

# An Introduction to Mechanized Reasoning\*

Manfred Kerber<sup>†</sup>    Christoph Lange<sup>‡</sup>    Colin Rowat<sup>§</sup>

August 12, 2016

## Abstract

Mechanized reasoning uses computers to verify proofs and to help discover new theorems. Computer scientists have applied mechanized reasoning to economic problems but – to date – this work has not yet been properly presented in economics journals. We introduce mechanized reasoning to economists in three ways. First, we introduce mechanized reasoning in general, describing both the techniques and their successful applications. Second, we explain how mechanized reasoning has been applied to economic problems, concentrating on the two domains that have attracted the most attention: social choice theory and auction theory. Finally, we present a detailed example of mechanized reasoning in practice by means of a proof of Vickrey’s familiar theorem on second-price auctions.

*Key words:* mechanized reasoning, formal methods, social choice theory, auction theory

*JEL classification numbers:* B41; C63; C88; D44

---

\*We are grateful to Makarius Wenzel for help refining our code, to Marco Caminati for research assistance, to Peter Cramton, Paul Klemperer, Peter Postl, Indra Ray, Rajiv Sarin, Arunava Sen and Ron Smith for comments, and to the EPSRC for funding (EP/J007498/1). Rowat thanks Birkbeck for its hospitality. The presentation of the formal proof of Vickrey’s theorem is based on Kerber, Lange, and Rowat (2014). Finally, we are grateful to two anonymous referees and the co-editor for working with us to improve this paper.

<sup>†</sup>School of Computer Science, University of Birmingham, UK

<sup>‡</sup>Fraunhofer IAIS and University of Bonn, Germany

<sup>§</sup>Department of Economics, University of Birmingham, Edgbaston B15 2TT, UK, c.rowat@bham.ac.uk, +44 121 414 3754. Corresponding author

# 1 Introduction

Mechanized reasoners automate logical operations, extending the scope of mechanical support for human reasoning beyond numerical computations (such as those carried out by a calculator) and symbolic calculations (such as those carried out by a computer algebra system). Such reasoners may be used to formulate new conjectures, check existing proofs, formally encode knowledge, or even prove new results. The idea of mechanizing reasoning dates back at least to Leibniz (1686), who envisaged a machine which could compute the validity of arguments and the truth of mathematical statements. The development of formal logic from 1850 to 1930, the advent of the computer, and the inception of *artificial intelligence* (AI) as a research field at the Dartmouth Workshop in 1956 all paved the way for the first mechanized reasoners in the 1950s and 1960s.<sup>1</sup>

Since then, mechanized reasoning has been both less and more successful than anticipated. In pure maths, mechanized reasoning has helped prove only a few high-profile theorems. Perhaps surprisingly – although consistent with the greater success of applied AI over ‘pure’ AI – mechanized reasoning and formal methods<sup>2</sup> have enjoyed greater success in industrial applications, as applied to both hardware and software design. In the past decade or so, computer scientists have also begun to apply formal methods to economics.

A central inspiration for this recent work are Geanakoplos’ three brief proofs of Arrow’s impossibility theorem (Geanakoplos, 2005).<sup>3</sup> Initially, Nipkow (2009), Wiedijk (2007), and Wiedijk (2009) used theorem provers to encode and verify two of Geanakoplos’ proofs. A subsequent generation of work, drawing on the inductive proof of Arrow’s theorem in Suzumura (2000), used formal methods to discover new theorems. Tang and Lin (2009) introduced a hybrid technique, using computational exhaustion to show that Arrow holds on a small base case

---

<sup>1</sup>Perhaps unsurprisingly, Gardner was ahead of his time in mechanized reasoning as well: four years before his regular columns with *Scientific American* began, his first article for them included a template allowing readers to make their own mechanized reasoners – out of paper.

<sup>2</sup>The term *formal methods* is used here to denote approaches to establishing the correctness of mathematical statements to a precision that they can be meticulously checked by a computer. Rather than being seen as distinct from other mathematical methods, researchers in the area see them as the next step in mathematics’ march towards greater precision and rigor (Wiedijk, 2008). Consider: “A Mathematical proof is rigorous when it is (or could be) written out in the first-order predicate language  $L(\epsilon)$  as a sequence of inferences from the axioms ZFC” (MacLane, 1986). The advantages of taking this next step with computers include: a computer system is never tired or intimidated by authority, it does not make hidden assumptions, and can easily be rerun. A pioneer of mechanized reasoning – who saw himself building on Bourbaki’s formalism – referred to computers as “slaves which are such persistent plodders” (Wang, 1960).

<sup>3</sup>All three use Barberà’s replacement of Arrow’s *decisive voter* with a *pivotal voter* (Barberà, 1980). Barberà (1983) also used this approach to find a direct proof of the Gibbard-Satterthwaite theorem.

of two agents and three alternatives, and then manual induction to extend that to the full theorem. By inspecting the results of the computational step, they were able to discover a new theorem subsuming Arrow's. Tang and Lin (2011a) used this approach – exhaustively generating and evaluating base cases, and then using a manual induction proof to generalize the results – to establish uniqueness conditions for pure strategy Nash equilibrium payoffs in two player static games; they published manual proofs of two of the most significant theorems discovered this way in Tang and Lin (2011b). Geist and Endriss (2011) used the approach to generate 84 impossibility theorems in the 'ranking sets of objects' problem (Barberà, Bossert, and Pattanaik, 2004).

To date, the economics literature remains almost untouched by research applying mechanized reasoning to economic problems.<sup>4</sup> The one exception that we are aware of is Tang and Lin (2011b), whose two theorems were discovered computationally, but proved manually.<sup>5</sup> As it is our view that these tools will become increasingly capable, this paper aims to introduce economists to mechanized reasoning.<sup>6</sup> It does so by means of three analytical lenses, each with narrower scope but greater magnification than its predecessor.

First, Section 2 presents an overview of mechanized reasoning in general. We do so by setting out a classificatory scheme, with the caveat that it should not be seen as implying a partition on the field: interesting research will straddle boundaries, perhaps even forcing them to be redefined.<sup>7</sup>

Second, Section 3 surveys the emerging literature applying mechanized reasoning to economics. We structure this survey primarily according to the problem domain within economics, referring only secondarily to our classificatory scheme. We do this to focus on the economic insights – primarily within social choice and auction theory – made possible by these techniques, rather than on the techniques *per se*.

Finally, to make this introduction more concrete, Section 4 provides an example of what mechanized reasoning looks like in practice, presenting a blueprint of a mechanized proof of Vickrey's theorem on second-price auctions. We present

---

<sup>4</sup>A recent symposium on economics and computer science, involving central figures at the interface between the disciplines, made no mention of mechanized reasoning (q.v. Blume et al., 2015).

<sup>5</sup>The process by which the theorems were discovered is described in Tang and Lin (2011a); Tang and Lin (2011b) itself is all but silent on its mechanized origins.

<sup>6</sup>For more general introductions, see Wiedijk (2008) and Avigad and Harrison (2014). Harrison (2007) introduces mechanized reasoning alongside computer algebra, presenting something of a unified view.

<sup>7</sup>For example, we shall see that mechanized theorem discovery is usually associated with inductive reasoning. However – in economic examples – the most fruitful examples of theorem discovery (Tang and Lin, 2009; Tang and Lin, 2011a; Tang and Lin, 2011b; Geist and Endriss, 2011) have combined very simple deductive reasoning systems with human intelligence.

such an established theorem to focus attention on its implementation.

Section 5 concludes, and suggests some possible next steps for mechanized reasoning in economics.

## 2 Mechanized reasoning

Our overview of mechanized reasoning distinguishes between deductive and inductive systems. While the distinction has been recognized at least since Aristotle, deductive reasoning – which allows reliable inference of unknown facts from established facts – has been in the focus of the mechanized reasoning community. Inductive reasoning also generalizes from individual cases, but does not restrict itself to reliable inferences; the cost of this additional freedom is that its conjectures must then be tested.

### 2.1 Deductive reasoning

Historically, deductive reasoning systems were among the first AI systems, dating back to the 1950s. While the origins of deductive reasoning date to at least Aristotle, modern advances in this area built on the work of logicians in the second half of the 19th century and the start of the 20th (e.g. Whitehead and Russell, 1910). At the Dartmouth Workshop in 1956, Newell and Simon introduced the Logic Theorist, an automated reasoner which re-proved 38 of the 52 theorems in Whitehead and Russell’s *Principia Mathematica* (Whitehead and Russell, 1910).<sup>8</sup>

Abstractly, a deductive reasoner implements a *logic* – which is comprised of a *syntax* defining well-formed formulae and a *semantics* assigning meaning to formulae – and a *calculus* for deriving formulae (called theorems) from formulae (called premises or axioms). Historically, subfields of mechanized reasoning have been defined by choice of logic, calculus and problem domain. This section provides a classificatory scheme based, first, on the choice of calculus. Following the choice of calculus, a logic is chosen to balance expressiveness and tractability. Finally, the problem domain itself will dictate some of the specialized features of a mechanized reasoner.

When a mechanized reasoner applies the calculus’ permissible operations to the axioms to obtain new, syntactically-correct formulae it does not make use of the semantics: the semantics, or ascribed meanings, yield models that may assist human intuition, but which are not necessary to the formal process of reasoning

---

<sup>8</sup>According to McCorduck (2004), Russell himself “responded with delight” when shown the Logic Theorist’s proof of the isosceles triangle theorem, whose proof was more elegant than their manual one.

itself.<sup>9</sup> Crucially, mechanized reasoning involves manipulating symbols.<sup>10</sup>

Thus, mechanized deductive reasoning since the Logic Theorist has seen reasoning as a search task for a syntactically well-defined goal.<sup>11</sup> Further, as the spaces through which search occurred was potentially large, successful reasoning would use *heuristics* to avoid unprofitable sequences of operations. From this point of view, mechanized reasoning operates as chess computers do.<sup>12</sup> For a chess computer, the premises' intended semantic interpretations are the board, its pieces and their positions; the calculus specified permissible moves. A chess computer could then test manually discovered solutions to chess puzzles by verifying that each move satisfies the requirements of its calculus, with the final operation yielding the goal-formula. More ambitiously, and interestingly, chess programs discover solutions (e.g. sequences of winning moves) by searching through permissible operations, with the benefit of heuristics (e.g. regarding relative values of pieces).

A set of premises and a formula may be related in two different ways. First, the *semantic consequence relation* describes situations in which the formula *follows from* the premises: if the symbols in the premises are interpreted in such a way that the formulae in the premises are all true, then the formula is also true when the symbols in it are interpreted in the same way. Second, the *syntactic derivability relation* describes situations in which the formula *can be derived from* the premises: it is possible to generate the formula from the premises by applying a fixed set of so-called calculus rules. (An example of such a rule is *modus ponens*: From  $A$  and  $A \rightarrow B$  it is possible to derive  $B$ , where  $A$  and  $B$  may match any formal expression). A proof that applies such rules, without any appeals to intuition or to the reader filling in steps on her own, is called a *formal proof* of the formula using the premises.

A calculus is called *sound* if only formulae can be derived from the premises that actually follow from them. Deductive reasoning is sound; inductive reasoning, considered below, is not.

A calculus is *complete* if it allows derivation of any formula that follows from the set of premises. A calculus is *decidable* if, for any set of premises and any formula, there is a procedure that either derives the formula from the premises or

---

<sup>9</sup>Beginning with Euclid's efforts to axiomatize geometry, logicians have produced syntactical descriptions that make semantic references obsolete: Hilbert allegedly said that we would still have an axiomatization of geometry if we replaced the words 'point', 'line', and 'plane' by 'beer mug', 'bench', and 'table' (Hoffmann, 2013, p.6).

<sup>10</sup>That this was an insight at one point may be inferred from Turing's famous explanation that, "computing is normally done by writing certain symbols on paper" (Turing, 1936)

<sup>11</sup>As noted by Harrison (2007), specialist provers have also been developed for particular problems for which more structured approaches than general search are appropriate.

<sup>12</sup>Indeed, Newell's collaboration with Simon began after the latter became aware of the former's work on a chess machine.

proves that no such derivation exists; a calculus is *semi-decidable* if a procedure exists that derives the formula from premises, whenever the formula follows from them (but may not terminate if it does not).

Decidability typically depends on the expressiveness of the logic used: more expressive logics model a richer set of concepts, but are generally harder to manipulate. While ambitious exercises in mechanized reasoning often begin by specifying a suitably tailored logic<sup>13</sup>, we largely restrict our attention to some of the best known *classical logics*.<sup>14</sup>

**Propositional (Boolean) logic:** *Propositional* or *Boolean logic*, the simplest classical logic, only uses propositional variables – which are either true or false – and *connectives* such as  $\wedge$  (and),  $\vee$  (or),  $\neg$  (not), and  $\rightarrow$  (implies). An example of a propositional formula is

$$first\_bidder\_bids\_highest \wedge second\_bidder\_bids\_lowest.$$

Propositional logic can only make concrete, finite statements, but has a sound, complete and decidable calculus.

An advantage of this decidability is that it may allow *push-button* technology, which does not require specialist knowledge in order to use. Once a problem is adequately represented a corresponding system solves the problem fully automatically.

**First-order logic:** *First-order logic* (FOL) is more expressive. First, it can speak about objects (e.g. “bidder  $b_1$ ”) and their properties (e.g. “bidder  $b_1$  wins auction”,  $bidder(b_1) \wedge wins(b_1)$ ). Second,  $\exists$  and  $\forall$  allow quantification over objects. For example, “every losing bidder pays nothing” may be expressed as

$$\forall i. bidder(i) \rightarrow (\neg wins(i) \rightarrow pay(i) = 0). \quad (1)$$

Expressions like *wins* are called *predicates*, Boolean functions which – when applied to their arguments – evaluate to either true or false. Gödel’s completeness theorem proves that FOL has a sound and complete calculus, but FOL has only semi-decidable calculi. Furthermore, FOL is not expressive enough to express the finitude or (per negation) infinity of the non-empty sets of objects.<sup>15</sup>

*Many-sorted FOL* uses *sorts* to extend first-order logic, not to add to its expressiveness, but to allow more concise representations, and – therefore – more

---

<sup>13</sup>See, for example, the *judgement aggregation logic* (JAL) of Ågotnes, Hoek, and Wooldridge (2011).

<sup>14</sup>The 17 volumes in the second edition of Gabbay and Guenther (2001/2014) make clear that the classical logics are a small subset of all logics.

<sup>15</sup>Thus, FOL could not express that only finitely many bidders participate in an auction.

efficient proving. Sorts restrict the instantiation of variables to expressions of a certain sort. For instance, sorts allow us to specify that variable  $i$  is a bidder, and variable  $x$  a good. Formula (1) is then more precisely stated as:

$$\forall i_{bidder} . \neg wins(i) \rightarrow pay(i) = 0. \quad (2)$$

$i$  (with the sort *bidder* mentioned only at the first occurrence) can be instantiated now by terms of sort *bidder*, but not by those of sort *good*, thus reducing the search space for a proof. Sorted formulae can be translated to unsorted formulae by converting the sorts to unary predicates (which take a single argument).

**Higher-order logic:** *Higher-order logic* (HOL) enriches the expressiveness of FOL by extending quantification to predicates and functions. It also allows predicates and functions to take certain<sup>16</sup> other predicates and functions as arguments. For example, bids,  $b$ , are both a function from bidders to prices and an argument (along with  $N, v$  and  $A$ ) in the predicate

$$equilibrium\_weakly\_dominant\_strategy\ N\ v\ b\ A.$$

Against this, HOL's calculi are not decidable, and are – by Gödel's incompleteness theorem – incomplete.

Two common ways in which the classical logics (in particular, FOL) are augmented are, first, by the addition of set theoretical axioms and, second, by the addition of modal operators. The first allows the approximation of higher order logic while maintaining advantages of first order logic; the second allows logic to be applied to modalities, such as knowledge, belief, or time.

Set theoretical axioms allow the definition of new symbols and operations on both predicates (e.g.  $\in$  and  $\subseteq$ ) and functions (e.g.  $\cup$ ,  $\cap$  and  $\emptyset$ ).<sup>17</sup> They also allow the specification of properties of sets (e.g.  $a \notin X$ ). Adding set theoretical axioms to FOL allows it to weakly simulate HOL: functions can be expressed as relations over  $X \times X$  that are left-total and right-unique; predicates are expressed as sets. While HOL is still more expressive than FOL augmented by set theory (e.g., FOL cannot express inductive arguments), HOL's incompleteness means that there are true statements that can be expressed in HOL but which may not have finite proofs. As FOL augmented by set theory uses FOL, it remains complete by using FOL's complete calculus.

---

<sup>16</sup>Unrestricted formula building leads to antinomies as discovered by Russell. The introduction of types imposes a hierarchy on logical objects, including predicates. This disables circular constructs such as  $X(Y) := \neg Y(Y)$ , which – when  $Y$  is instantiated with  $X$  – produces the set of all  $X$  for which  $X \notin X$ , Russell's famous antinomy.

<sup>17</sup>Constants such as  $\emptyset$  are considered as a special case of functions, nullary functions – functions that do not take any argument.

Modal operators – such as ‘next’ and ‘until’ – allow the consideration of *modes* (or *states* in economic parlance). *Linear temporal logic* (LTL) is a popular simple modal logic, modelling states in a linear fashion, thus excluding the consideration of multiple possible future states. Kamp’s theorem established the equivalence of LTL with a first-order logic. Another first-order approach to modelling states is the *situation calculus* (McCarthy and Hayes, 1969), which allows expression of states and the temporal development of systems in first-order logic by representing the state as an extra argument of the formulae (e.g., that agent  $i$  has £10 in state  $s_0$  can be expressed as  $has(i, 10, s_0)$ ). By referring to the state absolutely, rather than in relation to other states, the problem can be expressed in standard FOL without recourse to specialized modal relations.

Our final level of distinction is the domain of the problem; this level will allow us to present concrete examples of the preceding. Table 1 depicts these dimensions within deductive reasoning systems.

	decidable	undecidable
logic	SAT, CSP; description logic	ITP, ATP
computer system	model checking	program verification

Table 1: Mechanized reasoning using deductive logics

**Decidable logic:** In Table 1, the *decidable logic* cell refers to decidable calculi as applied to logical problems.

*Boolean satisfiability problems (SAT)* are among the simplest canonical problems in propositional logic. They specify a (finite) set of statements about a (finite) set of propositional variables, and ask whether there exists an assignment of values (i.e. true and false) to each of those variables that simultaneously satisfies all of the statements.

In SAT problems, clauses of Boolean variables are typically expressed in *conjunctive normal form*, conjunctions ( $\wedge$ ) of disjunctions ( $\vee$ ) such as

$$(\neg p \vee q) \wedge (p \vee \neg q); \quad (3)$$

where  $p$  and  $q$  are Boolean variables, evaluating either to true or false.<sup>18</sup> Revisiting the example that in auctions the non-winning player pays nothing, equation (1) can be translated for a finite number of bidders (here, three) to a propositional logic formula,

$$(wins1 \vee \neg pays1) \wedge (wins2 \vee \neg pays2) \wedge (wins3 \vee \neg pays3); \quad (4)$$

<sup>18</sup>The sentence given here is logically equivalent to  $p \equiv q$ , an equivalence exploited by Tang and Lin (2011a) in their search for uniqueness conditions in bimatrix games.



stating for each of the three players separately that they win or pay nothing.

Any formula in propositional logic can be expressed in this form, as can any formula in first-order logic when the domain is restricted to a concrete finite domain (such as three bidders in an auction). A SAT solver is used to try to assign the variables such that all of the clauses are true. For instance, assigning *wins1* and *pays1* to *true* and the other predicates to *false* shows that the single formula (4) is satisfiable.

SAT problems are  $\mathcal{NP}$ -hard (Karp, 1972), requiring – in the worst case – trial of every possible input. Thus, while the logic and calculi involved are simple, SAT problems may not be computable in practice except in small cases. However, techniques have been developed so that SAT solvers are able to solve typical cases very quickly. One application area of SAT solvers are model checkers, as described below.

*Constraint satisfaction problems (CSP)* are triples,  $\langle V, D, C \rangle$ , where  $V$  is a set of variables,  $D$  their domain, and  $C$  the constraint set. In CSPs, the variables may take on more values than in Boolean satisfiability's binary assignments. For example, an *hours* variable might take one of twelve values. While apparently richer, CSPs can be reduced to SATs by suitable definition of additional auxiliary variables.<sup>19</sup>

The third example of decidable calculi applied to logical problems that we consider are *description logics*. These are central to automated reasoning about concept hierarchies in classification (or ontological) tasks. One of their most important applications is to the *semantic web*, which allows computers to extract semantic information from web pages. As a simple example, semantically enabled web searches could recognize that  $x^2 + y^2 = z^2$  and  $a = \sqrt{c^2 - b^2}$  were both statements of Pythagoras' theorem.<sup>20</sup>

**Model checking:** Model checking (Clarke, Emerson, and Sistla, 1986; Clarke, Grumberg, and Long, 1994) builds finite models to describe computer hardware systems or simple software systems and then tests their properties. Typical questions include whether certain states of the system can be reached, or whether information is flowing properly through a circuit design.

Such models are typically expressed as *finite automata*. A finite automaton can model either a finite system or an infinite system if abstraction allows the infinite state space to be simplified to a finite one.<sup>21</sup> Then the model is systemat-

---

<sup>19</sup>See Bordeaux, Hamadi, and Zhang (2006) for a comparison of SAT and constraint programming.

<sup>20</sup>See Lange (2013) for a more in-depth discussion of applications of semantic web technology to mathematics.

<sup>21</sup>For example, in proofs involving real numbers, it may suffice to reduce an infinite number of possible values – which cannot be handled by a decidable calculus – to a trinary partition defined

ically checked for desired properties, e.g. by using SAT solvers. Viewing digital computer chips as a set of Boolean statements allows them to be modeled as *decidable computer systems* allowing, in turn, SAT solvers to automatically verify their properties. Since the mid-1990s, Intel has used formal methods to formally prove properties like ‘this chip implements the IEEE division standard’ following an embarrassing and costly recall of a Pentium chip that was discovered not to properly implement IEEE floating point division (Harrison, 2006). No further such problems have been reported since then.<sup>22</sup>

**Undecidable logic:** The upper right cell in Table 1 refers to the application of undecidable calculi to logical problems. The two types of mechanized reasoning mentioned here, *interactive theorem proving (ITP)* and *automated theorem proving (ATP)* have traditionally been equated with theorem proving, but seen as distinct, with the former involving more steering from a human user than the latter. Stereotypically, an ITP system could check an existing proof, while an ATP system could suggest steps in a proof or, in some cases, a whole proof. In practice, the distinction between the two has decreased, with ITP systems implementing ATP procedures.<sup>23</sup>

The traditional identification of theorem proving with work in these areas owes partly to some high profile successes in pure mathematics, the focus of the most hope in mechanized reasoning’s early days. The earliest major success was – as might be expected in an emerging field – not even a clear example of mechanized reasoning: in the 1970s, computers were used to carry out the exhaustive computations required to prove the four-color map theorem (q.v. Appel and Haken, 1977; Appel, Haken, and Koch, 1977). Here, the computers were used to perform simple (algebraic) calculations, rather than to (logically) ‘reason’. More recently, mechanized proof checkers have confirmed these results formally (q.v. Gonthier, 2008).<sup>24</sup>

The first major mathematical result to be established by mechanized reasoning – rather than ‘mere’ calculation – was Robbins’ conjecture that two bases for Boolean algebras are equivalent. While appearing to be a beguilingly simple problem, it remained unresolved for 60 years, becoming a favourite of Tarski, who set it as an open problem (q.v. Henkin, Monk, and Tarski, 1971, p. 245). One of

---

by  $>$ ,  $<$  and  $=$ . See Burch et al. (1990) for an application to large, complex microprocessor circuits.

<sup>22</sup>With chip design becoming more and more sophisticated, the reasoning in the verification needed to become also more sophisticated. Thus, HOL theorem provers such as HOL-Light are now also used for hardware verification.

<sup>23</sup>Harrison (2007) noted that ITP may be preferred to ATP, as – in working more closely alongside human reasoning – it may be better at developing human understanding.

<sup>24</sup>Gonthier’s team has now also formally checked the Feit-Thompson Odd Order Theorem (Gonthier et al., 2013).

the complicating factors of the conjecture was that the only known example of a Robbins algebra was also a Boolean algebra, reducing the evidence base that mathematicians could use to form intuitions about the problem. Nonetheless, in the late 1990s, McCune (1997) was able to pose the problem in a way that allowed EQP, an automated theorem prover related to his well-known Otter prover, to generate – not just check – a 17-step proof, later reduced to eight steps (McCune, 1997).<sup>25</sup>

Perhaps the highest profile success of mechanized reasoning in pure mathematics is the solution to Kepler’s conjecture that there is no denser packing of spheres in  $\mathbb{R}^3$  than the face-centred cubic. Hales’ original proof was 120 pages long (excluding computer code that exceeded 500MB), requiring a team of 12 referees five years to become “99% certain” that it was correct. Unsatisfied with this standard, Hales founded *Project Flyspeck* to establish a fully formal proof of the conjecture (Hales, 2012). In August 2014, the project was completed (Hales et al., 2015), close to Hales’s original estimate of 20 person-years (Avigad and Harrison, 2014).

More mundanely, ITP has been used to translate existing human proofs into formal proofs that are sufficiently detailed that a computer can mechanically verify them: as of January 2016, 91 of the ‘top 100’ mathematical theorems on a list maintained by Wiedijk (2014) had been formalized.<sup>26</sup> While most of these are considerably less spectacular than the examples cited above – in which theorem provers have been used to help convince mathematicians as to the validity of major, new results – the gradual accretion of small proof libraries builds a foundation for applying ATPs more widely.

The distinction between high-profile, major theorems and lower-profile bodies of theory has been suggested as a reason that ATP has yet to fulfil its early hopes: Buchberger (2006) noted that human mathematicians typically do not try to prove isolated theorems but explore a whole theory, thereby building up valuable intuition which helps them in proving related theorems. Additionally, Newell (1981) stated that standard theorem proving techniques – while often highly efficient – do not make use of advanced human approaches (as described in Pólya’s books) such as simplifying a problem to one they can solve; applying the simplified solution to the original problem may still be very hard, but the intuition gained by solving the simplified problem may help solve the original problem.<sup>27</sup>

---

<sup>25</sup>Dahn (1998) manually reworked EQP’s proof to provide a more human-readable proof.

<sup>26</sup>Exceptions include Fermat’s last theorem.

<sup>27</sup>Conversely, Dick (2011) observed that the ‘resolution’ inference rule (Robinson, 1965), central to mechanized reasoning, “was not based on any known human practice and was in fact difficult and counterintuitive for humans to understand”. Indeed, reviewing mechanized reasoning since resolution, Robinson lamented that it may have harmed mechanized reasoning by contributing to a parting of ways between human mathematicians and mechanized reasoners (Dick, 2015).

**Program verification** Table 1’s lower right cell corresponds to software engineering’s *program verification*, reasoning about software systems. This can be highly complex in the case of complex programs. Within program verification, traditional proof approaches have sought to prove that the software correctly implements properties specified in the design brief. As such proofs are very costly, full correctness proofs that seek to verify all desired properties of the code, are done only for ‘mission critical’ systems (D’Silva, Kroening, and Weissenbacher, 2008).

Some well known examples of program verification have come from transport and finance: in code controlling automated commuter rail systems, theorems that no two trains occupy the same location at the same time have been proved; within financial transactions software, theorems that transactions do not create or destroy value, but merely transfer it, have also been proved (Woodcock et al., 2009). More recently, a compiler for the C programming language has been formally verified (Boldo et al., 2013). These techniques are becoming more mainstream: in 2013, Facebook acquired Monoidics, a start-up firm applying theorem proving to software code analysis; in 2015, another start-up, Aesthetic Integration beat 600 competitors to win first prize in UBS’ Future of Finance Challenge for its ability to automatically prove failure or compliance in financial algorithms.<sup>28</sup>

Historically, program verification has been conducted as a *post mortem*: given existing code, program verification determines whether or not it is correct. More recently, *code extraction* techniques have been developed to generate code that provably implements the desired properties.

## 2.2 Inductive reasoning

As noted above, both inductive and deductive reasoning date back at least to Aristotle, but the former is not sound, while the latter has been the focus of the mechanized reasoning community. The distinction between the two – as well as the utility of each – was expressed by Pólya (1954, p. vi), who referred to deductive reasoning as *demonstrative reasoning*, and inductive reasoning as *plausible reasoning*:

We secure our mathematical knowledge by *demonstrative reasoning*, but we support our conjectures by *plausible reasoning* . . . Demonstrative reasoning is safe, beyond controversy, and final. Plausible reasoning is hazardous, controversial, and provisional. . . .

---

<sup>28</sup>Their entry formally defined a UBS ‘dark pool’ and a set of SEC regulations which the SEC had found the dark pool in breach of. Aesthetic Integration was able not only to verify the dark pool failure found by the SEC, but discovered that its order prioritization failed to satisfy transitivity (Ignatovich and Passmore, 2015).

In strict reasoning the principal thing is to distinguish a proof from a guess, a valid demonstration from an invalid attempt. In plausible reasoning the principal thing is to distinguish a guess from a guess, a more reasonable guess from a less reasonable guess. . . . [plausible reasoning] is the kind of reasoning on which [a mathematician's] creative work will depend.

Inductive systems seek to derive general statements based on a finite number of statements (e.g. if  $A_1$  is true, and  $A_2$  is true, and so on up to  $A_N$  for some finite  $N$ , then  $A_n$  is true for all natural numbers  $n$ ).<sup>29</sup> This sort of reasoning is immediately familiar to us when we reflect on how we form conjectures: we expect the sun to rise tomorrow without any understanding of astrophysics; this expectation, though, may lead to the formation of conjectures about astrophysics. However compelling the weight of evidence, inductive reasoning is not sound – as may be demonstrated by single counterexamples. In number theory, Euler's attempted generalization of Fermat's last theorem remained open for two centuries until a computer found a counterexample.<sup>30</sup> In game theory, Neumann and Morgenstern conjectured that stable sets ('solutions' in their parlance) always existed; it took almost a quarter-century for counterexamples to be found (Lucas, 1968).

Inductive reasoning may be used for *theorem discovery*, whereby regularities in observed data are used to form conjectures to test.<sup>31</sup>

Mechanized inductive reasoning dates back to two systems built in the 1970s and 1980s to discover new conjectures, AM (Automated Mathematician) (Lenat, 1976) and Eurisko (Lenat, 1983). These were able to detect conjectures such as the unique prime factorization theorem and Goldbach's conjecture.<sup>32</sup> The systems use certain measures of interestingness for concepts. For instance, concepts that are always true or always false are not interesting. However, if a concept is true for a significant proportion of examples (such as divisibility by only 1 and the number

---

<sup>29</sup>Inductive reasoning is distinct from mathematical induction, which involves proving  $A_0$  and that  $A_{n+1}$  is true given  $A_n$ . Mathematical induction is a sound *deductive* method.

<sup>30</sup>Euler's conjecture states: let  $n$  and  $k$  be integers greater than one, and let  $a_1, \dots, a_n$  and  $b$  be non-zero integers; then  $(\sum_{i=1}^n a_i^k = b^k) \Rightarrow (n \geq k)$ . The first known counterexample, found by computer, is  $27^5 + 84^5 + 110^5 + 133^5 = 144^5$  (Lander and Parkin, 1966).

<sup>31</sup>One of the most dynamic subfields of AI currently is *machine learning*. Some definitions are agnostic as to how the machines learn – e.g. whether deductively or inductively – while, perhaps more typically, others link machine learning more closely to inductive reasoning. Some of the highest profile applications of machine learning are statistical, positing rules that fit the existing data well, rather than perfectly.

<sup>32</sup>The prime factorization theorem states that any positive integer has a unique decomposition as the product of primes. Goldbach's conjecture states that every even integer beyond two can be expressed as the sum of two primes.

itself) then this is considered as an interesting concept ('primality' for divisibility by only 1 and the number itself).<sup>33</sup>

Lenat's work was continued by Colton in the HR (Hardy-Ramanujan) system (Colton, Bundy, and Walsh, 1999), where more advanced measures for interestingness were developed. For instance,

The novelty measure of a concept calculates how many times the categorisation produced by the concept has been seen. For example, square numbers categorise integers into two sets:  $\{1, 4, 9, \dots\}$  and  $\{2, 3, 5, \dots\}$ . If this categorisation had been seen often, square numbers would score poorly for novelty, and vice-versa. (Colton, Bundy, and Walsh, 2000).

Another important advance in Colton's work is that the HR system weeds out simple conjectures, namely those that can be easily verified or falsified by automated theorem provers.<sup>34</sup> One of the successes of HR was that it invented the concept of 'integers with a square number of divisors' which was added to Sloane's Encyclopedia of Integer Sequences.<sup>35</sup>

### 3 Mechanized reasoning for economic problems

Over the past decade, computer scientists have become interested in economic problems – often publishing economically novel and interesting results, but almost entirely within the computer science literature. This section reviews that literature, focusing on the applications to social choice and auction theory. We structure this survey primarily according to the problem domain within economics, and only secondarily according to our classificatory scheme, in order to focus on the insights into economic problems made possible by these techniques, rather than the techniques themselves.

Table 2 places the papers reviewed in this section into our original classificatory scheme. This classification is imperfect. For example, Tang and Lin (2009) and Geist and Endriss (2011) both used propositional logic solvers (and, therefore, deductive reasoning), but used them to discover new results – which we have associated, above, with inductive reasoning. Papers like this therefore span historical distinctions.

---

<sup>33</sup>Dick's case study of the Argonne National Laboratory's AURA system noted that, while "the capacity to identify what was 'promising' or 'interesting' was precisely one of those unautomatable human abilities ... the Argonne practitioners decided what was important on the basis of extensive experimenting with AURA."

<sup>34</sup>See also the introduction of Tang and Lin (2011a) for a brief review of the history of mechanized theorem discovery; a lengthier review is available in Tang (2010).

<sup>35</sup><https://oeis.org/>

Social choice has been mechanized reasoning’s main point of contact with economics, making it a convenient lens for illustrating mechanized reasoning. Auction theory is, we feel, promising as a new point of contact between mechanized reasoning and economics, due both to the technical parallels between social choice (where mechanized reasoning has proved fruitful) and mechanism design (q.v. Reny (2001)), and to auctions’ importance as allocation mechanisms.

	decidable	undecidable
logic	Geist and Endriss (2011), Brandt and Geist (2016): SAT Tang and Lin (2009): SAT, CSP Bai, Tadjouddine, and Guo (2014): description logic	Nipkow (2009), Wiedijk (2007), Wiedijk (2009), Lange et al. (2013): ITP Grandi and Endriss (2012): ATP
computer system	Xu and Cheng (2007), Arcos et al. (2005), Tadjouddine, Guerin, and Vasconcelos (2009) : model checking	Caminati et al. (2015): code extraction

Table 2: Some applications of mechanized reasoning to economic problems

### 3.1 Social choice

Geanakoplos’ three brief and distinct proofs of Arrow’s impossibility theorem – that, for three or more alternatives and a finite set of agents, there is no social choice rule satisfying unanimity (*UA*), independence of irrelevant alternatives (*IIA*) and non-dictatorship (*ND*) – served as the mechanized reasoning community’s entrée to economic problems: social choice was novel to this community, yet used familiar structures – particularly linear orders – and the three proofs by Geanakoplos (2005) gave the mechanized reasoning community an opportunity to attempt to compare the relative difficulty of encoding those proofs for computers.

One primitive measure of the relative difficulty of formal proofs is to compare their size to that of human proofs.<sup>36</sup> Table 3 reports on the relative sizes of Nipkow’s proofs in Isabelle – a higher-order logic theorem prover – and Wiedijk’s proof<sup>37</sup> in Mizar – a set theoretic proof checker, which augments first-order logic

<sup>36</sup>The easiest way of determining the size of a formal proof is by counting lines of source code. In Section 4 we discuss a less biased measure, the de Bruijn factor.

<sup>37</sup>Wiedijk justified his decision to formalize only Geanakoplos’ first proof by noting that they

by the axioms of Tarski–Grothendieck set theory.<sup>38</sup> Nipkow (2009) attributed the greater length of the Mizar proofs to Isabelle’s “higher level of automation” – something to which we return in our Isabelle proof of Vickrey’s theorem.

	1 <sup>st</sup> proof	3 <sup>rd</sup> proof
Paper (Geanakoplos, 2005)	1 page	1 page
Isabelle (Nipkow, 2009)	350 lines (6 pages)	300 lines
Mizar (Wiedijk, 2007; Wiedijk, 2009)	1100 lines	

Table 3: Relative lengths of human and machine proofs of Arrow’s theorem

Nipkow’s formalization attempts began with Geanakoplos (2001), a working paper that preceded the published version (Geanakoplos, 2005). In seeking to formalize the first proof, he discovered a statement in one of the lemmas that required a 20 line auxiliary proof to properly establish. Further, a relationship between a pivotal voter and a dictator only “hinted at” in the original text required elaboration. Nipkow did not discover any errors in this first proof. Similarly, Wiedijk (2009) reported on missing cases, but no “real errors”.

As to the third proof, Nipkow found two instances of omitted material in its central lemma, preventing him from formalizing the proof. Nipkow presented these concerns to Geanakoplos by e-mail; both concerns were resolved in Geanakoplos (2005).<sup>39</sup>

Both Nipkow and Wiedijk’s proofs were written by the authors themselves, and are therefore examples of ITP. By contrast, Grandi and Endriss (2012) sought to, first, restate Arrow’s theory in FOL and, then, to automatically generate a proof for it.<sup>40</sup> Expressing Arrow’s theory in FOL presented the challenge that quantifying over all possible linear orders of agents’ preference profiles appears to be a second-order quantification as it involves quantifying over agents, alternatives, and the agents’ preference profiles. Grandi and Endriss addressed this by adopting

---

became successively more abstract, making the first the most challenging as, generally “abstract mathematics is easier to formalize than concrete mathematics” (Wiedijk, 2009).

<sup>38</sup>The advantage of Tarski–Grothendieck set theory over Zermelo–Fraenkel is that the former only requires finitely many axioms to axiomatize sets.

<sup>39</sup>Mechanized reasoning can identify omissions by forcing close scrutiny. This, of course, is also possible without mechanical support. For example, in the matching literature, Aygün and Sönmez (2013) identified a hidden assumption in Hatfield and Milgrom (2005) – which they view as “widely considered to be one of the most important advances of the last two decades in matching theory” – without which many of their results fail to hold. The oversight arose from “an ambiguity in setting the primitives of the model”. This ambiguity would likely have been detected by a mechanized reasoner as well.

<sup>40</sup>Grandi and Endriss (2012) is also a good guide to related work on formalizing results in social choice.



the approach taken in Tang and Lin (2009), namely to apply the situation calculus (mentioned in section 2.1) for the representation. Thus, they could present a first-order formalization of the requisite axioms,  $T_{ARROW}$ , allowing them to restate Arrow’s theorem as:

**Theorem 1** (Arrow à la Grandi and Endriss (2012)).  $T_{ARROW}$  has no finite models.

A model in this sense is an instantiation (or example) of the variables used in the theory. For Arrow’s theorem, the variables include  $N$  (the set of agents),  $A$  (the set of alternatives), the set of the agents’ preference profiles, and the set of social welfare functions (SWFs) mapping from such profiles to a social preference. In the two-agent, three-alternative case, that  $T_{ARROW}$  “has no finite models” means that none of the  $6^{36}$  possible SWFs satisfy the theory’s axioms.<sup>41</sup> The theorem claims this property for any finite number of agents, and any finite number of alternatives in excess of three.

FOL’s completeness allows any property of the system to be explicitly derived. However, the second problem with FOL encountered by Grandi and Endriss is that FOL is unable to express finitude, for the same reason that it cannot express induction: intuitively, HOL defines finitude by considering the complement of the infinite, which it can define by induction on the natural numbers. Thus, formulating Arrow’s Theorem in FOL requires a separate formulation for each  $|N|$ . Similarly, proofs of Arrow’s theorem in FOL may differ for each  $|N|$ . Thus, Grandi and Endriss’ attempts to use a first-order theorem prover to automatically generate proofs of Arrow’s theorem failed outside of minimal cases.<sup>42</sup>

Independently of Geanakoplos’ proofs, Suzumura (2000) had presented an induction proof of Arrow’s impossibility theorem for a base case of two agents and  $|A|$  alternatives; an induction result then demonstrated its truth in general. This motivated Tang and Lin (2009) to manually derive a second induction result in the number of agents. Proving the impossibility in a two-agent, three-alternative base case, would – by their two induction lemmas – cause it to hold in general. They computationally exhausted this base case in two different ways.

First, they expressed the problem as a Boolean SAT problem. Tang and Lin then used the situation calculus, which allows many of the problem’s symmetries to be efficiently dealt with by the action of swapping arguments, to reduce the number of variables needed in the base case to 35,973 in 106,354 clauses. These are too many cases to check manually. However, using the SAT solver Chaff2 they could show the inconsistency between the three basic axioms in less than a second on a desktop computer.

---

<sup>41</sup>There are a total of 36 preference profiles in the domain, and six orders in the range, yielding a total of  $\prod_{i=1}^{36} 6$ .

<sup>42</sup>They used Prover9, a successor to Otter, and – therefore – a close relative of the system that found the proof of Robbins’ conjecture (McCune, 1997).

Second, Tang and Lin expressed the problem as a CSP, in which  $V$ , the set of variables, consists – in their base case – of 36 preference profiles;  $D$ , their domain, of six linear orderings for each profile; and  $C$ , their constraint set, of the  $UN$  and  $IJA$  axioms. As the base case implies  $6^{36} \approx 10^{28}$  possible SWFs – far too many to be feasibly generated – the authors used the (first-order) logical programming language Prolog to generate all SWF satisfying the constraints of  $UN$  and  $IJA$ . Running in less than a second on a desktop computer, their Prolog code generated two SWFs, both of which were also dictatorial.

A similar approach yielded the Muller-Satterthwaite theorem, and Sen’s Paretian liberal result, among others.<sup>43</sup>

When implementing the CSP, the authors noticed that imposing even just the  $IJA$  constraint reduced the set of SWFs from  $6^{36}$  to 94. By inspecting these manually, Tang and Lin (2009) posited a new theorem that implies both Arrow’s and Wilson’s. Before stating it, note that a social order is *inversely dictatorial* if it ranks elements in the opposite way to at least one agent; the *Kendall tau* distance between two orderings is the number of pairs on which they disagree. Then:

**Theorem 2** (Tang and Lin (2009)). *If a social welfare function  $W$  on  $(N, A)$  satisfies  $IJA$ , then for every subset  $Y$  of  $A$  such that  $|Y| = 3$ ,*

1.  $W_Y$  is dictatorial, or
2.  $W_Y$  is inversely dictatorial, or
3. The range of  $W_Y$  has at most 2 elements, whose [Kendall tau] distance is at most 1.

As an example of an SWF accepted under condition 3 of theorem, consider the function that always prefers the first alternative to the second, always prefers the first to the third, and prefers the second to the third alternative unless both agents prefer the third to the second. This is neither dictatorial nor inversely dictatorial: the agents’ preferences for the first item are ignored; there are only two elements in its range (e.g.  $a > b > c$  and  $a > c > b$ ), the distance between which is one.<sup>44</sup> As Tang and Lin noted, the third case of their result violates Arrow’s original non-imposition axiom, which requires that the SWF be surjective, mapping to every possible value in its range.

---

<sup>43</sup>See Geist (2010) for a more complete list.

<sup>44</sup>Represent preferences over three objects as a three-digit binary character, the first indicating whether  $a > b$ , the second whether  $a > c$  and the third whether  $b > c$ . There are six permissible three digit numbers, 000, 001, 011, 100, 110 and 111, after eliminating the two cyclical ones.  $IJA$  then requires that each digit in the social preference is a function of the corresponding digits in the individual preferences alone. The 1-distance condition then allows only one of those digits to vary.

Of the 94 SWFs satisfying *IIA*, there are 84 of the sort described above, 6 constant SWFs (one for each ordering), two dictatorial functions, and two inversely dictatorial functions.

As before, the theorem is established by exhaustive computation on the two-agent, three-alternative base case, and then extended to arbitrary finite domains by the manually-derived induction lemmas. Chatterjee and Sen (2014) observed that, as far as they were aware, this is the “only Arrow-type result in the literature that does not use an axiom other than *IIA*”, an achievement that they believe “could not have been conjectured without computational aid”.<sup>45</sup>

Social choice is replete with characterization and impossibility results. Geist and Endriss (2011) applied the Tang and Lin (2009) approach to the problem of ranking sets of objects (Kannai and Peleg, 1984), for which Barberà, Bossert, and Pattanaik (2004) supplied almost 50 possibly desirable axioms.<sup>46</sup>

Rather than deriving an induction lemma for every base case of interest, they derived a broadly applicable induction theorem based on model theory’s Łoś–Tarski preservation theorem which describes when properties ( $\varphi$ , below) are retained in substructures, namely essentially when the theory can be expressed using universal quantifiers in the form  $\forall x . \varphi$ .<sup>47</sup>

Furthermore, as they wished to distinguish between individual alternatives, sets of preferences, and preference orders the authors used a many-sorted FOL. Many-sorted FOL also allows relations (including set inclusion or union) to be defined on one domain that do not hold on the other.

Geist and Endriss then encoded 20 axioms drawn from Barberà, Bossert, and Pattanaik (2004) in their many-sorted FOL. As their induction result translated impossibilities generated on small, finite domains to full-blown impossibility results, they took advantage of these concrete, finite base cases to re-write the axioms in propositional logic (using the kind of rewriting that transformed formula (1) to formula (4) in section 2.1). This, in turn, allowed them to use SAT solvers to search for subsets of axioms which generate impossibility results in these base cases; once found, the induction theorem generalized them to full impossibility results. Doing so for all base cases up to sets of eight items yielded 84 impossibility theorems from about one million combinations.<sup>48</sup>

---

<sup>45</sup>In private correspondence, Sen has conjectured that the result of Malawski and Zhou (1994) linking Wilson’s and Arrow’s theorems may be an immediate consequence of Tang and Lin’s.

<sup>46</sup>Geist (2010) had initially attempted an approach more akin to Grandi and Endriss (2012), seeking to derive an automated proof of the Kannai and Peleg theorem using three different first-order theorem provers; none of them was able to derive a proof after 120 hours of CPU time on 2.26 GHz machines with 24 GB RAM.

<sup>47</sup>As a trivial example, the property that a structure contains three distinct elements cannot be preserved in substructures with fewer than three elements.

<sup>48</sup>Resource constraints limited them to eight items and 20 axioms. They derived their results in about one day.

Their results included known results (e.g. those of Kannai and Peleg (1984) and Barberà and Pattanaik (1984)); variations on known results, typically formed by strengthening axioms to reduce the impossibility's minimal domain; direct consequences of other results (as they did not prune implications of existing impossibilities); a trivial contradiction between the axioms of uncertainty aversion and uncertainty appeal; and – perhaps most interestingly – new theorems. These last resolved an open question in the literature, which we now describe.

Letting  $>$  (resp.  $\succeq$ ) denote strict (resp. weak) preference on individual choice objects (denoted by lower case letters), and  $\triangleright$  (resp.  $\supseteq$ ) strict (resp. weak) preference on sets of objects (denoted by capital letters), Bossert, Pattanaik, and Xu (2000) presented a theorem characterizing the min-max ordering in terms of four axioms. The min-max ordering is defined as

$$A \supseteq_{\text{mx}} B \Leftrightarrow [\min\{A\} > \min\{B\} \vee (\min\{A\} = \min\{B\} \wedge \max\{A\} \succeq \max\{B\})];$$

where  $\min\{A\}$  is the minimal element of  $A$  with respect to  $\succeq$  and  $\max\{A\}$  the maximal element. Thus, a set  $A$  is weakly preferred under the min-max ordering to set  $B$  iff either the worst element of  $A$  is strictly preferred to that of  $B$ , or (when the worst elements are equally preferred) the best element of  $A$  is weakly preferred to that of  $B$ .

The four axioms were:

1. *simple dominance*,

$$x > y \Rightarrow (\{x\} \triangleright \{x, y\} \wedge \{x, y\} \triangleright \{y\})$$

for all  $x$  and  $y$ , so that a set consisting of a strictly preferred object is preferred to a set containing it as well as a strictly less preferred object, which – in turn – is preferred to a set consisting only of that less preferred object.

2. *independence*,

$$A \triangleright B \Rightarrow A \cup \{x\} \supseteq B \cup \{x\}$$

for all  $A$  and  $B$  and  $x$  not contained in  $A$  or  $B$ . Thus, adding a single object to two sets ranked by strict preference does not reverse that ranking (but it may weaken it).

3. *uncertainty aversion*,

$$(x > y > z) \Rightarrow \{y\} \triangleright \{x, z\}$$

for all  $x, y$  and  $z$ , so that a set consisting only of an intermediately preferred object is strictly preferred to a set consisting of a strictly more favourable and a strictly less favourable object.

#### 4. *simple top monotonicity*,

$$x > y \Rightarrow \{x, z\} \triangleright \{y, z\}$$

for all  $x, y$  and  $z$  such that  $x > z$  and  $y > z$ , so that – if an object is strictly preferred to another – a set containing it and a third object is strictly preferred to a set containing the less preferred object and the third object.

Arlegi (2003) showed that the min-max ordering was, in fact, inconsistent with the independence axiom, and presented an alternative axiomatic basis for it. Geist and Endriss (2011) presented a complementary result to Arlegi’s, finding a contradiction between the four original axioms at even four choice objects, thus establishing that the original four axioms are inconsistent, so cannot form the basis of any transitive binary relationship.

Geist and Endriss (2011) also presented the first impossibility result in this literature not to use any dominance axiom.

In cases of interest, the authors were able to quickly derive manual proofs for the computationally discovered results.<sup>49</sup>

Finally, the large set of impossibility results allowed the authors to statistically consider the role of the various axioms. For example, the linear order axiom appeared in all theorems; the ‘even-numbered extension of equivalence’ and reflexivity occurred in none; ‘intermediate independence’ occurred in all results for seven or eight choice items, but never for fewer than five choice items.

Brandt and Geist (2016) extended the methodology of Geist and Endriss (2011) by performing an initial encoding in HOL, and then deriving implications capable of expression in propositional logic for small base cases. This allowed expression of more properties than was possible in the many-sorted FOL of Geist and Endriss (2011). Thus, Brandt and Geist (2016) could encode a neutrality axiom that Geist and Endriss (2011) could not, but at the cost of generating exponentially many new variables, restricting the size of cases that could be computed.

## 3.2 Auctions

Applications of mechanized reasoning to auction design and implementation are less sophisticated than those to social choice. Nevertheless, given auctions’ practical importance, we expect that these will ultimately become more widespread. This section surveys work in two separate areas – applying mechanized reasoning to checking results in auction theory, and checking implementations of auction designs.

---

<sup>49</sup>For the min-max ordering inconsistency, the manual proof is about a half-page long.

On the former, Vickrey’s theorem has provided a basic testbed result. Section 4 illustrates in detail our Isabelle implementation. It therefore complements Lange et al. (2013), which compared implementations of Vickrey’s theorem in four different mechanized reasoners.

Conceptually, as higher-order logic is sufficient to express all concepts in auction theory, it is not challenging to represent basic results in auction theory using a higher-order logic theorem prover like Isabelle. Doing so in more basic logics is both more conceptually challenging, and may offer more promise of automation.

In simpler logics, model checking can automatically establish properties of systems by exhaustively inspecting the system’s state space. Tadjouddine, Guerin, and Vasconcelos (2009) used SPIN, a widely-used commercial model checker based on a linear temporal logic (LTL), to verify Vickrey auctions’ strategy-proofness property that bidders cannot do better than to bid their valuations.<sup>50</sup> They implemented two techniques to reduce the search space while verifying strategy-proofness for arbitrary bid ranges and numbers of agents: *program slicing* removed variables irrelevant to the property; *abstraction* discretized the domain of bids into a three-element domain, depending on whether a bid exceeded, equalled, or was less than an agent’s valuation. A manual proof was required to establish the abstraction’s soundness. Together, the two simplifications allowed strategy-proofness to be verified for any number of agents in a Vickrey auction in a quarter of a second.

The second branch of applications of mechanized reasoning to auctions has sought to establish properties of auction designs as implemented. This is of interest for at least two reasons: first, even if theoretical properties of an auction are known, errors may be introduced when translating the auction from a design to an operational auction. Second, and more commonly for modern auctions, practice may simply outstrip theory. In both cases, mechanized reasoning can be used to reduce the likelihood that an auction will fail when run.

Caminati et al. (2015) used Isabelle to prove that a combinatorial Vickrey auction is soundly specified, in the sense of guaranteeing that – whatever the bids received as input – the output allocated only the available goods, at non-negative prices, and assigned a unique output to each input. Furthermore, it implemented two parallel specifications of the auction, the first close to its standard paper specification, and the second a constructive one. Constructive definitions are essentially algorithmic descriptions. By contrast, definitions in classical logics need only state properties of the defined object. For instance, a classical definition of the maximum of a (non-empty) list of bids identifies an element of the list that is

---

<sup>50</sup>Tadjouddine, Guerin, and Vasconcelos (2009) did not seem to use the modal capabilities of SPIN; instead, the authors seemed to adopt SPIN as they wished – in future work – to be able to accept C code as input, and to reason about it; reasoning about computer programs in which variables can be set does require modal capability.

greater than or equal to every other element in the list. A constructive definition would begin by noting that – for a one-element list – the maximum is the single element of the list; it would then proceed recursively by computing the maximum of the remainder of the list. It would then return the larger of the two: the initial element, or the maximum of the remaining elements.

Isabelle was used to formally prove the equivalence of the two specifications. While the constructive specification is less intuitive, its algorithmic nature allows Isabelle to automatically generate verified executable code from it.

Model checking has also been used to examine auctions for evidence of shill bidding. Xu and Cheng (2007) used SPIN to define predicates corresponding to suspicious behaviour, including pushing prices to a reserve price before dropping out, and bidding on the higher priced of two identical goods. The model checker was then used to see whether the predicates were present in a finite dataset of actual bidding behaviour.

Arcos et al. (2005) developed a toolkit to verify properties of multi-agent environments, with a traditional open outcry auction as their leading example. Their toolkit implemented liveness checks to ensure that agents are not blocked (i.e. can bid in every round), that each bidding round can be reached, and that the final bidding round is reachable from any other, as well as correctness of the bidding language (that is, that by following the rules, the system always remains in a defined state). Their toolkit also includes a simulation tool that conducts a ‘what-if’ analysis by performing a complete check of all cases. While the authors themselves do not refer to what they do as model checking, that is what it most closely resembles.

Finally, Bai, Tadjouddine, and Guo (2014) consider the question of how potential users of online auctions can trust the auctions’ protocols. They develop a protocol for specifying auction designs that can be read by Coq, a mechanized reasoner. Future work building on this should eventually allow Coq to verify properties claimed for the auction.

## **4 Blueprint of a formal proof of Vickrey’s theorem**

The preceding has provided an overview of mechanized reasoning, both in general, and as applied to economic problems. This section provides a detailed description of how a mechanized reasoner is used in practice, in this case to verify a formal proof of Vickrey’s theorem. We use Vickrey’s familiar theorem to focus attention on the formal proof’s implementation, rather than the details of the result or proof.

We begin with a standard statement of Vickrey’s theorem and proof, in this case from Maskin (2004):

**Theorem 3** (Vickrey 1961). *In a second-price auction, it is (weakly) dominant for each buyer  $i$  to bid its valuation  $v_i$ . Furthermore, the auction is efficient.*

*Proof #1.* Suppose that buyer  $i$  bids  $b_i < v_i$ . The only circumstance in which the outcome for  $i$  is changed by its bidding  $b_i$  rather than  $v_i$  is when the highest bid  $b$  by other bidders satisfies  $v_i > b > b_i$ . In that event, buyer  $i$  loses by bidding  $b_i$  (for which its net payoff is 0) but wins by bidding  $v_i$  (for which its net payoff is  $v_i - b$ ). Thus, it is *worse* off bidding  $b_i < v_i$ . By symmetric argument, it can only be worse off bidding  $b_i > v_i$ . We conclude that bidding its valuation (truthful bidding) is weakly dominant. Because it is optimal for buyers to bid truthfully and the high bidder wins, the second-price auction is efficient.  $\square$

However intelligible to humans, Maskin’s proof is too stylized for computers: that there is only one circumstance in which changing bids changes the outcome is merely asserted; the “symmetric argument” is not explicitly elaborated. Before formalizing it, we therefore elaborated the paper proof, and restructured it to four cases, rather than the original nine:

*Proof #2.* Let  $N$  be the set of bidders, and suppose bidder  $i$  bids  $b_i = v_i$ , whatever  $b_j$  each other bidder  $j \neq i$  bids. There are two cases:

1.  $i$  wins. This implies  $b_i = v_i = \max_{j \in N} \{b_j\}$ ,  $p_i = \max_{j \in N \setminus \{i\}} \{b_j\}$ , and  $u_i(\mathbf{b}) = v_i - p_i \geq 0$ . Now consider  $i$  submitting an arbitrary bid  $\hat{b}_i \neq b_i$  so that the bid vector is  $(b_1, \dots, b_{i-1}, \hat{b}_i, b_{i+1}, \dots, b_n)$ . This has two sub-cases:
  - (a)  $i$  wins with  $\hat{b}_i$ , so that  $u_i(b_1, \dots, b_{i-1}, \hat{b}_i, b_{i+1}, \dots, b_n) = u_i(\mathbf{b})$ :  $i$  receives the same utility from winning the item, and pays the same price as the second highest bid has not changed.
  - (b)  $i$  loses with  $\hat{b}_i$ , so that  $u_i(b_1, \dots, b_{i-1}, \hat{b}_i, b_{i+1}, \dots, b_n) = 0 \leq u_i(\mathbf{b})$ .
2.  $i$  loses. This implies  $p_i = 0$ ,  $u_i(\mathbf{b}) = 0$ , and  $b_i \leq \max_{j \in N \setminus \{i\}} \{b_j\}$  as, otherwise,  $i$  would have won. This yields again two cases for  $i$ ’s alternative bid  $\hat{b}_i$ :
  - (a)  $i$  wins, so that  $u_i(b_1, \dots, b_{i-1}, \hat{b}_i, b_{i+1}, \dots, b_n) = v_i - \max_{j \in N \setminus \{i\}} \{b_j\} = b_i - \max_{j \in N \setminus \{i\}} \{b_j\} \leq 0 = u_i(\mathbf{b})$ .
  - (b)  $i$  loses, so that  $u_i(b_1, \dots, b_{i-1}, \hat{b}_i, b_{i+1}, \dots, b_n) = 0 = u_i(\mathbf{b})$ .

By analogy for all  $i$ ,  $\mathbf{b} = \mathbf{v}$  supports an equilibrium in weakly dominant strategies. Efficiency is immediate: the highest bidder has the highest valuation.  $\square$

To formally prove Vickrey’s theorem, we used Isabelle, whose higher-order logic allows our formalization to remain close to paper mathematics.

Our proof, *Vickrey.thy*, is a 9 KB, 185 line file that draws on five ancillary files written for this project.<sup>51</sup> All six files amount to 17 KB and 404 lines – much

<sup>51</sup>See <https://github.com/formare/auctions/tree/master/isabelle/Auction> for the code.



longer than their paper counterparts. A more reliable estimate of the additional effort involved in formal proofs, the *de Bruijn factor* (Wiedijk, 2012), cleans and compresses files before dividing the size of the code by the size of an informal  $\text{\TeX}$  source. It thus avoids bias by semantically irrelevant differences in the syntaxes of formalisations such as languages or code styles using different lengths of lines or of identifiers. The de Bruijn factor relating Proof #2 and its definitions (including *max*) to our Isabelle code is 1.1; as our  $\text{\TeX}$  source is more elaborate than usual, this is lower than the typically observed factors of around four.

Figure 1 depicts the files used in the proof. Those already in Isabelle’s library are marked by ellipses. Dotted ellipses denote files containing general definitions and lemmas that we have added to Isabelle’s library. Rectangles denote this paper’s auction-specific files. Directed edges denote dependence, with the source code being imported into the target code.

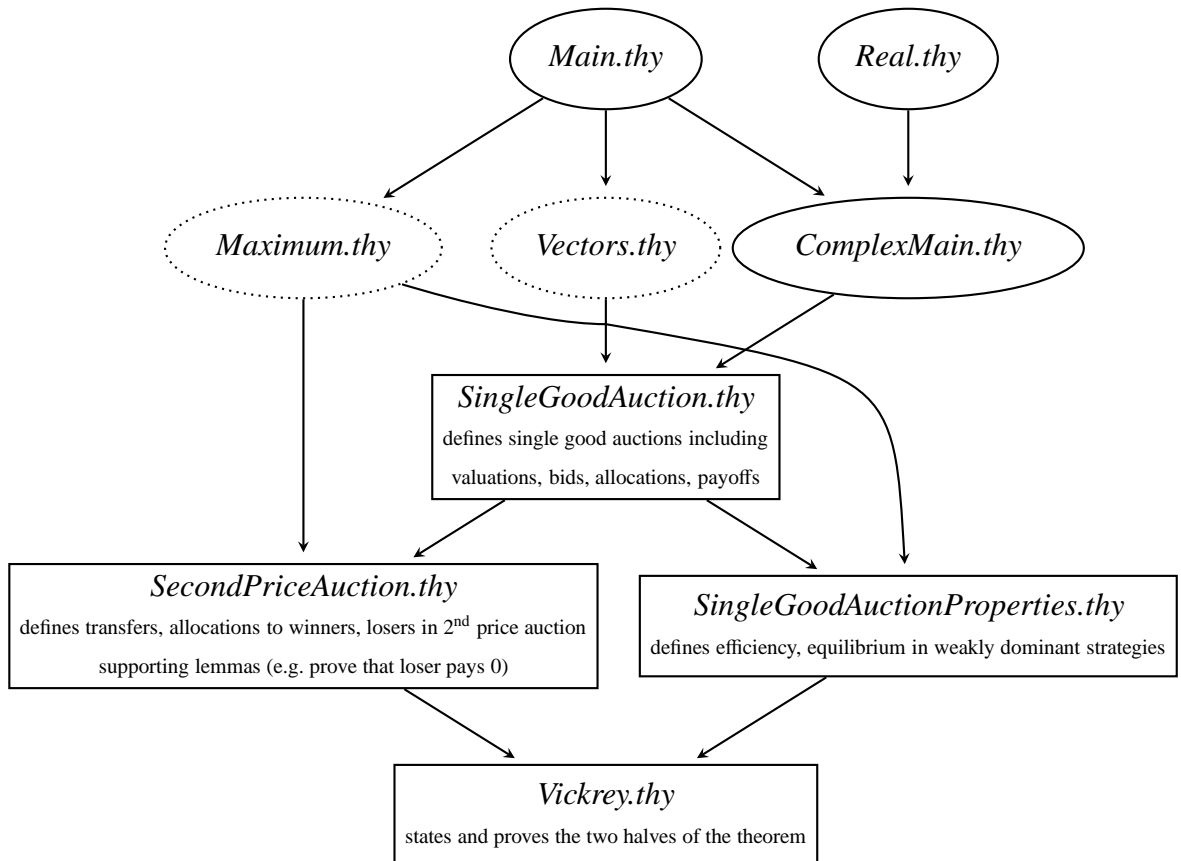


Figure 1: High level theory graph for the formal proof of Vickrey’s theorem

*Vickrey.thy* begins with *vickreyA*, which proves that truth telling is weakly dominant in Vickrey auctions:

**theorem** *vickreyA* :

**fixes**  $N$  :: “*participant set*” **and**  $v$  :: *valuations* **and**  $A$  :: *single\_good\_auction*  
**assumes**  $val$  : “*valuations N v*”  
**defines** “ $b \equiv v$ ”  
**assumes**  $spa$  : “*second\_price\_auction A*” **and**  $card\_N$  : “*card N > 1*”  
**shows** “*equilibrium\_weakly\_dominant\_strategy N v b A*”

The **fixes** keyword applies the theorem to any  $N$ ,  $v$  and  $A$  of the given types. The type *single\_good\_auction* is defined as an *input*  $\times$  *output* relation, with the bidders and their bids as input, and a Boolean allocation vector and a vector of transfers as outcome.<sup>52</sup> The *valuations* type is defined elsewhere to be a vector of real numbers. The **assumes** keyword on the next line states that the theorem holds under an assumption labeled  $val$ , namely that in the vector  $v$  of  $N$  real numbers, all numbers are non-negative (this defined at another place as the definition of ‘valuations’).

Next, the **defines** declaration equates bids and valuations. The following **assumes** keyword introduces and labels further assumptions (e.g.  $A$  is a second-price auction;  $N$  contains more than one bidder). The **shows** keyword states the theorem:  $N$  agents participating in auction  $A$ , with valuations  $v$  and bids  $b$  (equated to valuations) yields an equilibrium in weakly dominant strategies.

*SingleGoodAuctionProperties.thy* defines the equilibrium concept:

**definition** *equilibrium\_weakly\_dominant\_strategy* ::

“*participant set*  $\Rightarrow$  *valuations*  $\Rightarrow$  *bids*  $\Rightarrow$  *single\_good\_auction*  $\Rightarrow$  *bool*” **where**  
“*equilibrium\_weakly\_dominant\_strategy N v b A*  $\longleftrightarrow$   
*valuations N v*  $\wedge$  *bids N b*  $\wedge$  *single\_good\_auction A*  $\wedge$   
 $(\forall i \in N .$   
 $(\forall \text{whatever\_bid} . \text{bids } N \text{ whatever\_bid} \longrightarrow$   
 $(\text{let } b' = \text{whatever\_bid}(i := b \ i)$   
 $\text{in } (\forall x \ p \ x' \ p' . ((N, \text{whatever\_bid}), (x, p)) \in A \wedge ((N, b'), (x', p')) \in A$   
 $\longrightarrow \text{payoff } (v \ i) (x' \ i) (p' \ i) \geq \text{payoff } (v \ i) (x \ i) (p \ i))))$ ”

The definition’s second line declares the type of the *equilibrium\_weakly\_dominant\_strategy* to be a (Boolean) predicate whose arguments are a set of participants, a valuation vector, a bid vector, and an auction.<sup>53</sup> The definition’s body states that the predicate, given arguments  $N$ ,  $v$ ,  $b$  and  $A$ , evaluates to true if and only if the remaining expression does. The expressions in the subsequent line

<sup>52</sup>This can be seen from expressions such as  $((N, b'), (x', p')) \in A$ .

<sup>53</sup>The  $A \Rightarrow B \Rightarrow C$  notation, referred to as *currying*, is equivalent to  $A \times B \rightarrow C$ , but is conceptually simpler as it does not require definition of a  $\times$  operation.

ensure that all arguments have admissible values. Similarly, our first step when introducing *whatever\_bid* is to ensure that it is an admissible bid vector. The *whatever\_bid(i := b i)* notation then takes an arbitrary vector and replaces its *i*th component with *i*'s bid *b i* (which the theorem equates to *i*'s valuation).<sup>54</sup>

We denote the outcome of an arbitrary bid (*whatever\_bid*) by  $(x, p)$ , while  $(x', p')$  denotes that of *i*'s original bid and arbitrary bids by agents  $j \neq i$ . To satisfy the definition of an equilibrium in weakly dominant strategies, the outcome  $(x', p')$  of *i*'s truthful bid must yield a payoff no less than that resulting from an arbitrary bid. The **let** ... **in** ... notation<sup>55</sup> introduces local abbreviations, which can only be accessed within the **in** block; here, this makes the expression  $((N, b'), (x', p')) \in A$  more readable.

The code snippet below formalizes case 2b of Proof #2. It is *declarative*, resembling a textbook proof. *Procedural* proofs, by contrast, prescribe *tactics* to apply, thus more resembling the *process* humans use to find proofs. In either case, each theorem creates a *proof obligation*, or a *goal*; these may be broken into *subgoals* (e.g. by case distinction); the set of local proof obligations implied by these subgoals are stored on a *goal stack*.

*Proof #3.*

```

1 proof –
2   (*...*)
3   {
4     fix i :: participant
5     assume i_range : “i ∈ N”
6     (*...*)
7     let ?b = “whatever_bid(i := b i)”
8     (*...*)
9     have weak_dominance : “payoff (v i) (x' i) (p' i) ≥ payoff (v i) (x i) (p i)”
10    proof cases
11      assume non_alloc : “x' i ≠ 1”
12      with spa_pred' i_range have “x' i = 0” using spa_allocates_binary by blast
13      with spa_pred' i_range have loser_payoff : “payoff (v i) (x' i) (p' i) = 0”
14      by (rule second_price_auction_loser_payoff)
15      have i_bid_at_most_second : “?b i ≤ ?b_max'”
16      proof (rule ccontr)
17      assume “¬?thesis”
18      then have “?b i > ?b_max'” by simp

```

<sup>54</sup>The code snippet contains various instances of “.”: these are separators that improve readability.

<sup>55</sup>We use “...” to distinguish the standard use of ellipses from Isabelle’s “...” notation, whose meaning we introduce when explaining line 30 of the following code snippet.

```

19     with defined spa_pred' i_range have “second_price_auction_winner N ?b x' p' i”
20     by (simp add : only_max_bidder_wins)
21     with non_alloc show False
22     unfolding second_price_auction_winner_def
23     second_price_auction_winner_outcome_def by blast
24     qed
25     show ?thesis
26     proof cases
27     assume “x i ≠ 1”
28     then have “x i = 0” by (rule spa_allocates_binary')
29     with spa_pred i_range have “payoff (v i) (x i) (p i) = 0”
30     by (rule second_price_auction_loser_payoff')
31     also have “... = payoff (v i) (x' i) (p' i)” using loser_payoff ..
32     finally show ?thesis by (rule eq_refl)
33     next
34     (*...*)
35     qed
36     next
37     (*...*)
38     qed
39 }
40 (*...*)
41 qed

```

□

The **proof** keyword starts the proof. Invoked alone, Isabelle would automatically select inference rules to apply. **proof** – performs manual inference. Alternatively, one can specify existing inference rules:

- **proof** *cases* (lines 10 and 26) makes a case distinction; analysis of each case concludes by **showing** that the desired thesis holds; **qed** clears the goal stack; **next** begins the next case.
- **proof** (*rule ccontr*) (line 16) undertakes proof by contradiction, culminating in **show** *False*.

The proof considers an arbitrary but fixed participant  $i$ , which is introduced locally with the **fix** keyword, and assumed to be in the admissible range  $N$  for bidders.<sup>56</sup>

---

<sup>56</sup>In Isabelle, the descriptive form of a verb (e.g. **fixes**, **assumes** or **shows**) are often used when stating theorems, while their imperative counterparts (e.g. **fix**, **assume** or **show**) are used locally in proofs.

The **have** statements establish local facts, generating local proof obligations, which have to be discharged by corresponding **proofs**. Here, the *cases* proof establishes that  $u_i(\dots, v_i, \dots) \geq u_i(\dots, b_i, \dots)$ . This proof makes use of further facts, omitted to keep the snippet readable: *spa\_pred* and *spa\_pred'* state that  $((N, \textit{whatever\_bid}), (x, p))$  and  $((N, ?b), (x', p'))$  respectively are in an *(input, outcome)* relationship of a second price auction with each other.<sup>57</sup> *defined* states that a vector with one component per element of the (finite) set  $N$  has a well-defined maximum component.

Both **from** and **using** introduce facts to discharge the **have** obligations. The **by** keyword invokes an automated proof method, instead of discharging proof obligations by explicit declarative means. Isabelle thus combines ATP and ITP methods.

1. *simp* (lines 18 and 20) simplifies (e.g.  $x \wedge x = x$ ) the statement to be proved. Line 20 supplies a simplification rule of our own, *only\_max\_bidder\_wins*.
2. *blast* (lines 12 and 23) “is (in principle) a complete proof procedure for first-order formulas” (Nipkow, 2015). In practice, *blast* either succeeds, fails, or – giving a practical example of semi-decidability – runs until the user cancels it.
3. *rule* (lines 14, 16, 28, 30 and 32) applies the given lemma as an inference rule. In line 31, “..” abbreviates **by rule**, which automatically applies a matching inference rule.

While interactively developing the proof, we employed the **try** and **try0** commands, which apply a range of automated methods, to find the most appropriate proof methods. Automated calls can always be replaced by explicit declarative steps; Isabelle’s Sledgehammer tool (Blanchette and Paulson, 2015) can sometimes provide them automatically.

The **assume**  $\dots$  **then have** constructions (lines 17 and 18, and 27 and 28) list assumptions **then** state the proof obligations. Line 17’s identifier *?thesis* refers to the proof obligation at the proof’s current level of reasoning.

Lines 22 – 23’s **unfolding** also performs substitutions, replacing stated concepts’ names with the bodies of their definitions. Unlike abbreviations with *?*, the latter are semantic definitions, of which the reasoner make use (e.g. *second\_price\_auction\_winner\_def* is restated in terms of  $i \in N, i \in \arg \max \mathbf{b}, \dots$ ).

---

<sup>57</sup>Isabelle syntactically substitutes identifiers starting with *?* by other, usually more complex expressions before checking a proof step. Syntactic substitution is performed, for example, by the preprocessor of many programming languages, allowing the programmer to use shorthand designations rather than writing complicated expressions in full. It is distinct from the semantic equation of two variables, as in “ $b \equiv v$ ”.

Lines 29–32’s **have** . . . **also have** . . . **finally show** construction allows chains of reasoning with equality before discharging a proof obligation: the “...” following the **also have** are replaced by the right hand side of the previous **have** statement. In line 31, this establishes that  $i$  receives zero given valuation  $v_i$  and either  $(x, p)$ , or  $(x', p')$ .

## 5 Discussion

The decade since the mechanized reasoning community became interested in economic applications has seen rapid progress. When Nipkow reported on his formalization of Arrow’s theorem, he agreed that “[s]ocial choice theory turns out to be perfectly suitable for mechanical theorem proving”, but felt that it was “unclear if [it] will lead to new insights into either social choice theory or theorem proving” (Nipkow, 2009). However, that very year Tang and Lin (2009) used mechanized reasoning to discover a new theorem that subsumes Arrow’s, which Chatterjee and Sen (2014) believed to be novel, and unlikely to have been found with traditional methods. Shortly thereafter, Geist and Endriss (2011) contributed their 84 impossibility theorems.

If mechanized reasoning is to make further inroads into economics it must be sensitive to a number of concerns. First, economics has no proofs of comparable complexity or length to significant results in modern mathematics. Thus, the question of whether a proof will exceed the capability of human theorists to verify is less of a concern than in mathematics. Further, it is unclear that there have been any disastrous cases of mistaken proofs within economics; instead, our greater errors likely result from poor modelling in the first place, and coding or data errors in econometrics.

Second, even when mechanized reasoners have helped identify new results, economic theorists may dismiss them as unmotivated, non-transparent or lacking insight.<sup>58</sup> Even, however, in the worst case, we believe that a stock of poorly-motivated, non-transparent theorems generated blindly by computer provide cases for us to think about and reason with: the presence of the intermediate independence axiom in all of the larger impossibility theorems found by Geist and Endriss (2011) should provide precisely the sort of hunch that sets us sharpening our pencils.

We close by suggesting some further possible applications of mechanized reasoning to economic problems.

First, there are open problems in auction theory that seem amenable to solution by computation (rather than ‘reasoning’). For example, the simplest formulation

---

<sup>58</sup>See Avigad and Harrison (2014, p.73) for a discussion of the tension between rigour and insight in pure mathematics.

of optimal multi-object auctions (q.v. Armstrong, 2000) defines a linear programming problem that quickly becomes too large to solve manually as the number of items increases.<sup>59</sup> As efficient algorithms exist for solving linear programming problems, *automated mechanism design* (q.v. Conitzer and Sandholm, 2003) has already begun to address the purely computational aspects of optimal mechanism design. As formal methods can be used to verify the results of computations (q.v. Gonthier, 2008; Hales et al., 2015), proofs in automated mechanism design could also be verified by formal methods.

Second, we believe that the exhaust-then-induct technique pioneered by Tang and Lin (2009), and developed by Geist and Endriss (2011), offers the promise of automating search for theorems in other areas of economic theory. The formal similarities between social choice and matching theory – including a reliance on discrete objects – suggests that this technique could be applied directly to the latter. Although auction theory appears richer in its use of continuous objects (prices), there is a small literature establishing results by induction (Chew and Serizawa, 2007; Morimoto and Serizawa, 2015; Adachi, 2014; Kato, Ohseto, and Tamura, 2015); the possibility of coupling their induction steps with computational exhaustion has not been explored.

However these tools are applied within economics, it is hard to imagine them not becoming more important, as the tools themselves become faster and easier to use, as they gain acceptance within the pure mathematics community, and as the mechanized reasoning community seeks more applications for them.

## References

- Adachi, T. (2014). “Equity and the Vickrey allocation rule on general preference domains”. *Social Choice and Welfare* 42.4, pp. 813–830.
- Ågotnes, T., W. van der Hoek, and M. Wooldridge (Jan. 2011). “On the logic of preference and judgment aggregation”. *Autonomous Agents and Multi-Agent Systems* 22.1, pp. 4–30.
- Appel, K. and W. Haken (1977). “Every Planar Map is Four Colorable Part I: Discharging”. *Illinois Journal of Mathematics* 21.3, pp. 429–490.
- Appel, K., W. Haken, and J. Koch (1977). “Every Planar Map is Four Colorable Part II: Reducibility”. *Illinois Journal of Mathematics* 21.3, pp. 491–567.
- Arcos, J. L., M. Esteva, P. Noreiga, J. A. Rodríguez-Aguilar, and C. Sierra (2005). “Engineering open environments with electronic institutions”. *Engineering Applications of Artificial Intelligence* 18, pp. 191–204.

---

<sup>59</sup>See Armstrong and Rochet (1999) for the equivalent multi-dimensional screening problem for a monopolist.

- Arlegi, R. (Aug. 2003). “A note on Bossert, Pattanaik and Xu’s “Choice under complete uncertainty: axiomatic characterization of some decision rules””. *Economic Theory* 22.1, pp. 219–225.
- Armstrong, M. (2000). “Optimal multi-object auctions”. *Review of Economic Studies* 67.3, pp. 455–481.
- Armstrong, M. and J.-C. Rochet (1999). “Multi-dimensional screening: a user’s guide”. *European Economic Review* 43.4-6, pp. 959–979.
- Avigad, J. and J. Harrison (2014). “Formally verified mathematics”. *Communications of the ACM* 57.4, pp. 66–75.
- Aygin, O. and T. Sönmez (2013). “Matching with contracts: comment”. *American Economic Review* 103.5, pp. 2050–2051.
- Bai, W., E. M. Tadjouddine, and Y. Guo (2014). “Enabling Automatic Certification of Online Auctions”. In: *Proceedings 11th International Workshop on Formal Engineering Approaches to Software Components and Architectures*. (EPTCS 147, Apr. 2, 2014). Ed. by J. K. B. Buhnova L. Happe, pp. 123–132. doi: 10.4204/EPTCS.147.9.
- Barberà, S. (1980). “Pivotal voters: A new proof of Arrow’s theorem”. *Econom. Lett.* 6.1, pp. 13–16.
- (1983). “Strategy-proofness and pivotal voters: A direct proof of the Gibbard-Satterthwaite theorem”. *Internat. Econom. Rev.* 24.2, pp. 413–417.
- Barberà, S., W. Bossert, and P. K. Pattanaik (2004). “Ranking sets of objects”. In: *Handbook of Utility Theory*. Ed. by S. Barberà, P. J. Hammond, and C. Seidl. Vol. II. Dordrecht: Kluwer Academic Publishers, pp. 893–977.
- Barberà, S. and P. K. Pattanaik (Feb. 1984). “Extending an order on the set to the power set: some remarks on Kannai and Peleg’s approach”. *Journal of Economic Theory* 32.1, pp. 185–191.
- Blanchette, J. C. and L. C. Paulson (May 25, 2015). *Hammering Away. A User’s Guide to Sledgehammer for Isabelle/HOL*. URL: <http://isabelle.in.tum.de/dist/doc/sledgehammer.html>
- Blume, L., D. Easley, J. Kleinberg, R. Kleinberg, and Éva Tardos (2015). “Introduction to computer science and economic theory”. *Journal of Economic Theory* 156, pp. 1–13.
- Boldo, S., J.-H. Jourdan, X. Leroy, and G. Melquiond (Apr. 2013). “A Formally-Verified C Compiler Supporting Floating-Point Arithmetic”. In: *Arith - 21st IEEE Symposium on Computer Arithmetic*. Ed. by A. Nannarelli, P.-M. Seidel, and P. T. P. Tang. Austin, United States: IEEE, pp. 107–115. URL: <https://hal.inria.fr/hal-00761111>
- Bordeaux, L., Y. Hamadi, and L. Zhang (2006). “Propositional satisfiability and constraint programming: a comparative survey”. *ACM Computing Surveys* 38.4.
- Bossert, W., P. Pattanaik, and Y. Xu (Sept. 2000). “Choice under complete uncertainty: axiomatic characterizations of some decision rules”. *Economic Theory* 16.2, pp. 295–312.



- Brandt, F. and C. Geist (2016). “Finding strategy proof social choice functions via SAT solving”. *Journal of Artificial Intelligence Research*.
- Buchberger, B. (2006). “Mathematical Theory Exploration”. In: *Automated Reasoning, Third International Joint Conference, IJCAR 2006, Seattle, WA, USA, August 17-20, 2006, Proceedings*, pp. 1–2. DOI: 10.1007/11814771\_1. URL: [http://dx.doi.org/10.1007/11814771\\_1](http://dx.doi.org/10.1007/11814771_1).
- Burch, J. R., E. M. Clarke, K. L. McMillan, D. L. Dill, and J. Hwang (1990). “Symbolic model checking:  $10^{20}$  states and beyond”. In: *Proceedings of the 5th Annual Symposium on Logic in Computer Science*. IEEE Computer Society Press.
- Caminati, M. B., M. Kerber, C. Lange, and C. Rowat (2015). “Sound Auction Specification and Implementation”. In: *Economics and Computation*. 16<sup>th</sup> ACM Conference, EC’15 (Portland, Oregon, USA, June 15–19, 2015). Ed. by M. Feldman, M. Schwarz, and T. Roughgarden.
- Chatterjee, S. and A. Sen (2014). “Automated Reasoning In Social Choice Theory – Some Remarks”. *Mathematics in Computer Science* 8.1, pp. 5–10.
- Chew, S. H. and S. Serizawa (2007). “Characterizing the Vickrey combinatorial auction by induction”. *Economic Theory* 33.2, pp. 393–406.
- Clarke, E. M., E. A. Emerson, and A. P. Sistla (Apr. 1986). “Automatic verification of finite-state concurrent systems using temporal logic specifications”. *ACM Transactions on Programming Languages and Systems* 8.2, pp. 244–263. doi: 10.1145/5397.5399.
- Clarke, E. M., O. Grumberg, and D. E. Long (1994). “Model Checking and abstraction”. *ACM Transactions on Programming Languages and Systems* 16.5. doi:10.1145/186025.186051, pp. 1512–1542.
- Colton, S., A. Bundy, and T. Walsh (1999). “Automatic Concept Formation in Pure Mathematics”. In: *Proceedings of the 16th International Joint Conference on Artificial Intelligence - IJCAI ’99*. Morgan Kaufmann Pub Inc, pp. 786–791.
- (2000). *Automatic Invention of Integer Sequences*. [http://www.doc.ic.ac.uk/~sgc/html\\_pape](http://www.doc.ic.ac.uk/~sgc/html_pape)
- Conitzer, V. and T. Sandholm (2003). “Applications of Automated Mechanism Design”. In: *UAI-03 workshop on Bayesian Modeling Applications*. Acapulco, Mexico.
- Dahn, B. I. (1998). “Robbins algebras are Boolean: a revision of McCune’s computer-generated solution of Robbins’ problem”. *Journal of Algebra* 208.2, pp. 526–532.
- Dick, S. (2011). “AfterMath: The Work of Proof in the Age of Human-Machine Collaboration”. *Isis* 102.3, pp. 494–505.
- (2015). *After Math: Following Mathematics into the Digital*. Presentation to Microsoft Research New England.

- D’Silva, V., D. Kroening, and G. Weissenbacher (2008). “A Survey of Automated Techniques for Formal Software Verification”. *IEEE Trans. on CAD of Integrated Circuits and Systems* 27.7, pp. 1165–1178.
- Gabbay, D. M. and F. Guenther, eds. (2001/2014). *Handbook of Philosophical Logic*. 2nd ed. Vol. 1–17. Springer-Verlag.
- Gardner, M. (Mar. 1952). “Logic Machines”. *Scientific American* 186.3, pp. 68–73.
- Geanakoplos, J. D. (2001). *Three brief proofs of Arrow’s impossibility theorem*. Discussion Paper 1123RRR. New Haven: Cowles Foundation.
- (2005). “Three brief proofs of Arrow’s impossibility theorem”. *Economic Theory* 26.1, pp. 211–215.
- Geist, C. (July 2, 2010). “Automated Search for Impossibility Theorems in Choice Theory: Ranking Sets of Objects”. MSc Thesis. Institute for Logic, Language and Computation: Universiteit van Amsterdam.
- Geist, C. and U. Endriss (2011). “Automated search for impossibility theorems in social choice theory: ranking sets of objects”. *Journal of Artificial Intelligence Research* 40 (January–April), pp. 143–174.
- Gonthier, G. (2008). “Formal proof – the four color theorem”. *Notices of the AMS* 55.11, pp. 1382–1393.
- Gonthier, G., A. Asperti, J. Avigad, Y. Bertot, C. Cohen, F. Garillot, S. Le Roux, A. Mahboubi, R. O’Connor, S. Ould Biha, I. Pasca, L. Rideau, A. Solovyev, E. Tassi, and L. Théry (2013). “A Machine-Checked Proof of the Odd Order Theorem”. English. In: *ITP 2013, 4th Conference on Interactive Theorem Proving*. Ed. by S. Blazy, C. Paulin, and D. Pichardie. Vol. 7998. LNCS. Rennes, France: Springer, pp. 163–179. doi: 10.1007/978-3-642-39634-2\_14. URL: <http://hal.inria.fr/hal-00816699>.
- Grandi, U. and U. Endriss (2012). “First-Order Logic Formalisation of Impossibility Theorems in Preference Aggregation”. English. *Journal of Philosophical Logic*, pp. 1–24. doi: 10.1007/s10992-012-9240-8.
- Hales, T., M. Adams, G. Bauer, D. T. Dat, J. Harrison, H. L. Truong, C. Kaliszyk, V. Magron, S. McLaughlin, N. T. Thang, N. Q. Truong, T. Nipkow, S. Obua, J. Pleso, J. Rute, T. T. H. A. Alexey Solovyev, T. N. Trung, T. T. Diep, J. Urban, V. K. Ky, and R. Zumkeller (2015). “A formal proof of the Kepler conjecture”. *arXiv preprint arXiv:1501.02155*.
- Hales, T. C. (Nov. 2005). “A proof of the Kepler conjecture”. *Annals of Mathematics* 162.3, pp. 1063–1185.
- (Sept. 6, 2012). *Dense Sphere Packings. A Blueprint for Formal Proofs*. London Mathematical Society Lecture Note Series 400. Cambridge University Press.
- Harrison, J. (May 2006). “Floating-Point Verification using Theorem Proving”. In: *Formal Methods for Hardware Verification*. 6th International School on

- Formal Methods for the Design of Computer, Communication, and Software Systems, SFM 2006 (Bertinoro, Italy). Ed. by M. Bernardo and A. Cimatti. Lecture Notes in Computer Science 3965. Springer Verlag, pp. 211–242.
- Harrison, J. (2007). “A short survey of automated reasoning”. In: *Proceedings of the Second International Conference on Algebraic Biology, AB 2007*. Ed. by H. Anai, K. Horimoto, and T. Kutsia. Vol. 4545. Lecture Notes in Computer Science. Castle of Hagenberg, Austria: Springer-Verlag, pp. 334–349.
- Hatfield, J. W. and P. R. Milgrom (2005). “Matching with contracts”. *American Economic Review* 95.4, pp. 913–935.
- Henkin, L., J. D. Monk, and A. Tarski (1971). *Cylindric algebras, Part I*. Vol. 64. Studies in Logic. North Holland.
- Hoffmann, D. W. (2013). *Die Grenzen der Mathematik – Die Gödel’schen Unvollständigkeitssätze*. Springer-Verlag.
- Ignatovich, D. A. and G. O. Passmore (2015). *Case Study: 2015 SEC Fine Against UBS ATS*. white paper 1503. London: Aesthetic Integration.
- Kannai, Y. and B. Peleg (Feb. 1984). “A note on the extension of an order on a set to the power set”. *Journal of Economic Theory* 32.1, pp. 172–175.
- Karp, R. M. (1972). “Reducibility among Combinatorial Problems”. In: *Complexity of Computer Computations*. Ed. by R. E. Miller and J. W. Thatcher. New York: Plenum, pp. 85–103.
- Kato, M., S. Ohseto, and S. Tamura (2015). “Strategy-proofness versus symmetry in economies with an indivisible good and money”. *International Journal of Game Theory* 44.1, pp. 195–207.
- Kerber, M., C. Lange, and C. Rowat (Jan. 2014). *A formal proof of Vickrey’s theorem by blast, simp, and rule*. Working Paper 14-01. University of Birmingham, Department of Economics. URL: <http://ssrn.com/abstract=2376205>.
- Lander, L. J. and T. R. Parkin (1966). “Counterexample to Euler’s conjecture on sums of like powers”. *Bulletin of the American Mathematical Society* 72.6, p. 1079.
- Lange, C. (2013). “Ontologies and Languages for Representing Mathematical Knowledge on the Semantic Web”. *Semantic Web Journal* 4.2, pp. 119–158. doi: 10.3233/SW-2012-0059.
- Lange, C., M. B. Caminati, M. Kerber, T. Mossakowski, C. Rowat, M. Wenzel, and W. Windsteiger (2013). “A Qualitative Comparison of the Suitability of Four Theorem Provers for Basic Auction Theory”. In: *Intelligent Computer Mathematics*. Conferences on Intelligent Computer Mathematics (Bath, UK, July 8–12, 2013). Ed. by J. Carette, D. Aspinall, C. Lange, P. Sojka, and W. Windsteiger. Lecture Notes in Computer Science 7961. Springer, pp. 200–215. doi: 10.1007/978-3-642-39320-4. arXiv:1303.4193 [cs.LG].
- Leibniz, G. W. (1686). “Projet et Essais pour arriver à quelque certitude pour finir une bonne partie des disputes et pour avancer l’art d’inventer”. In: *Logik-*

- Texte: Kommentierte Auswahl zur Geschichte der modernen Logik.* Ed. by K. Berka and L. Kreisler. Deutsche Übersetzung aus G. W. Leibniz, *Fragmente zur Logik*, Akademie-Verlag, Berlin, 1960. Berlin, Deutschland: Akademie-Verlag. Chap. I.1, pp. 15–17.
- Lenat, D. B. (1976). “AM: An Artificial Intelligence Approach to Discovery in Mathematics as Heuristic Search”. AIM-286, STAN-CS-76-570, and Heuristic Programming Project Report HPP-76-8. PhD thesis. Stanford, California, USA: AI Lab, Stanford University.
- (1983). “EURISKO: A Program That Learns New Heuristics and Domain Concepts”. *Artificial Intelligence* 21, pp. 61–98.
- Lucas, W. F. (1968). “A game with no solution”. *Bulletin of the American Mathematical Society* 74.2, pp. 237–239.
- MacLane, S. (1986). *Mathematics: Form and Function*. Springer-Verlag.
- Malawski, M. and L. Zhou (1994). “A note on social choice theory without the Pareto principle”. *Social Choice and Welfare* 11.2, pp. 103–107.
- Maskin, E. (2004). “The unity of auction theory: Milgrom’s master class”. *Journal of Economic Literature* 42.4, pp. 1102–1115.
- McCarthy, J. and P. Hayes (1969). “Some Philosophical Problems from the Standpoint of Artificial Intelligence”. *Machine Intelligence* 4, pp. 463–502.
- McCorduck, P. (2004). *Machines who think*. 2nd ed. AK Peters.
- McCune, W. (Dec. 1997). “Solution of the Robbins problem”. *Journal of Automated Reasoning* 19.3, pp. 263–276.
- Morimoto, S. and S. Serizawa (2015). “Strategy-proofness and efficiency with non-quasi-linear preferences: a characterization of minimum price Walrasian rule”. *Theoretical Economics* 10.2, pp. 445–487.
- Neumann, J. von and O. Morgenstern (1953). *Theory of Games and Economic Behavior*. 2nd. Princeton University Press.
- Newell, A. (1981). *The Heuristic of George Polya and its Relation to Artificial Intelligence*. Tech. rep. CMU-CS-81-133. also in Rudolf Groner, Marina Groner and Walter F. Bishoof, eds., *Methods of Heuristics*, Lawrence Erlbaum, Hillsdale, New Jersey, USA, p. 195–243. Pittsburgh, Pennsylvania, USA: Department of Computer Science, Carnegie-Mellon University.
- Newell, A. and H. A. Simon (1956). *The logic theory machine: a complex information processing system*. Technical Report P-868. The RAND Corporation.
- Nipkow, T. (2009). “Social choice theory in HOL: Arrow and Gibbard-Satterthwaite”. *Journal of Automated Reasoning* 43.3, pp. 289–304.
- (May 25, 2015). *Programming and Proving in Isabelle/HOL*. URL: <http://isabelle.in.tum.de/>
- Pólya, G. (1945). *How to Solve It*. Princeton, New Jersey, USA: Princeton University Press.
- (1954). *Mathematics and Plausible Reasoning – Induction and Analogy in Mathematics*, Princeton, New Jersey, USA: Princeton University Press.

- Reny, P. J. (Jan. 2001). “Arrow’s theorem and the Gibbard-Satterthwaite theorem: a unified approach”. *Economics Letters* 70.1, pp. 99–105.
- Robinson, J. A. (Jan. 1965). “A Machine-Oriented Logic Based on the Resolution Principle”. *Journal of the Association for Computing Machinery* 12.1, pp. 23–41.
- Suzumura, K. (Mar. 2000). “Welfare economics beyond welfarist-consequentialism”. *Japanese Economic Review* 51.1, pp. 1–32.
- Tadjouddine, E. M., F. Guerin, and W. Vasconcelos (2009). “Abstracting and Verifying Strategy-Proofness for Auction Mechanisms”. In: *Declarative Agent Languages and Technologies VI*. Ed. by M. Baldoni, T. C. Son, M. B. van Riemsdijk, and M. Winikoff. 5397. Springer Verlag, pp. 197–214.
- Tang, P. (2010). “Computer-aided theorem discovery – a new adventure and its application to economic theory”. PhD dissertation. Hong Kong University of Science and Technology.
- Tang, P. and F. Lin (July 2009). “Computer-aided proofs of Arrow’s and other impossibility theorems”. *Artificial Intelligence* 173.11, pp. 1041–1053.
- (Sept. 2011a). “Discovering theorems in game theory: two-person games with unique pure Nash equilibrium payoffs”. *Artificial Intelligence* 175.14–15, pp. 2010–2020.
- (Mar. 2011b). “Two equivalence results for two-person strict games”. *Games and Economic Behavior* 71.2, pp. 479–486.
- Turing, A. M. (1936). “On Computable Numbers, with an Application to the Entscheidungsproblem”. *Proceedings of the London Mathematical Society. Second Series* 42, pp. 230–265.
- Wang, H. (1960). “Toward Mechanical Mathematics”. *IBM Journal of Research and Development* 4.1, pp. 2–22.
- Whitehead, A. N. and B. Russell (1910). *Principia Mathematica*. Vol. I. Cambridge, UK: Cambridge University Press.
- Wiedijk, F. (2007). “Arrow’s impossibility theorem”. *Journal of Formalized Mathematics* 15.4, pp. 171–174.
- (2008). “Formal proof: getting started”. *Notices of the AMS* 55.11, pp. 1408–1414.
- (Feb. 2009). “Formalizing Arrow’s theorem”. *Sādhanā* 34.1, pp. 193–220.
- (Mar. 1, 2012). *The “de Bruijn factor”*. URL: <http://www.cs.ru.nl/~freek/factor/> (visited on 01/24/2016).
- (2014). *Formalizing 100 Theorems*. URL: <http://www.cs.ru.nl/~freek/100/>.
- Woodcock, J., P. G. Larsen, J. Bicarregui, and J. Fitzgerald (Oct. 2009). “Formal method: practice and experience”. *ACM Computing Surveys* 41.4, pp. 1–40.
- Xu, H. and Y.-T. Cheng (2007). “Model checking bidding behaviors in internet concurrent auctions”. *International Journal of Computer Systems Science and Engineering* 22.4, pp. 179–191.