, h_m(x, \vec{y})$ be semialgebraic. We need to show that

$$\exists x \in \mathbb{R} \setminus F (E(x, \vec{y}) \wedge h_1(x, \vec{y}) \notin F \wedge \dots \wedge h_m(x, \vec{y}) \notin F) \quad (5)$$

is equivalent to a positive boolean combination of assertions $D(\vec{y})$ and $g(\vec{y}) \notin F$, for semialgebraic D and g .

By the theory of definability in real closed fields [6, 26], for each fixed \vec{y} , the set $\{x \mid E(x, \vec{y})\}$ is a finite union of disjoint intervals (including intervals of the form $(-\infty, a)$, $(-\infty, a]$, (a, ∞) , and $[a, \infty)$) with endpoints that are definable in the parameters \vec{y} . Similarly, fixing \vec{y} , for all but finitely many points x of \mathbb{R} all the functions h_i are either locally increasing or locally decreasing or locally constant at x . A bound p on the number of such intervals and exceptional points, independent of \vec{y} , can be determined effectively from the presentations of E, h_1, \dots, h_m . Furthermore, for fixed n , terms like “the left endpoint of the n th interval (in increasing order) in the decomposition of $\{x \mid E(x, \vec{y})\}$, if there is one, or 0 otherwise” and “the n th point at which one of the h_i ’s is neither locally monotone nor locally constant, if there is one, or 0 otherwise” are semialgebraic functions of \vec{y} .

As a result, for each fixed \vec{y} , there is a sequence of at most p disjoint nonempty open intervals J_1, \dots, J_q and at most p exceptional points u_1, \dots, u_r such that

- $\{x \mid E(x, \vec{y})\} = J_1 \cup \dots \cup J_q \cup \{u_1, \dots, u_r\}$, and
- on each interval J_n , all the functions h_i are either monotone or constant.

Furthermore, all the following are semialgebraic in \vec{y} :

- the predicates $D_{q,r}(\vec{y})$, where $q, r \leq p$, which assert that there are exactly q intervals in the decomposition of $\{x \mid E(x, \vec{y})\}$ and r exceptional points;
- the predicate $G_{i,n}(\vec{y})$ which asserts that h_i (as a function of x), is constant on J_n ; and
- the functions $k_{i,n}(\vec{y})$ which return the value of h_i on J_n , if h_i is constant on J_n , or 0 otherwise.

Given \vec{y} , assuming that there are q intervals J_n and r exceptional points, we claim that (5) is equivalent to the following disjunction:

1. there is an interval J_n , $n = 1, \dots, q$, such that for each function h_i , if h_i is constant on J_n , then the value of h_i on J_n is not in F ; or
2. for one of the exceptional points u_n , $n = 1, \dots, r$, we have $u_n \notin F$, and $h_i(u_n) \notin F$ for each i .

By the preceding paragraph, this can be expressed as a positive boolean combination $\psi_{q,r}(\vec{y})$ of assertions of the form $H(\vec{y})$ and $l(\vec{y}) \notin F$, where H and l are semialgebraic. This means that the expression

$$\bigvee_{q,r \leq p} (D_{q,r}(\vec{y}) \wedge \psi_{q,r}(\vec{y}))$$

is of the requisite form. Thus, to complete the proof of Theorem 7.5, it suffices to establish the equivalence of (5) with the disjunction of 1 and 2.

Suppose (5) holds, and, given \vec{y} , let $x \notin F$ witness the existential quantifier. Since $E(x, \vec{y})$ holds, either x is in J_n for some n , in which case clause 1 holds, or x is one of the exceptional points u_n , in which case clause 2 holds.

Conversely, given \vec{y} , suppose either 1 or 2 holds. If 2 holds, then that exceptional value u_n witnesses the existential quantifier in (5). So assume 1 holds, and let J be an interval on which all the functions that are constant take a value not in F . Renumbering, let h_1, \dots, h_l be functions that are not constant on J . It suffices to show that there is an $x \in J \setminus F$ such that $h_1(x, \vec{y}), \dots, h_l(x, \vec{y})$ are not in F .

We consider two cases. First, suppose R properly contains the real algebraic closure of $F(\vec{y})$ in R . Then one can choose an x transcendental over $F(\vec{y})$ in the interval J . This x has the desired property: if $h_i(x, \vec{y}) = a$ for some $i = 1 \dots l$, then $h_i(x, \vec{y}) - a = 0$ is a nontrivial algebraic identity in \vec{y} and elements of F , contradiction. Otherwise, R is equal to the real algebraic closure of $F(\vec{y})$ in R . Since F is properly contained in R , we can choose an x with sufficiently high algebraic degree over $F(\vec{y})$, in which case an equality $h_i(x, \vec{y}) = a$ for some $i = 1 \dots l$ again yields a contradiction. \square

Note that in the instance of Theorem 7.5 needed for Theorem 7.4, $R = \mathbb{R}$ and F is a countable subfield, in which case the implication from 1 to (5) in the last paragraph of the preceding proof follows more easily from cardinality considerations.

8 Normal forms

When dealing with an associative and commutative operation like addition, it is common to put terms in an appropriate normal form. For example, one can always rearrange a sum $t_1 + \dots + t_n$ so that parentheses are associated, say, to the left, and t_1, \dots, t_n are ordered according to a fixed ordering of terms; this makes it easy to tell whether or not two such sums agree up to the associativity and commutativity of addition. In the theories $T[F]$, not only do we have addition and multiplication (as well as subtraction and division), but also multiplication by constants from F . In this section, we will show that one can still, fruitfully, put terms in $T[F]$ into a normal form. This provides an algorithm for testing whether two terms are provably equal: just put them in normal form, and compare.

In fact, to show that normal forms are unique, we will take care to define an ordering on these terms that is compatible with the axioms for $<$ in $T[F]$. This will enable us to construct a term model of $T[F]$ in which different terms in normal form denote different elements. It will also enable us to show that any equality between terms that can be established in $T[F]$ can be proved without using the ordering.

We define a set of *preterms* inductively, each with an associated rank, as follows. For each n , a preterm of rank $2n + 1$ is called an “additive preterm,” and a preterm of rank $2n + 2$ is called a “multiplicative preterm.” A preterm of rank 0 is called a “basic preterm.”

- Each variable, x, y, z, \dots is a preterm of rank 0, as well as the constant, 1.
- For n greater than 0 and odd, if t_1, \dots, t_k are multiplicative or basic preterms of rank at most $n - 1$, $k \geq 2$, a_1, \dots, a_k are nonzero elements of F , and at least one t_i has rank $n - 1$, then $a_1 t_1 + a_2 t_2 + \dots + a_k t_k$ is a preterm of rank n .
- For n greater than 0 and even, if t_1, \dots, t_k are additive or basic preterms of rank at most $n - 1$ and other than 1, i_1, \dots, i_k are nonzero integers, either $k \geq 2$ or $i_1 \neq 1$, and at least one t_i has rank at least $n - 2$, then $t_1^{i_1} t_2^{i_2} \dots t_k^{i_k}$ is a preterm of rank n .

Here parentheses in products and sums are assumed to associate to the left, and for an integer i , t^i is the i -fold product of t with itself if i is positive, or 1 divided by the $-i$ -fold product of t with itself if i is negative. Note that there is no constant multiplier for multiplicative preterms. The condition “ $k \geq 2$ or $i_1 \neq 1$ ” in the third clause allows x^2 , for example, but rules out x^1 .

We now define, simultaneously, a normal form for preterms together with an ordering $s \prec t$ on preterms in normal form. We assume that variables have been indexed x_1, x_2, \dots . For each n , we define the notion of normal form, as well as the ordering, for terms of rank at most n , as follows:

1. $n = 0$: Each basic preterm is in normal form. These are ordered $1 \succ x_1 \succ x_2 \succ \dots$
2. $n > 0$, odd: An additive preterm $a_1 t_1 + a_2 t_2 + \dots + a_k t_k$ is in normal form if and only if each t_i is in normal form, $t_1 \succ t_2 \succ \dots \succ t_k$, and $a_1 = 1$.

To define $s \prec t$ when at least one of s and t has rank n and the other has rank at most n , write

$$s = a_1 u_1 + a_2 u_2 + \dots + a_k u_k$$

and

$$t = b_1 u_1 + b_2 u_2 + \dots + b_k u_k$$

where $u_1 \succ u_2 \succ \dots \succ u_k$ are preterms of rank at most $n - 1$, and now the a_i 's and b_i 's are allowed to be 0. Then use lexicographic order: $s \prec t$ if and only if $a_i \neq b_i$ for some i and $a_i < b_i$ for the least such i .

3. $n > 0$, even: A multiplicative preterm $t_1^{i_1} t_2^{i_2} \dots t_k^{i_k}$ is in normal form if and only if each t_m is in normal form, and $t_1 \succ t_2 \succ \dots \succ t_k$. To compare two multiplicative preterms of rank at most n , the procedure is slightly more complicated now, since we now consider the standing of the subterms in relation to the basic preterm 1. Write the subterms s_i occurring in s and

the subterms t_j occurring t , together with the preterm 1, in \succ -decreasing order as $u_1, \dots, u_m, 1, u_{m+1}, \dots, u_k$. Then express

$$s = u_1^{i_1} u_2^{i_2} \cdots u_m^{i_m} \cdot 1 \cdot u_{m+1}^{i_{m+1}} \cdots u_k^{i_k} \quad (6)$$

and

$$t = u_1^{j_1} u_2^{j_2} \cdots u_m^{j_m} \cdot 1 \cdot u_{m+1}^{j_{m+1}} \cdots u_k^{j_k} \quad (7)$$

where now the i_n 's and j_n 's may be 0. We now say $s \prec t$ if and only if

- there is an $n \leq m$ such that $i_n \neq j_n$, and, for the least such n , $i_n < j_n$; or
- For every $n \leq m$, $i_n = j_n$, but there is some $n > m$ such that $i_n \neq j_n$, and $i_n > j_n$ for the *largest* such n .

Note that the clause 1 of the definition of \succ makes sense if we think of the variables as being positive values, with each x_{i+1} infinitesimally small compared to x_i and 1. Clause 2, which treats the case where the term of highest rank is additive, is also intuitively consistent with an interpretation of \prec as denoting a relation, “is infinitely smaller than,” on positive numbers. Clause 3, which treats the case where the term of highest rank is multiplicative, has similarly been designed to admit such an interpretation. The main constraint there was to ensure that the ordering cohere, in the following sense:

Lemma 8.1. *Let $n > 0$ be even, and let s and t be preterms of rank less than or equal to n . Then the ordering of s and t is equivalent to the order obtained under clause 3, when s and t are put in the form (6) and (7), respectively.*

Lemma 8.1 is needed to prove Lemma 8.6. The proof proceeds by running through the cases where each of s and t is a variable, the constant 1, an additive term, or a multiplicative term. For example, if s and t are additive and $1 \succ s \succ t$, one easily verifies that $1s^1t^0 \succ 1s^0t^1$ under Clause 3. The other cases are similarly straightforward.

Say that a term is in *normal form* if it is either 0 or of the form at , where t is a preterm in normal form and a is a nonzero element of F . Let $T'_{add}[F]$ be the restriction of $T_{add}[F]$ to the language without the ordering $<$. Let $T'_{mult}[F]$ be corresponding restriction of $T_{mult}[F]$. Let $T'[F] = T'_{add}[F] \cup T'_{mult}[F]$. It is straightforward to verify the following:

Theorem 8.2. *For every term t , there is a term \hat{t} in normal form, such that $T'[F]$ proves $t = \hat{t}$.*

Our main goal, in this section, is to prove the following:

Theorem 8.3. *If \hat{s} and \hat{t} are terms in normal form, and $T[F]$ proves $\hat{s} = \hat{t}$, then $\hat{s} = \hat{t}$.*

Note that the last equality is *syntactic* equality; in other words, T proves that two terms in normal form are equal if and only if they are the same term.

As corollaries, we obtain the following:

Corollary 8.4. *There is an efficient procedure for determining whether $T[F]$ proves $s = t$.*

Proof. Just put s and t in normal form, and compare. □

Corollary 8.5. *$T[F]$ and $T'[F]$ have the same provable equalities.*

Proof. If $T[F]$ proves $s = t$, then s and t have the same normal form u . Since $T'[F]$ proves $s = u$ and $t = u$, it proves $s = t$. □

To prove Theorem 8.3, first let us extend the ordering \prec from preterms in normal form to terms in normal form, as follows: if s and t are preterms in normal form, then

- $0 \prec at$ if and only if $a > 0$.
- $at \prec 0$ if and only if $a < 0$.
- $0 \not\prec 0$
- $as \prec bt$ if and only if:
 - a is negative, and b is positive
 - a and b are both positive, and either $s \prec t$ or $s = t$ and $a < b$
 - a and b are both negative, and either $s \succ t$ or $s = t$ and $a < b$

It suffices to show

Lemma 8.6. *There is a model \mathcal{M} of $T[F]$ such that if \hat{s} and \hat{t} are terms in normal form and $\hat{s} \prec \hat{t}$, then $\hat{s} < \hat{t}$ holds in \mathcal{M} .*

Proof. Note that operations of addition, subtraction, multiplication, and division are naturally defined on terms in normal form. For example, suppose $a(a_1s_1 + a_2s_2 + \dots + a_ks_k)$ and $b(b_1t_1 + b_2t_2 + \dots + b_lt_l)$. To express their sum as a term in normal form, multiply through by a and b , respectively, combine terms, and express the sum as $c_1u_1 + c_2u_2 + \dots + c_mu_m$, where $u_1 \succ u_2 \succ \dots \succ u_m$ and each $c_i \neq 0$, or 0. In the former case, the desired normal-form term is $c_1(u_1 + (c_2/c_1)u_2 + \dots + (c_m/c_1)u_m)$. This term model *almost* satisfies the claim of Lemma 8.6; it satisfies all the axioms of $T[F]$ indicated in Section 4, except for the axiom that asserts that the multiplicative group of positive elements is divisible. That is, all that is missing are n th roots of positive elements. To remedy the situation, we embed this term model in an expanded set of formal terms, defined as follows.

Let F' be the smallest subfield of \mathbb{R} that includes F and is closed under n th roots of positive elements, for positive n . Define the set of *extended preterms* inductively, as above, with the following changes:

- in the additive extended preterms $a_1t_1 + \dots + a_kt_k$, the coefficients a_i are taken from F' ; and

- multiplicative extended preterms are taken to be formal products $t_1^{i_1} t_2^{i_2} \dots t_k^{i_k}$ where now the exponents i_j are *rational numbers*.

Define the set of extended preterms in normal form, the ordering on these, the set of extended terms in normal form, and the ordering on these, exactly as before. Once again, operations of addition and multiplication can be defined on extended terms in normal form. Lemma 8.1, as well as the analogue for additive preterms, carry over to extended preterms as well.

Let \mathcal{M} be the model whose universe is the set of extended terms in normal form, with the associated ordering and operations of addition and multiplication. Clearly there is an embedding of the set of terms in normal form into the set of extended terms in normal form which preserves all the operations. So it suffices to show that \mathcal{M} satisfies $T[F]$.

We simply run through the axioms given in Section 4. Verifying the axioms of $T_{comm}[F]$ is straightforward, as well as the fact that the terms form an abelian group under addition, and the positive terms form an abelian group under multiplication.

To show that the ordering is compatible with multiplication of positive elements, we need to show that $s \prec t \rightarrow su \prec tu$ holds of positive terms s, t, u in normal form. Let $s = as'$, $t = bt'$, and $u = cu'$ where s' , t' , and u' are preterms in normal form, and a, b , and c are positive. Then $su = (ac)s'u'$ and $tu = (bc)t'u'$. Since $s \prec t$, we have either $s' \prec t'$, or $s' = t'$ and $a < b$. In the first case, Lemma 8.1 and Clause 3 of the definition of \prec guarantees that $s'u' \prec t'u'$, and hence $su \prec tu$. In the second case, $s'u' = t'u'$ and $ac < bc$, so, again, $su \prec tu$.

Showing that the ordering is compatible with addition is similarly straightforward. So we only need to show that the multiplicative group of positive elements is divisible. Let at be an extended term in normal form satisfying $at \succ 0$. Then $a > 0$, and we can view t as a multiplicative preterm $t_1^{i_1} t_2^{i_2} \dots t_k^{i_k}$, possibly with $k = 1$ and $i_1 = 1$. But this has n th root $\sqrt[n]{at_1^{i_1/n} t_2^{i_2/n} \dots t_k^{i_k/n}}$, where this is identified with $\sqrt[n]{at_1}$ if $k = 1$ and $i_1/n = 1$. \square

We note that the complicated definition of \prec in the multiplicative clause of the ordering of preterms was designed to ensure that \prec is compatible with the axioms of $T[F]$. This, in turn, was used to construct the term model in the proof of Theorem 8.3. Theorem 8.3 remains true, however, for a simpler version of \prec , in which we simply use a lexicographic ordering at the multiplicative stage. This simpler ordering, and the associated normal forms, are more amenable to implementation. (Indeed, it may also be natural to order terms of lower rank before terms of higher rank.) To derive the variant of Theorem 8.3 for these normal forms, it suffices to show that the map from terms in the simpler normal form to the normal form we have used here is injective. In other words, it suffices to show that if s and t are in the simpler normal form, u a term in the normal form we have used here, and $T[F]$ proves both $s = u$ and $t = u$, then s and t are syntactically identical. This can be done by a careful induction on the maximum rank of s and t .

Note also that it is harmless, and again useful from an implementation point of view, to extend the language of $T[F]$ to include exponentiation to arbitrary integers. Since n th roots of positive elements can be defined in $T[F]$, one can similarly expand the language of $T[F]$ to allow n th root functions for positive n , or even exponentiation to any rational power. One has to be careful, however, to provide a consistent interpretation of the n th root function on negative elements, and natural simplifications may depend on knowing the sign of the relevant terms. For example, $\sqrt{x^2}$ can be simplified to x if x is positive and $-x$ if x is negative. For that reason, determining an appropriate normal form representation for terms involving n th roots is more complicated. Similar complications arise in obtaining an adequate handling of absolute value, max, and min. The issue of obtaining useful canonical representations for such extensions is of practical importance, and is discussed further in Section 14 below.

Finally, we note that the method of computing normal forms only gives a decision procedure for provable equations in the absence of hypotheses. For example, $T[F]$ proves $1+x^2+y^2 \neq 0$ (or, equivalently, $1+x^2+y^2 = 0 \rightarrow 0 = 1$), but this is not provable in $T'[F]$.

9 Building models of $T[F]$

In Sections 10 and 11, our goal will be to prove undecidability results (and conditional undecidability results) for the theories $T[F]$. Recall the alternative formulations $T[F]^*$ introduced in Section 4, in the language with symbols $0, 1, +, \times, <$ and constants c_a for each $a \in F$. In light of Theorem 4.3, we will work exclusively with the theories $T[F]^*$. Our strategy will be to build models of $T[F]^*$ in which F and \mathbb{Z} are, respectively, definable. In this section, we will develop techniques for building such models.

Let $\mathcal{R} = \langle R, <, +, -, \times \rangle$ be an ordered real closed field extending the countable ordered subfield $F \subseteq \mathbb{R}$. More specifically, we assume that F is a subfield of \mathcal{R} , where the ordering on F agrees with the ordering in \mathcal{R} .

Definition 9.1. We say that h is an F -bijection of \mathcal{R} if and only if

1. $h : R \rightarrow R$ is an order preserving bijection.
2. $h(0) = 0$ and $h(1) = 1$.
3. For all $x \in R$ and $a \in F$, we have $h(ax) = ah(x)$.

Given an F -bijection h , we define the structure $h^{-1}[\mathcal{R}]$ in the language of $T[F]^*$ as follows. The domain of $h^{-1}[\mathcal{R}]$ is R . The symbols $0, 1, +$, and $<$ are interpreted as in \mathcal{R} . For $a \in F$, c_a is interpreted as a . The symbol \times is interpreted in $h^{-1}[\mathcal{R}]$ as \otimes , defined by the equation

$$x \otimes y = h^{-1}(h(x)h(y)).$$

It follows from the definition that $x \otimes y = z$ if and only if $h(x)h(y) = h(z)$. Hence h is an isomorphism from $\langle R, \otimes, < \rangle$ onto $\langle R, \times, < \rangle$.

Theorem 9.2. *Let h be an F -bijection of \mathcal{R} . The model $h^{-1}[\mathcal{R}]$ satisfies $T[F]^*$.*

Proof. Recall the axiomatization of $T[F]^*$ given in Section 4. We first verify axioms 1,2 in $h^{-1}[\mathcal{R}]$. The group given by $0, +, <$ is obviously an ordered commutative group. Since h is an isomorphism from $\langle R, \otimes, < \rangle$ onto $\langle R, \times, < \rangle$, we have that $1, \times, <$ is a divisible ordered commutative group on the positive elements of R .

Axioms 3a-3c obviously hold in $h^{-1}[\mathcal{R}]$. For axioms 4a,4b, note that for all $a \in F$,

$$a \otimes x = h^{-1}(h(a)h(x)) = h^{-1}(ah(x)) = ah^{-1}(h(x)) = ax.$$

Hence

$$(a + b) \otimes x = (a + b)x = ax + bx = a \otimes x + b \otimes x$$

and

$$a \otimes (x + y) = a(x + y) = ax + ay = (a \otimes x) + (a \otimes y).$$

□

So far, we have only assumed that \mathcal{R} is an ordered real closed field extending the countable ordered subfield $F \subseteq \mathbb{R}$. We will now need to assume that \mathcal{R} obeys some additional conditions. Note that R is a densely ordered set. An *interval* in R is a $J \subseteq R$ such that for all $x < y < z$, $x, z \in J$, $y \in R$, we have $y \in J$. J is said to be *nontrivial* if and only if J has infinitely many elements. This is the same as saying that J has at least two elements.

By a standard saturation argument, we will fix an ordered real closed field \mathcal{R} , such that the following hold:

1. R is countable.
2. \mathcal{R} extends F in the sense above.
3. Let $n \geq 1$. Suppose that for all $i \geq 1$, $g_i, h_i : R^n \rightarrow R$ are \mathcal{R} -definable, where n may depend on i . Then $\cup_i g_i[F^n]$ has an upper bound. Furthermore, suppose each $g_i[F^n]$ lies strictly below each $h_j[F^n]$. Then the interval strictly above each $g_i[F^n]$ and strictly below each $h_j[F^n]$ is nontrivial.

Here, as always, \mathcal{R} -definability allows the use of parameters from R , and the notation $f[S]$ denotes the forward image of f on S . The existence of such a field can be proved by starting with a countable ordered real closed subfield R_0 of \mathbb{R} containing F , and then building a countably infinite chain of elementary extensions. At each stage, use compactness to ensure that the required upper bounds in 3 exist, and also that there are $x < y$ forming the required nontrivial intervals. (For similar constructions see, for example, [9, Chapter 5].)

Below, we will refer to condition 3 as the “saturation condition on F, \mathcal{R} .” We will use the terms “lower bound” and “upper bound” in the weak sense (\leq, \geq), and we will use the terms “strict lower bound” and “strict upper bound” in the strong sense ($<, >$). For $x_1, \dots, x_n \in R$, we write $F[x_1, \dots, x_n]$ for the subfield of R obtained by adjoining x_1, \dots, x_n to F .

Lemma 9.3. *Let $x_1, \dots, x_n, y, z \in R$, where $y < z$. There exists $y < w < z$ such that w is not algebraic over $F[x_1, \dots, x_n]$.*

Proof. Let x_1, \dots, x_n, y, z be as given. Let g_1, g_2, \dots be \mathcal{R} -definable functions where the union of their images over appropriate Cartesian powers of F consists of all elements $1/(u - y)$, where $u > y$ is algebraic over $F[x_1, \dots, x_n]$. By the saturation property of F, \mathcal{R} , these elements have a strict upper bound b . Hence $y + 1/b$ is a strict lower bound on these elements. Set $w = y + 1/b$. \square

Our goal in the next two sections will be to construct F -bijections of R such that properties of \mathbb{Q} or \mathbb{Z} are coded into $h^{-1}[\mathcal{R}]$. Our strategy will be to iteratively extend partial F -homomorphisms until they become total and onto. The following definitions and lemmas will support our constructions.

Definition 9.4. Let $V[F, \mathcal{R}]$ be the family of all sets $E \subseteq R$ such that for some $x_1, \dots, x_n \in R$, $n \geq 0$,

$$E = \{ax_i \mid 1 \leq i \leq n \wedge a \in F\}.$$

Let $W[F, \mathcal{R}]$ be the set of all partial one-one functions h from R into R such that the following hold:

1. $\text{dom}(h) \in V[F, \mathcal{R}]$.
2. h is order preserving.
3. $h(0) = 0$ and $h(1) = 1$.
4. For all $x \in \text{dom}(h)$ and $a \in F$, we have $h(ax) = ah(x)$.

Note that for all $h \in W[F, \mathcal{R}]$, $\text{rng}(h) \in V[F, \mathcal{R}]$.

Lemma 9.5. *Every $E \in V[F, R]$ is the image of an \mathcal{R} -definable function on some F^n . Every $h \in W[F, R]$ is the restriction of an \mathcal{R} -definable function to its domain.*

Proof. The first claim follows immediately from the definition. For the second claim, fix $x_1, \dots, x_n \in R$ such that $\text{dom}(h) = \{ax_i \mid 1 \leq i \leq n, a \in F\}$. Then $h = h_1 \cup \dots \cup h_n$, where each $h_i : \{ax_i \mid a \in F\} \rightarrow \{ah(x_i) \mid a \in F\}$ is given by $h_i(ax_i) = ah_i(x_i)$. \square

Lemma 9.6. *For all $h \in W[F, \mathcal{R}]$, $h^{-1} \in W[F, \mathcal{R}]$.*

Proof. Let $h \in W[F, \mathcal{R}]$. For all $x, y \in \text{rng}(h) = \text{dom}(h^{-1})$, if $x < y$, then $h(h^{-1}(x)) < h(h^{-1}(y))$, and so $h^{-1}(x) < h^{-1}(y)$. Similarly, $h^{-1}(0) = h^{-1}(h(0)) = 0$ and $h^{-1}(1) = h^{-1}(h(1)) = 1$. For any a in F and x in $\text{rng}(h)$, $h^{-1}(ax) = h^{-1}(h(a)h(h^{-1}(x))) = h^{-1}(h(ah^{-1}(x))) = ah^{-1}(x)$, as required. \square

The following proposition provides a connection between types over $T_{\text{comm}}[F]$, which were discussed in Section 6, and the elements of $W[F, \mathcal{R}]$.

Proposition 9.7. *Let $x_1, \dots, x_n, y_1, \dots, y_n$ be elements of R . Then there is an $h \in W[F, \mathcal{R}]$ satisfying $h(x_i) = y_i$ for every i if and only if \vec{x} and \vec{y} have the same types over $T_{\text{comm}}[F]$.*

We will not use Proposition 9.7 below, and so we omit the proof, which is straightforward.

We now determine ways in which elements of $W[F, \mathcal{R}]$ can be extended. We write $\text{fld}(h)$ for $\text{dom}(h) \cup \text{rng}(h)$. The F -multiples of $x \in R$ are the elements ax , for $a \in F$. We write $h \subseteq_1 h'$ if and only if the following hold:

1. $h, h' \in W[F, \mathcal{R}]$.
2. $h \subseteq h'$.
3. There exists $x \in \text{dom}(h') \setminus \text{dom}(h)$ such that $\text{dom}(h') = \text{dom}(h) \uplus \{ax \mid a \in F \setminus \{0\}\}$.

Here, \uplus denotes a disjoint union. Then $h \subseteq_1 h'$ is equivalent to the following assertions.

1. $h, h' \in W[F, \mathcal{R}]$.
2. $h \subseteq h'$.
- 3' There exists $y \in \text{rng}(h') \setminus \text{rng}(h)$ such that $\text{rng}(h') = \text{rng}(h) \uplus \{ay \mid a \in F \setminus \{0\}\}$.

Note that $h \subseteq_1 h'$ if and only if $h^{-1} \subseteq_1 h'^{-1}$. Note also that in 3,3' above, x and y are not unique, but they are unique up to multiplication by an element of F .

Lemma 9.8. *Let $h \in W[F, \mathcal{R}]$ and $x \in R \setminus \text{dom}(h)$, $x > 0$. There exists a nontrivial interval J such that the following holds: for all $y \in J$, there exists $h \subseteq_1 h'$ such that $h'(x) = y$.*

Proof. Let h, x be as given. Obviously $\text{rng}(h) = h[\text{dom}(h) \upharpoonright_{<x}] \uplus h[\text{dom}(h) \upharpoonright_{>x}]$, where $h[\text{dom}(h) \upharpoonright_{<x}]$ lies strictly below $h[\text{dom}(h) \upharpoonright_{>x}]$.

Case 1. $\text{dom}(h) \upharpoonright_{>x}$ is empty. Let J be the interval of elements of R strictly above $\text{rng}(h)$. By Lemma 9.5 and the saturation property of F, \mathcal{R} , $\text{fld}(h)$ has a strict upper bound. Hence J is nontrivial. Let $y \in J$, and define $h'(ax) = ay$, for all $a \in F$. We have only to verify that $h' \in W[F, \mathcal{R}]$.

It suffices to show that h' is order preserving. First, suppose $ax < a'x$, $a, a' \in F \setminus \{0\}$. Then $a < a'$, and so $h'(ax) = ah'(x) < a'h'(x) = h'(a'x)$.

Next, suppose $v < ax$, $a \in F \setminus \{0\}$, $v \in \text{dom}(h)$. If $a < 0$ then $v/a > x$, which is impossible. Hence $a > 0$. Now $h(v/a) < h'(x)$. Hence $h(v) < ah'(x) = h'(ax)$.

Finally, suppose $ax < v$, $a \in F \setminus \{0\}$, $v \in \text{dom}(h)$. If $a > 0$ then $x < v/a$, which is impossible. Hence $a < 0$. Now $h(v/a) < h'(x)$, so $h(v)/a < h'(x)$, $h(v) > ah'(x) = h'(ax)$, and $h'(ax) < h(v)$.

Case 2. $dom(h)|_{<x}$ and $dom(h)|_{>x}$ are nonempty. Let J be the interval lying strictly above $h[dom(h)|_{<x}]$ and strictly below $h[dom(h)|_{>x}]$. By Lemma 9.5, these two sets are each images of an \mathcal{R} -definable function on some F^n . Hence by the saturation condition on F, \mathcal{R} , J is nontrivial. Let $y \in J$, and define $h'(ax) = ay$, for all $a \in F$. We have only to verify that $h' \in W[F]$.

It suffices to show that h' is order preserving. Suppose $ax < a'x$, $a, a' \in F \setminus \{0\}$. Then $a < a'$, and so $h'(ax) = ah'(x) < a'h'(x) = h'(a'x)$.

Suppose $v < ax$, $a \in F \setminus \{0\}$, $v \in dom(h)$. First assume $a > 0$. Then $v/a < x$, and so $h(v/a) < h'(x)$, $h(v)/a < h'(x)$, and $h(v) < ah'(x) = h'(ax)$. Now assume $a < 0$. Then $v/a > x$, and so $h(v/a) > h'(x)$, $h(v)/a > h'(x)$, and $h(v) < ah'(x) = h'(ax)$.

Finally, suppose $ax < v$, $a \in F \setminus \{0\}$, $v \in dom(h)$. First assume $a > 0$. Then $x < v/a$, and so $h'(x) < h(v/a) = h(v)/a$, $ah'(x) < h(v)$, $h'(ax) < h(v)$. Now assume $a < 0$. Then $x > v/a$, and so $h'(x) > h(v/a) = h(v)/a$, $ah'(x) < h(v)$, and $h'(ax) < h(v)$. \square

Lemma 9.9 (First Extension Lemma). *Let $h \in W[F]$ and $x \notin dom(h)$. Then there exists a nontrivial interval J such that the following holds: for all $y \in J$, there exists $h \subseteq_1 h'$ such that $h'(x) = y$.*

Proof. Let h, x be as given. The case $x > 0$ is given by Lemma 9.8. So, suppose $x < 0$. Apply Lemma 9.8 to the case $-x > 0$, obtaining a nontrivial J such that for all $y \in J$, there exists $h \subseteq_1 h'$ such that $h'(-x) = y$.

We claim that $-J$ is a nontrivial interval such that for all $y \in -J$, there exists $h \subseteq_1 h'$ such that $h'(x) = y$. To see this, let $y \in -J$. Then $-y \in J$, and hence there exists $h \subseteq_1 h'$ such that $h'(-x) = -y$. But $h'(-x) = -y$ implies $h'(x) = y$, as required. \square

Lemma 9.10 (Second Extension Lemma). *Let $h \in W[F]$ and $x \notin rng(h)$. There exists a nontrivial interval J such that the following holds: for all $y \in J$, there exists $h \subseteq_1 h'$ such that $h'(y) = x$.*

Proof. We obtain this from Lemma 9.9 as follows. Let h, x be as given. Then $h^{-1} \in W[F]$ and $x \notin dom(h^{-1})$. By Lemma 9.8, let J be a nontrivial interval such that for all $y \in J$, there exists $h^{-1} \subseteq_1 h'$ such that $h'(x) = y$.

We claim that for all $y \in J$, there exists $h \subseteq_1 h''$ such that $h''(y) = x$. To see this, let $h \subseteq_1 h'$ be such that $h'(x) = y$. Then $h^{-1} \subseteq_1 h'^{-1}$ and $h'^{-1}(y) = x$. That is, we can set $h'' = h'^{-1}$. \square

10 Existential consequences of $T[F]$

The existential theory of F consists of all sentences

$$\exists x_1, \dots, x_n \in F \varphi(x_1, \dots, x_n)$$

where φ is a quantifier free formula involving $+, \times, <$, and is interpreted in \mathbb{R} . Here we show that the existential theory of F can be effectively reduced to

the existential consequences of $T[F]$ without auxiliary functions. This yields, in particular, a conditional undecidability result for $T[\mathbb{Q}]$; see Corollary 10.6 below.

We adhere strictly to the convention that if an equation holds, then both sides must be defined. Also, a term is defined if and only if each subterm is defined. For example,

$$h^{-1}(h(x)h(1+x)) = x + h^{-1}(h(x)^2)$$

implies that both sides of this equation are defined. In particular, the above equation implies that $x, 1+x \in \text{dom}(h)$.

Let $h \in W[F, \mathcal{R}]$. We write $\text{alg}(F, h)$ for the elements that are algebraic over some $F[x_1, \dots, x_n]$, $x_1, \dots, x_n \in \text{fld}(h)$. We write $\text{trans}(F, h)$ for $R \setminus \text{alg}(F, h)$.

Note that by Lemma 9.5, there exists $x_1, \dots, x_n \in R$ such that every element of $\text{alg}(F, h)$ is algebraic over $F[x_1, \dots, x_n]$. This allows us to use Lemma 9.3 to obtain an element of $\text{trans}(F, h)$ in every nontrivial interval.

Lemma 10.1. *Let $h \in W[F, \mathcal{R}]$ be such that for every x , if*

$$h^{-1}(h(x)h(1+x)) = x + h^{-1}(h(x)^2),$$

then $x \in F$. Let $b \notin \text{dom}(h)$. There exists $h \subseteq_1 h' \in W[F, \mathcal{R}]$, $h'(b)$ defined, such that h' has the same property; i.e. for every x , if

$$h'^{-1}(h'(x)h'(1+x)) = x + h'^{-1}(h'(x)^2),$$

then $x \in F$.

Proof. Let h, b be as given. By Lemmas 9.9 and 9.3, define $h \subseteq_1 h'$ such that $h'(b) \in \text{trans}(F, h)$. We first show the conclusion for all $ab, a \in F$. We assume

$$h'^{-1}(h'(ab)h'(1+ab)) = ab + h'^{-1}(h'(ab)^2).$$

and derive a contradiction. Clearly $h'(ab)^2 = (ah'(b))^2 = a^2h'(b)^2 \in \text{rng}(h')$. Since $a \in F$ and $h'(b) \in \text{trans}(F, h)$, $a^2h'(b)^2 \in \text{rng}(h') \setminus \text{rng}(h)$, which consists of the nonzero F -multiples of $h'(b)$. This contradicts that $h'(b) \in \text{trans}(F, h)$.

Finally, we show the conclusion for all $x \in \text{dom}(h) \setminus F$. We assume

$$h'^{-1}(h'(x)h'(1+x)) = x + h'^{-1}(h(x)^2) \tag{8}$$

and derive a contradiction. By the hypothesis on h , (8) does not hold with h' replaced by h . Hence if we replace h' by h , at least one side of (8) is undefined.

Case 1. $h(1+x)$ is undefined. Let $1+x = ab$, $a \in F \setminus \{0\}$. Hence

$$h'^{-1}(h(x)ah'(b)) = x + h'^{-1}(h'(ab^{-1})^2).$$

Hence $h(x)h'(b) \in \text{rng}(h')$. Since $h'(b) \in \text{trans}(F, h)$, $h(x)h'(b) \in \text{rng}(h') \setminus \text{rng}(h)$. Hence $h(x)h'(b)$ is a nonzero F -multiple of $h'(b)$. This contradicts that $h(x) \notin F$.

Case 2. $h(1+x)$ is defined, but $h^{-1}(h(x)h(1+x))$ is not defined. Then $h(x)h(1+x)$ is a nonzero F -multiple of $h'(b)$. Since $x \neq -1$, this product is nonzero. This contradicts that $h'(b) \in \text{trans}(F, h)$.

Case 3. $h^{-1}(h(x)h(1+x))$ is defined, but $h^{-1}(h(x)^2)$ is undefined. Then $h(x)^2$ is a nonzero F -multiple of $h'(b)$. This contradicts that $h'(b) \in \text{trans}(F, h)$. \square

Lemma 10.2. *Let $h \in W[F, \mathcal{R}]$ be such that for every x , if*

$$h^{-1}(h(x)h(1+x)) = x + h^{-1}(h(x)^2)$$

then $x \in F$. Let $b \notin \text{rng}(h)$. Then there exists $h \subseteq_1 h' \in W$ such that $h'^{-1}(b)$ is defined, and for every x , if

$$h'^{-1}(h'(x)h'(1+x)) = x + h'^{-1}(h'(x)^2)$$

then $x \in F$.

Proof. Let h, b be as given. By Lemmas 9.10 and 9.3, let $h \subseteq_1 h'$, where $h'^{-1}(b) \in \text{trans}(F, h)$. Write $c = h'^{-1}(b)$.

We first show the conclusion for all $ac, a \in F \setminus \{0\}$. We assume

$$h'^{-1}(h'(ac)h'(1+ac)) = ac + h'^{-1}(h'(ac)^2)$$

and derive a contradiction. From the assumption, we have $1+ac \in \text{dom}(h')$. Since $c \in \text{trans}(F, h)$, $1+ac \in \text{dom}(h') \setminus \text{dom}(h)$. Hence $1+ac$ is a nonzero F -multiple of c . This contradicts $c \in \text{trans}(F, h)$.

Finally, we show the conclusion for all $x \in \text{dom}(h) \setminus F$. We assume

$$h'^{-1}(h(x)h'(1+x)) = x + h'^{-1}(h(x)^2) \tag{9}$$

and derive a contradiction. By the hypothesis on h , (9) does not hold with h' replaced by h . Hence if we replace h' by h , at least one side of (9) is undefined.

Case 1. $h^{-1}(h(x)^2)$ is undefined. Then $h(x)^2$ is a nonzero F -multiple of b and $h'^{-1}(h(x)^2)$ is a nonzero F -multiple ac of c . Clearly the left side of (9) either lies in $\text{dom}(h)$ or is a nonzero F -multiple ac of c . Both possibilities contradict that $c \in \text{trans}(F, h)$.

Case 2. $h^{-1}(h(x)^2)$ is defined and $h(1+x)$ is undefined. Then $h(1+x)$ is a nonzero F -multiple of b and $1+x$ is a nonzero F -multiple of c . This contradicts that $c \in \text{trans}(F, h)$.

Case 3. $h^{-1}(h(x)^2)$ and $h(1+x)$ are defined, but $h^{-1}(h(x)h(1+x))$ is undefined. Hence $h(x)h(1+x)$ is a nonzero F -multiple of b and $h'^{-1}(h(x)h(1+x))$ is a nonzero F -multiple of c . But the right side of (9) is algebraic in $\text{fld}(h)$. This is a contradiction. \square

Theorem 10.3. *There is a model \mathcal{M} of $T[F]^*$ with domain R , with the same $0, 1, +, <$ of R , in which for all b , $b(1+b) = b + b^2$ holds if and only if $b \in F$. In this equation, we use the multiplication of \mathcal{M} to multiply b and $1+b$.*

Proof. Let h be the identity function on F . Then $h \in W[F, \mathcal{R}]$, and trivially we have that

- for every x , if $h^{-1}(h(x)h(1+x)) = x + h^{-1}(h(x)^2)$ then $x \in F$; and
- for every $x \in F$, $h^{-1}(h(x)h(1+x)) = x + h^{-1}(h(x)^2)$.

Thus we can iterate Lemmas 10.1 and 10.2, starting with the identity function on F , diagonalizing over the countably many elements of R . We then obtain $h \in W[F, \mathcal{R}]$ with domain R , such that

for every x in R , $h^{-1}(h(x)h(1+x)) = x + h^{-1}(h(x)^2)$ if and only if $x \in F$.

The required model \mathcal{M} of $T[F]^*$ is $h^{-1}[\mathcal{R}]$. Calculating in \mathcal{M} , we have

$$x \otimes (1+x) = h^{-1}(h(x)h(1+x))$$

and

$$x + (x \otimes x) = x + h^{-1}(h(x)^2)$$

Hence, for every x in R , we have $x \otimes (1+x) = x + (x \otimes x)$ if and only if $x \in F$, as required. \square

Corollary 10.4. *An existential sentence φ over F in the language of ordered fields is true if and only if in any model of $T[F]^*$, φ has witnesses among the b with $b(1+b) = b + b^2$.*

Proof. Suppose φ has the form $\exists x_1, \dots, x_n \psi(x_1, \dots, x_n)$ with ψ quantifier-free, and suppose $\psi(a_1, \dots, a_n)$ holds with $a_1, \dots, a_n \in F$. Let \mathcal{M} be a model of $T[F]^*$. Then for all $1 \leq i \leq n$, $T[F]^*$ proves $\varphi(c_{a_1}, \dots, c_{a_n})$ and $c_{a_i}(1+c_{a_i}) = c_{a_i} + c_{a_i}^2$.

For the converse, Let \mathcal{M} be a model of $T[F]^*$ given by Theorem 10.3. Then the witnesses must lie in F . \square

Corollary 10.5. *The existential theory over F is effectively reducible to the existential consequences of $T[F]^*$ without auxiliary constants, and to the existential consequences of $T[F]$ without auxiliary functions. The reduction can be accomplished in linear time.*

Proof. From Theorem 4.1 and Corollary 10.4. By Theorem 4.3, the we can use $T[F]$ in place of $T[F]^*$. \square

Corollary 10.6. *If Hilbert's 10th Problem over the rationals is undecidable (as expected), then the existential consequences of $T[\mathbb{Q}]$ and $T[\mathbb{Q}]^*$, not mentioning auxiliary constants or auxiliary functions, respectively, are each undecidable. The former can be reduced to the latter by a linear time reduction.*

Proof. Immediate from Corollary 10.5. \square

11 $\forall\forall\forall\exists^*$ consequences of $T[F]$

We use \mathbb{Z}^+ for the set of all positive integers, and \mathbb{N} for the set of all nonnegative integers.

Lemma 11.1. *There exists $\mu, \kappa, \lambda \in R$ such that*

1. *For every n in \mathbb{N} , we have $n < \mu$, $\mu^n < \kappa$, and $\kappa^n < \lambda$.*
2. $[\mu, \infty) \cap F = \emptyset$.

Proof. By the saturation condition on F, \mathcal{R} . □

We fix μ, κ, λ given by Lemma 11.1. Let $K[F, \mathcal{R}]$ be the set of all functions h such that

1. $h \in W[F, \mathcal{R}]$.
2. h is the identity on $\{\mu, \kappa, \lambda, \mu\kappa, \mu\lambda, \kappa\lambda\}$.

We will build a bijection $h \in K[F, \mathcal{R}]$, $h : R \rightarrow R$, such that for all $x \in R$, $1 \leq x \leq \mu$, the equation

$$(\kappa + x)(\lambda + x) = \kappa\lambda + \kappa x + \lambda x + x^2$$

holds in $h^{-1}[\mathcal{R}]$ if and only if $x \in \mathbb{N}$. That is, for all $x \in R$, $1 \leq x \leq \mu$,

$$\begin{aligned} f^{-1}(f(\kappa + x)f(\lambda + x)) = \\ f^{-1}(f(\kappa)f(\lambda)) + f^{-1}(f(\kappa)f(x)) + f^{-1}(f(\lambda)f(x)) + f^{-1}(f(x)^2) \end{aligned}$$

if and only if $x \in \mathbb{Z}^+$. In other words, for all $x \in R$, $1 \leq x \leq \mu$,

$$f^{-1}(f(\kappa + x)f(\lambda + x)) = \kappa\lambda + f^{-1}(\kappa f(x)) + f^{-1}(\lambda f(x)) + f^{-1}(f(x)^2)$$

if and only if $x \in \mathbb{Z}^+$.

Lemma 11.2. *Let $h \in K[F, \mathcal{R}]$, where for every x in $[1, \mu]$, if*

$$h^{-1}(h(\kappa + x)h(\lambda + x)) = \kappa\lambda + h^{-1}(\kappa h(x)) + h^{-1}(\lambda h(x)) + h^{-1}(h(x)^2),$$

then x is in \mathbb{Z}^+ . Let $b \notin \text{dom}(h)$. Then there exists $h \subseteq_1 h'$ such that $h'(b)$ is defined and for every x in $[1, \mu]$, if

$$h'^{-1}(h'(\kappa + x)h'(\lambda + x)) = \kappa\lambda + h'^{-1}(\kappa h'(x)) + h'^{-1}(\lambda h'(x)) + h'^{-1}(h'(x)^2),$$

then x is \mathbb{Z}^+ .

Proof. Let h, b be as given. By Lemmas 9.9 and 9.3, let $h \subseteq_1 h'$, where $h'(b) \in \text{trans}(F, h)$. Note that $\text{rng}(h') \setminus \text{rng}(h)$ consists of the nonzero F -multiples of $h'(b)$.

We first show the conclusion for all $ab, a \in F \setminus \{0\}$. We assume

$$h'^{-1}(h'(\kappa + ab)h'(\lambda + ab)) = \kappa\lambda + h^{-1}(\kappa h(ab)) + h^{-1}(\lambda h(ab)) + h'^{-1}(h'(ab)^2)$$

and derive a contradiction.

Clearly $h'^{-1}(h'(ab)^2) = h'^{-1}(a^2h'(b)^2)$ is defined. Since $h'(b) \in \text{trans}(F, h)$, $a^2h'(b)^2 \in \text{rng}(h') \setminus \text{rng}(h)$. Hence $a^2h'(b)^2$ is an F -multiple of $h'(b)$. This contradicts that $h'(b) \in \text{trans}(F, h)$.

Finally, we show the conclusion for all $x \in \text{dom}(h) \setminus \mathbb{Z}^+, 1 \leq x \leq \mu$. We assume

$$h'^{-1}(h'(\kappa + x)h'(\lambda + x)) = \kappa\lambda + h'^{-1}(\kappa h(x)) + h'^{-1}(\lambda h(x)) + h'^{-1}(h(x)^2) \quad (10)$$

and derive a contradiction. By the hypothesis on h , (10) does not hold with h' replaced by h . Hence if we replace h' by h , at least one side of (10) is undefined.

First, we claim that $h^{-1}(h(x)^2)$ is defined. Otherwise, $h(x)^2$ is a nonzero F -multiple of $h'(b)$. This contradicts that $h'(b) \in \text{trans}(F, h)$.

Second, we claim that $h^{-1}(\kappa h(x))$ is defined. Otherwise, $\kappa h(x)$ is a nonzero F -multiple of $h'(b)$.

Third, we claim that $h^{-1}(\lambda h(x))$ is defined. Otherwise, $\lambda h(x)$ is a nonzero F -multiple of $h'(b)$.

From these three claims, we see that the right side of (10) is defined if we replace h' by h . Therefore $h^{-1}(h(\kappa + x)h(\lambda + x))$ is undefined.

Case 1. $h(\kappa + x)$ and $h(\lambda + x)$ are undefined. Then $h'(\kappa + x), h'(\lambda + x)$ are nonzero F -multiples of $h'(b)$. Since $h'(b) \in \text{trans}(F, h)$, the product $h'(\kappa + x)h'(\lambda + x) \in \text{rng}(h') \setminus \text{rng}(h)$. Hence $h'(\kappa + x)h'(\lambda + x)$ is a nonzero F -multiple of $h'(b)$. Also $h'(\kappa + x)h'(\lambda + x)$ is a nonzero F -multiple of $h'(b)^2$. This contradicts that $h'(b) \in \text{trans}(F, h)$.

Case 2. $h(\kappa + x)$ is undefined, but $h(\lambda + x)$ is defined. Since $\lambda + x \neq 0$, we have $h(\lambda + x) \neq 0$. Now $h'(\kappa + x)$ is a nonzero F -multiple of $h'(b)$. Since $h'(b) \in \text{trans}(F, \text{fld}(h))$, $h'(\kappa + x)h(\lambda + x) \in \text{rng}(h') \setminus \text{rng}(h)$. Hence $h'(\kappa + x)h(\lambda + x)$ is a nonzero F -multiple of $h'(b)$. Therefore $h(\lambda + x) \in F$, and hence $h(\lambda + x) = \lambda + x \in F$. In particular, $\lambda + x \in F$ and $x \geq 0$. This contradicts Lemma 11.1.

Case 3. $h(\kappa + x)$ is defined, $h(\lambda + x)$ is undefined. Since $\kappa + x \neq 0$, we have $h(\kappa + x) \neq 0$. Now $h'(\lambda + x)$ is a nonzero F -multiple of $h'(b)$. Since $h'(b) \in \text{trans}(F, \text{fld}(h))$, $h(\kappa + x)h'(\lambda + x) \in \text{rng}(h') \setminus \text{rng}(h)$. Hence $h(\kappa + x)h'(\lambda + x)$ is a nonzero F -multiple of $h'(b)$. Therefore $h(\kappa + x) \in F$, and hence $h(\kappa + x) = \kappa + x \in F$. In particular, $\kappa + x \in F$ and $x \geq 0$. This contradicts Lemma 11.1.

Case 4. $h(\kappa + x)$ and $h(\lambda + x)$ are defined. Since $h^{-1}(h(\kappa + x)h(\lambda + x))$ is undefined, $h(\kappa + x)h(\lambda + x)$ is a nonzero F -multiple of $h'(b)$. This contradicts that $h'(b) \in \text{trans}(F, h)$. \square

Lemma 11.3. *Let $h \in K[F, \mathcal{R}]$ be such that for every x in $[1, \mu]$, if*

$$h^{-1}(h(\kappa + x)h(\lambda + x)) = \kappa\lambda + h^{-1}(\kappa h(x)) + h^{-1}(\lambda h(x)) + h^{-1}(h(x)^2),$$

then x is in \mathbb{Z}^+ . Let $b \notin \text{rng}(h)$. Then there exists $h \subseteq_1 h'$ such that $h'^{-1}(b)$ defined and for every x in $[1, \mu]$, if

$$h'^{-1}(h'(\kappa + x)h'(\lambda + x)) = \kappa\lambda + h'^{-1}(\kappa h'(x)) + h'^{-1}(\lambda h'(x)) + h'^{-1}(h'(x)^2),$$

then x is in \mathbb{Z}^+ .

Proof. Let h, b be as given. By Lemmas 9.10 and 9.3, let $h \subseteq_1 h'$, where $h'^{-1}(b) \in \text{trans}(F, \text{fld}(h))$. Write $c = h'^{-1}(b)$. Note that $\text{dom}(h') \setminus \text{dom}(h)$ consists of the nonzero F -multiples of c .

We first show the conclusion for all $ac, a \in F \setminus \{0\}$. We assume

$$h'^{-1}(h'(\kappa + ac)h'(\lambda + ac)) = \kappa\lambda + h'^{-1}(\kappa h(ac)) + h'^{-1}(\lambda h(ac)) + h'^{-1}(h'(ac)^2)$$

and derive a contradiction. In particular, the assumption implies that $h'(\kappa + ac)$ is defined, and so $\kappa + ac \in \text{dom}(h)$ or $\kappa + ac$ is an F -multiple of c . Both alternatives contradict that $c \in \text{trans}(F, \text{fld}(h))$.

Finally, we show the conclusion for all $x \in \text{dom}(h) \setminus \mathbb{Z}^+, 1 \leq x \leq \mu$. We assume

$$h'^{-1}(h'(\kappa + x)h'(\lambda + x)) = \kappa\lambda + h'^{-1}(\kappa h(x)) + h'^{-1}(\lambda h(x)) + h'^{-1}(h(x)^2) \quad (11)$$

and derive a contradiction.

There are five terms in (11). The four terms other than $\kappa\lambda$ are each either a nonzero F -multiple of c or an element of $\text{fld}(h)$. Since $c \in \text{trans}(F, \text{fld}(h))$, the ones that are nonzero F -multiples of c must cancel.

We now use the inequalities on x, μ, κ , and λ . Note that

- $(\kappa + x)(\lambda + x) > \kappa\lambda$.
- $h'((\kappa + x)(\lambda + x)) > h'(\kappa\lambda) = \kappa\lambda$.
- $h'^{-1}(h'(\kappa + x)h'(\lambda + x)) > h'^{-1}(\kappa\lambda) = \kappa\lambda$.
- $x \leq \mu$.
- $h(x) \leq h(\mu) = \mu$.
- $\kappa h(x) \leq \mu\kappa$.
- $h'^{-1}(\kappa h(x)) \leq h'^{-1}(\mu\kappa) = \mu\kappa$.
- $\lambda h(x) \leq \mu\lambda$.
- $h'^{-1}(\lambda h(x)) \leq h'^{-1}(\mu\lambda) = \mu\lambda$.
- $h(x)^2 \leq \mu^2 < \kappa$.
- $h'^{-1}(h(x)^2) \leq h'^{-1}(\kappa) < \kappa$.
- $h'^{-1}(h'(\kappa + x)h'(\lambda + x)) > \kappa\lambda > \mu\kappa + \mu\lambda + \kappa \geq h'^{-1}(\kappa h(x)) + h'^{-1}(\lambda h(x)) + h'^{-1}(h(x)^2)$.

It is now obvious that the terms that are nonzero F -multiples of c cannot include $h'^{-1}(h'(\kappa+x)h'(\lambda+x))$.

This leaves $h'^{-1}(\kappa h(x))$, $h'^{-1}(\lambda h(x))$, $h'^{-1}(h(x)^2)$ as the terms that might be nonzero F -multiples of c . Using the above, we have

- $h'^{-1}(h(x)^2) < \kappa$.
- $h'^{-1}(\kappa h(x)) \leq \mu\kappa$.
- $h'^{-1}(\lambda h(x)) \leq \mu\lambda$.
- $x \geq 1$.
- $h(x) \geq h(1) = 1$.
- $\kappa h(x) \geq \kappa$.
- $h'^{-1}(\kappa h(x)) \geq h'^{-1}(\kappa) = \kappa$.
- $h(x) \geq h(1) = 1$.
- $\lambda h(x) \geq \lambda$.
- $h'^{-1}(\lambda h(x)) \geq h'^{-1}(\lambda) = \lambda$.

Hence

- $h'^{-1}(h(x)^2) < \kappa$.
- $\kappa \leq h'^{-1}(\kappa h(x)) \leq \mu\kappa$.
- $\lambda \leq h'^{-1}(\lambda h(x))$.

It is now clear that none of $h'^{-1}(\kappa h(x))$, $h'^{-1}(\lambda h(x))$, $h'^{-1}(h(x)^2)$ can be a nonzero F -multiple of c . Hence

$$h'^{-1}(h'(\kappa+x)h'(\lambda+x)), \quad h'^{-1}(\kappa h(x)), \quad \text{and} \quad h'^{-1}(\lambda h(x)), h'^{-1}(h(x)^2)$$

all lie in $\text{dom}(h)$. Therefore

$$h'(\kappa+x)h'(\lambda+x), \quad \kappa h(x), \quad \lambda h(x), \quad \text{and} \quad h(x)^2$$

lie in $\text{rng}(h)$. We claim that $h'(\kappa+x), h'(\lambda+x) \in \text{rng}(h)$. To see this, first suppose both are not in $\text{rng}(h)$. Then $\kappa+x$ and $\lambda+x$ are F -multiples of c , and so $(\kappa+x)(\lambda+x)$ is of the form $aa'c^2$, where $a, a' \in F$. This contradicts the fact that c is in $\text{trans}(F, h)$.

Now suppose one of them, say, by symmetry, $h'(\kappa+x)$, is an F -multiple of c , and the other, $h'(\lambda+x)$, lies in $\text{rng}(h)$. Since $\lambda+x \neq 0$, we have $h'(\lambda+x) \neq 0$. Then $h'(\kappa+x)h'(\lambda+x)$ is of the form acu , where $a \in F \setminus \{0\}$ and $u \in \text{rng}(h)$. But $h'(\kappa+x)h'(\lambda+x) \in \text{rng}(h)$. Hence $acu \in \text{rng}(h) \setminus \{0\}$. This contradicts the fact that c is in $\text{trans}(F, h)$.

From $h'(\kappa+x), h'(\lambda+x) \in \text{rng}(h)$, we obtain that $\kappa+x, \lambda+x \in \text{dom}(h)$. Thus we see that both sides of (11) are defined if we replace h' by h . Hence (11) holds with h' replaced by h . This is a contradiction. \square

We want to iterate Lemmas 11.2 and 11.3, but we first need to deal with the base case. Let

$$S = \{\kappa + x : x \in \mathbb{Z}^+\} \cup \{\lambda + x : x \in \mathbb{Z}^+\} \cup \{\kappa\lambda + \kappa x + \lambda x + x^2 \mid x \in \mathbb{Z}^+\} \cup \{1, \mu, \kappa, \lambda, \mu\kappa, \mu\lambda, \kappa\lambda\}.$$

Let S' be the set of all F -multiples of elements of S .

Lemma 11.4. *Let $x \in S'$, $1 \leq x \leq \mu$. If $\kappa + x \in S'$ then $x \in \mathbb{Z}^+$. If $\lambda + x \in S'$ then $x \in \mathbb{Z}^+$.*

Proof. Let x be as given. Suppose $\kappa + x \in S'$. Since $\kappa + x < 2\kappa$, clearly $\kappa + x$ is not a nonzero F -multiple of any element of

$$\{\lambda + x \mid x \in \mathbb{Z}^+\} \cup \{\kappa\lambda + \kappa x + \lambda x + x^2 \mid x \in \mathbb{Z}^+\} \cup \{\lambda, \mu\kappa, \mu\lambda, \kappa\lambda\}.$$

Since $\kappa + x$ is greater than every μ^n , $n \in \mathbb{Z}^+$, $\kappa + x$ is not a nonzero F -multiple of any element of $\{1, \mu\}$.

Now suppose $\kappa + x$ is an F -multiple of $\kappa + y$, $y \in \mathbb{N}$. Write $\kappa + x = a(\kappa + y)$, $a \in F$. Then $\kappa = (ay - x)/(1 - a)$ or $a = 1$. Now $|ay - x| \leq |ay| + |x| \leq \mu + \mu = 2\mu$. Also $1/|1 - a| \leq \mu$ or $a = 1$. Hence $\kappa \leq 2\mu^2$ or $a = 1$. Therefore $a = 1$. Hence $\kappa + x = \kappa + y$, and $x = y$. Therefore $x \in \mathbb{Z}^+$.

Suppose $\lambda + x \in S'$. Since $\lambda + x < 2\lambda$, clearly $\lambda + x$ is not a nonzero F -multiple of any element of $\{\mu\lambda, \kappa\lambda\} \cup \{\kappa\lambda + \kappa x + \lambda x + x^2 \mid x \in \mathbb{N}\}$. Since $\lambda + x$ is greater than every κ^n , $n \in \mathbb{Z}^+$, $\lambda + x$ is not a nonzero F -multiple of any element of $\{\kappa + x : x \in \mathbb{Z}^+\} \cup \{1, \mu, \kappa, \mu\kappa\}$.

Now suppose $\lambda + x$ is a nonzero F -multiple of $\lambda + y$, $y \in \mathbb{Z}^+$. Argue as above that $x \in \mathbb{Z}^+$. \square

Lemma 11.5. *There exists a bijection $h \in K[F, \mathcal{R}]$, $h : R \rightarrow R$, such that the following holds. For all $x \in \text{dom}(h)$ with $1 \leq x \leq \mu$, we have*

$$h^{-1}(h(\kappa + x)h(\lambda + x)) = \kappa\lambda + h^{-1}(\kappa h(x)) + h^{-1}(\lambda h(x)) + h^{-1}(h(x)^2)$$

if and only if x is in \mathbb{Z}^+ .

Proof. Let h be the identity function on S' . Obviously $h \in K[F, \mathcal{R}]$. By Lemma 11.4, for all $x \in \text{dom}(h)$ such that $1 \leq x \leq \mu$, if

$$h^{-1}(h(\kappa + x)h(\lambda + x)) = \kappa\lambda + h^{-1}(\kappa h(x)) + h^{-1}(\lambda h(x)) + h^{-1}(h(x)^2)$$

then $x \in \mathbb{Z}^+$. This is because for the relevant x , if $h(\kappa + x)$ is defined then $x \in \mathbb{Z}^+$.

For the reverse, let $x \in \mathbb{Z}^+$, and note that

$$h^{-1}(h(\kappa + x)h(\lambda + x)) = h^{-1}((\kappa + x)(\lambda + x)) = h^{-1}(\kappa\lambda + \kappa x + \lambda x + x^2) = \kappa\lambda + \kappa x + \lambda x + x^2.$$

So

$$\begin{aligned}\kappa\lambda + h^{-1}(\kappa h(x)) + h^{-1}(\lambda h(x)) + h^{-1}(h(x)^2) &= \\ \kappa\lambda + h^{-1}(\kappa x) + h^{-1}(\lambda x) + h^{-1}(x^2) &= \kappa\lambda + \kappa x + \lambda x + x^2.\end{aligned}$$

□

Lemma 11.6. *There exists a bijection $h \in K[F, \mathcal{R}]$, $h : R \rightarrow R$, such that the following holds. For all $x \in R$ with $1 \leq x \leq \mu$, we have*

$$h^{-1}(h(\kappa + x)h(\lambda + x)) = \kappa\lambda + h^{-1}(\kappa h(x)) + h^{-1}(\lambda h(x)) + h^{-1}(h(x)^2)$$

if and only if x is in \mathbb{Z}^+ .

Proof. Start with the h given by Lemma 11.5, and iterate Lemmas 11.2 and 11.3, diagonalizing over the countably many elements of R . □

Theorem 11.7. *There is a model \mathcal{M} of $T[F]^*$ with domain R , with the same $0, 1, +, <$ as \mathcal{R} , with three elements μ, κ, λ such that the following holds. For all $x \in R$ with $1 \leq x \leq \mu$, we have $(\kappa + x)(\lambda + x) = \kappa\lambda + \kappa x + \lambda x + x^2$ if and only if x is in \mathbb{Z}^+ . In this equation, we use the multiplication of \mathcal{M} .*

Proof. By Theorem 9.2 and Lemma 11.6. □

We say that a quadruple $\langle M, \mu, \kappa, \lambda \rangle$ has property (*) if and only if

1. \mathcal{M} is a model of $T[F]^*$.
2. $\mu, \kappa, \lambda \in \text{dom}(\mathcal{M})$.
3. The $x \in \text{dom}(M)$ for which $1 \leq x \leq \mu$ and $(\kappa + x)(\lambda + x) = \kappa\lambda + \kappa x + \lambda x + x^2$ contain 1 and are closed under $+1$.

There is the stronger property (**) of $\langle \mathcal{M}, \mu, \kappa, \lambda \rangle$ that asserts the following.

1. \mathcal{M} is a model of $T[F]^*$.
2. $\mu, \kappa, \lambda \in \text{dom}(\mathcal{M})$.
3. The $x \in \text{dom}(M)$ for which $1 \leq x \leq \mu$ and $(\kappa + x)(\lambda + x) = \kappa\lambda + \kappa x + \lambda x + x^2$ are exactly the positive integers in \mathcal{M} .

Corollary 11.8. *Let D be a Diophantine equation over the positive integers. Then D has a solution in nonnegative integers if and only if the following holds. For all quadruples $\langle \mathcal{M}, \mu, \kappa, \lambda \rangle$ with property (*), D has a solution over the x such that $1 \leq x \leq \mu$ and $(\kappa + x)(\lambda + x) = \kappa\lambda + \kappa x + \lambda x + x^2$.*

Proof. Let D be as given. Suppose D has a solution in the positive integers. Let $\langle \mathcal{M}, \mu, \kappa, \lambda \rangle$ have property (*). Then the x such that $1 \leq x \leq \mu$ and $(\kappa + x)(\lambda + x) = \kappa\lambda + \kappa x + \lambda x + x^2$ must contain the positive integers.

Conversely, suppose that for all quadruples $\langle \mathcal{M}, \mu, \kappa, \lambda \rangle$ with property (*), D has a solution over the x such that $1 \leq x \leq \mu$ and $(\kappa + x)(\lambda + x) = \kappa\lambda + \kappa x + \lambda x + x^2$. By Theorem 11.8, there exists $\langle \mathcal{M}, \mu, \kappa, \lambda \rangle$ with property (**). Hence D has a solution over the positive integers. □

Theorem 11.9. *The set of consequences of $T[F]^*$ without auxiliary constants, and of $T[F]$ without auxiliary functions, is undecidable. In fact, the set of $\forall\forall\forall\exists^*$ consequences of $T[F]^*$ without auxiliary constants, and of $T[F]$ without auxiliary functions, is complete r.e.*

Proof. We use Corollary 11.8 and that Hilbert's 10th problem over \mathbb{Z}^+ is complete r.e. We can express

$$\langle \mathcal{M}, \mu, \kappa, \lambda \rangle \text{ has property } (*)$$

as the formula $\varphi(\mu, \kappa, \lambda)$ given by

$$\begin{aligned} (\kappa + 1)(\lambda + 1) = \kappa\lambda + \kappa + \lambda + 1 \wedge \\ \forall x ((1 \leq x \leq \mu \wedge (\kappa + x)(\lambda + x) = \kappa\lambda + \kappa x + \lambda x + x^2) \rightarrow \\ (1 \leq x + 1 \leq \mu \wedge (\kappa + (x + 1))(\lambda + (x + 1)) = \kappa\lambda + \kappa(x + 1) + \lambda(x + 1) + (x + 1)^2)). \end{aligned}$$

Then we can write

$$\text{for all quadruples } \langle \mathcal{M}, \mu, \kappa, \lambda \rangle \text{ with property } (*), D \text{ has a solution} \\ \text{over the } x \text{ such that } 1 \leq x \leq \mu \text{ and } (\kappa + x)(\lambda + x) = \kappa\lambda + \kappa x + \lambda x + x^2$$

as the assertion that

$$\begin{aligned} \forall \mu, \kappa, \lambda (\varphi(\mu, \kappa, \lambda) \rightarrow \\ D \text{ has a solution over the } x \text{ such that } 1 \leq x \leq \mu \text{ and} \\ (\kappa + x)(\lambda + x) = \kappa\lambda + \kappa x + \lambda x + x^2) \end{aligned}$$

is provable in $T[F]^*$. Note that the sentence above is in the form $\forall\forall\forall\exists^*$. By Theorem 4.3, we can replace $T[F]$ by $T[F]^*$. \square

12 Avoiding disjunctions

In Section 7, we saw that the universal fragment of $T[\mathbb{Q}]$ is decidable. The proof, however, involves a complex reduction to the language of real closed fields. As a result, the procedure is of little practical importance: $T[\mathbb{Q}]$ is weaker than the theory of real closed fields, our decision procedure works for only the universal fragment of the language, and it does so less efficiently than procedures for the corresponding fragment of real closed fields. The procedure we describe is in no sense more extensible to larger languages than procedures for real closed fields. It may therefore seem as though we have taken a step in the wrong direction.

We maintain, however, that the analysis provides guidance in designing heuristic procedures for the reals that address the aims outlined in Section 1. An obvious strategy for capturing inferences like the ones described there is to work backwards from the desired conclusion, using the obvious monotonicity laws. For example, when the terms s , t , and u are known to be positive, one can prove $st \leq uv$ by proving $s \leq u$ and $t \leq v$. The examples presented in Section 1 can be verified by iteratively applying such rules.

There are drawbacks to such an approach, however. For one thing, excessive case splits can lead to exponential blowup; e.g. one can show $st > 0$ by showing that s and t are either both strictly positive or both strictly negative. And the relevant monotonicity inferences are generally nondeterministic: one can show $r + s + t > 0$ by showing that two of the terms are nonnegative and the third is strictly positive, and one can show $r + s < t + u + v + w$, say, by showing $r < u$, $s \leq t + v$, and $0 \leq w$.

In “straightforward” inferences that arise in practice, however, sign information is typically available. This is the case with the examples in Section 1, where all the relevant terms are easily seen to be positive. It is also the case with the following representative example, taken from the first author’s formalization of the prime number theorem [2]: verify

$$\left(1 + \frac{\varepsilon}{3(C+3)}\right) \cdot n < Kx$$

using the hypotheses

$$\begin{aligned} n &\leq (K/2)x \\ 0 &< C \\ 0 &< \varepsilon < 1. \end{aligned}$$

This is easily verified by noting that $1 + \frac{\varepsilon}{3(C+3)}$ is strictly less than 2, and so the product with n is strictly less than $2(K/2)x = Kx$. In this case, backchaining does not work, unless one thinks of replacing Kx by $2((K/2)x)$ in the goal inequality.

This example suggests that some form of forward search may be more fruitful: starting from the hypotheses, iteratively derive useful consequences, until the goal is obtained. Alternatively, we negate the conclusion and add it to the list of hypotheses, and then iteratively derive consequences until we obtain a contradiction. Our analysis shows that if we separate terms, we can in fact use $T_{add}[F]$ and $T_{mult}[F]$ independently to derive consequences, and that we only have to consider consequences in the language of $T_{comm}[F]$. This procedure is complete for the universal consequences of $T[F]$, and works equally well if we combine other local decision procedures for languages that are disjoint except for $=$ and \leq .

But what consequences shall we look for? Once again, our analysis shows us that a single well-chosen interpolant suffices: if we pick the right θ , $T_{add}[F]$ will be able to derive θ from our initial set of hypotheses, while $T_{mult}[F]$ will be able to prove $\neg\theta$. According to Proposition 6.3 and the discussion after it, we can assume, without loss of generality, that θ is a conjunction of disjunctions of literals of the form $x_i < ax_j$, $x_i \leq ax_j$, $x_i > ax_j$, $x_i \geq ax_j$, and comparisons between variables and constants in F . As a result, if the initial sequence of hypotheses can be refuted, there is a sequence $\theta_1, \theta_2, \dots, \theta_n$ of disjunctions of atomic formulas of the form above, such that $T_{add}[F]$ proves each formula θ_i from the initial set of hypotheses, and $T_{mult}[F]$ proves a contradiction from

these hypotheses and $\theta_1, \dots, \theta_n$. Of course, the situation is symmetric, so we can just as well switch $T_{add}[F]$ and $T_{mult}[F]$ in the previous assertion.

This reduces the task to that of deriving appropriate disjunctions θ_i of atomic formulas $x_i \leq ax_j$ from the initial hypotheses. The problem is that there are always infinitely many disjunctions that one can prove, and it may not be clear which ones are likely to be useful. For example, from $x + y \geq 0$, $T_{add}[F]$ can prove $x \geq a \vee y \geq -a$ for any a , and, a priori, any of these may be useful to $T_{mult}[F]$.

One solution is simply to ignore disjunctions. By Proposition 2.2, with some initial case splits we can reduce the problem of proving a universal formula to refuting a finite number of sets of formulas of the form $\Delta_{add} \cup \Delta_{mult} \cup \Delta_{comm}$, where

- Δ_{add} is a set of formulas of the form $x_i = t$, where t is a term in the language of $T_{add}[F]$;
- Δ_{mult} is a set of formulas of the form $x_i = t$, where t is a term in the language of $T_{mult}[F]$;
- Δ_{comm} is a set of formulas of the form $x_i < ax_j, x_i \leq ax_j, x_i > ax_j, x_i \geq ax_j$, or a comparison between a variable and a constant.

Definition 12.1. Let $\Delta = \Delta_{add} \cup \Delta_{mult} \cup \Delta_{comm}$ be as above. Say $T[F]$ *refutes* Δ *without case splits* if there is a sequence of atomic formulas $\theta_0, \dots, \theta_{2n}$ such that the following hold:

- for $m < 2n$, θ_m has the same form as the formulas in Δ_{comm} ;
- θ_{2n} is \perp ;
- for each $m < n$,

$$T_{add}[F] \cup \Delta_{add} \cup \Delta_{comm} \cup \{\theta_0, \dots, \theta_{2m-1}\} \vdash \theta_{2m};$$

- for each $m < n$,

$$T_{mult}[F] \cup \Delta_{mult} \cup \Delta_{comm} \cup \{\theta_0, \dots, \theta_{2m}\} \vdash \theta_{2m+1}.$$

In other words, $T[F]$ refutes Δ without case splits if $T_{add}[F]$ and $T_{mult}[F]$ can iteratively augment a database of derivable atomic formulas in the common language until a contradiction is reached. This is a proper restriction on the theories $T[F]$, which is to say, there are sets Δ that can be refuted by $T[F]$, but not without case splits. It takes some effort, though, to cook up an example. Here is one. Let

$$\Delta_{add} = \{x + y \geq 2, w + z \geq 2\}$$

From this, $T_{add}[F]$ proves $(x \geq 1 \vee y \geq 1) \wedge (w \geq 1 \vee z \geq 1)$. Let

$$\Delta_{mult} = \{ux^2 < ux, uy^2 < uy, uw^2 > uw, uz^2 > uz\}.$$

From this, $T_{mult}[F]$ proves $u > 0 \vee u < 0$, and hence $(x < 1 \wedge y < 1) \vee (w < 1 \wedge z < 1)$. As a result, $T[F]$ refutes $\Delta_{add} \cup \Delta_{mult}$. But one can check that there are no atomic consequences involving the common variables, x, y, z and w , that follow from either set. (Strictly speaking, our characterization of Δ has us using new variables to name the additive and multiplicative terms in Δ_{add} and Δ_{mult} , respectively, and then putting the comparisons in Δ_{comm} . But the net effect is the same.)

Situations like this are contrived, however, and we expect that focusing on atomic consequences will be effective in many ordinary situations. The following proposition provides some encouragement.

Proposition 12.2. *Let Δ be a set of atomic formulas in the language of $T_{add}[F]$. Let u and v be any two variables. Then there is a consequence, θ , of $T_{add}[F] \cup \Delta$ in the language of $T_{comm}[F]$, involving only u and v , that implies all the consequences of the form $u < av$, $u \leq av$, $v < au$, or $v \leq au$ that can be derived from $T_{add}[F] \cup \Delta$. In fact, θ can be expressed as a conjunction of at most two formulas of the form $u < av$, $u \leq av$, $u > av$, $u \geq av$, $v < 0$, $v \leq 0$, $v > 0$, or $v \geq 0$.*

Proof. Use a linear elimination procedure to eliminate all variables except for u and v from Δ . The result is a set of linear inequalities involving u and v , which implies every other relation between u and v that is derivable from $T_{add}[F] \cup \Delta$. (If a relation is not a consequence of the resulting set of linear inequalities, its negation is consistent with them, and hence with $T_{add}[F] \cup \Delta$.) This set of linear inequalities determines a convex subset of the cartesian plane. Considering extremal points, one can determine the minimal intersection of at most two half planes through the origin that includes this convex subset. \square

An efficient algorithm for determining the convex polygon determined by a sequence of half-planes can be found in [12, Section 4.2]. Keep in mind that there may be *no* nontrivial consequences of Δ , in which case we can take θ to be the empty conjunction, \top . Or Δ may be contradictory, in which case we can take θ to be \perp , or $v < 0 \wedge v > 0$. Furthermore, θ may not be strong enough to determine whether u and v are positive, negative, etc. In that case, as in the discussion after Proposition 6.3, determining whether one inequality is stronger than another can be confusing. For example, θ may be $u > 2v \wedge u > 3v$; in the absence of sign information, neither conjunct is stronger. If one adds the information $v > 0$, θ becomes $v > 0 \wedge u > 3v$.

On the multiplicative side, we have to assume we know the signs of the variables, and that F is closed under n th roots.

Proposition 12.3. *Let Δ be a set of atomic formulas in the language of $T_{mult}[F]$. Assume that for each variable x occurring in Δ , Δ contains either the formula $x > 0$ or the formula $x < 0$. Assume also that F is closed under n th roots of positive numbers for positive integers n . Let u and v be any two variables. Then there is a consequence, θ , of $T_{mult}[F] \cup \Delta$ in the language of $T_{comm}[F]$, involving only u and v , that implies all the consequences of the form*

$u < av$, $u \leq av$, $v < au$, or $v \leq au$ that can be derived from $T_{\text{mult}}[F] \cup \Delta$. In fact, θ can be expressed as a conjunction of at most two formulas of the form $u < av$, $u \leq av$, $u > av$, or $u \geq av$.

The good news is that the proof is even easier in this case.

Proof. Introduce a new variable w , and the equation $w = u/v$. Eliminate all variables except for w . The result is a set of inequalities of the form $w < a$, $w \leq a$, $w > a$, and $w \geq a$, of which we can choose the strongest and then replace w by u/v . \square

The requirement that we have sign information on the variables is generally needed to carry out the elimination procedure for $T_{\text{mult}}[F]$. We can always ensure that this information is present using case splits, though this can be computationally expensive. The requirement that F is closed under taking roots is also needed for the conclusion; for example, from $\{u > 0, u^2 > 2v^2\}$ we would like to conclude $u > \sqrt{2}v$. For practical purposes, however, we will suggest, in the next section, that one should choose \mathbb{Q} for F in an implementation, and avoid case splits. In that case, we can only hope for an approximation to Proposition 12.3. For example, when trying to put a multiplicative equation in pivot form, if we do not have sufficient sign information to determine the appropriate direction of an inequality, we can simply ignore this equation. And when required to take n th roots at the very end of the procedure, we can rely on crude approximations, such as $\sqrt[n]{a} > 1$ whenever $a > 1$. Once again, we expect that even with these concessions, the resulting procedure will be helpful in verifying commonplace inferences.

This strategy, then, will form the basis for the heuristic procedure that we will suggest in the next section. We leave open one interesting theoretical question, though: is it decidable whether a theory $T[F]$ can refute a set Δ without case splits? The proof of Theorem 5.2 shows that trying to refute the set Δ corresponding to $x^2 + 2x - 1 < 0$ leads to an infinite iteration, so the obvious search procedure is not guaranteed to terminate.

13 Towards a heuristic procedure

In this section, we discuss some possible avenues towards developing heuristic decision procedures, based on the analysis we have provided here. We are, of course, sensitive to the tremendous gap between neat decidability results and heuristic procedures that work well in practice. But we expect that the former can serve as a useful guide in the development of the latter, by clarifying the inherent possibilities and limitations of the method, and separating heuristic issues from theoretical ones. Of course, different heuristic approaches will have distinct advantages and disadvantages, and so different procedures can be expected to work better in different domains. We expect the type of algorithm we propose here to be fruitful for the kinds of examples discussed in Section 1.

Given a quantifier-free sequent in the language of $T[\mathbb{Q}]$, first, put all terms in normal form, as described in Section 8. This will make it possible to identify subterms that are provably equal. For that purpose, one can use the simpler normal form described at the end of Section 8.

Next, use new variables, recursively, to name additive and multiplicative subterms. These will form the sets Δ_{add} and Δ_{mult} . With these renamings, the original sequent will be equivalent to one in the language of $T_{comm}[\mathbb{Q}]$.

Convert the resulting sequent to a finite sequence of sets Δ_{comm} of inequalities $x < ay$, $x \leq ay$, $x > ay$, $x \geq ay$, to be refuted. For example, proving the sequent

$$x = y, w < z \Rightarrow u < v$$

amounts to refuting the set Δ_{comm} of formulas

$$\{x \geq y, x \leq y, w < z, v \leq u\}.$$

Note that the equality in the hypothesis is replaced by two inequalities. This seems to be a reasonable move, since with Δ_{add} and Δ_{mult} , x and y may name complex terms; we imagine that this procedure will be called after obvious simplifications and rewriting have been performed. Also note that the task of proving an equality $u = v$ splits into two tasks, namely, refuting $u > v$ and refuting $u < v$. Again, this seems reasonable, since we envision this procedure being called when direct methods for proving equalities have failed.

Now, try to refute each set Δ_{comm} , with the following iterative procedure. First, for each pair of variables x, y in Δ_{comm} , use $T_{add}[\mathbb{Q}] \cup \Delta_{comm}$ to derive new or stronger inequalities of the form $x < ay$, $x \leq ay$, $x > ay$, or $x \geq ay$, as well as comparisons between x and constants for each variable x . Add the new inequalities to Δ_{comm} , removing ones that are subsumed by the new information. Δ_{comm} can be represented as a table of comparisons for each pair $\{x, y\}$ (for each pair, at most two formulas need to be stored), as well as a table of comparisons with constants for each variable x . Even though the procedure implicit in Proposition 12.2 invokes a linear elimination procedure (see the discussion and references in Section 3), the work can be shared when cycling through all possible pairs. For example, to determine all inequalities obtainable from a set with n variables, eliminate the first variable, x , and recursively determine all the inequalities obtainable from the resulting set with n variables; then determine all the inequalities that can be obtained with x and one other variable. Furthermore, at least initially, for most pairs no information will be available at all, and so will be eliminated quickly. We expect that for the types of problems that arise in ordinary practice, the number of variables and named subterms will be small enough to make the procedure manageable. If not, heuristics can be used to focus attention on pairs that are likely to provide useful information.

Do the same with $T_{mult}[\mathbb{Q}] \cup \Delta_{mult}$. First, use the information in Δ_{mult} to determine the variables for which one has comparisons with 0. For a defining equation such as $u = x^2y^4$, the multiplicative procedure can infer $u \geq 0$ at the start, and add it to Δ_{comm} for possible use by the additive procedure. With limited sign information on the variables, let the procedure for $T_{mult}[\mathbb{Q}] \cup \Delta_{mult}$

do the best it can to eliminate variables. If it cannot make use of an inequality $x^k s < t$ to eliminate x because the sign of s is not known, simply ignore the inequality at this stage. It may become useful later on, if the sign of s becomes known.

Iterate the additive and multiplicative steps, until one of $\Delta_{add} \cup \Delta_{comm}$ or $\Delta_{mult} \cup \Delta_{comm}$ yields a contradiction. Of course, there is the question as to when to give up. One can certainly report failure when no new inequalities have been derived. But as noted at the end of Section 12, nonterminating iterations are possible; in that case, the procedure can simply give up after a fixed amount of time, or rely on the user to halt the procedure.

14 Extending the heuristic

There are many ways that one may extend the proposal in the previous section. These fall into general classes.

Improvements to the heuristic. There are likely to be better ways of searching for useful comparisons between terms. For example, one can have a list of “focus” formulas – initially, one wants to include the goal formula as a focus formula – and search for inequalities between subterms of those. Also, one does not need to search for comparisons between two variables unless information has been added to Δ_{comm} since the last such search that could potentially yield new information. Thus, a wise choice of data structures and representations of information in the database may yield significant improvements.

Extensions to stronger fragments of $T[\mathbb{Q}]$. The procedure we have described does not try to derive disjunctions, which requires potentially costly case splits. Are there situations in which it makes sense to introduce such splits? For example, it may be useful to split on the sign of a variable, $x \geq 0 \vee x < 0$; or to split on a comparison between two variables, $x \geq y \vee x < y$, where x and y name terms in the search.

Conservative extensions of $T[\mathbb{Q}]$. The functions which return n th roots, absolute value, minimums, and maximums can all be defined in $T[\mathbb{Q}]$, and it would be useful to extend the heuristic to languages that include these. But, as discussed at the end of Section 8, one has to either introduce case splits at the outset to simplify terms appropriately, or simplify a term like $\sqrt{x^2}$ to x when $x \geq 0$ is determined in the course of the search. What is the best way to handle such extensions?

Nonconservative extensions of $T[\mathbb{Q}]$, in the same language. An obvious shortcoming of $T[\mathbb{Q}]$ is that it fails to capture straightforward inferences that are easily obtained using distributivity. On the other hand, using distributivity to simplify an expression before calling a decision procedure for $T[\mathbb{Q}]$ can erase valuable information; for example, after simplification, $T[\mathbb{Q}]$ can no longer verify $(x + 1)^2 \geq 0$. A better strategy is to perform such simplifications as the search proceeds, when occasion seems to warrant it, perhaps retaining the factored versions as well.

As noted in Section 1, it is reasonable to claim that any validity that requires complex factoring falls outside the range of the “obvious,” and hence outside the scope of the problem we are concerned with here. But one would expect a good procedure to multiply through in at least some contexts, i.e. only use distributivity in the “left-to-right” direction to simplify expressions at hand. The question is how to work this in to the procedures described below in a principled way. It would also be nice to have a better theoretical framework to discuss provability with equalities “applied only in the left-to-right direction.”

Amalgamating other decision and heuristic procedures. A major advantage of the method described in Section 13 is that it can easily be scaled to allow other procedures to add facts to the common database. For example, one can easily make use of the equivalence $x < y \leftrightarrow f(x) < f(y)$ for a strictly monotone function f . One can similarly add procedures that make use of straightforward properties of transcendental functions like exp , ln , sin , cos , and so on.

Extending the overlap. Just as one might make use of limited forms of distributivity, one can add restricted uses of laws like $e^{x+y} = e^x e^y$, for the exponential function.

Handling subdomains, like \mathbb{Z} and \mathbb{Q} , and extended domains, like \mathbb{C} . For example, it is known that the linear theory of the reals with a predicate for the integers is decidable (see, for example, [28]). Handling mixed domains involving \mathbb{N} , \mathbb{Z} , \mathbb{Q} , \mathbb{R} , and/or \mathbb{C} is an important challenge for heuristic procedures.

15 Conclusions

In order to obtain useful methods for verifying inferences in nontrivial mathematical situations, undecidability and infeasibility should encourage one to search for novel ways of delimiting manageable, restricted classes of inferences that include the ones that come up in ordinary mathematical practice. We hope our study of inferences involving inequalities between real-valued expressions that can be verified without using distributivity is an interesting and fruitful investigation along these lines. We also feel that the paradigm of amalgamating decision or heuristic procedures when there is nontrivial overlap between the theories is an important one for automated reasoning.

However, we expect that similar investigations can be carried out in almost any mathematical domain. This yields both theoretical and practical challenges. On the theoretical side, for example, there are questions of decidability and complexity. On the practical side, there is always the question of how to implement proof searches that work well in practice. As a result, we feel that this type of research represents a promising interaction between theory and practice.

References

- [1] Krzysztof Apt. *Principles of Constraint Programming*. Cambridge University Press, Cambridge, 2003.
- [2] Jeremy Avigad, Kevin Donnelly, David Gray, and Paul Raff. A formally verified proof of the prime number theorem. To appear in *ACM Transactions on Computational Logic*.
- [3] Franz Baader, Silvio Ghilardi, and Cesare Tinelli. A new combination procedure for the word problem that generalizes fusion decidability results in modal logics. In David A. Basin and Michaël Rusinowitch, editors, *IJCAR '04*, pages 183–197. Springer-Verlag, Berlin, 2004.
- [4] Clark Wayne Barrett. *Checking Validity of Quantifier-free Formulas in Combinations of First-order Theories*. PhD thesis, Stanford University, 2002.
- [5] Saugata Basu. New results on quantifier elimination over real closed fields and applications to constraint databases. *Journal of the ACM*, 46:537–555, 1999.
- [6] Saugata Basu, Richard Pollack, and Marie-Françoise Roy. *Algorithms in Real Algebraic Geometry*. Springer-Verlag, Berlin, 2003.
- [7] Michael Beeson. Design principles of Mathpert: software to support education in algebra and calculus. In N. Kajler, editor, *Computer-Human Interaction in Symbolic Computation*, pages 89–115. Springer-Verlag, Berlin, 1998.
- [8] B. F. Caviness and J. R. Johnson, editors. *Quantifier Elimination and Cylindrical Algebraic Decomposition*. Springer-Verlag, Vienna, 1998.
- [9] C. C. Chang and H. Jerome Keisler. *Model theory*. North-Holland, Amsterdam, third edition, 1990.
- [10] George E. Collins. Quantifier elimination for real closed fields by cylindrical algebraic decomposition. In *Automata Theory and Formal Languages (Second GI Conf., Kaiserslautern, 1975)*, pages 134–183. Springer-Verlag, Berlin, 1975. Reprinted in [8].
- [11] S. Conchon and S. Krstić. Strategies for combining decision procedures. In P. Narendran and M. Rusinowitch, editors, *Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*, pages 537–553. Springer-Verlag, Berlin, 2003.
- [12] Mark de Berg, Marc van Kreveld, Mark Overmars, and Otfried Schwartzkopf. *Computational Geometry: Algorithms and Applications*. Second edition. Springer-Verlag, Berlin, 2000.

- [13] David Detlefs and Greg Nelson and James B. Saxe. Simplify: a theorem prover for program checking. *Journal of the ACM*, 52:365–473, 2005.
- [14] Andreas Dolzmann, Thomas Sturm, and Volker Weispfenning. Real quantifier elimination in practice. In B. H. Matzat, G.-M. Greuel, and G. Hiss, editors, *Algorithmic Algebra and Number Theory*, pages 221–248. Springer-Verlag, Berlin, 1998.
- [15] Silvio Ghilardi. Model-theoretic methods in combined constraint satisfiability. *Journal of Automated Reasoning*, 33:221–249, 2004.
- [16] John Harrison. *Introduction to Logic and Automated Theorem Proving*. In preparation.
- [17] Warren A. Hunt, Robert Bellarmine Krug, and J. Moore. Linear and non-linear arithmetic in ACL2. In Daniel Geist and Enrico Tronci, editors, *Correct Hardware Design and Verification Methods, Proceedings of CHARME 2003*, pages 319–333. Springer-Verlag, Berlin, 2003.
- [18] Predrag Janičić and Alan Bundy. A general setting for flexibly combining and augmenting decision procedures. *Journal of Automated Reasoning*, 28:257–305, 2002.
- [19] Rüdiger Loos and Volker Weispfenning. Applying linear quantifier elimination. *The Computer Journal*, 36:450-461, 1993.
- [20] Sean McLaughlin and John Harrison. A proof producing decision procedure for real arithmetic. In Robert Nieuwenhuis, editor, *Automated Deduction – CADE-20. 20th International Conference on Automated Deduction, Tallinn, Estonia, July 22-27, 2005, Proceedings*, pages 295–314, 2005.
- [21] Greg Nelson and Derek C. Oppen. Simplification by cooperating decision procedures. *ACM Transactions on Programming Languages and Systems*, 1:245–257, 1979.
- [22] Silvio Ranise, Christophe Ringeissen, and Duc-Khanh Tran. Nelson-Oppen, Shostak and the extended canonizer: A family picture with a newborn. In Zhiming Liu and Keijiro Araki, editors, *ICTAC*, pages 372–386. Springer-Verlag, Berlin, 2004.
- [23] Alfred Tarski. *A Decision Procedure for Elementary Algebra and Geometry*. Prepared for publication by J. C. C. McKinsey. University of California Press, second edition edition, 1951. Reprinted in [8].
- [24] Cesare Tinelli and Mehdi T. Harandi. A new correctness proof of the Nelson-Oppen combination procedure. In Franz Baader and Klaus U. Schulz, editors, *Frontiers of Combining Systems (FroCos)*, pages 103–119. Kluwer Academic Publishers, 1996.

- [25] A. Tiwari. Abstraction based theorem proving: An example from the theory of reals. In C. Tinelli and S. Ranise, editors, *Proceedings of the CADE-19 Workshop on Pragmatics of Decision Procedures in Automated Deduction, PDPAR 2003*, pages 40–52. INRIA, Nancy, 2003.
- [26] Lou van den Dries. *Tame Topology and O-minimal Structures*. Cambridge University Press, Cambridge, 1998.
- [27] Volker Weispfenning. The complexity of linear problems in fields. *Journal of Symbolic Computation*, 5:3-27, 1988.
- [28] Volker Weispfenning. Mixed real-integer linear quantifier elimination. In S. Dooley, editor, *Proceedings of the 1999 International Symposium on Symbolic and Algebraic Computation, ISSAC '99*, pages 129–136. ACM Press, 1999.