

Modularly Programmable Syntax and Type Structure

Cyrus Omar

July 30, 2015

School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213

Thesis Committee:

Jonathan Aldrich, Chair

TODO: confirm rest of committee

*Submitted in partial fulfillment of the requirements
for the degree of Doctor of Philosophy.*

Copyright © 2015 Cyrus Omar

July 30, 2015
DRAFT

Abstract

Functional programming languages like ML descend conceptually from minimal lambda calculi, but to be pragmatic, expose a concrete syntax and type structure to programmers of a more elaborate design. Language designers have many viable choices along these dimensions, as evidenced by the diversity of dialects that continue to proliferate around these languages. But such language dialects cannot be modularly combined, limiting the choices available to programmers. We describe and formally specify new language primitives designed to decrease the need for dialects by giving library providers the ability to safely and modularly control syntactic expansion, typechecking and translation to a minimal type-theoretic internal language.

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Contributions	3
2	Language Overview	7
2.1	External Language	7
2.2	Internal Language	7
2.3	Static Language	7
2.4	Module Language	7
I	Modularly Programmable Syntax	9
3	Motivation	11
3.1	Motivating Examples	11
3.1.1	Lists	11
3.1.2	HTML	11
3.1.3	Regular Expressions	11
3.1.4	Monadic Commands	11
3.1.5	Quasiquotation	11
3.2	Existing Approaches	11
3.2.1	Dynamic String Parsing	11
3.2.2	Direct Syntax Extension	11
3.2.3	Term Rewriting	11
4	Typed Syntax Macros	13
4.1	Examples	13
4.2	Minimal Formalization	13
4.3	Parameterized TSMs	13
5	Type-Specific Languages	15
5.1	Examples	15
5.2	Minimal Formalization	15
5.3	Parameterized TSLs	15

5.4	Conclusion & Future Work	15
II	Modularly Programmable Type Structure	17
6	Motivation	19
6.1	Motivating Examples	19
6.2	Existing Approaches	19
7	Metamodules	21
8	Conclusion & Future Work	23
	Bibliography	25

List of Figures

Chapter 1

Introduction

1.1 Motivation

Functional programming language designers often turn to minimal typed lambda calculi to develop a principled understanding of fundamental metatheoretic issues, like type safety, and to examine the essential character of various primitives of interest. To design a “full-scale”¹ language, language designers must carefully “combine” such primitives and develop various generalizations and primitive “embellishments” motivated by human factors. For example, major functional programming languages like Standard ML (SML) [11, 16], OCaml [13] and Haskell [12] all primitively build in record types (generalizing the nullary and binary product types that suffice in simpler settings) because explicitly labeled components are cognitively useful to human programmers. Similarly, they all build in derived syntax (colloquially, “syntactic sugar”) that decreases the syntactic cost of common idioms, like list construction.

One might hope that a limited number of semantic primitives and primitive “embellishments” like these will suffice to go from minimal calculi to a broadly useful (or “general-purpose”) programming language. But a stable language design that fully achieves this ideal has yet to emerge, as evidenced by the diverse array of “dialects” that continue to proliferate around all major contemporary languages. Indeed, tools that aid in the construction of so-called “domain-specific” language dialects (DSLs)² are becoming increasingly prominent. This calls for an investigation: why is it that programmers and researchers are still so often unable to satisfyingly express the constructs that they need in libraries, as modes of use of the “general-purpose” primitives already available in major languages today?

Why are there so many dialects?

Perhaps the most common reason for this ongoing proliferation of dialects may simply be that the *syntactic cost* of expressing a construct of interest using contemporary general-purpose primitives is not always ideal. In response, library providers construct *syntactic dialects*, i.e. dialects

¹Throughout this work, colloquialisms that should be read as having an intuitive meaning, rather than a strict mathematical meaning, will be first introduced with quotation marks.

²Such dialects are also sometimes called “external DSLs”, to distinguish them from “internal” or “embedded DSLs”, which are actually library interfaces that “resemble” distinct dialects [8].

that can be specified by purely syntactic (i.e. context-independent) elaboration to the existing language. In other words, they introduce only new derived forms. For example, Ur/Web is a syntactic dialect of Ur (a language that itself descends from ML) that builds in derived syntax for SQL queries, HTML elements and other datatypes used in the domain of web programming [2]. These are far from the only types of data that could similarly benefit from the availability of specialized derived syntax. As another example, we will consider regular expression patterns expressed using abstract data types (as a mode of use of the ML module system) in Sec. ???. Tools like Camlp4 [13], Sugar* [3, 4] and Racket [6] have lowered the engineering costs of constructing syntactic dialects in such situations, contributing to their proliferation.

More advanced dialects introduce new type structure, going beyond what is possible with derived syntax. As a simple example, the static and dynamic semantics of records cannot be expressed by context-independent elaboration to a language with only nullary and binary products. Various languages have explored “record-like” primitives that go still further, supporting “functional update” operators, width and depth coercions (sometimes implicit), methods, prototypic dispatch and other such “semantic embellishments” that cannot be expressed by context-independent expansion to a language with only standard record types (we will detail an example in Sec. ???). OCaml primitively builds in the type structure of polymorphic variants, open datatypes and operations that use format strings like `sprintf` [13]. ReactiveML builds in primitives for functional reactive programming [14]. ML5 builds in high-level primitives for distributed programming based on a modal lambda calculus [17]. Manticore [7] and AliceML [19] build in parallel programming primitives with a more elaborate type structure than is found in simpler accounts of parallelism. MLj builds in the type structure of the Java object system (motivated by a desire to interface safely and naturally with Java libraries) [1]. Other dialects do the same for other foreign languages, e.g. Furr and Foster describe a dialect of OCaml that builds in the type structure of C [9]. Tools like proof assistants and logical frameworks are used to specify and reason metatheoretically about dialects like these, and tools like compiler generators and language frameworks [5] lower their implementation cost, again contributing to their proliferation.

Dialects Considered Harmful

The reason why this proliferation of language dialects should be considered alarming is that it is, in an important sense, anti-modular: a library written in one dialect cannot, in general, safely and idiomatically interface with a library written in another dialect. As MLj demonstrates, addressing this interoperability problem requires somehow “combining” the dialects into a single language. However, in the most general setting where the dialects in question might be specified by judgements of arbitrary form, this is not a well-defined notion. Even if we restrict our interest to dialects specified using formalisms that do operationalize the notion of dialect combination, there is generally no guarantee that the combined dialect will conserve important syntactic and semantic properties that can be established about the dialects in isolation. For example, consider two syntactic dialects, one specifying derived syntax for finite mappings, the other specifying a similar syntax for ordered finite mappings. Though each dialect can be shown to have an unambiguous concrete syntax in isolation, when their grammars are naïvely combined by, for example, Camlp4, ambiguities arise. Due to this paucity of modular reasoning principles, the

“dialect-oriented” approach (also called the “language-oriented approach” [21]) is not sensible for software development “in the large”.

Dialect designers must instead take a less direct approach to have an impact on large-scale software development: they must convince the designers in control of comparatively popular languages, like OCaml and Scala, to include some suitable variant of the primitives they espouse into backwards compatible language revisions. This *ad hoc* approach is not sustainable, for three main reasons. First, as suggested by the diversity of examples given above, there are simply too many potentially useful such primitives, and many of these are only relevant in relatively narrow application domains (for derived syntax, our group has gathered initial data speaking to this [18]). Second, primitives introduced earlier in a language’s lifespan end up monopolizing finite “syntactic resources”, forcing subsequent primitives to use ever more esoteric forms. And third, primitives that prove to be flawed in some way cannot be removed or changed without breaking backwards compatibility.

This leaves the subset of the language design community interested in keeping general-purpose languages small and free of *ad hoc* primitives with two possible paths forward. One, exemplified (arguably) by SML, is to simply eschew the introduction of specialized syntax and type structure and settle on the existing primitives, which can be said to sit at a “sweet spot” in the overall language design space (accepting that in some circumstances, this trades away expressive power or leads to high syntactic cost). The other is to search for more general primitives that reduce the primitives found in dialects today to modularly composable library constructs. Encouragingly, primitives of this sort do occasionally arise. For example, a recent revision of OCaml added support for “generalized algebraic data types” (GADTs), based on research on guarded recursive datatype constructors [22]. Using GADTs, OCaml was able to move some of the *ad hoc* machinery for typechecking operations that use format strings, like `sprintf`, out of the language and into a library (however, syntactic machinery remains primitively built in).

1.2 Contributions

Our aim in the work being proposed is to introduce primitive language constructs that take further steps down the second path just described. In particular, we plan to introduce the following constructs:

1. **Typed syntax macros** (TSMs), introduced in Sec. ??, reduce the need to primitively build in derived concrete syntax specific to library constructs, e.g. `list` syntax as in SML or `XML` syntax as in Scala and `Ur/Web`, by giving library providers static control over the parsing and expansion of delimited segments of concrete syntax (at a specified type or parameterized family of types).
2. **Metamodules**, introduced in Sec. ??, reduce the need to primitively build in the type structure of constructs like records (and variants thereof), labeled sums and other more esoteric constructs that we will introduce later by giving library providers programmatic hooks directly into the semantics. We will see direct analogies between ML-style modules and metamodules in Sec. ??.

Both TSMs and metamodules involve *static code generation* (also called *static* or *compile-time metaprogramming*), meaning that the relevant rules in the static semantics of the language call for the evaluation of *static functions* that generate static representations of expressions and types. This design has conceptual roots in earlier work on *active libraries*, which similarly envisioned using compile-time computation to give library providers more control over aspects of the language and compilation process [20].

The main challenge in the design of these primitives will come in ensuring that they are metatheoretically well-behaved. If we are not careful, many of the problems that arise when combining language dialects, discussed earlier, could simply shift into the semantics of these primitives.³ Our main technical contributions will be in rigorously showing how to address these problems in a principled manner. In particular, syntactic conflicts will be impossible by construction and the semantics will validate code statically generated by TSMs and metamodules to maintain a strong *hygienic type discipline* and, most uniquely, powerful *modular reasoning principles*. In other words, library providers will have the ability to reason about the constructs that they have defined in isolation, and clients will be able to use them safely in any program context and in any combination, without the possibility of conflict.⁴ We will make these notions completely precise as we continue.

As vehicles for this work, we plan to formally specify typed lambda calculi that capture each of the novel primitives that we introduce “minimally”. We will also describe (but not formally specify) a new “full-scale” functional language called Verse.⁵ The reason we will not follow Standard ML [16] in giving a complete formal specification of Verse is both to emphasize that the primitives we introduce can be considered for inclusion in a variety of language designs, and to avoid distracting the reader with specifications for “orthogonal” primitives that are already well-understood in the literature. We will give a brief overview covering how these languages are organized in Sec. ??.

Thesis Statement

In summary, we propose a thesis defending the following statement:

A functional programming language can give library providers the ability to express new syntactic expansions and new type structure atop a minimal type-theoretic internal language while maintaining a hygienic type discipline and modular reasoning principles.

Disclaimers

Before we continue, it may be useful to explicitly acknowledge that completely eliminating the need for dialects would indeed be asking for too much: certain design decisions are fundamen-

³This is why languages like Verse are often called *extensible languages*, though this is somewhat of a misnomer. The defining characteristic of an extensible language is that it *doesn't* need to be extended in situations where other languages would need to be extended. We will avoid this somewhat confusing terminology.

⁴This is not quite true – simple naming conflicts can arise. We will tacitly assume that they are being avoided extrinsically, e.g. by using a URI-based naming scheme as in the Java ecosystem.

⁵We distinguish Verse from Wyvern, which is the language referred to in prior publications about some of the work being proposed here, because Wyvern is a group effort evolving independently in some important ways.

tally incompatible with others or require coordination across a language design. We aim only to decrease the need for dialects.

It may also be useful to explicitly acknowledge that library providers could leverage the primitives we introduce to define constructs that are in “poor taste”. We expect that in practice, Verse will come with a standard library defining a carefully curated collection of standard constructs, as well as guidelines for advanced users regarding when it would be sensible to use the mechanisms we introduce (following the example of languages that support operator overloading or type classes [10], which also have the potential for such “abuse”). For most programmers, using Verse should not be substantially different from using a language like ML or one of its dialects.

Finally, Verse intentionally is not a dependently-typed language like Coq, Agda or Idris, because these languages do not maintain a phase separation between “compile-time” and “run-time.” This phase separation is useful for programming tasks (where one would like to be able to discover errors before running a program, particularly programs that may have an effect) but less so for theorem proving tasks (where it is mainly the fact that a pure expression is well-typed that is of interest, by the propositions-as-types principle). Verse is designed to be used for programming tasks where SML, OCaml, Haskell or Scala would be used today, not for advanced theorem proving tasks. That said, we conjecture that the primitives we describe could be added to languages like Gallina (the “external language” of the Coq proof assistant [15]) or to the program extraction mechanisms of proof assistants like Coq with modifications, but do not plan to pursue this line of research in this dissertation.

Chapter 2

Language Overview

2.1 External Language

2.2 Internal Language

2.3 Static Language

2.4 Module Language

Part I

Modularly Programmable Syntax

Chapter 3

Motivation

3.1 Motivating Examples

3.1.1 Lists

3.1.2 HTML

3.1.3 Regular Expressions

3.1.4 Monadic Commands

3.1.5 Quasiquotation

3.2 Existing Approaches

3.2.1 Dynamic String Parsing

3.2.2 Direct Syntax Extension

Related work I haven't mentioned yet:

- Fan: <http://zhanghongbo.me/fan/start.html>
- Well-Typed Islands Parse Faster:
<http://www.ccs.neu.edu/home/ejs/papers/tfp12-island.pdf>

3.2.3 Term Rewriting

Chapter 4

Typed Syntax Macros

4.1 Examples

4.2 Minimal Formalization

4.3 Parameterized TSMs

Chapter 5

Type-Specific Languages

5.1 Examples

5.2 Minimal Formalization

5.3 Parameterized TSLs

5.4 Conclusion & Future Work

Part II

Modularly Programmable Type Structure

Chapter 6

Motivation

6.1 Motivating Examples

6.2 Existing Approaches

Chapter 7

Metamodules

Chapter 8

Conclusion & Future Work

Bibliography

- [1] Nick Benton and Andrew Kennedy. Interlanguage Working Without Tears: Blending SML with Java. In *ICFP '99*, pages 126–137, 1999. ISBN 1-58113-111-9. doi: 10.1145/317636.317791. URL <http://doi.acm.org/10.1145/317636.317791>. 1.1
- [2] Adam Chlipala. Ur/Web: A simple model for programming the web. In *POPL '15*, pages 153–165, 2015. ISBN 978-1-4503-3300-9. URL <http://dl.acm.org/citation.cfm?id=2676726>. 1.1
- [3] Sebastian Erdweg and Felix Rieger. A framework for extensible languages. In *GPCE '13*, pages 3–12, 2013. 1.1
- [4] Sebastian Erdweg, Tillmann Rendel, Christian Kastner, and Klaus Ostermann. SugarJ: Library-based syntactic language extensibility. In *OOPSLA '11*, pages 187–188, 2011. 1.1
- [5] Sebastian Erdweg, Tijs van der Storm, Markus Völter, Meinte Boersma, Remi Bosman, William R Cook, Albert Gerritsen, Angelo Hulshout, Steven Kelly, Alex Loh, Gabriël D. P. Konat, Pedro J. Molina, Martin Palatnik, Risto Pohjonen, Eugen Schindler, Klemens Schindler, Riccardo Solmi, Vlad A. Vergu, Eelco Visser, Kevin van der Vlist, Guido H. Wachsmuth, and Jimi van der Woning. The state of the art in language workbenches. In *Software Language Engineering (SLE '13)*. 2013. 1.1
- [6] Matthew Flatt. Creating languages in Racket. *Commun. ACM*, 55(1):48–56, January 2012. ISSN 0001-0782. doi: 10.1145/2063176.2063195. URL <http://doi.acm.org/10.1145/2063176.2063195>. 1.1
- [7] Matthew Fluet, Mike Rainey, John H. Reppy, Adam Shaw, and Yingqi Xiao. Manticore: a heterogeneous parallel language. In *Workshop on Declarative Aspects of Multicore Programming (DAMP '07)*, pages 37–44. ACM, 2007. ISBN 978-1-59593-690-5. URL <http://doi.acm.org/10.1145/1248648.1248656>. 1.1
- [8] M. Fowler and R. Parsons. *Domain-Specific Languages*. Addison-Wesley Professional, 2010. 2
- [9] Michael Furr and Jeffrey S. Foster. Checking type safety of foreign function calls. In *PLDI '05*, pages 62–72, 2005. ISBN 1-59593-056-6. doi: 10.1145/1065010.1065019. URL <http://doi.acm.org/10.1145/1065010.1065019>. 1.1
- [10] Cordelia V. Hall, Kevin Hammond, Simon L. Peyton Jones, and Philip L. Wadler. Type classes in Haskell. *ACM Trans. Program. Lang. Syst.*, 18(2):109–138, March 1996. ISSN 0164-0925. doi: 10.1145/227699.227700. URL <http://doi.acm.org/10.1145/227699.227700>. 1.2

- [11] Robert Harper. Programming in Standard ML. <http://www.cs.cmu.edu/~rwh/smlbook/book.pdf>. Retrieved June 21, 2015., 1997. 1.1
- [12] Simon L Peyton Jones. *Haskell 98 language and libraries: the revised report*. Cambridge University Press, 2003. 1.1
- [13] Xavier Leroy, Damien Doligez, Alain Frisch, Jacques Garrigue, Didier Rémy, and Jérôme Vouillon. *The OCaml system release 4.01 Documentation and user's manual*. Institut National de Recherche en Informatique et en Automatique, September 2013. 1.1, 1.1
- [14] Louis Mandel and Marc Pouzet. ReactiveML: a reactive extension to ML. In *PPDP '05*, pages 82–93. ACM, 2005. 1.1
- [15] The Coq development team. *The Coq proof assistant reference manual*. LogiCal Project, 2004. URL <http://coq.inria.fr>. Version 8.0. 1.2
- [16] Robin Milner, Mads Tofte, Robert Harper, and David MacQueen. *The Definition of Standard ML (Revised)*. The MIT Press, 1997. 1.1, 1.2
- [17] Tom Murphy, VII., Karl Crary, and Robert Harper. Type-safe Distributed Programming with ML5. In *Proceedings of the 3rd Conference on Trustworthy Global Computing, TGC'07*, pages 108–123, Berlin, Heidelberg, 2008. Springer-Verlag. ISBN 3-540-78662-7, 978-3-540-78662-7. URL <http://dl.acm.org/citation.cfm?id=1793574.1793585>. 1.1
- [18] Cyrus Omar, Darya Kurilova, Ligia Nistor, Benjamin Chung, Alex Potanin, and Jonathan Aldrich. Safely composable type-specific languages. In *ECOOP '14*, 2014. 1.1
- [19] Andreas Rossberg, Didier Le Botlan, Guido Tack, Thorsten Brunklaus, and Gert Smolka. *Alice Through the Looking Glass*, volume 5 of *Trends in Functional Programming*, pages 79–96. Intellect Books, Bristol, UK, ISBN 1-84150144-1, Munich, Germany, February 2006. 1.1
- [20] Todd L. Veldhuizen. *Active Libraries and Universal Languages*. PhD thesis, Indiana University, 2004. 1.2
- [21] Martin P. Ward. Language-oriented programming. *Software - Concepts and Tools*, 15(4): 147–161, 1994. 1.1
- [22] Xi, Chen, and Chen. Guarded recursive datatype constructors. In *POPL '03*, 2003. 1.1