

Nitpick: A Counterexample Generator for Higher-Order Logic Based on a Relational Model Finder[★]

Jasmin Christian Blanchette and Tobias Nipkow

Institut für Informatik, Technische Universität München, Germany
{blanchette,nipkow}@in.tum.de

Abstract. Nitpick is a counterexample generator for Isabelle/HOL that builds on Kodkod, a SAT-based first-order relational model finder. Nitpick supports unbounded quantification, (co)inductive predicates and datatypes, and (co)recursive functions. Fundamentally a finite model finder, it approximates infinite types by finite subsets. As case studies, we consider a security type system and a hotel key card system. Our experimental results on Isabelle theories and the TPTP library indicate that Nitpick generates more counterexamples than other model finders for higher-order logic, without restrictions on the form of the formulas to falsify.

1 Introduction

Anecdotal evidence suggests that most “theorems” initially given to an interactive theorem prover do not hold, typically because of a typo or a missing assumption, but sometimes because of a fundamental flaw. Modern proof assistants often include counterexample generators that can be run on putative theorems or on specific subgoals in a proof to spare users the Sisyphean task of trying to prove non-theorems.

Isabelle/HOL [17] includes two such tools: Quickcheck [4] generates functional code for the higher-order logic (HOL) formula and evaluates it for random values of the free variables, and Refute [23] searches for finite countermodels of a formula through a reduction to SAT (Boolean satisfiability). Their areas of applicability are almost disjoint: Quickcheck excels at inductive datatypes but is restricted to the executable fragment of HOL (which excludes unbounded quantifiers) and may loop endlessly on inductive predicates. Refute copes well with logical symbols, but inductive datatypes and predicates are mostly out of reach due to combinatorial explosion.

In the first-order world, the Alloy Analyzer [13], a testing tool for first-order relational logic (FORL), has enjoyed considerable success lately. Alloy’s backend, the relational model finder Kodkod [21], is available as a stand-alone Java library and is used in many projects.

Alloy’s success inspired us to develop a new counterexample generator for Isabelle, called Nitpick.¹ It uses Kodkod as its backend, thereby benefiting from Kodkod’s optimizations (notably its symmetry breaking) and its rich relational logic. The basic translation from HOL to FORL is conceptually simple (Section 3); however, common HOL

[★] This work is supported by the DFG grant Ni 491/11-1.

¹ The name Nitpick is appropriated from Alloy’s venerable precursor.

idioms such as (co)inductive datatypes and (co)inductive predicates necessitate a translation scheme tailored for SAT solving (Section 4). In addition, Nitpick benefits from many novel optimizations that greatly improve its performance, especially in the presence of higher-order constructs (Section 5).

As case studies, we consider the Isabelle formalizations of a hotel key card system and a security type system (Section 6), both of which are currently beyond the reach of Quickcheck and Refute. Our evaluation indicates that Nitpick falsifies more formulas than Quickcheck and Refute (Section 7), to a large extent because it imposes no syntactic restrictions on the formulas to falsify. Nitpick is integrated with the TPTP benchmark suite [20] and exposed three bugs in the higher-order provers TPS [1] and LEO-II [3].

2 Background

2.1 Higher-Order Logic (HOL)

The types and terms of HOL [12] are that of the simply typed λ -calculus extended with type constructors and constants:

<i>Types:</i>		<i>Terms:</i>	
$\sigma ::= \alpha$	(type variable)	$t ::= x^\sigma$	(variable)
$ (\sigma, \dots, \sigma) \kappa$	(type construction)	$ c^\sigma$	(constant)
		$ t t$	(application)
		$ \lambda x^\sigma. t$	(abstraction)

We write κ for $() \kappa$, $\sigma \kappa$ for $(\sigma) \kappa$, and $\sigma \kappa \tau$ for $(\sigma, \tau) \kappa$. HOL's standard semantics interprets the Boolean type o and the function space $\sigma \rightarrow \tau$. Other types are defined, notably the product type $\sigma \times \tau$. The function arrow associates to the right, reflecting the left-associativity of application. We assume throughout that terms are well-typed using the standard typing rules and write x and c instead of x^σ and c^σ when the type σ is irrelevant or can be inferred from the context. A formula is a term of type o .

Type variables occurring in the type of a constant can be instantiated, offering a restricted form of polymorphism. Standard models interpret the constant $\simeq^{\alpha \rightarrow \alpha \rightarrow o}$ as equality on α for any instance of α . Logical connectives and quantifiers can be defined in terms of \simeq ; for example, $\text{True}^o = (\lambda x^o. x) \simeq (\lambda x. x)$ and $\forall^{(\alpha \rightarrow o) \rightarrow o} = (\lambda P^{\alpha \rightarrow o}. P \simeq (\lambda x. \text{True}))$. The traditional binder notation $Qx. t$ abbreviates $Q (\lambda x. t)$.

2.2 First-Order Relational Logic (FORL)

Kodkod's idiosyncratic logic, FORL, combines elements from first-order logic and relational calculus, to which it adds the transitive closure operator [21]. Its formulas involve variables and terms ranging over relations (sets of tuples drawn from a universe of uninterpreted atoms) of arbitrary arities. The logic is unsorted, but each term denotes a relation of a fixed arity that can be inferred from the arities of its variables. Our translation relies on the following FORL fragment.

Formulas:

$\varphi ::= \text{false}$	(falsity)
$\mid \text{true}$	(truth)
$\mid m\ r$	(multiplicity constraint)
$\mid r \simeq r$	(equality)
$\mid r \subseteq r$	(inclusion)
$\mid \neg \varphi$	(negation)
$\mid \varphi \wedge \varphi$	(conjunction)
$\mid \forall d: \varphi$	(universal quantification)

$d ::= x \in r$

$m ::= \text{no} \mid \text{lone} \mid \text{one}$

$n ::= 1 \mid 2 \mid \dots$

Terms:

$r ::= \text{none}$	(empty set)
$\mid \text{iden}$	(identity relation)
$\mid a_n$	(atom)
$\mid x$	(variable)
$\mid \{ \langle d, \dots, d \rangle \mid \varphi \}$	(comprehension)
$\mid \pi_n^n(r)$	(projection)
$\mid r^+$	(transitive closure)
$\mid r.r$	(dot-join)
$\mid r \times r$	(Cartesian product)
$\mid r \cup r$	(union)
$\mid r - r$	(difference)
$\mid \text{if } \varphi \text{ then } r \text{ else } r$	(conditional)

FORL syntactically distinguishes between terms and formulas. The universe of discourse is $\mathcal{A} = \{a_1, \dots, a_k\}$, where each a_i is an uninterpreted atom. Atoms and n -tuples are identified with singleton sets and singleton n -ary relations, respectively. Bound variables in quantifications and comprehensions range over the tuples in a relation; thus, $\forall x \in (a_1 \cup a_2) \times a_3: \varphi(x)$ is equivalent to $\varphi(a_1 \times a_3) \wedge \varphi(a_2 \times a_3)$.

Although they are not listed above, we will sometimes make use of \vee , \longrightarrow , \exists , $*$, and \cap in examples. The constraint $\text{no } r$ expresses that r is the empty relation, $\text{one } r$ expresses that r is a singleton, and $\text{lone } r \iff \text{no } r \vee \text{one } r$. The projection and dot-join operators are unconventional; their semantics is given by the equations

$$\llbracket \pi_i^k(r) \rrbracket = \{ (r_i, \dots, r_{i+k-1}) \mid (r_1, \dots, r_m) \in \llbracket r \rrbracket \}$$

$$\llbracket r.s \rrbracket = \{ (r_1, \dots, r_{m-1}, s_2, \dots, s_n) \mid \exists t. (r_1, \dots, r_{m-1}, t) \in \llbracket r \rrbracket \wedge (t, s_2, \dots, s_n) \in \llbracket s \rrbracket \}.$$

The dot-join operator admits three important special cases. Let s be unary and r, r' be binary relations. The expression $s.r$ gives the direct image of the set s under r ; if s is a singleton and r a function, it coincides with the function application $r(s)$. Analogously, $r.s$ gives the inverse image of s under r . Finally, $r.r'$ expresses relational composition.

To pass an n -tuple s to a function r , we write $\langle s \rangle.r$, which stands for the n -fold dot-join $\pi_n(s).(\dots(\pi_1(s).r)\dots)$. We write $\pi_i(r)$ for $\pi_i^1(r)$.

The relational operators often make it possible to express first-order problems concisely. The following Kodkod specification attempts to fit 30 pigeons in 29 holes:

```
vars pigeons = {a1, ..., a30}, holes = {a31, ..., a59}
var  $\emptyset \subseteq \text{nest} \subseteq \{a_1, \dots, a_{30}\} \times \{a_{31}, \dots, a_{59}\}$ 
solve  $(\forall p \in \text{pigeons}: \text{one } p.\text{nest}) \wedge (\forall h \in \text{holes}: \text{lone } \text{nest}.h)$ 
```

The example declares three free variables: *pigeons* and *holes* are given fixed values, whereas *nest* is specified with a lower and an upper bound. Variable declarations are an extralogical way of specifying sort constraints and partial solutions.

The constraint $\text{one } p.\text{nest}$ states that pigeon p is in relation with exactly one hole, and $\text{lone } \text{nest}.h$ that hole h is in relation with at most one pigeon. Taken as a whole, the formula states that *nest* is a one-to-one function. It is, of course, not satisfiable, a fact that Kodkod can establish in less than a second.

When reducing FORL to SAT, each n -ary relational variable y is in principle translated to an $|\mathcal{A}|^n$ array of propositional variables $V[i_1, \dots, i_n]$, with $V[i_1, \dots, i_n] \iff \langle a_{i_1}, \dots, a_{i_n} \rangle \in y$. Most relational operations can be coded efficiently; for example, \cup is simply \vee . The quantified formula $\forall r \in s: \varphi(r)$ is treated as $\bigwedge_{j=1}^n t_j \subseteq s \longrightarrow \varphi(t_j)$, where the t_j 's are the tuples that may belong to s . Transitive closure is unrolled to saturation.

3 The Basic Translation

Nitpick employs Kodkod to find a finite model (a satisfying assignment to the free variables and constants) of $\neg P$, where P is the formula to refute. The translation of a formula from HOL to FORL is parameterized by the cardinalities of the types occurring in it, provided as a function $|\sigma|$ from types to positive integers obeying

$$|\sigma| \geq 1 \quad |\sigma| = 2 \quad |\sigma \rightarrow \tau| = |\tau|^{|\sigma|} \quad |\sigma \times \tau| = |\sigma| \cdot |\tau|.$$

Following Jackson [13], we call such a function a *scope*. Like other SAT-based model finders, Nitpick enumerates the possible scopes for each basic type, so that if a formula has a finite counterexample, the tool eventually finds it, unless it runs out of resources.

The basic translation presented in this section handles the following HOL constants:

$False^o$	(falsity)	$insert^{\alpha \rightarrow (\alpha \rightarrow o) \rightarrow \alpha \rightarrow o}$	(element insertion)
$True^o$	(truth)	$UNIV^{\alpha \rightarrow o}$	(universal set)
$\simeq^{\alpha \rightarrow \alpha \rightarrow o}$	(equality)	$\cup^{\alpha \rightarrow o \rightarrow (\alpha \rightarrow o) \rightarrow \alpha \rightarrow o}$	(union)
$\subseteq^{(\alpha \rightarrow o) \rightarrow (\alpha \rightarrow o) \rightarrow o}$	(subset)	$-^{\alpha \rightarrow o \rightarrow (\alpha \rightarrow o) \rightarrow \alpha \rightarrow o}$	(set difference)
\neg^o	(negation)	$Pair^{\alpha \rightarrow \beta \rightarrow \alpha \times \beta}$	(pair constructor)
\wedge^o	(conjunction)	$fst^{\alpha \times \beta \rightarrow \alpha}$	(first projection)
$\forall^{(\alpha \rightarrow o) \rightarrow o}$	(universal quantifier)	$snd^{\alpha \times \beta \rightarrow \beta}$	(second projection)
$\emptyset^{\alpha \rightarrow o}$	(empty set)	$(\)^+_{(\alpha \times \alpha \rightarrow o) \rightarrow \alpha \times \alpha \rightarrow o}$	(transitive closure)

SAT solvers are particularly sensitive to the encoding of problems, so special care is needed when translating HOL formulas. Whenever practicable, HOL constants should be mapped to their FORL equivalents, rather than expanded to their definitions. This is especially true for the transitive closure r^+ , which is defined as the least fixed point of $\lambda R(x, y). (\exists a b. x \simeq a \wedge y \simeq b \wedge r(a, b)) \vee (\exists a b c. x \simeq a \wedge y \simeq c \wedge R(a, b) \wedge r(b, c))$.

As a rule, HOL functions should be mapped to FORL relations accompanied by a constraint. For example, assuming the scope $|\alpha| = 2$ and $|\beta| = 3$, the presumptive theorem $\forall x^\alpha. \exists y^\beta. f x \simeq y$ corresponds to the Kodkod problem

$$\begin{aligned} \text{var } \emptyset &\subseteq f \subseteq \{a_1, a_2\} \times \{a_3, a_4, a_5\} \\ \text{solve } &(\forall x \in a_1 \cup a_2: \text{one } x.f) \wedge \neg(\forall x \in a_1 \cup a_2: \exists y \in a_3 \cup a_4 \cup a_5: x.f \simeq y) \end{aligned}$$

The first conjunct ensures that f is a function, and the second conjunct is the negation of the HOL formula translated to FORL.

An n -ary first-order function (curried or not) can be coded as an $(n+1)$ -ary relation accompanied by a constraint. However, if the return type is o , the function is more efficiently coded as an unconstrained n -ary relation. This allows formulas such as $A^+ \cup B^+ \simeq (A \cup B)^+$ to be translated without taking a detour through ternary relations.

Higher-order quantification and functions bring complications of their own. For example, we would like to translate $\forall g^{\beta \rightarrow \alpha}. g \ x \not\approx y$ into something like

$$\forall g \subseteq (a_3 \cup a_4 \cup a_5) \times (a_1 \cup a_2): (\forall x \in a_3 \cup a_4 \cup a_5: \text{one } x.g) \longrightarrow x.g \not\approx y,$$

but the \subseteq symbol is not allowed at the binding site; only \in is. Skolemization solves half of the problem (Section 5.1), but for the remaining quantifiers we are forced to adopt an unwieldy n -tuple singleton representation of functions, where n is the cardinality of the domain. For the formula above, this gives

$$\forall G \in (a_1 \cup a_2) \times (a_1 \cup a_2) \times (a_1 \cup a_2): x.(\overbrace{a_3 \times \pi_1(G) \cup a_4 \times \pi_2(G) \cup a_5 \times \pi_3(G)}^g) \not\approx y,$$

where G is the triple corresponding to g . In the body, we convert the singleton G to the relational representation, then we apply x on it using dot-join. The singleton encoding is also used for passing functions to functions; fortunately, two optimizations, function specialization and boxing (Section 5.1), make this rarely necessary.

We are now ready to look at the basic translation in more detail. The translation distinguishes between formulas (F), singletons (S), and relations (R). We start by mapping HOL types to sets of FORL atom tuples. For each type σ , we provide two codings, a singleton representation $S\langle\sigma\rangle$ and a relational representation $R\langle\sigma\rangle$:²

$$\begin{aligned} S\langle\sigma \rightarrow \tau\rangle &= S\langle\tau\rangle^{|\sigma|} & R\langle\sigma \rightarrow o\rangle &= S\langle\sigma\rangle \\ S\langle\sigma \times \tau\rangle &= S\langle\sigma\rangle \times S\langle\tau\rangle & R\langle\sigma \rightarrow \tau\rangle &= S\langle\sigma\rangle \times R\langle\tau\rangle \\ S\langle\sigma\rangle &= \{a_1, \dots, a_{|\sigma|}\} & R\langle\sigma\rangle &= S\langle\sigma\rangle. \end{aligned}$$

In the S representation, an element of type σ is mapped to a single tuple $\in S\langle\sigma\rangle$. In the R representation, an element of type $\sigma \rightarrow o$ is mapped to a subset of $S\langle\sigma\rangle$ consisting of the points at which the predicate is *True*; an element of $\sigma \rightarrow \tau$ (where $\tau \neq o$) is mapped to a relation $\subseteq S\langle\sigma\rangle \times R\langle\tau\rangle$; any other element is coded as a singleton. For simplicity, we reuse the same atoms for distinct types. Doing so is sound for well-typed terms.

For each free variable y^σ , we generate the declaration $\mathbf{var} \ \emptyset \subseteq y \subseteq R\langle\sigma\rangle$ as well as a constraint $\Phi^\sigma(y)$ to ensure that functions are functions and single values are singletons:

$$\Phi^{\sigma_1 \rightarrow \dots \rightarrow \sigma_n \rightarrow o}(r) = \text{true} \quad \Phi^{\sigma \rightarrow \tau}(r) = \forall b_f \in S\langle\sigma\rangle: \Phi^\tau(\langle b_f \rangle.r) \quad \Phi^\sigma(r) = \text{one } r.$$

We assume that free and bound variables are syntactically distinguishable, and use the letter y for the former and b for the latter. The symbol b_f denotes a fresh bound variable.

We assume a total order on n -tuples of atoms and let $S_i\langle\sigma\rangle$ denote the i th tuple from $S\langle\sigma\rangle$ according to that order. Furthermore, we define $s(\sigma)$ and $r(\sigma)$ as the arity of the tuples in $S\langle\sigma\rangle$ and $R\langle\sigma\rangle$, respectively. The translation of terms requires the following rather technical conversions between singletons (S), relations (R), and formulas (F):

$$\begin{aligned} s2r^{\sigma \rightarrow o}(r) &= \bigcup_{i=1}^{|\sigma|} \pi_i(r). (a_2 \times S_i\langle\sigma\rangle) & f2s(\varphi) &= \text{if } \varphi \text{ then } a_2 \text{ else } a_1 \\ s2r^{\sigma \rightarrow \tau}(r) &= \bigcup_{i=1}^{|\sigma|} S_i\langle\sigma\rangle \times s2r^\tau\left(\pi_{(i-1) \cdot s(\tau)+1}^{s(\tau)}(r)\right) & s2f(r) &= r \simeq a_2 \\ r2s^{\sigma \rightarrow \tau}(r) &= \{ \langle b_f \in S\langle\sigma \rightarrow \tau\rangle \mid s2r^{\sigma \rightarrow \tau}(b_f) \simeq r \} & s2r^\sigma(r) &= r \quad r2s^\sigma(r) = r. \end{aligned}$$

The Boolean values false and true are arbitrarily coded as a_1 and a_2 , respectively.

² Metatheoretic functions here and elsewhere are defined using sequential pattern matching, eliminating the need for side conditions such as “if $\tau \neq o$ ” and “otherwise.”

The translation of HOL terms is performed by three functions, $F\langle t \rangle$, $S\langle t \rangle$, and $R\langle t \rangle$. Their defining equations are to be matched modulo η -equivalence:

$$\begin{array}{lll}
F\langle y \rangle = s2f(y) & F\langle t \subseteq u \rangle = R\langle t \rangle \subseteq R\langle u \rangle & S\langle b \rangle = b \\
F\langle b \rangle = s2f(b) & F\langle \neg t \rangle = \neg F\langle t \rangle & S\langle \text{Pair } t \ u \rangle = S\langle t \rangle \times S\langle u \rangle \\
F\langle \text{False} \rangle = \text{false} & F\langle t \wedge u \rangle = F\langle t \rangle \wedge F\langle u \rangle & S\langle \text{fst } t^{\sigma \times \tau} \rangle = \pi_1^{s(\sigma)}(S\langle t \rangle) \\
F\langle \text{True} \rangle = \text{true} & F\langle \forall b^\sigma. t \rangle = \forall b \in S\langle \sigma \rangle: F\langle t \rangle & S\langle \text{snd } t^{\sigma \times \tau} \rangle = \pi_{s(\sigma)+1}^{s(\tau)}(S\langle t \rangle) \\
F\langle t \simeq u \rangle = R\langle t \rangle \simeq R\langle u \rangle & F\langle t \ u \rangle = S\langle u \rangle \subseteq R\langle t \rangle & S\langle t \rangle = r2s^\sigma(R\langle t \rangle) \\
R\langle c^\sigma \rangle = f2s(c) & R\langle \text{insert } t \ u \rangle = S\langle t \rangle \cup R\langle u \rangle & \\
R\langle y \rangle = y & R\langle t \cup u \rangle = R\langle t \rangle \cup R\langle u \rangle & \\
R\langle b^\sigma \rangle = s2r^\sigma(b) & R\langle t - u \rangle = R\langle t \rangle - R\langle u \rangle & \\
R\langle \text{Pair } t \ u \rangle = S\langle \text{Pair } t \ u \rangle & R\langle (t^{\sigma \times \sigma \rightarrow o})^+ \rangle = R\langle t \rangle^+ \quad \text{if } r(\sigma) = 1 & \\
R\langle \text{fst } t^{\sigma \times \tau} \rangle = s2r^\sigma(S\langle \text{fst } t \rangle) & R\langle t^{\sigma \rightarrow o} \ u \rangle = f2s(F\langle t \ u \rangle) & \\
R\langle \text{snd } t^{\sigma \times \tau} \rangle = s2r^\tau(S\langle \text{snd } t \rangle) & R\langle t \ u \rangle = \langle S\langle u \rangle \rangle \cdot R\langle t \rangle & \\
R\langle \emptyset^\sigma \rangle = \text{none}^{r(\sigma)} & R\langle \lambda b^\sigma. t^o \rangle = \{ \langle b \in S\langle \sigma \rangle \rangle \mid F\langle t \rangle \} & \\
R\langle \text{UNIV}^\sigma \rangle = R\langle \sigma \rangle & R\langle \lambda b^\sigma. t^\tau \rangle = \{ \langle b \in S\langle \sigma \rangle, b_f \in R\langle \tau \rangle \rangle \mid b_f \subseteq R\langle t \rangle \}. &
\end{array}$$

Annoyingly, the translation of transitive closure is defined only if $r(\sigma) = 1$. We will see ways to lift this restriction in the next two sections.

Theorem 1 (Soundness). *Given a putative theorem P with free variables $y_1^{\sigma_1}, \dots, y_n^{\sigma_n}$ within our HOL fragment and a scope S , P admits a counterexample if there exists a valuation V with $V(y_j) \subseteq R\langle \sigma_j \rangle$ that satisfies the FORL formula $F\langle \neg P \rangle \wedge \bigwedge_{j=1}^n \Phi^{\sigma_j}(y_j)$.*

Proof sketch. Let $\llbracket t \rrbracket_A$ denote the set-theoretic semantics of the HOL term t w.r.t. a variable assignment A and the scope S . Let $\llbracket \rho \rrbracket_V$ denote the semantics of the FORL term or formula ρ w.r.t. a variable valuation V and the scope S . Furthermore, let $\lfloor v \rfloor_X$ denote the X -encoded FORL value corresponding to the HOL value v , for $X \in \{F, S, R\}$. Using recursion induction, it is straightforward to prove that $\llbracket X\langle t \rangle \rrbracket_V = \llbracket \llbracket t \rrbracket_A \rrbracket_X$ if $V(y_i) = \lfloor A(y_i) \rfloor_R$ for all free variables y_i and $V(b_i) = \lfloor A(b_i) \rfloor_S$ for all locally free bound variables b_i occurring in t . Moreover, from the satisfying valuation V of the free variables y_i , we can construct a type-correct HOL assignment A such that $\lfloor A(y_i) \rfloor_R = V(y_i)$; the $\Phi^{\sigma_j}(y_j)$ constraints and the variable bounds $V(y_j) \subseteq R\langle \sigma_j \rangle$ ensure that such an assignment exists. Hence, $\llbracket F\langle \neg P \rangle \rrbracket_V = \text{true} = \llbracket \llbracket \neg P \rrbracket_A \rrbracket_F$, which shows that A falsifies P .

A very thorough soundness proof of a translation from HOL to SAT can be found in Tjark Weber's Ph.D. thesis [23].

4 Refinements to the Basic Translation

4.1 Approximation of Infinite Types and Partiality

Because of the axiom of infinity, the type *nat* of natural numbers does not admit any finite models. To work around this, Nitpick considers finite subsets $\{0, 1, \dots, K-1\}$

of *nat* and maps numbers $\geq K$ to the undefined value (\perp), coded as *none*, the empty set. Formulas of the form $\forall n^{nat}. P(n)$ are treated in essence as $(\forall n < K. P(n)) \wedge P(\perp)$, which usually evaluates to either *False* (if $P(i)$ gives *False* for some $i < K$) or \perp , but not to *True*, since we do not know whether $P(K)$, $P(K+1)$, \dots (collectively represented by $P(\perp)$) are true. In view of this, Nitpick generally cannot soundly disprove conjectures that contain an infinite existential quantifier in their conclusion or an infinite universal quantifier in their assumptions. As a fallback, the tool enters an unsound mode in which the quantifiers are artificially bounded. Counterexamples obtained under these conditions are marked as “potential.”

Functions from *nat* to α are abstracted by relations $\subseteq \{a_1, \dots, a_K\} \times \{a_1, \dots, a_{|\alpha|}\}$ constrained to be partial functions. Partiality makes it possible to encode the successor function *Suc* as the relation $S = (a_1 \times a_2) \cup \dots \cup (a_{K-1} \times a_K)$, which associates no value with a_K . Conveniently, the dot-join $a_K.S$ yields *none*, and so does *none.S*. This is desirable because *Suc* ($K-1$) is unrepresentable and *Suc* \perp is unknown.

Partiality leads to a Kleene three-valued logic, which is expressed in terms of Kodkod’s two-valued logic as follows. At the outermost level, we let the FORL truth value *false* stand for both *False* (no counterexample) and \perp (potential counterexample), and reserve *true* for *True* (genuine counterexample). The same convention is obeyed in other positive contexts within the formula. In negative contexts, *false* codes *False* and *true* codes *True* or \perp . Finally, in unpolarized contexts (for example, as argument to a function), the atom a_1 codes *False*, a_2 codes *True*, and *none* codes \perp . Unlike similar approximation approaches [18, p. 164; 23], Nitpick’s logic is sound, although the tool also has an unsound mode as noted above.

4.2 Nonuniform Representation of HOL Terms

FORL gives Nitpick a lot of flexibility when encoding terms. A value of type $\alpha \times \beta$, for example, can be translated as before to a pair $\in \{a_1, \dots, a_{|\alpha|}\} \times \{a_1, \dots, a_{|\beta|}\}$, but it can also be mapped to a single atom $\in \{a_1, \dots, a_{|\alpha \times \beta|}\}$. Predicates on α (or functions from α to σ with $|\sigma| = 2$) can be coded as single atoms, sets of atoms, relations from atoms to $\{a_1, a_2\}$, or $|\alpha|$ -tuples over $\{a_1, a_2\}$.

Nitpick uses FORL’s flexibility to a larger extent than was hinted at in Section 3. For example, it ensures that the operand of transitive closure is always a binary relation (a set of pairs), no matter what the HOL type is, lifting an annoying limitation in the basic translation described earlier. It also keeps track of whether a term can evaluate to \perp , which makes many optimizations possible in FORL. For example, because free variables never yield \perp , we encode $x \simeq y$ as $x \subseteq y$, which is more efficient.

The current representation selection scheme proceeds in a straightforward bottom-up fashion, inserting conversions as appropriate. More sophisticated schemes that would minimize the number of conversions have yet to be tried.

4.3 Encoding of (Co)Inductive Predicates

Isabelle lets users specify (co)inductive predicates p by their introduction rules and synthesizes a fixed point definition $p \simeq \text{lfp } F$ or $p \simeq \text{gfp } F$. For performance reasons,

Nitpick avoids expanding *lfp* and *gfp* to their definitions and translates (co)inductive predicates directly, using appropriate FORL concepts.

A first intuition is that an inductive predicate p is a fixed point, so we could use the equation $p \simeq F p$ as the axiomatic specification of p . In general, this is unsound since it underspecifies p , but there are two important cases for which this method is sound. First, if the recursion in F is well-founded, the fixed point equation $p \simeq F p$ admits exactly one solution and we can safely use it as p 's specification. Second, if p occurs negatively in the formula, we can replace these occurrences by a fresh constant q satisfying the axiom $q \simeq F q$; this transformation preserves equisatisfiability.

To deal with positive occurrences of p , we adapt a technique from bounded model checking [5]: We replace p by a fresh predicate r_k defined by

$$r_0 \simeq (\lambda \bar{x}. \perp) \qquad r_{n+1} \simeq F r_n,$$

which corresponds to p unrolled k times. For unpolarized occurrences, we use $q \cap r_k$. In essence, we have made p well-founded by adding a counter that decreases by one with each recursive call. This unrolling comes at a price: The search space and the size of the propositional formula for r_k is k times that of q . Hence, it makes sense to look for a counterexample with a small value of k first and increment it gradually if needed.

The situation is mirrored for coinductive predicates: Negative occurrences of p become r_k , positive occurrences become q , and unpolarized occurrences become $q \cup r_k$.

To determine whether a predicate is well-founded, Nitpick generates a wellfoundedness goal and invokes Isabelle's termination prover [7] with a time limit. Given introduction rules of the form

$$\frac{p \bar{t}_{i1} \quad \cdots \quad p \bar{t}_{in_i} \quad Q_i}{p \bar{u}_i}$$

for $i \in \{1, \dots, m\}$, the termination prover must exhibit a well-founded relation R such that $\bigwedge_{i=1}^m \bigwedge_{j=1}^{n_i} Q_i \longrightarrow \langle \bar{t}_{ij}, \bar{u}_i \rangle \in R$ holds.

In our experience, about half of the inductive predicates occurring in practice are well-founded—this includes most type systems and other compositional formalisms, but generally excludes state transition systems.

4.4 Encoding of (Co)Inductive Datatypes and (Co)Recursive Functions

In contrast to Isabelle's constructor-oriented treatment of inductive datatypes, Nitpick's FORL axiomatization revolves around selectors and discriminators, inspired by Kuncak and Jackson's modeling of lists and trees in Alloy [14]. The selector and discriminator view is usually more efficient than the constructor view because it breaks high-arity constructors into several low-arity selectors.

Consider the type $\alpha \text{ list}$ generated from $\text{Nil}^{\alpha \text{ list}}$ and $\text{Cons}^{\alpha \rightarrow \alpha \text{ list} \rightarrow \alpha \text{ list}}$. The FORL axiomatization is done in terms of the discriminators $\text{isNil}^{\alpha \text{ list} \rightarrow o}$ and $\text{isCons}^{\alpha \text{ list} \rightarrow o}$ and the selectors $\text{get1Cons}^{\alpha \text{ list} \rightarrow \alpha}$ and $\text{get2Cons}^{\alpha \text{ list} \rightarrow \alpha \text{ list}}$, which give access to a nonempty list's head and tail. Following Dunets et al. [10], Nil and $\text{Cons } x \text{ xs}$ are translated as isNil and $\text{get1Cons}.x \cap \text{get2Cons}.xs$, respectively.

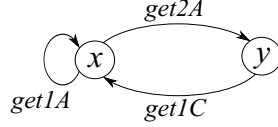
The following axioms, with $N = 1, 2$, specify a subterm-closed finite universe of lists using the atoms $\mathcal{A}_{\alpha \text{ list}}$:

DISJ: $\text{no } \text{isNil} \cap \text{isCons}$
 EXH: $\text{isNil} \cup \text{isCons} \simeq \mathcal{A}_{\text{list}}$
 SEL_N: $\forall xs \in \mathcal{A}_{\text{list}}: \text{if } xs \subseteq \text{isCons} \text{ then one } xs.\text{getNCons} \text{ else no } xs.\text{getNCons}$
 UNIQ: $\text{lone isNil} \wedge (\forall x \in \mathcal{A}_{\alpha}, xs \in \mathcal{A}_{\text{list}}: \text{lone } \text{get1Cons}.x \cap \text{get2Cons}.xs)$
 ACYCL: $\text{no } \text{get2Cons}^+ \cap \text{iden}.$

Examples of subterm-closed list universes using traditional list notation are $\{\[], [a_1], [a_2], [a_3]\}$ and $\{\[], [a_2], [a_3, a_2], [a_1, a_3, a_2]\}$. For recursive functions, Nitpick ignores the construction synthesized by Isabelle and relies instead on the user-specified equations.

The approach can be generalized to mutually recursive datatypes. To generate the ACYCL axioms for the mutually recursive datatypes x with constructors $A^{x \rightarrow y \rightarrow x}$ and B^x and y with constructor $C^{x \rightarrow y}$, we compute their datatype dependency graph, in which vertices are labeled with datatypes and arcs are labeled with selectors. Then we compute for each datatype a regular expression capturing the nontrivial paths from the datatype to itself, with \cdot standing for concatenation, \cup for alternative, and $*$ and $^+$ for repetition. We require the paths to be disjoint from identity:

$$\begin{aligned}
 &\text{no } (\text{get1A} \cup \text{get2A}.\text{get1C})^+ \cap \text{iden} \\
 &\text{no } (\text{get1C}.\text{get1A}^*.\text{get2A})^+ \cap \text{iden}.
 \end{aligned}$$



Nitpick supports coinductive datatypes, even though Isabelle does not provide a high-level mechanism for defining them. Users can define custom coinductive datatypes from first principles and tell Nitpick to substitute its efficient FORL axiomatization for their definitions. Nitpick also knows about Isabelle’s coinductive “lazy list” datatype, $\alpha \text{ llist}$, with the constructors $LNil^{\alpha \text{ llist}}$ and $LCons^{\alpha \rightarrow \alpha \text{ llist} \rightarrow \alpha \text{ llist}}$. The FORL axiomatization is similar to that used for $\alpha \text{ list}$, but the ACYCL axiom is omitted to allow cyclic (ω -regular) lists. Infinite lists are presented to the user as lassos, with a finite stem and cycle. The following coinductive bisimulation principle is translated along with the HOL formula, to ensure that distinct atoms correspond to observably distinct lists:

$$\frac{}{LNil \simeq LNil} \text{BISIM}_1 \qquad \frac{x \simeq y \quad xs \simeq ys}{LCons \ x \ xs \simeq LCons \ y \ ys} \text{BISIM}_2.$$

5 Optimization Steps

5.1 HOL Preprocessing

Function Specialization. A function argument is said to be static if it is passed unaltered to all recursive calls. A typical example is f in the definition of map :

$$\text{map } f \ [] \simeq [] \qquad \text{map } f \ (x \cdot xs) \simeq f \ x \cdot \text{map } f \ xs.$$

An optimization reminiscent of the static argument transformation or lambda-dropping [9, pp. 148–156] is to specialize the function for each eligible call site, thereby avoiding passing the static argument altogether. At the call site, any term whose free variables

are all globally free is eligible for this optimization. Following this scheme, map Suc ns would become $\text{map}_{\text{Suc}} ns$, where map_{Suc} is defined as follows:

$$\text{map}_{\text{Suc}} [] \simeq [] \quad \text{map}_{\text{Suc}} (x \cdot xs) \simeq \text{Suc } x \cdot \text{map}_{\text{Suc}} xs.$$

For this example, specialization reduces the number of propositional variables needed to encode the function by a factor of $|nat|^{|nat|}$.

Boxing. Nitpick normally translates function and product types directly to the homologous Kodkod concepts. This is not always desirable; for example, a transition relation on states represented as n -tuples leads to a $2n$ -ary relation, which gives rise to a combinatorial explosion and precludes the use of FORL’s binary transitive closure.

Our experience suggests that it is almost always advantageous to approximate n -tuples where $n \geq 3$ as well as higher-order arguments. This is achieved by wrapping them in an isomorphic type $\alpha \text{ box}$ with the single constructor $\text{Box}^{\alpha \rightarrow \alpha \text{ box}}$, inserting constructors and selectors as appropriate. Assuming that specialization is not in use, the second equation for map would then become

$$\text{map } f^{(\text{nat} \rightarrow \text{nat}) \text{ box}} (x \cdot xs) \simeq \text{get1Box } f \ x \cdot \text{map } f \ xs,$$

with $\text{map } (\text{Box Suc}) ns$ at the call site. Notice that for function types, boxing is similar to defunctionalization [2], with selectors playing the role of “apply” functions. Further opportunities for boxing are created by uncurrying high-arity constants beforehand.

Quantifier Massaging. (Co)inductive definitions are marred by existential quantifiers, which blow up the size of the resulting propositional formula. The following steps are applied to eliminate quantifiers or reduce their binding range: (1) Replace quantifications of the forms $\forall x. x \simeq t \longrightarrow P(x)$ and $\exists x. x \simeq t \wedge P(x)$ by $P(t)$ if x does not occur free in t . (2) Skolemize. (3) Distribute quantifiers over congenial connectives (\forall over \wedge , \exists over \vee and \longrightarrow). (4) For any remaining subformula $Qx_1 \dots x_n. p_1 \otimes \dots \otimes p_m$, where Q is a quantifier and \otimes is a connective, move the p_i ’s out of as many quantifiers as possible by rebuilding the formula using $\text{qfy}(\{x_1, \dots, x_n\}, \{p_1, \dots, p_m\})$, defined as

$$\text{qfy}(\emptyset, P) = \otimes P \quad \text{qfy}(x \uplus X, P) = \text{qfy}(X, P - P_x \cup \{Qx. \otimes P_x\}),$$

where $P_x = \{p \in P \mid x \text{ occurs free in } p\}$.

The order in which individual variables x are removed from the first argument is crucial because it affects which p_i ’s can be moved out. For clusters of up to 7 quantifiers, Nitpick considers all permutations of the bound variables and chooses the one that minimizes the sum $\sum_{i=1}^m |\tau_{i1}| \cdot \dots \cdot |\tau_{ik_i}| \cdot \text{size}(p_i)$, where $\tau_{i1}, \dots, \tau_{ik_i}$ are the types of the variables that have p_i in their binding range, and $\text{size}(p_i)$ is a rough syntactic measure of p_i ’s size; for larger clusters, it falls back on a heuristic inspired by Paradox’s clause splitting procedure [8]. Thus, the formula $\exists x^\alpha y^\alpha. p \ x \wedge q \ x \ y \wedge r \ y \ (f \ y \ y)$ is transformed into $\exists y^\alpha. r \ y \ (f \ y \ y) \wedge (\exists x^\alpha. p \ x \wedge q \ x \ y)$. Processing y before x in qfy would instead give $\exists x^\alpha. p \ x \wedge (\exists y^\alpha. q \ x \ y \wedge r \ y \ (f \ y \ y))$, which is more expensive because $r \ y \ (f \ y \ y)$, the most complex conjunct, is doubly quantified and hence $|\alpha|^2$ copies of it are needed in the resulting propositional formula.

Constructor Elimination. Since datatype constructors may return \perp in our encoding, we can increase precision by eliminating them. A formula such as $[x, y] \simeq [a, b]$ can

easily be rewritten into $x \simeq a \wedge y \simeq b$, which evaluates to either *True* or *False* even if $[x, y]$ or $[a, b]$ would yield \perp .

For multiple-argument constructors, eliminating constructors helps reduce the number of nested quantifiers. Consider a datatype of AVL trees with two constructors, $Null^{\alpha \text{ tree}}$ and $Node^{\alpha \rightarrow \alpha \text{ tree} \rightarrow \alpha \text{ tree} \rightarrow \text{nat} \rightarrow \alpha \text{ tree}}$, and a *data* constant defined by the equations

$$\text{data } Null \simeq \emptyset \quad \forall a \, t_1 \, t_2 \, h. \text{data } (Node \, a \, t_1 \, t_2 \, h) \simeq \{a\} \cup \text{data } t_1 \cup \text{data } t_2.$$

Our target is the constructor application $Node \, a \, t_1 \, t_2 \, h$ in the second equation's left-hand side. We first pull it out and assign it to a fresh bound variable y :

$$\forall a \, t_1 \, t_2 \, h \, y. \, y \simeq Node \, a \, t_1 \, t_2 \, h \longrightarrow \text{data } y \simeq \{a\} \cup \text{data } t_1 \cup \text{data } t_2.$$

Then we express the constructor arguments in terms of selectors in the conclusion, rewrite the assumption to use a discriminator, and omit the obsolete variables:

$$\forall y. \text{isNode } y \longrightarrow \text{data } y \simeq \{\text{getNode } y\} \cup \text{data } (\text{get2Node } y) \cup \text{data } (\text{get3Node } y).$$

By quantifying over a single variable, we reduce the number of copies of the body from $|\alpha| \cdot |\alpha \text{ tree}|^2 \cdot |\text{nat}|$ to $|\alpha \text{ tree}|$ in the SAT problem, without losing counterexamples. This technique is also useful for constructors taking a single higher-order argument, such as those inserted by the boxing optimization described above.

5.2 Monotonicity Inference

Many formulas occurring in practice are monotonic in the sense that if the formula is falsifiable for a given scope, it is also falsifiable for all larger scopes [13, p. 165]. That not all formulas are monotonic will become clear after considering $|UNIV| = 3$.

Monotonicity can be exploited to prune the search space. For a formula involving n uninterpreted types, a model finder must a priori consider k^n scopes to exhaust all models up to the cardinality bound k . With monotonicity, it is sufficient to consider the single scope in which all types have cardinality k .

We developed and implemented two calculi for inferring monotonicity, and proved them sound [6]. The first calculus, on which we focus here, has limited support for sets encoded as predicates. The second, more powerful calculus addresses this problem by annotating function arrows and relying on a SAT solver to ensure consistent annotations.

For simplicity of exposition, the first calculus is defined for a HOL fragment in which the only constants are \simeq and \longrightarrow . We let *True* abbreviate $(\lambda x^o. x) \simeq (\lambda x. x)$ and $\forall x. p$ abbreviate $(\lambda x. p) \simeq (\lambda x. \text{True})$. We assume a distinguished type variable α with respect to which monotonicity is inferred. The calculus is defined below:

$$\mathbf{TV}^s(o) = \emptyset \quad \mathbf{TV}^+(\beta) = \{\beta\} \quad \mathbf{TV}^-(\beta) = \emptyset \quad \mathbf{TV}^s(\sigma \rightarrow \tau) = \mathbf{TV}^{-s}(\sigma) \cup \mathbf{TV}^s(\tau)$$

$$\frac{}{\mathbf{K}(x)} \quad \frac{}{\mathbf{K}(\longrightarrow)} \quad \frac{\alpha \notin \mathbf{TV}^-(\sigma)}{\mathbf{K}(\simeq^{\sigma \rightarrow \sigma \rightarrow o})} \quad \frac{\mathbf{K}(t) \quad \mathbf{K}(u)}{\mathbf{K}(t \, u)} \quad \frac{\mathbf{K}(t)}{\mathbf{K}(\lambda x. t)}$$

$$\frac{\mathbf{K}(t)}{\mathbf{M}^s(t)} \quad \frac{\mathbf{M}^{-s}(t) \quad \mathbf{M}^s(u)}{\mathbf{M}^s(t \longrightarrow u)} \quad \frac{\mathbf{M}^+(t) \quad \alpha \notin \mathbf{TV}^+(\sigma)}{\mathbf{M}^+(\forall x^\sigma. t)} \quad \frac{\mathbf{M}^-(t)}{\mathbf{M}^-(\forall x. t)} \quad \frac{\mathbf{K}(t) \quad \mathbf{K}(u)}{\mathbf{M}^-(t \simeq u)}.$$

The $\mathbf{TV}^s(\sigma)$ function gives the set of type variables occurring positively (if s is $+$) or negatively (if s is $-$) with respect to \rightarrow . The judgment $\mathbf{K}(t)$ expresses that t 's value remains essentially the same when α 's cardinality is increased, assuming that the free variables also stay the same. A formula P is monotonic if $\mathbf{M}^-(P)$ is derivable.

We evaluated both calculi on the theorems from six highly polymorphic Isabelle theories (*AVL2*, *Fun*, *Huffman*, *List*, *Map*, and *Relation*). We found that the simple calculus inferred monotonicity for 41% to 97% of the theorems depending on the theory, while the more sophisticated calculus achieved 65% to 100% [6].

6 Case Studies

6.1 Volpano–Smith–Irvine Security Type System

Assuming a partition of program variables into public and private ones, Volpano, Smith, and Irvine [22] provide typing rules guaranteeing that the contents of private variables stay private. They define two types, *High* (private) and *Low* (public). An expression is *High* if it involves private variables; otherwise it is *Low*. A command is *High* if it modifies private variables only; commands that could alter public variables are *Low*.

As our first case study, we consider a fragment of the formal soundness proof by Snelting and Wasserrab [19]. Given a variable partition Γ , the inductive predicate $\Gamma \vdash e : \sigma$ tells whether e has type σ , whereas $\Gamma, \sigma \vdash c$ tells whether command c has type σ . Below is a flawed definition of $\Gamma, \sigma \vdash c$:

$$\frac{}{\Gamma, \sigma \vdash \text{skip}} \quad \frac{\Gamma v \simeq [\text{High}]}{\Gamma, \sigma \vdash v := e} \quad \frac{\Gamma \vdash e : \text{Low} \quad \Gamma v \simeq [\text{Low}]}{\Gamma, \text{Low} \vdash v := e} \quad \frac{\Gamma, \sigma \vdash c_1}{\Gamma, \sigma \vdash c_1 ; c_2}$$

$$\frac{\Gamma \vdash b : \sigma \quad \Gamma, \sigma \vdash c_1 \quad \Gamma, \sigma \vdash c_2}{\Gamma, \sigma \vdash \text{if } (b) \ c_1 \text{ else } c_2} \quad \frac{\Gamma \vdash b : \sigma \quad \Gamma, \sigma \vdash c}{\Gamma, \sigma \vdash \text{while } (b) \ c} \quad \frac{\Gamma, \text{High} \vdash c}{\Gamma, \text{Low} \vdash c}.$$

The following theorem constitutes a key step in the soundness proof:

$$\Gamma, \text{High} \vdash c \wedge \langle c, s \rangle \rightsquigarrow^* \langle \text{skip}, s' \rangle \longrightarrow \forall v. \Gamma v \simeq [\text{Low}] \longrightarrow s v \simeq s' v.$$

Informally, it asserts that if executing the *High* command c in state s terminates in state s' , then the public variables of s and s' must agree. This is consistent with our intuition that *High* commands should only modify private variables. However, because we planted a bug in the definition of $\Gamma, \sigma \vdash c$, Nitpick finds a counterexample:

$$\begin{array}{ll} \Gamma = [v_1 \mapsto \text{Low}] & s = [v_1 \mapsto \text{false}] \\ c = \text{skip} ; v_1 := (\text{Var } v_1 == \text{Var } v_1) & s' = [v_1 \mapsto \text{true}]. \end{array}$$

Even though the command c has type *High*, it assigns *true* to the *Low* variable v_1 . The bug is a missing assumption $\Gamma, \sigma \vdash c_2$ in the typing rule for sequential composition.

6.2 Hotel Key Card System

We consider a state-based model of a vulnerable hotel key card system with recordable locks [16], inspired by an Alloy specification due to Jackson [13, pp. 299–306]. The formalization relies on three opaque types, *room*, *guest*, and *key*. A key card, of

type $card = key \times key$, combines an old key and a new key. A state is a 7-field record $\langle \text{owns} :: room \rightarrow guest \rightarrow option, \text{curr} :: room \rightarrow key, \text{issued} :: key \rightarrow o, \text{cards} :: guest \rightarrow card \rightarrow o, \text{roomk} :: room \rightarrow key, \text{isin} :: room \rightarrow guest \rightarrow o, \text{safe} :: room \rightarrow o \rangle$. The set $reach$ of reachable states is defined inductively by the following rules:

$$\begin{array}{c}
 \text{inj init} \\
 \hline
 \langle \text{owns} = (\lambda r. \perp), \text{curr} = \text{init}, \text{issued} = \text{range init}, \text{cards} = (\lambda g. \emptyset), \\
 \text{roomk} = \text{init}, \text{isin} = (\lambda r. \emptyset), \text{safe} = (\lambda r. \text{True}) \rangle \in reach \quad \text{INIT} \\
 \\
 s \in reach \quad k \notin \text{issued } s \\
 \hline
 s(\text{curr} := (\text{curr } s)(r := k), \text{issued} := \text{issued } s \cup k, \\
 \text{cards} := (\text{cards } s)(g := \text{cards } s \cup \langle \text{curr } s \ r, k \rangle), \\
 \text{owns} := (\text{owns } s)(r := \lfloor g \rfloor), \text{safe} := (\text{safe } s)(r := \text{False})) \in reach \quad \text{CHECK-IN} \\
 \\
 s \in reach \quad \langle k, k' \rangle \in \text{cards } s \quad \text{roomk } s \ r \in \{k, k'\} \\
 \hline
 s(\text{isin} := (\text{isin } s)(r := \text{isin } s \ r \cup g), \text{roomk} := (\text{roomk } s)(r := k'), \\
 \text{safe} := (\text{safe } s)(r := \text{owns } s \ r \simeq \lfloor g \rfloor \wedge \text{isin } s \ r \simeq \emptyset \vee \text{safe } s \ r)) \in reach \quad \text{ENTRY} \\
 \\
 s \in reach \quad g \in \text{isin } s \ r \\
 \hline
 s(\text{isin} := (\text{isin } s)(r := \text{isin } s \ r - \{g\})) \in reach \quad \text{EXIT.}
 \end{array}$$

A desirable property of the system is that it should prevent unauthorized access:

$$s \in reach \wedge \text{safe } s \ r \wedge g \in \text{isin } s \ r \longrightarrow \text{owns } s \ r \simeq \lfloor g \rfloor.$$

Nitpick needs some help to contain the state space explosion: We restrict the search to one room and two guests. Within seconds, we get the counterexample

$$\begin{aligned}
 s = \langle \text{owns} = (r_1 := \lfloor g_1 \rfloor), \text{curr} = (r_1 := k_1), \text{issued} = \{k_1, k_2, k_3, k_4\}, \\
 \text{cards} = (g_1 := \{\langle k_3, k_1 \rangle, \langle k_4, k_2 \rangle\}, g_2 := \{\langle k_2, k_3 \rangle\}), \text{roomk} = (r_1 := k_3), \\
 \text{isin} = (r_1 := \{g_1, g_2\}), \text{safe} = \{r_1\} \rangle
 \end{aligned}$$

with $g = g_2$ and $r = r_1$.

To retrace the steps from the initial state to the s , we can ask Nitpick to show the interpretation of $reach$ at each iteration. This reveals the following “guest in the middle” attack: (1) Guest g_1 checks in and gets a card $\langle k_4, k_2 \rangle$ for room r_1 , whose lock expects k_4 . Guest g_1 does not enter the room yet. (2) Guest g_2 checks in, gets a card $\langle k_2, k_3 \rangle$ for r_1 , and waits. (3) Guest g_1 checks in again, gets a card $\langle k_3, k_1 \rangle$, inadvertently unlocks room r_1 with her previous card, $\langle k_4, k_2 \rangle$, leaves a diamond on the nightstand, and exits. (4) Guest g_2 enters the room and “borrows” the diamond.

This flaw was already detected by Jackson using the Alloy Analyzer on his original specification and can be fixed by adding $k' \simeq \text{curr } s \ r$ to the conjunction in ENTRY.

7 Evaluation

An ideal way to assess Nitpick’s strength would be to run it against Refute and Quick-check on a representative database of Isabelle/HOL non-theorems. Lacking such a database, we chose instead to derive formulas from existing theorems by mutation, replacing constants with other constants and swapping arguments, as was done when

evaluating Quickcheck [4]. The vast majority of formulas obtained this way are invalid, and those few that are valid do not influence the ranking of the counterexample generators. For executable theorems, we made sure that the generated mutants are also executable to prevent a bias against Quickcheck.

The table below summarizes the results of running the tools on 3200 random mutants from 20 Isabelle theories (200 per theory), with a limit of 10 seconds per formula. Most counterexamples are found within a few seconds; giving the tool more time would have little impact on the results.

THEORY	QUICK.	REF.	NITP.	THEORY	QUICK.	REF.	NITP.
<i>Divides</i>	134/184	3+15	141 +2	<i>ArrowGS</i>	0/0	0+126	139 +2
<i>Fun</i>	5/9	162+1	163 +0	<i>Coinductive</i>	4/6	16+7	87 +13
<i>GCD</i>	119/162	1+17	124 +10	<i>CoreC++</i>	7/30	3+6	29 +1
<i>List</i>	78/130	3+113	117 +9	<i>FFT</i>	31/40	1+2	47 +15
<i>MacLaurin</i>	43 /62	0+0	26+7	<i>Huffman</i>	84/160	1+45	119 +2
<i>Map</i>	19/34	103+45	157 +0	<i>MiniML</i>	14/33	0+116	79 +49
<i>Predicate</i>	2/2	147+14	161 +0	<i>NBE</i>	41/62	0+18	81 +24
<i>Relation</i>	0/2	144+3	150 +1	<i>Ordinal</i>	0/66	10+3	12 +0
<i>Set</i>	17/25	149+0	151 +0	<i>POPLmark</i>	56/96	4+6	103 +15
<i>Wellfounded</i>	10/24	118+20	141 +1	<i>Topology</i>	0/0	124+4	139 +3

The table's entries have the form G/X for Quickcheck and $G+P$ for Refute and Nitpick, where G = number of genuine counterexamples found and reported as such, P = number of potential counterexamples found (in addition to G), and X = number of executable mutants (among 200).

Refute's three-valued logic is unsound, so all counterexamples for formulas that involve an infinite type are potentially spurious and reported as such to the user. Nitpick also has an unsound mode, which contributes some potential counterexamples. Unfortunately, there is no easy way to tell how many of these are actually genuine.

Quickcheck and Nitpick are comparable on executable formulas, but Nitpick also fares well on non-executable ones. Notable exceptions are formalizations involving real arithmetic (*MacLaurin* and *FFT*), complex set-theoretic constructions (*Ordinal*), or a large state space (*CoreC++*).

Independently, Nitpick competes against Refute in the higher-order model finding division of the TPTP [20]. In a preliminary run, it disproved 293 out of 2729 formulas (mostly theorems), compared with 214 for Refute. Much to our surprise, Nitpick exhibited counterexamples for five formulas that had previously been proved by TPS [1] or LEO-II [3], revealing two bugs in the former and one bug in the latter.

8 Related Work

The approaches for testing conjectures can be classified in three broad categories:

- *Random testing*. The formula is evaluated for random values of the free variables. This approach is embodied by Isabelle's Quickcheck [4] and similar tools for other proof assistants. It is restricted to executable formulas.

- *SAT solving.* The formula is translated to propositional logic and handed to a SAT solver. This procedure was pioneered by McCune in his first-order finder MACE [15]. Other first-order MACE-style finders include Paradox [8] and Kodkod [21]. The higher-order finders Refute [23] and Nitpick also belong to this category.
- *Direct search.* The search for a model is performed directly on the formula, without translation to propositional logic. This approach was introduced by SEM [24].

Some proof methods deliver sound or unsound counterexamples upon failure, notably model checking, semantic tableaux, and satisfiability modulo theory (SMT) solving. Also worth of mention is the Dynamite tool [11], which lets users prove Alloy formulas in the interactive theorem prover PVS. Weber [23, pp. 3–4] provides a more detailed discussion of related work.

9 Conclusion

Nitpick is to our knowledge the first higher-order model finder that supports both inductive and coinductive predicates and datatypes. It works by translating higher-order formulas to first-order relational logic (FORL) and invoking the highly-optimized SAT-based Kodkod model finder [21] to solve these. Compared with Quickcheck, which is restricted to executable formulas, Nitpick shines by its generality—the hallmark of SAT-based model finding.

The translation to FORL is designed to exploit Kodkod’s strengths. Datatypes are encoded following an Alloy idiom [10, 14] extended to mutually recursive and coinductive datatypes. FORL’s relational operators provide a natural encoding of partial application and λ -abstraction, and the transitive closure plays a crucial role in the encoding of inductive datatypes. Our main contributions have been to isolate three ways to translate (co)inductive predicates to FORL, based on wellfoundedness, polarity, and linearity, and to devise optimizations—notably function specialization, boxing, and monotonicity inference—that dramatically increase scalability in practical applications.

Nitpick is included with the latest version of Isabelle and is invoked automatically whenever users enter new formulas to prove, helping to catch errors early, thereby saving time and effort. But Nitpick’s real beauty is that it lets users experiment with formal specifications in the playful way championed by Alloy but with Isabelle’s higher-order syntax, definition principles, and theories at their fingertips.

Acknowledgment. Stefan Berghofer, Lukas Bulwahn, Marcelo Frias, Florian Haftmann, Alexander Krauss, Mark Summerfield, Emina Torlak, and several anonymous reviewers provided useful comments on drafts of this paper. Alexander Krauss also helped devise the monotonicity inference calculi, and Geoff Sutcliffe ran Nitpick against Refute on the TPTP benchmark suite. We thank them all.

References

1. Andrews, P.B., Bishop, M., Issar, S., Nesmith, D., Pfenning, F., Xi, H.: TPS: A theorem-proving system for classical type theory. *J. Auto. Reas.* 16(3), 321–353 (1996)
2. Bell, J.M., Bellegarde, F., Hook, J.: Type-driven defunctionalization. *ACM SIGPLAN Notices* 32(8), 25–37 (1997)

3. Benz Müller, C., Paulson, L., Theiss, F., Fietzke, A.: Progress report on LEO-II, an automatic theorem prover for higher-order logic. In: Schneider, K., Brandt, J. (eds.) TPHOLs: Emerging Trends. C.S. Dept., University of Kaiserslautern, Internal Report 364/07 (2007)
4. Berghofer, S., Nipkow, T.: Random testing in Isabelle/HOL. In: Cuellar, J., Liu, Z. (eds.) SEFM 2004, pp. 230–239. IEEE C.S., Los Alamitos (2004)
5. Biere, A., Cimatti, A., Clarke, E.M., Zhu, Y.: Symbolic model checking without BDDs. In: Cleaveland, W.R. (ed.) TACAS 1999. LNCS, vol. 1579, pp. 193–207. Springer, Heidelberg (1999)
6. Blanchette, J.C., Krauss, A.: Monotonicity inference for higher-order formulas. In: Giesl, J., Hähnle, R. (eds.) IJCAR 2010. LNCS, Springer, Heidelberg (to appear, 2010)
7. Bulwahn, L., Krauss, A., Nipkow, T.: Finding lexicographic orders for termination proofs in Isabelle/HOL. In: Schneider, K., Brandt, J. (eds.) TPHOLs 2007. LNCS, vol. 4732, pp. 38–53. Springer, Heidelberg (2007)
8. Claessen, K., Sörensson, N.: New techniques that improve MACE-style model finding. In: MODEL (2003)
9. de Medeiros Santos, A.L.: Compilation by Transformation in Non-Strict Functional Languages. Ph.D. thesis, C.S. Dept., University of Glasgow (1995)
10. Dunets, A., Schellhorn, G., Reif, W.: Bounded relational analysis of free datatypes. In: Beckert, B., Hähnle, R. (eds.) TAP 2008. LNCS, vol. 4966, pp. 99–115. Springer, Heidelberg (2008)
11. Frias, M.F., Pombo, C.G.L., Moscato, M.M.: Alloy Analyzer+PVS in the analysis and verification of Alloy specifications. In: Grumberg, O., Huth, M. (eds.) TACAS 2007. LNCS, vol. 4424, pp. 587–601. Springer, Heidelberg (2007)
12. Gordon, M.J.C., Melham, T.F. (eds.): Introduction to HOL: A Theorem Proving Environment for Higher Order Logic. Cambridge University Press, Cambridge (1993)
13. Jackson, D.: Software Abstractions: Logic, Language, and Analysis. MIT Press, Cambridge (2006)
14. Kuncak, V., Jackson, D.: Relational analysis of algebraic datatypes. In: Gall, H.C. (ed.) ESEC/FSE 2005 (2005)
15. McCune, W.: A Davis–Putnam program and its application to finite first-order model search: Quasigroup existence problems. Technical report, ANL (1994)
16. Nipkow, T.: Verifying a hotel key card system. In: Barkaoui, K., Cavalcanti, A., Cerone, A. (eds.) ICTAC 2006. LNCS, vol. 4281, pp. 1–14. Springer, Heidelberg (2006)
17. Nipkow, T., Paulson, L.C., Wenzel, M.: Isabelle/HOL. LNCS, vol. 2283. Springer, Heidelberg (2002)
18. Schumann, J.M.: Automated Theorem Proving in Software Engineering. Springer, Heidelberg (2001)
19. Snelting, G., Wasserrab, D.: A correctness proof for the Volpano/Smith security typing system. In: Klein, G., Nipkow, T., Paulson, L.C. (eds.) AFP (September 2008)
20. Sutcliffe, G., Suttner, C.: The TPTP problem library for automated theorem proving, <http://www.cs.miami.edu/~tptp/>
21. Torlak, E., Jackson, D.: Kodkod: A relational model finder. In: Grumberg, O., Huth, M. (eds.) TACAS 2007. LNCS, vol. 4424, pp. 632–647. Springer, Heidelberg (2007)
22. Volpano, D., Smith, G., Irvine, C.: A sound type system for secure flow analysis. *J. Comp. Sec.* 4(3), 167–187 (1996)
23. Weber, T.: SAT-Based Finite Model Generation for Higher-Order Logic. Ph.D. thesis, Dept. of Informatics, T.U. München (2008)
24. Zhang, J., Zhang, H.: SEM: A system for enumerating models. In: Kaufmann, M. (ed.) IJCAI 95, vol. 1, pp. 298–303 (1995)