

Mechanized Metatheory for the Masses: The POPLMARK Challenge

Brian E. Aydemir¹, Aaron Bohannon¹, Matthew Fairbairn², J. Nathan Foster¹,
Benjamin C. Pierce¹, Peter Sewell², Dimitrios Vytiniotis¹, Geoffrey
Washburn¹, Stephanie Weirich¹, and Steve Zdancewic¹

¹ Department of Computer and Information Science, University of Pennsylvania

² Computer Laboratory, University of Cambridge

Abstract. How close are we to a world where every paper on programming languages is accompanied by an electronic appendix with machine-checked proofs?

We propose an initial set of benchmarks for measuring progress in this area. Based on the metatheory of System F_<, a typed lambda-calculus with second-order polymorphism, subtyping, and records, these benchmarks embody many aspects of programming languages that are challenging to formalize: variable binding at both the term and type levels, syntactic forms with variable numbers of components (including binders), and proofs demanding complex induction principles. We hope that these benchmarks will help clarify the current state of the art, provide a basis for comparing competing technologies, and motivate further research.

1 Introduction

Many proofs about programming languages are long, straightforward, and tedious, with just a few interesting cases. Their complexity arises from the management of many details rather than from deep conceptual difficulties; yet small mistakes or overlooked cases can invalidate large amounts of work. These effects are amplified as languages scale: it becomes hard to keep definitions and proofs consistent, to reuse work, and to ensure tight relationships between theory and implementations. Automated proof assistants offer the hope of significantly easing these problems. However, despite much encouraging progress in recent years and the availability of several mature tools (ACL2, Coq, HOL, HOL Light, Isabelle, Lego, NuPRL, PVS, Twelf, etc.), their use is still not commonplace.

We believe that the time is ripe to join the efforts of the two communities, bringing developers of automated proof assistants together with a large pool of eager potential clients—programming language designers and researchers. In particular, we would like to answer two questions:

1. What is the current state of the art in formalizing language metatheory and semantics? What can be recommended as best practices for groups (typically not proof-assistant experts) embarking on formalizing language definitions, either small- or large-scale?

2. What improvements are needed to make the use of tool support commonplace? What can each community contribute?

Over the past several months, we have surveyed the landscape of proof assistants, language representation strategies, and related tools. Collectively, we have applied automated theorem proving technology to a number of problems, including proving transitivity of the algorithmic subtype relation in System $F_{<}$ [3], proving type soundness of Featherweight Java, proving type soundness of variants of the simply typed λ -calculus and $F_{<}$, and a substantial formalization of the behavior of TCP, UDP, and the Sockets API. We have carried out these case studies using a variety of object-language representation strategies, proof techniques, and proving environments. We have also experimented with lightweight tools designed to make it easier to define and typeset both formal and informal mathematics. Although experts in programming language theory, we are relative outsiders with respect to computer-aided proof.

Our conclusion from these experiments is that the relevant technology has developed *almost* to the point where it can be widely used by language researchers. We seek to push it over the threshold, making the use of proof tools common practice in programming language research—mechanized metatheory for the masses.

Tool support for formal reasoning about programming languages would be useful at many levels:

1. *Machine-checked metatheory.* These are the classic problems: type preservation and soundness theorems, unique decomposition properties for operational semantics, proofs of equivalence between algorithmic and declarative variants of type systems, etc. At present such results are typically proved by hand for small to medium-sized calculi, and are not proved at all for full language definitions. We envision a future in which the papers in conferences such as *Principles of Programming Languages (POPL)* are routinely accompanied by mechanically checkable proofs of the theorems they claim.
2. *Use of definitions as oracles for testing and animation.* When developing a language implementation together with a formal definition one would like to use the definition as an oracle for testing. This requires tools that can decide typing and evaluation relationships, and they might differ from the tools used for (1) or be embedded in the same proof assistant. In some cases one could use a definition directly as a prototype.
3. *Support for engineering large-scale definitions.* As we move to full language definitions—on the scale of Standard ML [19] or larger—pragmatic “software engineering” issues become increasingly important, as do the potential benefits of tool support. For large definitions, the need for elegant and concise notation becomes pressing, as witnessed by the care taken by present-day researchers using informal mathematics. Even lightweight tool support, without full mechanized proof, could be very useful in this domain, e.g., for sort checking and typesetting of definitions and of informal proofs, automatically instantiating definitions, performing substitutions, etc.

Large scale formalization of languages is already within reach of current technology. For examples, see the work on proofs of correctness of the Damas-Milner type inference algorithm for ML [6, 22], semantics for C [25], semantics for Standard ML [32, 34, 13], and semantics and proofs of correctness for substantial subsets of Java [24, 17, 23]. Some other significant existing applications of mechanized metatheory include Foundational Proof Carrying Code [1] and Typed Assembly Languages [4]. Inspired by these successes, we seek to make mechanized metatheory more accessible to programming languages researchers.

We hope to stimulate progress by providing a framework for comparing alternative technologies. We issue here an initial set of challenge problems, dubbed the POPLMARK Challenge, chosen to exercise some aspects of programming languages that are known to be difficult to formalize: variable binding at both term and type levels, syntactic forms with variable numbers of components (including binders), and proofs demanding complex induction principles. Such challenge problems have been used in the past within the theorem proving community to focus attention on specific areas and to evaluate the relative merits of different tools; these have ranged in scale from benchmark suites and small problems [31, 12, 5, 15, 9, 21] up to the grand challenges of Floyd, Hoare, and Moore [7, 14, 20]. We hope that our challenge will have a similarly stimulating effect.

Our problems are drawn from the basic metatheory of a call-by-value variant of System F_λ [3], enriched with records, record subtyping, and record patterns. We provide an informal definition of its type system and operational semantics. This language is of moderate scale—significantly more complex than simply typed lambda-calculus or “mini ML,” but much smaller than a full-blown programming language—to keep the work involved in attempting the challenges manageable. (Our challenges therefore explicitly address only points 1 and 2 above; we regard the pragmatic issues of point 3 as equally critical, but it is not yet clear to us how to formulate a useful challenge problem at this larger scale.) As these challenges are met, we anticipate that the set of problems will be extended to emphasize other language features and proof techniques.

We have begun collecting and publicizing solutions to these challenge problems on our web site.³ In the longer run, we hope that this site and accompanying e-mail discussion list will serve as a forum for promoting and advancing the current best practices in proof assistant technology and making this technology available to the broader programming languages community and beyond. We encourage researchers to try out the POPLMARK Challenge using their favorite tools and send us their solutions.

In the next section, we discuss in more detail our reasons for selecting this specific set of challenge problems. Section 3 describes the problems themselves, and Section 4 sketches some avenues for further development of the challenge problem set.

³ <http://www.cis.upenn.edu/proj/plclub/mmm/>

2 Design of the Challenge

This section motivates our choice of challenge problems and discusses the evaluation criteria for proposed solutions to the challenges. Since variable binding is a central aspect of the challenges, we briefly discuss relevant techniques and sketch some of our own experience in this area.

2.1 Problem Selection

The goal of the POPLMARK Challenge is to provide a small, well-defined set of problems that capture many of the most critical issues in formalizing programming language metatheory. By its nature, such a benchmark will not be able to reflect *all* important issues. Instead, the POPLMARK problems concentrate on a few important features:

- *Binding*. Most programming languages have some form of binding in their syntax and require a treatment of α -equivalence and substitution in their semantics. To adequately represent many languages, the representation strategy must support multiple kinds of binders (e.g. term and type), constructs introducing multiple binders over the same scope, and potentially unbounded lists of binders (e.g. for record patterns).
- *Complex inductions*. Programming language definitions often involve complex, mutually recursive definitions. Structural induction over such objects, mutual induction, and induction on heights or pairs of derivations are all commonplace in metatheory.
- *Experimentation*. Proofs about programming languages are just one aspect of formalization; for some applications, experimenting with formalized language designs is equally interesting. It should be easy for the language designer to execute typechecking algorithms, generate sample program behaviors, and—most importantly—test real language implementations against the formalized definitions.
- *Component reuse*. To further facilitate experimentation with and sharing of language designs, the infrastructure should support some way of reusing prior definitions and parts of proofs.

We have carefully constructed the POPLMARK Challenge to stress these features; a theorem-proving infrastructure that addresses the whole challenge should be applicable across a wide spectrum of programming language theory. While we believe that the features above are essential, our challenge does not address many other interesting and tricky-to-formalize constructs and reasoning principles. We discuss possible extensions to the challenge in Section 4.

2.2 Evaluation Criteria

A solution to the POPLMARK Challenge will consist of appropriate software tools, a language representation strategy, and a demonstration that this infrastructure is sufficient to formalize the problems described in Section 3. The long

version of this paper (available at our web site) includes an appendix with reasonably detailed informal proofs of the challenge properties. Solutions to the challenge should follow the overall structure of these proofs, though we expect that details will vary from prover to prover and across term representations. In all cases, there must be an argument for why the formalization is equivalent to the presentation in Section 3—i.e., an adequacy theorem—which should be as simple as possible.

The primary metric of success (beyond correctness, of course) is that a solution should give us confidence of future success of other formalizations carried out using similar techniques. In particular, this implies that:

- *The technology should impose reasonable overheads.* We accept that there is a cost to formalization, and our goal is *not* to be able to prove things more easily than by hand (although that would certainly be welcome). We are willing to spend more time and effort to use the proof infrastructure, but the overhead of doing so must not be prohibitive. (For example, as we discuss below, our experience is that explicit de Bruijn-indexed representations of variable binding structure fail this test.)
- *The technology should be transparent.* The representation strategy and proof assistant syntax should not depart too radically from the usual conventions familiar to the technical audience, and the content of the theorems themselves should be apparent to someone not deeply familiar with the theorem proving technology used or the representation strategy chosen.
- *The technology should have a reasonable cost of entry.* The infrastructure should be usable (after, say, one semester of training) by someone who is knowledgeable about programming language theory but not an expert in theorem prover technology.

2.3 Representing Binders

The problem of representing and reasoning about inductively-defined structures with binders is central to the POPLMARK challenges. Representing binders has been recognized as crucial by the theorem proving community, and many different solutions to this problem have been proposed. In our (still limited) experience, none emerge as clear winners. In this section we briefly summarize the main approaches and, where applicable, describe our own experiments using them. Our survey is far from complete and we refrain from drawing any hard conclusions, to give the proponents of each method a chance to try their hand at meeting the challenge.

A first-order, named approach very similar in flavor to standard informal presentations was used by Vestergaard and Brotherston to formalize some metatheory of untyped λ -calculus [35, 36]. Their representation requires that each binder initially be assigned a unique name—one aspect of the so-called Barendregt convention.

Another popular concrete representation is de Bruijn’s nameless representation. De Bruijn indices are easy to understand and support the full range of

induction principles needed to reason over terms. In our experience, however, de Bruijn representations have two major flaws. First, the *statements* of theorems require complicated clauses involving “shifted” terms and contexts. These extra clauses make it difficult to see the correspondence between informal and formal versions of the same theorem—there is no question of simply typesetting the formal statement and pasting it into a paper. Second, while the notational clutter is manageable for “toy” examples of the size of the simply-typed lambda calculus, we have found it becomes quite a heavy burden even for fairly small languages like $F_{<}$.

In their formalization of properties of pure type systems, McKinna and Pollack use a hybrid approach that combines the above two representation strategies. In this approach, free variables are ordinary names while bound variables are represented using de Bruijn indices [18].

A radically different approach to representing terms with binders is higher-order abstract syntax (HOAS) [28]. In HOAS representations, binders in the meta-language are used to represent binders in the object language. Our experience with HOAS encodings (mainly as realized in Twelf) is that they provide a conveniently high level of abstraction, encapsulating much of the complexity of reasoning about binders. However, the strengths and limitations of the approach are not yet clearly understood, and it can sometimes require significant ingenuity to encode particular language features or proof ideas in this style.

Gordon and Melham propose a way to axiomatize inductive reasoning over untyped lambda-terms [11] and suggest that other inductive structures with binding can be encoded by setting up a correspondence with the untyped lambda terms. Norrish has pursued this direction [26, 27], but observes that these axioms are cumbersome to use without some assistance from the theorem-proving tool. In particular, the axioms use universal quantification in inductive hypotheses where in informal practice “some/any” quantification is used. He has developed a library of lemmas about a system of permutations on top of the axioms that aids reasoning significantly.

Several recent approaches to binding take the concept of “swapping” as a primitive, and use it to build a nominal logic. Gabbay and Pitts proposed a method of reasoning about binders based upon a set theory extended with an intrinsic notion of permutation [8]. Pitts followed this up by proposing a new “nominal” logic based upon the idea of permutations [30]. More recent work by Urban proposes methods based on the same intuitions but carried out within a conventional logic [33]. Our own preliminary experiments with Urban’s methods have been encouraging.

3 The Challenge

Our challenge problems are taken from the basic metatheory of System $F_{<}$. This system is formed by enriching the types and terms of System F with a subtype relation, refining universal quantifiers to carry subtyping constraints,

and adding records, record subtyping, and record patterns. Our presentation is based on Pierce’s *Types and Programming Languages* [29].

The challenge comprises three distinct parts. The first deals just with the type language of $F_{<}$; the second considers terms, evaluation, and type soundness. Each of these is further subdivided into two parts, starting with definitions and properties for *pure* $F_{<}$ and then asking that the same properties be proved for $F_{<}$ enriched with records and patterns. This partitioning allows the development to start small, but also—and more importantly—focuses attention on issues of reuse: How much of the first sub-part can be re-used verbatim in the second sub-part? The third problem asks that useful algorithms be extracted from the earlier formal definitions and used to “animate” some simple properties.

Challenge 1A: Transitivity of Subtyping

The first part of this challenge problem deals purely with the type language of $F_{<}$. The syntax for this language is defined by the following grammar and inference rules. Although the grammar is simple—it has only four syntactic forms—some of its key properties require fairly sophisticated reasoning.

Syntax:

$T ::=$	<i>types</i>
X	<i>type variable</i>
Top	<i>maximum type</i>
$T \rightarrow T$	<i>type of functions</i>
$\forall X < : T. T$	<i>universal type</i>
$\Gamma ::=$	<i>type environments</i>
\emptyset	<i>empty type env.</i>
$\Gamma, X < : T$	<i>type variable binding</i>

In $\forall X < : T_1. T_2$, the variable X is a binding occurrence with scope T_2 (X is not bound in T_1). In $\Gamma, X < : T$, the X must not be in the domain of Γ , and the free variables of T must all be in the domain of Γ .

Following standard practice, issues such as the use of α -conversion, capture avoidance during substitution, etc. are left implicit in what follows. There are several ways in which these issues can be formalized: we might take Γ as a concrete structure (such as an association list of named variables and types) but quotient types and terms up to alpha equivalence, or we could take entire judgments up to alpha equivalence. We might axiomatize the well-formedness of typing environments and types using auxiliary $\vdash \Gamma \text{ ok}$ and $\vdash T \text{ ok}$ judgments. And so on. We leave these decisions to the individual formalization.

It is acceptable to make small changes to the rules below to reflect these decisions, such as adding well-formedness premises. Changing the presentation of the rules to a notationally different but “obviously equivalent” style such as HOAS is also acceptable, but there must be a clear argument that it is really equivalent. Also, whatever formalization is chosen should make clear that we are

only dealing with *well-scoped* terms. For example, it should not be possible to derive $X <: \text{Top}$ in the empty typing environment.

The subtyping relation captures the intuition “if S is a subtype of T (written $S <: T$) then an instance of S may safely be used wherever an instance of T is expected.” It is defined as the least relation closed under the following rules.

Subtyping

$$\boxed{\Gamma \vdash S <: T}$$

$$\Gamma \vdash S <: \text{Top} \quad (\text{SA-TOP})$$

$$\Gamma \vdash X <: X \quad (\text{SA-REFL-TVAR})$$

$$\frac{X <: U \in \Gamma \quad \Gamma \vdash U <: T}{\Gamma \vdash X <: T} \quad (\text{SA-TRANS-TVAR})$$

$$\frac{\Gamma \vdash T_1 <: S_1 \quad \Gamma \vdash S_2 <: T_2}{\Gamma \vdash S_1 \rightarrow S_2 <: T_1 \rightarrow T_2} \quad (\text{SA-ARROW})$$

$$\frac{\Gamma \vdash T_1 <: S_1 \quad \Gamma, X <: T_1 \vdash S_2 <: T_2}{\Gamma \vdash \forall X <: S_1. S_2 <: \forall X <: T_1. T_2} \quad (\text{SA-ALL})$$

These rules present an *algorithmic* version of the subtyping relation. In contrast to the more familiar *declarative* presentation, these rules are syntax-directed, as might be found in the implementation of a type checker; the algorithmic rules are also somewhat easier to reason with, having, for example, an obvious inversion property. The declarative rules differ from these by explicitly stating that subtyping is reflexive and transitive. However, reflexivity and transitivity also turn out to be derivable properties in the algorithmic system. A straightforward induction shows that the algorithmic rules are reflexive. The first challenge is to show that they are also transitive.

3.1 LEMMA [TRANSITIVITY OF ALGORITHMIC SUBTYPING]: If $\Gamma \vdash S <: Q$ and $\Gamma \vdash Q <: T$, then $\Gamma \vdash S <: T$. \square

The difficulty here lies in the reasoning needed to prove this lemma. Transitivity must be proven simultaneously with another property, called *narrowing*, by an inductive argument with case analyses on the final rules used in the given derivations.

3.2 LEMMA [NARROWING]: If $\Gamma, X <: Q, \Delta \vdash M <: N$ and $\Gamma \vdash P <: Q$ then $\Gamma, X <: P, \Delta \vdash M <: N$. \square

Challenge 1B: Transitivity of Subtyping with Records

We now extend this challenge by enriching the type language with record types. The new syntax and subtyping rule for record types are shown below. Implicit

in the syntax is the condition that the labels $\{l_i^{i \in 1..n}\}$ appearing in a record type $\{l_i : T_i^{i \in 1..n}\}$ are pairwise distinct.

New syntactic forms:

$T ::= \dots$	<i>types</i>
$\{l_i : T_i^{i \in 1..n}\}$	<i>type of records</i>

New subtyping rules

$\Gamma \vdash S <: T$

$$\frac{\{l_i^{i \in 1..n}\} \subseteq \{k_j^{j \in 1..m}\} \quad \text{if } k_j = l_i, \text{ then } \Gamma \vdash S_j <: T_i}{\Gamma \vdash \{k_j : S_j^{j \in 1..m}\} <: \{l_i : T_i^{i \in 1..n}\}} \quad (\text{SA-RCD})$$

Although it has been shown that records can actually be encoded in pure $F_{<}$ [2, 10], dealing with them directly is a worthwhile task since, unlike other syntactic forms, record types have an arbitrary (finite) number of fields. Also, the informal proof for Challenge 1A extends to record types by only adding the appropriate cases. A formal proof should reflect this.

Challenge 2A: Type Safety for Pure $F_{<}$

The next challenge considers the type soundness of pure $F_{<}$ (without record types, for the moment). Below, we complete the definition of $F_{<}$ by describing the syntax of terms, values, and typing environments with term binders and giving inference rules for the typing relation and a small-step operational semantics.

As usual in informal presentations, we elide the formal definition of substitution and simply assume that the substitutions of a type P for X in T (denoted $[X \mapsto P]T$) and of a term q for x in t (denoted $[x \mapsto q]t$) are capture-avoiding.

Syntax:

$t ::=$	<i>terms</i>
x	<i>variable</i>
$\lambda x : T. t$	<i>abstraction</i>
$t \ t$	<i>application</i>
$\lambda X < : T. t$	<i>type abstraction</i>
$t \ [T]$	<i>type application</i>
$v ::=$	<i>values</i>
$\lambda x : T. t$	<i>abstraction value</i>
$\lambda X < : T. t$	<i>type abstraction value</i>
$\Gamma ::=$	<i>type environments</i>
\emptyset	<i>empty type env.</i>
$\Gamma, x : T$	<i>term variable binding</i>
$\Gamma, X < : T$	<i>type variable binding</i>

Typing

$$\boxed{\Gamma \vdash \mathbf{t} : \mathbf{T}}$$

$$\frac{\mathbf{x} : \mathbf{T} \in \Gamma}{\Gamma \vdash \mathbf{x} : \mathbf{T}} \quad (\text{T-VAR})$$

$$\frac{\Gamma, \mathbf{x} : \mathbf{T}_1 \vdash \mathbf{t}_2 : \mathbf{T}_2}{\Gamma \vdash \lambda \mathbf{x} : \mathbf{T}_1. \mathbf{t}_2 : \mathbf{T}_1 \rightarrow \mathbf{T}_2} \quad (\text{T-ABS})$$

$$\frac{\Gamma \vdash \mathbf{t}_1 : \mathbf{T}_{11} \rightarrow \mathbf{T}_{12} \quad \Gamma \vdash \mathbf{t}_2 : \mathbf{T}_{11}}{\Gamma \vdash \mathbf{t}_1 \ \mathbf{t}_2 : \mathbf{T}_{12}} \quad (\text{T-APP})$$

$$\frac{\Gamma, \mathbf{X} < : \mathbf{T}_1 \vdash \mathbf{t}_2 : \mathbf{T}_2}{\Gamma \vdash \lambda \mathbf{X} < : \mathbf{T}_1. \mathbf{t}_2 : \forall \mathbf{X} < : \mathbf{T}_1. \mathbf{T}_2} \quad (\text{T-TABS})$$

$$\frac{\Gamma \vdash \mathbf{t}_1 : \forall \mathbf{X} < : \mathbf{T}_{11}. \mathbf{T}_{12} \quad \Gamma \vdash \mathbf{T}_2 < : \mathbf{T}_{11}}{\Gamma \vdash \mathbf{t}_1 \ [\mathbf{T}_2] : [\mathbf{X} \mapsto \mathbf{T}_2] \mathbf{T}_{12}} \quad (\text{T-TAPP})$$

$$\frac{\Gamma \vdash \mathbf{t} : \mathbf{S} \quad \Gamma \vdash \mathbf{S} < : \mathbf{T}}{\Gamma \vdash \mathbf{t} : \mathbf{T}} \quad (\text{T-SUB})$$

Evaluation

$$\boxed{\mathbf{t} \longrightarrow \mathbf{t}'}$$

$$(\lambda \mathbf{x} : \mathbf{T}_{11}. \mathbf{t}_{12}) \ \mathbf{v}_2 \longrightarrow [\mathbf{x} \mapsto \mathbf{v}_2] \mathbf{t}_{12} \quad (\text{E-APPABS})$$

$$(\lambda \mathbf{X} < : \mathbf{T}_{11}. \mathbf{t}_{12}) \ [\mathbf{T}_2] \longrightarrow [\mathbf{X} \mapsto \mathbf{T}_2] \mathbf{t}_{12} \quad (\text{E-TAPPTABS})$$

Evaluation contexts:

$E ::=$

$[-]$
 $E \ \mathbf{t}$
 $\mathbf{v} \ E$
 $E \ [\mathbf{T}]$

evaluation contexts

hole
 app fun
 app arg
 type fun

Evaluation in context:

$$\frac{\mathbf{t}_1 \longrightarrow \mathbf{t}'_1}{E[\mathbf{t}_1] \longrightarrow E[\mathbf{t}'_1]} \quad (\text{E-CTX})$$

The evaluation relation is presented in two parts: the rules E-APPABS and E-TAPPTABS capture the immediate reduction rules of the language, while E-CTX permits reduction under an arbitrary *evaluation context* E . For $F_{<}$, one would have an equally clear definition and slightly simpler proofs using explicit closure rules for the evaluation relation. We use evaluation contexts with an eye to larger languages and languages with non-local control operators such as

exceptions, for which (in informal mathematics) they are an important tool for reducing notational clutter in definitions.⁴ Evaluation contexts are also particularly interesting from the point of view of formalization when they include binders, though unfortunately there are no examples of this in call-by-value F_{\leq} .

Type soundness is usually proven in the style popularized by Wright and Felleisen [37], in terms of *preservation* and *progress* theorems. Challenge 2A is to prove these properties for pure F_{\leq} .

3.3 THEOREM [PRESERVATION]: If $\Gamma \vdash t : T$ and $t \longrightarrow t'$, then $\Gamma \vdash t' : T$. \square

3.4 THEOREM [PROGRESS]: If t is a closed, well-typed F_{\leq} term (i.e., if $\vdash t : T$ for some T), then either t is a value or else there is some t' with $t \longrightarrow t'$. \square

Unlike the proof of transitivity of subtyping, the inductive arguments required here are straightforward. However, variable binding becomes a more significant issue, since this language includes binding of both type and term variables. Several lemmas relating to both kinds of binding must also be shown, in particular lemmas about type and term substitutions. These lemmas, in turn, require reasoning about permuting, weakening, and strengthening typing environments.

Challenge 2B: Type Safety with Records and Pattern Matching

The next challenge is to extend the preservation and progress results to cover records and pattern matching. The new syntax and rules for this language appear below. As for record types, the labels $\{l_i \mid i \in 1..n\}$ appearing in a record $\{l_i = t_i \mid i \in 1..n\}$ are assumed to be pairwise distinct. Similarly, the variable patterns appearing in a pattern are assumed to bind pairwise distinct variables.

New syntactic forms:

$t ::=$...	<i>terms</i>
	$\{l_i = t_i \mid i \in 1..n\}$	<i>record</i>
	$t.l$	<i>projection</i>
	$\text{let } p = t \text{ in } t$	<i>pattern binding</i>
$p ::=$		<i>patterns</i>
	$x : T$	<i>variable pattern</i>
	$\{l_i = p_i \mid i \in 1..n\}$	<i>record pattern</i>

⁴ This design choice has generated a robust debate on the POPLMARK discussion list as to whether evaluation contexts *must* be used in order for a solution to count as valid, or whether an “obviously equivalent” presentation such as an evaluation relation with additional congruence rules is acceptable. We prefer evaluation contexts for the reasons we have given, but the consensus of the community appears to be that one should accept solutions in other styles. However, a good solution must be formulated in a style that provides similar clarity as the language scales.

$v ::= \dots$
 $\{l_i = v_i \mid i \in I..n\}$

values
record value

New typing rules

$\boxed{\Gamma \vdash t : T}$

$$\frac{\Gamma \vdash t_1 : T_1 \quad \vdash p : T_1 \Rightarrow \Delta \quad \Gamma, \Delta \vdash t_2 : T_2}{\Gamma \vdash \text{let } p = t_1 \text{ in } t_2 : T_2} \quad (\text{T-LET})$$

$$\frac{\text{for each } i \quad \Gamma \vdash t_i : T_i}{\Gamma \vdash \{l_i = t_i \mid i \in I..n\} : \{l_i : T_i \mid i \in I..n\}} \quad (\text{T-RCD})$$

$$\frac{\Gamma \vdash t_1 : \{l_i : T_i \mid i \in I..n\}}{\Gamma \vdash t_1.l_j : T_j} \quad (\text{T-PROJ})$$

Pattern typing rules:

$$\vdash (x : T) : T \Rightarrow x : T \quad (\text{P-VAR})$$

$$\frac{\text{for each } i \quad \vdash p_i : T_i \Rightarrow \Delta_i}{\vdash \{l_i = p_i \mid i \in I..n\} : \{l_i : T_i \mid i \in I..n\} \Rightarrow \Delta_n, \dots, \Delta_1} \quad (\text{P-RCD})$$

New evaluation rules

$\boxed{t \longrightarrow t'}$

$$\text{let } p = v_1 \text{ in } t_2 \longrightarrow \text{match}(p, v_1)t_2 \quad (\text{E-LETV})$$

$$\{l_i = v_i \mid i \in I..n\}.l_j \longrightarrow v_j \quad (\text{E-PROJRCD})$$

New evaluation contexts:

$E ::= \dots$
 $E.l$
 $\{l_i = v_i \mid i \in I..j-1, l_j = E, l_k = t_k \mid k \in j+1..n\}$
 $\text{let } p = E \text{ in } t_2$

evaluation contexts
projection
record
let binding

Matching rules:

$$\text{match}(x : T, v) = [x \mapsto v] \quad (\text{M-VAR})$$

$$\frac{\{l_i \mid i \in I..n\} \subseteq \{k_j \mid j \in I..m\} \quad \text{if } l_i = k_j, \text{ then } \text{match}(p_i, v_j) = \sigma_i}{\text{match}(\{l_i = p_i \mid i \in I..n\}, \{k_j = v_j \mid j \in I..m\}) = \sigma_n \circ \dots \circ \sigma_1} \quad (\text{M-RCD})$$

Compared to the language of Challenge 2A, the `let` construct is a fundamentally new binding form, since patterns may bind an arbitrary (finite) number of term variables.

Challenge 3: Testing and Animating with Respect to the Semantics

Given a complete formal definition of a language, there are at least two interesting ways in which it can be used (as opposed to being reasoned about). When implementing the language, it should be possible to use the formal definition as an oracle for *testing* the implementation—checking that it does conform to the definition by running test cases in the implementation and confirming formally that the outcome is as prescribed. Secondly, one would like to construct a prototype implementation from the definition and use it for *animating* the language, i.e., exploring the language’s properties on particular examples. In both cases, this should be done without any unverified (and thus error-prone) manual translation of definitions.

Our final challenge is to provide an implementation of this functionality, specifically for the following three tasks (using the language of Challenge 2B):

1. Given F_{\leq} terms t and t' , decide whether $t \longrightarrow t'$.
2. Given F_{\leq} terms t and t' , decide whether $t \longrightarrow^* t' \not\rightarrow$, where \longrightarrow^* is the reflexive-transitive closure of \longrightarrow .
3. Given an F_{\leq} term t , find a term t' such that $t \longrightarrow t'$.

The first two subtasks are useful for testing language implementations, while the last is useful for animating the definition. For all three subtasks, the system(s) should accept syntax that is “reasonably close” to that of informal (ASCII) mathematical notation, though it may be necessary to translate between the syntaxes of a formal environment and an implementation. We will provide an implementation of an interpreter for F_{\leq} with records and patterns at the challenge’s website in order to make this challenge concrete, together with a graded sequence of example terms. To make a rough performance comparison possible, solutions should indicate execution times for these terms.

A solution to this challenge might make use of decision procedures and tactics of a proof assistant or might extract stand-alone code. In general, it may be necessary to combine theorems (e.g. that a rule-based but algorithmic definition of typing coincides with a declarative definition) and proof search (e.g. deciding particular instances of the algorithmic definition).

4 Beyond the Challenge

The POPLMARK Challenge is not meant to be exhaustive: other aspects of programming language theory raise formalization difficulties that are interestingly different from the problems we have proposed—to name a few: more complex binding constructs such as mutually recursive definitions, logical relations proofs, coinductive simulation arguments, undecidability results, and linear handling of type environments. As time goes on, we will issue a small number of further challenges highlighting the most important of these issues; suggestions from the community would be welcome. However, we believe that a technology that provides a good solution to the POPLMARK challenge as we have formulated it

here will be sufficient to attract eager adopters in the programming languages community, beginning with the authors.

So what are you waiting for? It's time to bring mechanized metatheory to the masses!

Acknowledgments

A number of people have joined us in preliminary discussions of these challenge problems, including Andrew Appel, Karl Crary, Frank Pfenning, Bob Harper, Hugo Herbelin, Jason Hickey, Michael Norrish, Andrew Pitts, Randy Pollack, Carsten Schürmann, Phil Wadler, and Dan Wang. We are especially grateful to Randy Pollack for helping guide our earliest efforts at formalization and to Randy Pollack, Michael Norrish, and Carsten Schürmann for their own work on the challenge problems. Discussions on the POPLMARK mailing list have greatly deepened our understanding of the current state of the art.

Sewell and Fairbairn acknowledge support from a Royal Society University Research Fellowship and EPSRC grant GR/T11715.

References

1. Andrew W. Appel. Foundational proof-carrying code. In *IEEE Symposium on Logic in Computer Science (LICS)*, Boston, Massachusetts, pages 247–258, June 2001.
2. Luca Cardelli. Extensible records in a pure calculus of subtyping. Research report 81, DEC/Compaq Systems Research Center, January 1992. Also in C. A. Gunter and J. C. Mitchell, editors, *Theoretical Aspects of Object-Oriented Programming: Types, Semantics, and Language Design*, MIT Press, 1994.
3. Luca Cardelli, Simone Martini, John C. Mitchell, and Andre Scedrov. An extension of System F with subtyping. *Information and Computation*, 109(1–2):4–56, 1994. Summary in TACS '91 (Sendai, Japan, pp. 750–770).
4. Karl Crary. Toward a foundational typed assembly language. In *ACM SIGPLAN–SIGACT Symposium on Principles of Programming Languages (POPL)*, New Orleans, Louisiana, pages 198–212, January 2003.
5. Louise A. Dennis. Inductive challenge problems, 2000. <http://www.cs.nott.ac.uk/~lad/research/challenges>.
6. Catherine Dubois and Valerie Menissier-Morain. Certification of a type inference tool for ML: Damas-Milner within Coq. *Journal of Automated Reasoning*, 23(3–4):319–346, 1999.
7. Robert W. Floyd. Assigning meanings to programs. In J. T. Schwartz, editor, *Mathematical Aspects of Computer Science*, volume 19 of *Proceedings of Symposia in Applied Mathematics*, pages 19–32, Providence, Rhode Island, 1967. American Mathematical Society.
8. Murdoch Gabbay and Andrew Pitts. A new approach to abstract syntax involving binders. In *14th Symposium on Logic in Computer Science*, pages 214–224, 1999.
9. I. P. Gent and T. Walsh. CSPLib: a benchmark library for constraints. Technical Report APES-09-1999, APES, 1999. Available from <http://www.dcs.st-and.ac.uk/~apes/apesreports.html>. A shorter version appears in the Proceedings of the

- 5th International Conference on Principles and Practices of Constraint Programming (CP-99).
10. Giorgio Ghelli. *Proof Theoretic Studies about a Minimal Type System Integrating Inclusion and Parametric Polymorphism*. PhD thesis, Università di Pisa, 1990. Technical report TD-6/90, Dipartimento di Informatica, Università di Pisa.
 11. Andrew D. Gordon and Tom Melham. Five axioms of alpha-conversion. In J. von Wright, J. Grundy, and J. Harrison, editors, *Theorem Proving in Higher Order Logics: 9th International Conference, TPHOLs '96*, volume 1125 of *Lecture Notes in Computer Science*, pages 173–190. Springer-Verlag, 1996.
 12. Ian Green. The dream corpus of inductive conjectures, 1999. <http://dream.dai.ed.ac.uk/dc/lib.html>.
 13. Elsa Gunter and Savi Maharaj. Studying the ML module system in HOL. *The Computer Journal: Special Issue on Theorem Proving in Higher Order Logics*, 38(2):142–151, 1995.
 14. Tony Hoare. The verifying compiler: A grand challenge for computing research. *J. ACM*, 50(1):63–69, 2003.
 15. Holger Hoos and Thomas Stuetzle. SATLIB. <http://www.intellektik.informatik.tu-darmstadt.de/SATLIB/>.
 16. J. J. Joyce and C.-J. H. Seger, editors. *Proceedings of the 6th International Workshop on Higher Order Logic Theorem Proving and its Applications (HUG'93)*, volume 780 of *Lecture Notes in Computer Science*, Vancouver, B.C., Canada, August 1993. Springer-Verlag, 1994.
 17. Gerwin Klein and Tobias Nipkow. A machine-checked model for a Java-like language, virtual machine and compiler. Technical Report 0400001T.1, National ICT Australia, Sydney, March 2004.
 18. James McKinna and Robert Pollack. Some lambda calculus and type theory formalized. *Journal of Automated Reasoning*, 23(3–4), November 1999.
 19. Robin Milner, Mads Tofte, Robert Harper, and David MacQueen. *The Definition of Standard ML*, Revised edition. MIT Press, 1997.
 20. J Strother Moore. A grand challenge proposal for formal methods: A verified stack. In Bernhard K. Aichernig and T. S. E. Maibaum, editors, *Formal Methods at the Crossroads. From Panacea to Foundational Support, 10th Anniversary Colloquium of UNU/IIST, Lisbon, Portugal*, volume 2757 of *Lecture Notes in Computer Science*, pages 161–172. Springer, 2002.
 21. J Strother Moore and George Porter. The apprentice challenge. *ACM Trans. Program. Lang. Syst.*, 24(3):193–216, 2002.
 22. Wolfgang Naraschewski and Tobias Nipkow. Type inference verified: Algorithm W in Isabelle/HOL. *Journal of Automated Reasoning*, 23:299–318, 1999.
 23. Tobias Nipkow, David von Oheimb, and Cornelia Pusch. μ Java: Embedding a programming language in a theorem prover. In F.L. Bauer and R. Steinbrüggen, editors, *Foundations of Secure Computation. Proc. Int. Summer School Marktoberdorf 1999*, pages 117–144. IOS Press, 2000.
 24. Tobias Nipkow and David von Oheimb. Java_{light} is type-safe—definitely. In *POPL '98: Proceedings of the 25th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 161–170, New York, NY, USA, 1998. ACM Press.
 25. Michael Norrish. *Formalising C in HOL*. PhD thesis, Computer Laboratory, University of Cambridge, 1998.
 26. Michael Norrish. Mechanising Hankin and Barendregt using the Gordon-Melham axioms. In *MERLIN '03: Proceedings Of The 2003 Workshop On Mechanized Reasoning About Languages With Variable Binding*, pages 1–7. ACM Press, 2003.

27. Michael Norrish. Recursive function definition for types with binders. In Konrad Slind, Annette Bunker, and Ganesh Gopalakrishnan, editors, *Theorem Proving in Higher Order Logics: 17th International Conference, TPHOLs 2004*, volume 3223 of *Lecture Notes In Computer Science*, pages 241–256. Springer-Verlag, 2004.
28. Frank Pfenning and Conal Elliot. Higher-order abstract syntax. In *PLDI '88: Proceedings of the ACM SIGPLAN 1988 Conference on Programming Language Design and Implementation*, pages 199–208. ACM Press, 1988.
29. Benjamin C. Pierce. *Types and Programming Languages*. MIT Press, 2002.
30. Andrew M. Pitts. Nominal logic, a first order theory of names and binding. *Inf. Comput.*, 186(2):165–193, 2003.
31. Geoff Sutcliffe and Christian Suttner. The TPTP problem library. *Journal of Automated Reasoning*, 21(2):177–203, 1998.
32. Donald Syme. Reasoning with the formal definition of Standard ML in HOL. In Joyce and Seger [16], pages 43–60.
33. Christian Urban and Christine Tasson. Nominal techniques in Isabelle/HOL. Accepted at CADE-20 in Tallinn. See <http://www.mathematik.uni-muenchen.de/~urban/nominal/>.
34. Myra VanInwegen and Elsa Gunter. HOL-ML. In Joyce and Seger [16], pages 61–74.
35. René Vestergaard and James Brotherston. The mechanisation of Barendregt-style equational proofs (the residual perspective). In *Mechanized Reasoning about Languages with Variable Binding (MERLIN)*, volume 58 of *Electronic Notes in Theoretical Computer Science*. Elsevier, 2001.
36. René Vestergaard and James Brotherston. A formalised first-order confluence proof for the λ -calculus using one-sorted variable names. *Information and Computation*, 183(2):212 – 244, 2003. Special edition with selected papers from RTA01.
37. Andrew K. Wright and Matthias Felleisen. A syntactic approach to type soundness. *Information and Computation*, 115(1):38–94, November 1994.